

NUCLEAR: An Efficient Methods for Mining Frequent Itemsets and Generators from Closed Frequent Itemsets

Huy Quang Pham^{1,2}, Duc Tran³, Ninh Bao Duong², Philippe Fournier-Viger⁴, Alioune Ngom¹

Abstract—Frequent itemset (FI) mining is an interesting data mining task. Instead of directly mining the FIs from data it is preferred to mine only the closed frequent itemsets (CFIs) first and then extract the FIs for each CFI. However, some algorithms require the generators for each CFI in order to extract the FIs, leading to an extra cost. In this paper, we introduce an effective algorithm, called NUCLEAR, which can induce the FIs from the lattice of CFIs without the need of the generators. It can enumerate generators as well by similar fashion. Experimental results showed that NUCLEAR is effective as compared to previous studies, especially, the time for extracting the FIs is usually much smaller than that for mining the CFIs.

Keywords— Association rule, minimal association rule, kernel and extendable set, frequent itemset, closed frequent itemset, mining frequent itemset from closed frequent itemset, NUCLEAR.

I. INTRODUCTION

ASSOCIATION rule mining (ARM) is one of the most interesting and popular problems in data mining. It is widely used for decision making in retail, e-commerce, medicine, and many other domains. Mining frequent itemsets (FIs) is the first and the main step in the discovery of association rules (ARs). Since its first introduction in 1993 [1] it has attracted a lot of attention and has been extended and applied in various ways. For instance, some popular variations of the FI mining problem are to discover high utility patterns [2, 3], uncertain frequent patterns [2] and high utility association rules [2, 4]. Most algorithms for mining FIs partition the search space into subclasses in order to apply the parallel approaches to improve their performance. However, the performance of many parallel

FI mining algorithms is limited by the speed of disk accesses, as they repeatedly scan the input database, which can still lead to long execution [5], [6]. To address this issue, some

researchers proposed more efficient parallel algorithms, which compress the database in a frequent pattern tree and perform tree projections [7], [8]. Another approach which can speed up the FI mining process is to first mine all the closed frequent itemsets (CFIs) and then derive the FIs from them without the need of rescanning the data file. This approach is more efficient than mining FIs directly because the number of CFIs is usually much less than that of FIs (see Table IV). Charm [9], FPClose [7], DCI_PLUS [10] and NAFCP [6] are among the best algorithms for mining CFIs. In 2010, a parallel algorithm (PLCMQS) for mining CFIs has been proposed [11]. The authors of Charm proposed the CharmL algorithm [8], which builds the lattice of CFIs. Formal concept (i.e., lattice of CFIs) analysis is also another way of mining FIs as well as ARs [6], [12], [13].

Mining FIs from the lattice of CFIs has several advantages over mining FIs directly from data. First, the number of CFIs is often much smaller than the number of FIs; therefore, this requires less memory. Second, each CFI can stand for an equivalence class of FIs having the same closure (i.e., these FIs shares the same set of transactions containing them); thus, we can develop parallel algorithms or “divide-and-conquer” approaches to facilitate the process of deriving FIs from CFIs. Third, when users want to try different minimum support (*minsup*) thresholds to find an optimal set of FIs for a certain downstream procedure, the cost of updating the set of CFIs can be much lower than mining them from scratch in a database. However, to extract to FIs from CFIs, the current algorithms requires the generators for each CFI, and mining them might take a significant amount of time.

Contribution. In this paper, we present a method to mine the FIs from a lattice of CFIs without the need of the generators. We introduce the concepts of “kernels” and “extendable sets” which further partition the equivalence classes represented by the CFIs into smaller subclasses. Then, each pair of kernel and extendable set stands for a subclass of FIs which are supersets of the kernel and subsets of the union of the kernel and the extendable set. Thus, once a pair of kernel and extendable set are identified, enumerating FIs in the subclass is straightforward. Our proposed algorithm, called NUCLEAR, to generate kernels and extendable sets for each CFI is simple and efficient. Its inputs are just the largest CFIs

¹University of Windsor, Windsor, Ontario, Canada

²University of Dalat, Dalat, Vietnam

³Faculty of Information Technology, Ho Chi Minh City University of Food Industry, Vietnam

⁴School of Natural Sciences and Humanities, Harbin Institute of Technology Shenzhen Graduate School, Shenzhen, 518055, China

those are subsets of that CFI, and the time for it to induce all FIs from the lattice of CFIs is significantly shorter than the time to construct the lattice. As the generators play an important role in extracting the “minimal association rules”, we also provide the options for NUCLEAR to mine them by similar fashion. We also present a simple method, called NUC, that wraps CharmL and NUCLEAR to mine a lattice of CFIs from data and then infer the FIs from the lattice. The fact that NUCLEAR does not require the generators makes it more efficient than the approach infer FIs from CFIs and generators. In the comparison of NUC against dEclat [8], a well-known algorithm which mines FIs directly from data, NUC is faster when the number of FIs is much larger than the number of CFIs. When NUC is slower, the reason is that CharmL, the algorithm used to construct the lattice of CFIs from data before applying NUCLEAR, is already slower than dEclat.

The rest of this paper is organized as follows. Section 2 introduces some related concepts. Section 3 reviews the related works. Section 4 presents novel theoretical results that are the basis of the proposed algorithm, including a recurrent formula for generating kernels and extendable sets. Section 5 presents the proposed NUCLEAR algorithm. Section 6 reports experimental results that show the efficiency of the proposed algorithm. Finally, a conclusion is drawn and future work is discussed in section 7.

II. PRELIMINARIES

Let consider a context (T, I, R) where I is a set of items (or attributes), T is a set of transactions (or objects) and R is a binary relation on $T \times I$.

Definition 1 (Frequent itemset). For each non-empty subset A of I and non-empty subset O of T , the two functions λ and ρ below define a Galois connection between 2^T and 2^I (reader can refer to [28] for more details):

$$\lambda: 2^T \rightarrow 2^I: \lambda(O) = \{a \in I \mid (o, a) \in R, \forall o \in O\}, \lambda(\emptyset) = I$$

$$\rho: 2^I \rightarrow 2^T: \rho(A) = \{o \in T \mid (o, a) \in R, \forall a \in A\}, \rho(\emptyset) = T.$$

A is also called an itemset.

Given a user-specified minimum support threshold $minsup$, such that $0 < minsup \leq 1$, the support of an itemset A is denoted and defined as $supp(A) = |\rho(A)| / |T|$. A is said to be “frequent” if $supp(A) \geq minsup$.

Definition 2 (Association rule mining). Given T , $minsup$, and $minconf$, which are a transactional dataset, a minimum support threshold, and a minimum confidence threshold, respectively. The task of association rule mining (ARM) is to find all rules of the form $X \rightarrow Y$ such that $supp(X \cup Y) \geq minsup$ and $supp(X \cup Y) / supp(X) \geq minconf$. $supp(X \cup Y)$ is

called the support of the rule, and $supp(X \cup Y) / supp(X)$ is called the confidence of the rule.

As $(X \cup Y)$ must be frequent in order to $X \rightarrow Y$ be valid, identifying the frequent itemsets becomes the main task in ARM.

Definition 3 (Closed frequent itemset). $h = \lambda_o \rho$ in 2^I is called the closure operator [33], and $h(A) = \lambda_o \rho(A)$ is said to be the closure of A . Itemset C is “closed” if and only if $h(C) = C$. C is a closed frequent itemset if it is “closed” and it is “frequent”.

Let $[C] = \{A \subseteq I: h(C) = h(A)\}$ be the set of all itemsets having the same closure, which is $h(C)$, then, the itemsets in $[C]$ share the same set of transactions, which is $\rho(C)$, i.e., they have the same support.

Let CS and CFS denote the set of all closed itemsets and the set of all CFIs, respectively. Then, $L = (CFS, \preceq_A)$ is a lattice of CFIs, where \preceq_A is an order relation based on the \subseteq operator between subsets of I .

Definition 4 (Generator). An itemset G is called a “generator” of a closed itemset C if and only if $h(G) = h(C)$ and $\forall G': \emptyset \neq G' \subset G \Rightarrow h(G') \subset h(G)$.

For any itemset $A \subseteq I$, the equivalence class $[A]$ has only one closed itemset, and one or more generators.

Example 1. Given $I = \{1, 2, 3, 4, 5, 6\}$, $T = \{t_1, t_2, t_3, t_4, t_5, t_6\}$, and R as in the Table I. Let itemset $X = \{1, 4, 5\}$, then $\rho(X) = \{t_1, t_2\}$, $supp(X) = 2/4$, and $\lambda(\{t_1, t_2\}) = \{1, 4, 5, 6\}$. Then, X is not a CFI since $h(X) = \lambda_o \rho(X) = \{1, 4, 5, 6\} \neq X$. Now, let $C = \{1, 4, 5, 6\}$, we have $h(C) = C$. Thus, C is a CFI. And, we have, $[X] = [C] = \{\{1, 4\}, \{1, 5\}, \{1, 6\}, \{1, 4, 5\}, \{1, 4, 6\}, \{1, 4, 5, 6\}\}$, in which, the sets $\{1, 4\}$, $\{1, 5\}$, $\{1, 6\}$ are the generators of C .

The lattice of CFIs mined from R for $minsup = 0.25$ (i.e., absolute minimum support = 1) is shown in Fig.1. In this lattice, each node (rectangle) represents a CFI and its absolute support, separated by a colon (“:”) and each edge links a CFI to the CFIs which are its largest subsets.

TABLE I. THE RELATION $R = T \times I$, WHERE $I = \{1, 2, 3, 4, 5, 6\}$, $T = \{T_1, T_2, T_3, T_4, T_5, T_6\}$.

Transactions	Items					
t_1	1	2	3	4	5	6
t_2	1			4	5	6
t_3	1		3			
t_4		2	3	4	5	6

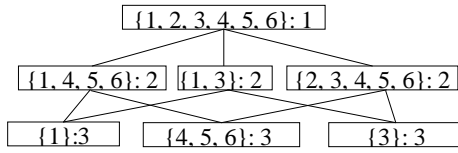


Fig. 1. The lattice of CFIs mined from R with minsup = 0.25.

Definition 5 (Minimal association rule). A rule $X \rightarrow Y$ is called “minimal” if there is no other rule $A \rightarrow B$ such that $(A \subseteq X)$ and $(Y \subseteq B)$ and $h(X \cup Y) = h(A \cup B)$ and $h(X) = h(A)$.

One can see that for any minimal association rule $X \rightarrow Y$, $X \cup Y$ is a closed (frequent) itemset and X is a generator. (X is either in $[X \cup Y]$ or in $[C]$, where $C \subset X \cup Y$.) [14, 15]

III. RELATED WORK

A. Mining Frequent (Closed) Itemsets

In recent years, several algorithms have been proposed for FI mining such as dEclat, and Node-list-based algorithms [5], [16]–[18]. The dEclat was one of the most effective algorithms according to [19]. It scans a database once to generate the transaction sets (tidsets) for all itemsets of 1 item (1-itemset). Then, it applies the “diffset” strategy to enumerate all FIs without repeatedly scanning the database. Deng et al. [16] proposed a novel structure named N-list for mining FIs. The proposed algorithm first compresses the dataset into a PPC-tree structure, and then, using that tree, the algorithm generates N-lists for each 1-itemset. Finally, the algorithm applies a divide-and-conquer approach to mine FIs using these lists. Experimental evaluation has shown that the N-list-based algorithm outperforms state-of-the-art FI mining algorithms on a variety of real and synthetic datasets. Recently, Deng and Lv [17] proposed an improved N-list based frequent itemset mining algorithm named PrePost+, which applies a novel pruning strategy called children-parent equivalence pruning to reduce the search space. Subsequently, Vo et al. [18] combined the N-list structure with the subsume concept to further increase the performance of FI mining. Recently, Deng [16] proposed an efficient algorithm relying on an improved Nodeset structure, named DiffNodesets.

As defined in section II, the closure of an itemset is the set of items that appear in all transactions containing the itemset. CFIs have attracted a lot of studies as they can be used to partition FIs into equivalence classes. This inspires the development of parallel or “divide-and-conquer” approaches to mine FIs from CFIs without scanning the database for the support. However, few approaches have been proposed to perform this efficiently. Several researchers have studied retrieving FIs using the generator itemsets and eliminable itemsets in the equivalence classes of their closures [10], [14],

[15], [19], [20]. For this purpose, algorithms were proposed that efficiently discover FIs using the lattice of CFIs, without performing duplicate checks, and by processing only one CFI at a time, that is, without considering its relationship to other CFIs. Generator itemsets can be mined independently or at the same time as CFIs. Zaki et al. [9] proposed the Minimal Generators algorithm to mine generators from the lattice of CFIs using a level-wise approach inspired by the Apriori algorithm. However, to identify the generators of a CFI, the algorithm had to scan all its subsets. Thus, the algorithm can be very slow. Szathmary et al. [21] proposed the Talky-G algorithm to mine generators from data, using an IT-tree structure. The algorithm uses the Charm algorithm [8], [9] to separately mine the CFIs and then, matches the generators with each CFI. Talky-G guarantees that when an itemset X is visited during the search, all its subsets have been already visited, and thus all generators that are subsets of X have already been found. Consequently, an itemset X is a generator if no already found generator is subset X and has the same support as X ’s support. To quickly select the generators for that check Talky-G stores the support of visited generators in a hash table using the number of transactions containing each itemset as the hash function. This hash function is also used to match each CFI to its generators. The algorithm is effective when minsup is high. However, the time required for finding generators is similar to the time for mining CFIs. GENCLOSE [22] is an algorithm that concurrently mines CFIs and generators. The authors introduced necessary and sufficient conditions to generate generators $(k + 1)$ -itemsets, i.e., itemset containing $k+1$ items, using the generators of k -itemsets. Using these conditions, the closure of each generator can be extended gradually to find generators. In 2005, the CHARM [9] and dCHARM [8] algorithms have been proposed for mining CFIs using the “diffset” structure introduced in the dEclat algorithm. In 2012, the DBV-Miner [23] algorithm improved this approach by compressing the tidsets of 1-itemsets using dynamic bit vectors. It was shown that this can greatly reduce the memory required for storing tidsets and compute the support of itemsets efficiently. Then, Sahoo et al. [10] proposed the DCI_PLUS algorithm for mining CFIs and their generators. The algorithm compresses the database using a BitTable structure, which is built using a single database scan. In [20], Tran et al. proposed the GEN_ITEMSETS algorithm to generate all itemsets from a lattice of CFIs and generators without repetitions. More recently, Le and Vo [12] proposed an N-list-based algorithm for mining CFIs, named NAFCP. The experimental evaluation of this work has shown that NAFCP outperforms state-of-art CFI mining algorithms in terms of runtime and memory usage in most cases.

B. Lattice-based Approaches for Mining Association Rules

In general, two types of lattices are considered for ARM, which are the frequent itemset lattice (FIL) and the closed

frequent itemset lattice (CFIL) [12]. Vo and Le [13] have presented a lattice-based approach for building the FIL, here called FIL-2009. Each node of the structure used in FIL-2009 represents an itemset X and stores a tuple $(X, Tidset, Children)$ where $Tidset$ is the list of transactions containing X and $Children$ are pointers to nodes representing supersets of X . To mine the minimal non-redundant association rules, Vo and Le [24] extended the structure used by FIL-2009 (here called FIL-2011) by adding two fields in each node indicating if a node is a minimal generator or a closed frequent itemset, respectively. These values are determined during lattice construction. The structure is then used by FIL-2011 to effectively mine minimal non-redundant association rules. Thereafter, an efficient approach named PFIL was proposed, which supports incremental mining using the pre-large concept. It was shown that this approach is especially efficient for huge databases containing a large number of FIs [25]. The PFIL algorithm uses the diffset structure to quickly build a FIL. Then it uses the pre-large concept and diffset structure for maintaining the pre-large FIL.

For a given dataset and a *minsup* threshold, building the CFIL is generally much faster than building the FIL because the number of CFIs is usually much less than that of FIs. CharmL [8] is an effective algorithm to build the lattice of CFIs. To update the lattice, researchers have proposed an algorithm [5] that runs efficiently in the case of large databases with a small number of inserted transactions.

For parallel algorithms or for the survey on algorithms for mining FIs and ARs we refer readers to [26]-[31].

IV. THEORETICAL RESULTS

In this section, we introduce theoretical results that are the basis of our proposed algorithm.

From now on, for convenience, whenever we use the variables C , and/or G without condition, it implicitly means that: $C \in CS$, $\emptyset \neq G \subseteq C$ respectively. And, let “.” stand for “such that”, and “,” stand for the logical operator “ \wedge ” in the logical propositions.

Definition 6 (the immediately closed subsets of a closed itemset). Let $S_C = \{Y \in CS: (Y \subset C) \wedge (\nexists Z \in CS: Y \subset Z \subset C)\}$ be the set of the largest closed itemsets that are subset of a given closed itemset C . These itemsets in S_C are called the immediately closed subsets of C .

Proposition 1. $\forall G \subseteq C$, $\forall Y \in S_C$ ($G \in [C] \Leftrightarrow \exists x \in G: x \notin Y$).

Proposition 1 points out a way to find $[C]$ by searching every itemset G satisfying the right-hand side of Proposition 1. However, it might be not efficient if we have to scan every

subset of C .

Proof. For all $G \subseteq C$ and $Y \in S_C$.

“ \Rightarrow ”: Since $G \in [C]$, we have $h(G) = C$. Now, assume that $\forall x \in G (x \in Y)$. Thus, $G \subseteq Y \subset C \Rightarrow h(G) \subseteq h(Y) = Y \subset C$. This leads to a contradiction: $h(G) \subset C$.

“ \Leftarrow ”: Assume that $\forall Y \in S_C (\exists x \in G: x \notin Y) \wedge (G \notin [C])$. We have: $\forall Y \in S_C ((h(G) = h(Y) = Y) \vee (h(G) \neq Y)) \wedge (h(G) \subset C)$. By the definition of S_C , we have: $\forall X \subset C (X \in CS \Rightarrow \exists Y \in S_C: X \subseteq Y)$. Then, $\forall Y \in S_C (h(G) \subseteq Y)$. Thus, we have: $\forall Y \in S_C, \forall x \in G (x \in Y)$, which is a contradiction to the hypothesis.

Corollary 1. $\forall Y_k \in S_C$, let $M_k = C \setminus Y_k$, $NS = |S_C|$, $1 \leq k \leq NS$, and $M = \{M_1, \dots, M_{NS}\}$. We have $G \in [C] \Leftrightarrow \forall M_k, \exists x \in G: x \in M_k$.

Proof. We have $G \in [C] \Leftrightarrow \forall Y \in S_C, \exists x \in G: x \notin Y \Leftrightarrow \forall M_k, \exists x \in G: x \in M_k$. Therefore, this corollary is proven.

Definition 7 (Kernel and extendable set). Given two itemsets G and E those are disjoint. Let the notation $[G, E]$ denote the class of all itemsets those are supersets of G and subsets of $G + E$, i.e., $[G, E] = \{X: G \subseteq X \subseteq G + E\}$.

G is called the kernel and E is called the extendable set of the class.

The following results provide an easy way to find the pairs of kernel and extendable set that can help partition $[C]$ into equivalence classes.

Definition 8. Let $M_k = C \setminus Y_k$, $\forall Y_k \in S_C$, $M = \{M_1, \dots, M_{NS}\}$. For $1 \leq k \leq |S_C|$, let $S_k = \{M_1, \dots, M_k\}$ denotes the set of k first elements of M , and $S_0 = \emptyset$.

An itemset G is said to “satisfy S_k ” if $\forall M_i \in S_k, \exists x \in G: x \in M_i$. Let $[S_k]$ denote the set of all FIs that satisfies S_k , and $[S_k] = 2^C$ if $k = 0$.

The following lemmas are obtained by Definition 8 and Corollary 1.

Lemma 1. $\forall G \subseteq C$, $E \subseteq C$, $G \cap E = \emptyset$, $1 \leq k \leq |S_C|$, we have:

$[S_k] \subseteq [S_{k-1}]$, (i.e., if G satisfies S_k , it also satisfies S_{k-1}).

$\forall G \in [S_k] \Rightarrow \forall X \in [G, E], X \in [S_k]$.

$[C] = [S_{NS}]$.

Lemma 2. For $1 \leq k \leq |S_C|$, $G^* \in [S_{k-1}]$, and let $G = G^* + \{x\}$, $x \in C$, we have: $(M_k \cap G^* \neq \emptyset) \vee (x \in M_k) \Rightarrow G \in [S_k]$

In the first condition case (i.e., $M_k \cap G^* \neq \emptyset$), x is not necessary for G to satisfy S_k since G^* already satisfy S_k , meanwhile, in latter it is.

From now on, we denote “+” (and “ Σ ”) the union operator for two (and many, respectively) disjoint sets. Definition 9 below lead to the idea of generating the kernel sets.

Definition 9 (k-minimal set). Given $G = G^* \cup X$, and $X \subseteq M_k$. G is “k-minimal” if one of the following conditions is satisfied:

G^* is “(k-1)-minimal”, $M_k \cap G^* \neq \emptyset$, $X = \emptyset$ (i.e., $G = G^*$, and no more item are needed for G^* to satisfy S_k).

G^* is “(k-1)-minimal” $M_k \cap G^* = \emptyset$, and $X = \{x\}$, $x \in M_k$ (i.e., x is the new item needed for G^* to satisfy S_k).

\emptyset is “0-minimal”.

We can see that if G is “k-minimal” then G satisfies S_k . Cases (a), and (b) are based on Lemma 2.

G is “k-minimal” does not imply that there is no subset of G satisfying S_k . It just implies that no prefix of G satisfies S_k if G is treated as a sequence of items.

Here after, we assume that there exists an order over the items in C (e.g., alphabetic order), and every M_i is sorted in the increasing order.

Definition 10 (A partition of $[S_k]$ by kernels and

extendable sets). Given C , and S_C . Let the set Q_k contain the pairs of kernel set and extendable set defined recurrently as follows:

$$Q_0 = \{(\emptyset, C)\}$$

$$Q_1 = \{(G_i, E_i): G_i = \{x_i\}, x_i \in M_1, E_i = C \setminus \{y \in M_1: y \leq x_i\}\}$$

$$\forall k > 1, Q_k = B_k + C_k \text{ where:}$$

$$B_k = \{(G, E): (G, E) \in Q_{k-1}, G \cap M_k \neq \emptyset\},$$

$$C_k = \{(G + \{x_i\}, E|E_i): (G, E) \in Q_{k-1}, G \cap M_k = \emptyset, N_k = M_k \cap E, x_i \in N_k, E_i = \{y \in N_k: y \leq x_i\}\}.$$

In details, B_k contains pairs (G, E) in Q_{k-1} where G is “(k-1)-minimal” and is also “k-minimal”, according to Definition 9.a. Thus, (G, E) belongs to Q_k also. Meanwhile, C_k contains the pairs (G, E) such that $G = G' + \{x_i\}$, where G' is “(k-1)-minimal” but not “k-minimal”, and x_i is necessary for G to become “k-minimal” (as Definition 9.b).

Lemma 3. For all $(G_i, E_i), (G_j, E_j) \in Q_k$, $i \neq j$, $1 \leq k \leq |S_C|$, we have:

$$a) X \in [G_i, E_i] \Rightarrow X \in [S_k].$$

$$b) G_i \cap E_i = \emptyset.$$

$$c) [G_i, E_i] \cap [G_j, E_j] = \emptyset.$$

In other words, Q_k induces a partition of all FIs in $[S_k]$, where each equivalence classes is defined by $[G, E]$ as in Definition 10, with $(G, E) \in Q_k$. Furthermore, for every $(G, E) \in Q_k$, G is “k-minimal”.

Theorem 1. For $0 \leq k \leq |S_C|$, $\{[G, E]: (G, E) \in Q_k\}$ is a partition of $[S_k]$, where each $[G, E]$ is an equivalence class, and G is “k-minimal”

Proof.

We'll first prove that Theorem 1 hold with $k = 0$.

Since $k = 0$, $Q_0 = \{(\emptyset, C)\}$ by Definition 10.a, and $[S_k] = 2^C$ by Definition 8. We have: $[\emptyset, C] = \{X: X \subseteq C\} = 2^C$. Since $\{[\emptyset, C]\}$ is a partition of 2^C and \emptyset is “0-minimal”, let $G = \emptyset$ and $E = C$ then Theorem 1 is proven for $k = 0$.

Assume that Theorem 1 holds for any $k-1$, $0 < k \leq |S_C|$ (i.e., the set $\{[G, E]: (G, E) \in Q_{k-1}\}$ is a partition of $[S_{k-1}]$, where each $[G, E]$ is an equivalence class, and G is “(k-1)-minimal”), we will prove that Theorem 1 holds for k .

By assumption, we have $[S_{k-1}] = \sum [G_i, E_i]$, where $(G_i, E_i) \in Q_{k-1}$. By Lemma 1.a, for an itemset X to be in $[S_k]$, it must be in $[S_{k-1}]$ ^(a). By Lemma 3.c, for all (G_i, E_i) and $(G_j, E_j) \in Q_{k-1}$, $i \neq j$, we have: $[G_i, E_i] \cap [G_j, E_j] = \emptyset$ ^(b). From ^(a) and ^(b), we only need to prove that given a pair $(G, E) \in Q_{k-1}$, we can partition $[G, E]$ into disjoint subclasses, where each subclass either is in the form of $[G', E']$ and G' is a “k-minimal” (i.e., $(G', E') \in Q_k$), or contains only the itemsets which do not satisfy S_k ^(*).

If $M_k \cap G \neq \emptyset$, then G satisfies S_k . Then, $(G, E) \in Q_k$. Then ^(*) is proved. In this case, (G, E) belongs to B_k as in Definition 10.c, thus, it belongs to Q_k ⁽ⁱ⁾. Now, let assume that $M_k \cap G = \emptyset$. If $M_k \cap E = \emptyset$, then for all $Y \in [G, E]$, Y does not satisfy S_k . Then ^(*) is proved ⁽ⁱⁱ⁾.

Now, let assume that $M_k \cap G = \emptyset$ and $M_k \cap E \neq \emptyset$. Denote $N_k = M_k \cap E = \{x_1, \dots, x_n\}$, and for $0 < i \leq n$, denote $G_i = G + \{x_i\}$ ($x_i \in N_k$), $E_i = E \setminus \{x_j: x_j \in N_k, j \leq i\}$. Let $U_1 = \{G_1 + T: T \subseteq E_1\} = [G_1, E_1]$, and $V_1 = \{G + Y: Y \subseteq E \setminus \{x_1\}\} = [G, E_1]$. Then, $\forall X \in U_1$ ($x_1 \in X$) and $\forall Y \in V_1$ ($x_1 \notin Y$). This means: U_1 and V_1 are disjoint. We can further divide V_1 into two disjoint sets: U_2 – the set of itemsets containing x_2 , and V_2 – the set of itemsets not containing x_2 . One can see that $E_i = E_{i-1} \setminus \{x_i\}$, then, we have: $U_2 = \{G_2 + T: T \subseteq E_2\} = [G_2, E_1 \setminus \{x_2\}] = [G_2, E_2]$, $V_2 = [G, E_2]$, and U_2 and $V_2 = \emptyset$. By the same way, the division process continues until we divide V_{n-1} into two disjoint sets: $U_n = [G_n, E_n]$ and $V_n = [G, E_n]$. One can see that, $[G, E] = V_n + \sum [G_i, E_i]$.

Since $M_k \cap G = \emptyset$ and $E_n = \emptyset$, for all $Y \in V_n$, Y does not satisfy S_k . Meanwhile, since $N_k \cap G_i = \{x_i\} \neq \emptyset$, for all $i \leq n$. This means $M_k \cap G_i \neq \emptyset$. Thus, G_i is “k-minimal” by Definition 9.b and for all $X \in [G_i, E_i]$, X satisfies S_k . This mean (G_i, E_i) is generated as Definition 10.c then $(G_i, E_i) \in Q_k$. Then ^(*) is proved. ⁽ⁱⁱⁱ⁾

By ⁽ⁱ⁾, ⁽ⁱⁱ⁾, and ⁽ⁱⁱⁱ⁾, ^(*) is proved, and Theorem 1 is proved.

Corollary 2. $\{(G, E) : (G, E) \in Q_{NS}\}$ is a partition of $[C]$, where each $[G, E]$ is an equivalence class, and G is “NS-minimal”.

Proof. This is result of Theorem 1, where $k = NS$, $[C] = [S_k]$.

Example 2. Let consider the relation R shown in Table I and the lattice of CFIs mined from R with $minsup = 0.25$ (i.e., absolute minimum support = 1) shown in Fig. 1. The following paragraphs explain how to find all FIs in $[C]$ for $C = \{1, 2, 3, 4, 5, 6\}$.

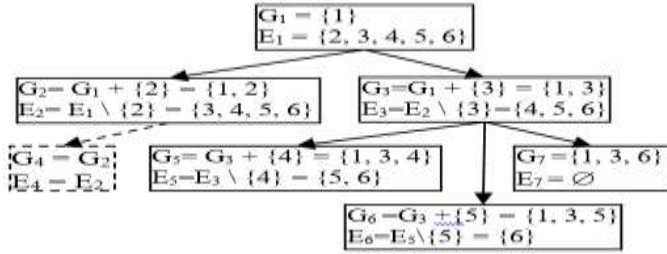


Fig. 2. Search tree for generating the kernels and extendable sets of Q_3 , where $Q_3 = \{(G_2, E_2), (G_5, E_5), (G_6, E_6), \{(G_6, E_7)\}$.

According to the lattice, the immediately closed subsets of C are $Y_1 = \{2, 3, 4, 5, 6\}$, $Y_2 = \{1, 4, 5, 6\}$ and $Y_3 = \{1, 3\}$. In other words, $S_C = \{Y_1, Y_2, Y_3\}$. Fig. 2 presents the search tree that can be built, implicitly, during the process of generating the kernels and extendable sets using a breadth-first search. Here, $Q_1 = \{(G_1, E_1)\}$, $Q_2 = \{(G_2, E_2), (G_3, E_3)\}$, and $Q_3 = \{(G_2, E_2), (G_5, E_5), (G_6, E_6), \{(G_6, E_7)\}$, where G_i is a kernel and E_i is its corresponding extendable set. (We do not have to compute the pair (G_4, E_4) , the rectangle with dash-lined border.) They are found by the following steps:

With $M_1 = C \setminus Y_1 = \{1\}$, by Definition 5.b, let $G_1 = \{1\}$, $E_1 = C \setminus \{1\} = \{2, 3, 4, 5, 6\}$ we have: $Q_1 = \{(G_1, E_1)\}$. (Then, $[S_1] = [G_1, E_1]$, but we do not need to compute it!)

Now, we will find Q_2 based on (G_1, E_1) and M_2 , where $M_2 = C \setminus Y_2 = \{2, 3\}$. Let $N_2 = M_2 \cap E = \{2, 3\}$. Using item 2 in N_2 we have $G_2 = G_1 + \{2\} = \{1, 2\}$, and $E_2 = E_1 \setminus \{2\} = \{3, 4, 5, 6\}$. Using item 3 in N_2 we have $G_3 = G_1 + \{3\} = \{1, 3\}$, and $E_3 = E_1 \setminus \{3\} = \{2, 4, 5, 6\}$. Then, $Q_2 = \{(G_2, E_2), (G_3, E_3)\}$ as nodes of level 2 in Fig. 2.

Now, we will find Q_3 based on $\{(G_2, E_2), (G_3, E_3)\}$ and M_3 , where $M_3 = C \setminus Y_3 = \{2, 4, 5, 6\}$. With (G_2, E_2) , one can see that G_2 satisfies S_3 since $M_3 \cap E_2 = \{2\} \neq \emptyset$. Thus, (G_2, E_2) belongs to Q_3 . With (G_3, E_3) , we have $N_3 = M_3 \cap E_3 = \{4, 5, 6\}$. With item 4 in N_3 , we have $G_5 = G_3 + \{4\} = \{1, 3, 4\}$, and $E_5 = E_3 \setminus \{4\} = \{2, 5, 6\}$. With item 5 in N_3 , we have $G_6 = G_3 + \{5\} = \{1, 3, 5\}$, and $E_6 = E_3 \setminus \{5\} = \{2, 4, 6\}$. With item 6 in N_3 , we have $G_7 = G_3 + \{6\} = \{1, 3, 6\}$, and $E_7 = E_3 \setminus \{6\} = \emptyset$. Then,

$Q_3 = \{(G_2, E_2), (G_5, E_5), (G_6, E_6), \{(G_6, E_7)\}$ as nodes of level 3 in Fig. 2.

TABLE II. THE PARTITION OF 23 FIs IN $\{1, 2, 3, 4, 5, 6\}$ BASED ON Q_3 AS IN FIG. 2.

	$[G_2, E_2]$	$[G_5, E_5]$	$[G_6, E_6]$	$[G_7, E_7]$
1	$\{1, 2\}$,	$\{1, 3, 4\}$,	$\{1, 3, 5\}$,	$\{1, 3, 6\}$
2	$\{1, 2, 3\}$,	$\{1, 3, 4, 5\}$,	$\{1, 3, 5, 6\}$	
3	$\{1, 2, 3, 4\}$,	$\{1, 3, 4, 5, 6\}$,		
4	$\{1, 2, 3, 4, 5\}$,	$\{1, 3, 4, 6\}$		
5	$\{1, 2, 3, 4, 5, 6\}$,			
6	$\{1, 2, 3, 4, 6\}$,			
7	$\{1, 2, 3, 5\}$,			
8	$\{1, 2, 3, 5, 6\}$,			
9	$\{1, 2, 3, 6\}$,			
10	$\{1, 2, 4\}$,			
11	$\{1, 2, 4, 5\}$,			
12	$\{1, 2, 4, 5, 6\}$,			
13	$\{1, 2, 4, 6\}$,			
14	$\{1, 2, 5\}$,			
15	$\{1, 2, 5, 6\}$,			
16	$\{1, 2, 6\}$			

Theorem 2. Let GS be the set of all generators of C , and $KS = \{G : (G, E) \in Q_{NS}\}$ be the set of all kernels. We have:

$$GS \subseteq KS$$

$$G \in GS \Leftrightarrow (G \in KS) \wedge (\nexists K \in KS: K \subset G).$$

Proof. By Corollary 2, Q_{NS} induce a partition of $[C]$, where each $[K, E]$ ($(K, E) \in Q_{NS}$) is an equivalence class. Then, $\forall G \in GS, \exists (K, E) \in Q_{NS}: G \in [K, E]$. This mean $h(G) = h(K) = C$. Then, by Definition 4, $G = K$ since $\forall X \in [K, E], K \subseteq X$. Thus, $\forall G \in GS, \exists K \in KS: G = K$. Or, $GS \subseteq KS$. This means Theorem 2.a is proven. As a consequence, a kernel is a generator if and only if it is not a superset of any other kernel, i.e., Theorem 2.a is proven.

V. OUR ALGORITHMS

In this section, based on Definition 10, Theorem 1, and Theorem 2 we present NUCLEAR (Fig. 4), the algorithm for enumerating the FIs/generators from a lattice of CFIs, for all closed frequent itemset. There are options for it to find only the generators, or both FIs and generators. The NUC algorithm, in Fig. 3, wraps NUCLEAR to mine FIs/generators from data (files). It uses a certain algorithm (e.g., CharmL) to construct the lattice of CFIs from data which will be used as the input for NUCLEAR. NUCLEAR will call FindKandE_BFS (Fig. 5) to enumerate the pairs of kernels and extendable sets for each closed frequent itemset, in breadth-first-search manner. Depending on the output type(s) required, the FIs and/or the generators will be enumerated accordingly.

Input: D : the transactional dataset,
 $minsup$: the minimum support threshold,
 $returnFIs$: “true” to return all FIs for each closed frequent itemset,
 $returnGenerators$: “true” to return all generators for each closed frequent itemset.
Output: all FIs/generators satisfying $minsup$.

- 1 $L =$ The lattice of CFIs mined from D for the given $minsup$ (using a certain algorithm, e.g., CharmL).
- 2 return NUCLEAR (L , $returnFIs$, $returnGenerators$)

Fig. 3. The NUC algorithm.

Input: L = the lattice of CFIs
 $minsup$: the minimum support threshold,
 $returnFIs$: “true” to return all FIs for each closed frequent itemset,
 $returnGenerators$: “true” to return all generators for each closed frequent itemset.
Output: all FIs satisfying $minsup$.

- 1 $\forall C \in L$,
- 2 $M = \{M_i \mid M_i = C \setminus Y_i, Y_i \in SC\}$;
- 3 $Q_0 = \{(\emptyset, C)\}$;
- 4 $Q_{[NS]} = FindKandE_BFS(1, Q_0, M)$;
- 5 for each (G, E) in $Q_{[NS]}$
- 6 if ($returnFIs$) then generate and save the class $[G, E]$ for C
- 7 if ($returnGenerators$)
- 8 $GS = \emptyset$
- 9 for each (G, E) in $Q_{[NS]}$
- 10 if $(\exists(K, E) \in Q_{[NS]} : K \subset G)$ then $GS = GS \cup \{G\}$
- 11 save GS for C

Fig. 4. The NUCLEAR algorithm.

Input: k : a positive integer,
 Q_{k-1} : The set of pairs of kernel and extendable set generated as Definition 10.
Output: $Q_{[NS]}$ all pairs of kernel and extendable of the current closed frequent itemset C .

- 1 if $(k > |SC|)$ return Q_{k-1}
- 2 else
- 3 $Q_k = \emptyset$;
- 4 for each (G, E) in Q_{k-1}
- 5 if $(M_k \cap G \neq \emptyset)$ then $Q_k = Q_k \cup \{(G, E)\}$;
- 6 else
- 7 $E_i = E$;
- 8 for each x in $M_k \cap E$
- 9 $E_i = E_i \setminus \{x\}$;
- 10 $Q_k = Q_k \cup \{(G \cup \{x\}, E_i)\}$;
- 11 FindKandE_BFS ($k+1, Q_k$);

Fig. 5 The FindKandE_BFS algorithm.

For each immediate frequent closed subset C in lattice L , the variables G and E are initialized as $G = \emptyset$ and $E = C$. If $SC = \emptyset$, the class $[C]$ is all non-empty subsets of C . Otherwise, the recursive FindKandE_BFS algorithm is called (Fig. 4).

VI. EXPERIMENTAL RESULTS AND DISCUSSION

TABLE III. CHARACTERISTICS OF THE DATASETS.

Database	Abbreviation	#Items	#Transactions	#Average
Chess	CH	75	3,196	37
Connect	CO	129	67,557	43
Mushroom	MU	119	8,124	23
Retail	RE	16,469	88,162	10.3
T40I10100K	T4	1,000	100,000	40
C20d10k	C2	192	10,000	20
C73d10k	C7	1,592	10,000	73

In this section, we compare the running time (in seconds) for mining all FIs from data of three frameworks: dEclat, GenIT, and NUC. dEclat mines FIs directly from data. GenIT is not an algorithm but a combination of algorithms of some previous studies, which is slightly different to our proposed approach. In GenIT, we first use CharmL algorithm to mine the lattice of CFIs from data, then, use Minimal Generator algorithm to mine the generators from the CFIs before applying GEN_ITEMSETS to extract FIs from the CFIs and the generators by FIs. In NUC, we also use CharmL algorithm to mine the lattice of CFIs from data before applying NUCLEAR to generate the kernels and extendable sets. The FIs are inferred during this process. Thus, to be fair, the total times reported for NUC includes the time of CharmL and the time of NUCLEAR; whereas, that for GenIT includes the time of CharmL, the time of Minimal Generator, and the time of GEN_ITEMSETS. These algorithms have been executed on a Pentium (R) Dual-Core CPU E6500 @ 2.93GHz, equipped with 1.94GB of RAM, running the Microsoft Windows XP Version 2002 operating system.

Seven datasets (available at [32], [33]) have been used to compare the frameworks under different $minsup$ threshold values. Information about the datasets is given in Table III.

The number of patterns can be mined from each dataset are shown in Table IV (Appendix). Table V in (Appendix) shows the overall runtime for NUC, GenIT, and dEclat in the columns of the corresponding names, and other details. The visual comparisons of the three approaches are also given in Fig. 5 – Fig.12.

In our experiment, NUC is faster than dEclat when testing on the Mushroom and Connect datasets (Fig. 9, Fig.10). The

reasons are: (1) the time for CharmL (column tCS) is smaller than that of dEclat because the number of CFIs in a dataset is much less than that of FIs, and (2) the time for NUCLEAR (column tNNI) is also smaller than that of dEclat. On the other datasets, NUC is slower than dEclat. However, the main reason is just CharmL is slower dEclat, whereas, the time for NUCLEAR is significantly small as compared to those of CharmL and dEclat.

From Table V, we can estimate that NUC is about 1.25 time faster than GenIT. To clearer see the advantages of NUCLEAR, we break down the runtime of GenIT and NUC into the times for different stages including:

- Constructing the lattice of CFIs by CharmL (column tCS),
- Mining generators by Minimal Generator (column tG),
- Extracting FIs from the CFIs and the generators by GEN_ITEMSETS (column tGI),
- Generating kernels and extendable sets from the lattice of CFIs by NUCLEAR (column tN),
- Enumerating the FIs from the kernels and extendable sets by NUCLEAR (column tNI).

In comparing tG against tN, and tGI against tNI, one can see that the runtime of Minimal Generator to mine the generators is about double that of NUCLEAR for generating the kernels and extendable sets (tN), and the time for GEN_ITEMSETS (tGI) is similarly about twice as that of NUCLEAR for extracting FIs. These make NUCLEAR more efficient than GenIT in extracting the intermediate and the final results.

The time for minings FIs from the lattice using NUCLEAR (tNNI) is mostly much smaller than that for constructing the lattice of CFIs (tCS) using the CharmL algorithm. Thus, in the applications where users have to try different *minsup* thresholds to find an optimal set of FIs for a certain downstream process, our approach might be more efficient than repeatedly mining FIs from scratch like dEclat. Because, the lattice can be constructed only once for a small enough *minsup*, the cost for updating/filtering the lattice is expected to be small, and after that NUCLEAR can be used to query FIs many times. For example, in bioinformatics, we can use NUC to conduct a feature selection which shrinks the high dimensional data to a smaller one before applying a machine learning algorithm. The feature selection might have to be conducted many times to obtain an optimal set of features.

From Table IV, we can see that the number of pairs of kernels and extendable sets ($\#N$) is almost similar to the number of generators ($\#G$) and just slightly bigger than the number of CFIs ($\#C$). The lowest value of $\#C/\#N$ is 0.55 means that there are no more than two kernels per CFI. Thus,

identifying them as demonstrated in Example 2 will be fast. As there are only a few kernels in each CFI and the number of generators is almost similar to the number of kernels ($\#N/\#G$ is almost equal to 1), extracting the generators among the kernels would be very quick. Furthermore, we do not need to store the kernels and extendable sets but just the CFIs. Because, for each CFI, we only need the largest CFIs those are its subsets to generate the kernels, extendable sets, and FIs, and this can be done quickly. Thus, using NUCLEAR can save a lot of memory as well.

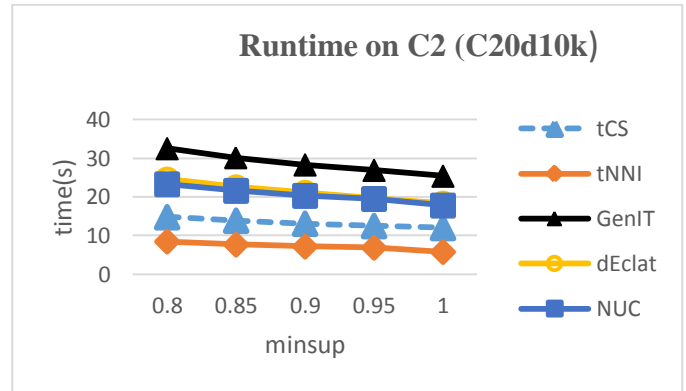


Fig. 6 Time execution comparison on C20d10k.

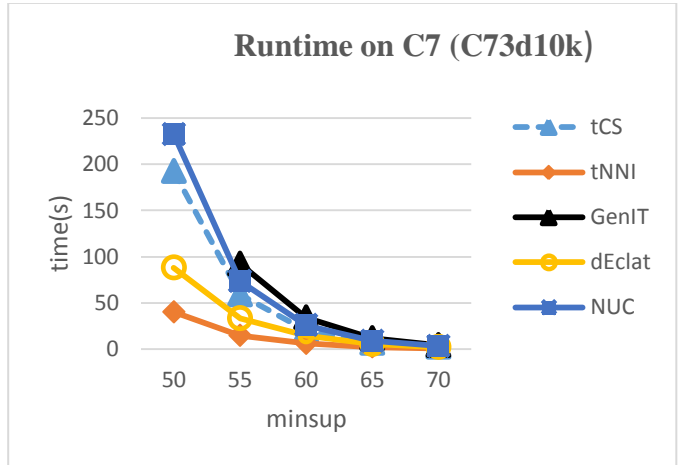


Fig. 7 Time execution comparison on C73d10k.

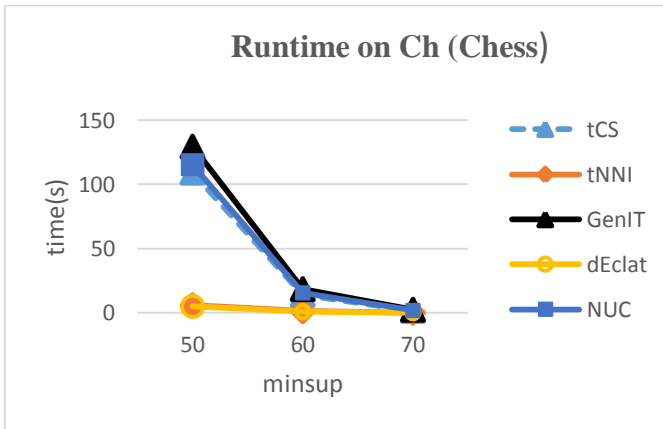


Fig. 8 Comparison of time on Chess.

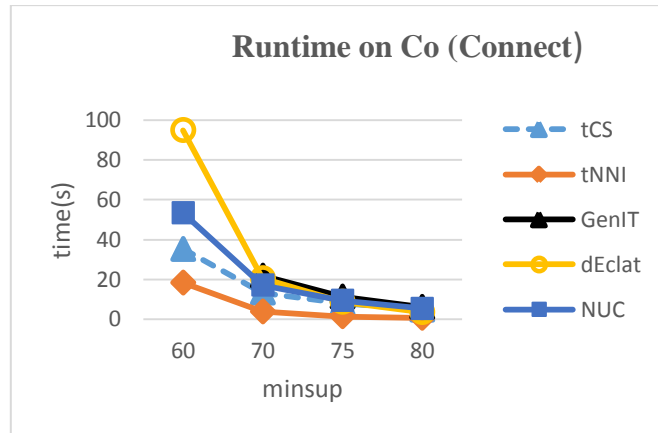


Fig. 9 Comparison of time on Connect.

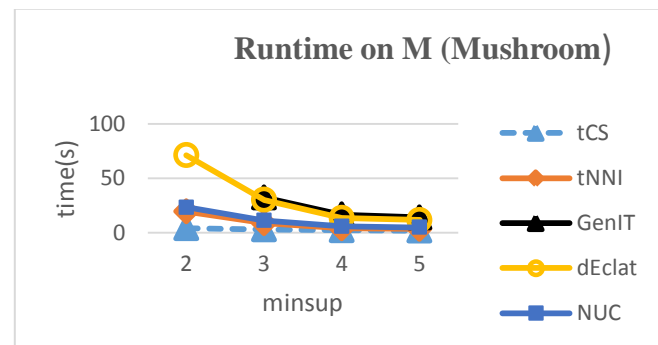


Fig. 10 Comparison of time on Mushroom.

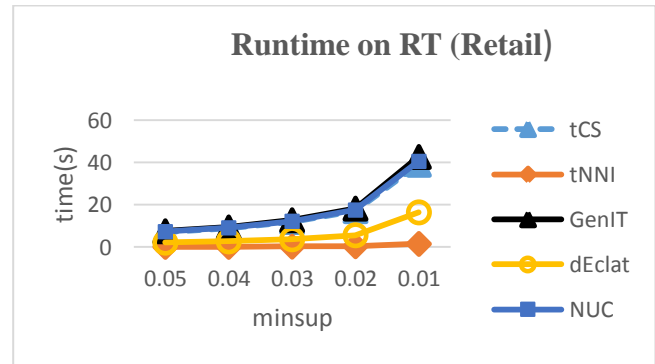


Fig. 11 Time execution comparison on Retail.

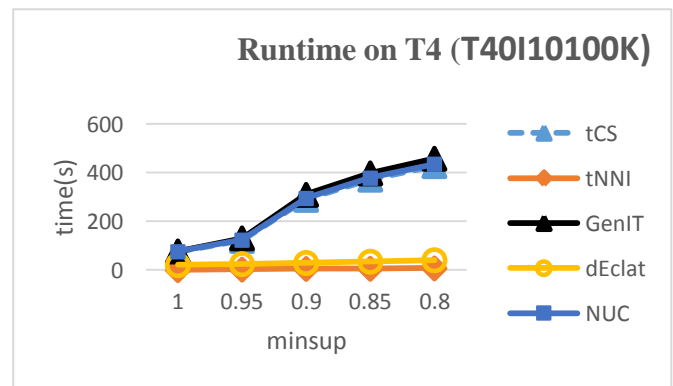


Fig. 12 Time execution comparison on T40I10100K.

VII. CONCLUSIONS AND FUTURE WORK

Mining FIs from the lattice of CFIs is a reasonable approach since the number of CFIs is often much smaller than the number of FIs. Thus, CFIs can be mined with a limited amount of memory. Especially, because there are parallel and distributed algorithms to mine CFIs for large or high dimensional data, the CFIs are easier to be available than the FIs. Besides, CFIs can be used to partition FIs into equivalence classes that can be used to efficiently process FIs in parallel. This approach is interesting as the lattice of CFIs can be mined once for a minimum support threshold that is small enough and used many times later to derive FIs for different minimum support thresholds.

In this paper, we presented a recurrent formula for generating the kernels and extendable sets for a closed frequent itemset, without the need of the generators. They are simple enough so that users can easily and quickly derive the FIs and/or the generators from them, and we even don't need to store them. Thank for that, NUC, the approach using NUCLEAR to mine the FIs from the lattice of CFIs is more efficient than GEN_IT, a similar approach that requires the

generators for mining FIs from the lattice of CFIs. NUC is slower than dEclat in the major cases, but it's just mainly because the construction of the lattice by CharmL takes more time than dEclat; whereas, the time for obtaining the FIs from the lattice by NUCLEAR is still considerably small.

In the future, the methods for updating the FIs when *minsup* is changed will be studied for the case that lattice can be

constructed only once and reuse many times. We would like to test our approach on the real data, such as bioinformatics data, where there are a lot more items/features than samples. In such case, probably, FIs cannot be mined directly from data within a reasonable amount of time while our approach with or without parallel implementation can.

APPENDIX

TABLE IV. NUMBER OF PATTERNS EXTRACTED FROM DATASETS.

Data	<i>minsup</i>	#FS	#CS	#G	#N	#FS/#CS	#CS/#N	#N/#G
C2	0.8	8165081	99785	122031	122359	81.83	0.82	1.00
C2	0.85	7525408	95533	116126	116416	78.77	0.82	1.00
C2	0.9	7017040	92087	111297	111564	76.2	0.83	1.00
C2	0.95	6525355	88695	106575	106813	73.57	0.83	1.00
C2	1	6092449	85608	102316	102519	71.17	0.84	1.00
C7	50	2569643 ₉	482902	765450	765449	53.21	0.63	1.00
C7	55	9698268	222253	346029	346028	43.64	0.64	1.00
C7	60	4188627	108428	166918	166917	38.63	0.65	1.00
C7	65	1472818	47491	71875	71874	31.01	0.66	1.00
C7	70	543081	19501	29008	29007	27.85	0.67	1.00
CH	50	900355	369450	372603	372603	2.44	0.99	1.00
CH	60	156551	98392	98418	98418	1.59	1.00	1.00
CH	70	24997	23991	23991	23991	1.04	1.00	1.00
CO	60	2118445 ₄	68349	68349	68349	309.95	1.00	1.00
CO	70	4093971	35875	35875	35875	114.12	1.00	1.00
CO	75	1561212	24346	24346	24346	64.13	1.00	1.00
CO	80	518875	15107	15107	15107	34.35	1.00	1.00
MU	2	2359664 ₉	31767	57728	82483	742.8	0.55	1.43
MU	3	9934877	22229	37972	52165	446.93	0.59	1.37
MU	4	4324745	16565	26984	35597	261.08	0.61	1.32
MU	5	3727905	12854	21160	27801	290.02	0.61	1.31
RE	0.006	975063	504142	532342	542565	1.93	0.95	1.02
RE	0.008	480620	286435	293235	294709	1.68	0.98	1.01
RE	0.01	240852	189077	191265	191650	1.27	0.99	1.00
RE	0.02	67186	65301	65329	65330	1.03	1.00	1.00
RE	0.03	40153	39552	39552	39552	1.02	1.00	1.00
RE	0.04	26925	26666	26666	26666	1.01	1.00	1.00
RE	0.05	19836	19698	19698	19698	1.01	1.00	1.00
T4	0.8	480531	480531	480531	480531	1	1.00	1.00

T4	0.85	432211	432211	432211	432211	1	1.00	1.00
T4	0.9	350323	350323	350323	350323	1	1.00	1.00
T4	0.95	210610	210610	210610	210610	1	1.00	1.00
T4	1	66278	66278	66278	66278	1	1.00	1.00
Note: #FS: the number of FIs; #CS: the number of CFIs; #G: the number of generators; #N: the number of pairs of kernels and extendable sets.								

TABLE V. RUNTIMES OF THE FRAMEWORKS AND OF THE BREAKDOWN PROCESS.

Data	$\text{min}_{su} p$	tCS	tG	tN	tNI	tNNI	tGI	GenIT	dEclat	NUC
C2	0.800	14.8	3.6	2.4	6.1	8.5	14.1	32.5	24.6	23.3
C2	0.850	13.9	3.4	2.3	5.3	7.7	12.8	30.1	22.6	21.6
C2	0.900	13.1	3.3	2.2	5.0	7.2	11.9	28.3	21.2	20.3
C2	0.950	12.5	3.2	2.3	4.6	6.9	11.2	26.9	19.6	19.5
C2	1.000	12.0	3.1	1.9	3.9	5.8	10.3	25.4	18.4	17.8
C7	50.000	192.5	33.3	23.6	16.3	40.0	outM	outM	88.1	232.4
C7	55.000	58.9	14.0	9.2	5.4	14.6	19.4	92.2	33.2	73.5
C7	60.000	20.0	6.3	3.9	2.0	5.9	7.8	34.1	14.6	25.9
C7	65.000	6.5	2.4	1.5	0.7	2.2	2.7	11.6	5.3	8.7
C7	70.000	2.3	0.9	0.4	0.4	0.8	0.9	4.1	2.1	3.1
Ch	50.000	109.3	17.4	5.5	0.6	6.1	3.0	129.6	5.0	115.4
Ch	60.000	14.3	3.9	1.3	0.1	1.5	0.6	18.7	1.1	15.8
Ch	70.000	1.8	0.8	0.3	0.0	0.3	0.1	2.7	0.3	2.1
Co	60.000	35.2	2.4	1.8	16.5	18.3	outM	outM	95.0	53.4
Co	70.000	13.3	1.2	0.9	2.9	3.8	7.3	21.8	20.5	17.1
Co	75.000	7.9	0.8	0.5	0.9	1.4	2.6	11.3	8.7	9.4
Co	80.000	4.7	0.6	0.3	0.3	0.6	0.8	6.1	3.5	5.3
M	2.000	3.8	2.3	1.1	18.3	19.3	outM	outM	71.2	23.1
M	3.000	2.7	1.3	0.4	8.0	8.4	28.2	32.1	30.1	11.1
M	4.000	1.9	0.8	0.5	3.2	3.7	13.8	16.5	13.3	5.6
M	5.000	1.6	0.9	0.3	2.9	3.2	11.1	13.7	11.3	4.8
RT	0.006	101.2	9.8	4.9	0.8	5.7	7.8	118.8	60.6	106.9
RT	0.008	57.2	5.2	2.4	0.3	2.7	1.2	63.6	27.2	59.8
RT	0.010	39.0	3.1	1.5	0.1	1.6	0.5	42.7	16.5	40.6
RT	0.020	17.1	1.1	0.4	0.0	0.5	0.1	18.3	5.7	17.6
RT	0.030	12.0	0.6	0.2	0.0	0.3	0.1	12.8	3.8	12.3
RT	0.040	9.0	0.4	0.2	0.0	0.2	0.1	9.4	2.8	9.2
RT	0.050	7.2	0.4	0.1	0.0	0.2	0.0	7.6	2.2	7.4

T4	0.800	427.7	28.9	7.4	0.3	7.7	1.4	458.0	41.2	435.4
T4	0.850	372.0	24.6	5.9	0.2	6.1	1.3	397.9	35.2	378.2
T4	0.900	289.2	19.3	5.3	0.2	5.5	1.3	309.8	29.9	294.7
T4	0.950	120.8	9.7	2.8	0.1	2.8	0.6	131.1	25.2	123.6
T4	1.000	75.4	1.9	0.6	0.0	0.6	0.1	77.4	21.1	76.0

Note: *tCS*: time to find CFIs using CharmL; *tG*: time to find the generators using Minimal Generator; *tN*: time to generate kernels and extendable sets by NUCLEAR; *tNI*: time to enumerate FIs from kernels and extendable sets by NUCLEAR; *tNNI*: time to find FIs from lattice of CFIs, which is the sum of *tN* and *tNI*; *tGI*: time to generate FIs based on CFIs and generators using GEN_ITEMSETS; *NUC*: time to find FIs from data using CharmL and NUCLEAR; *GenIT*: time to find FIs using CharmL, Minimal Generator, and GEN_ITEMSETS; *dEclat*: time to find FIs from data using dElat. *outM*: out of memory.

ACKNOWLEDGMENT

This work has been partially supported by the Natural Sciences and Engineering Research Council of Canada (NSERC).

The authors would like to thank Tran N.A., Duong V.H., and the co-authors of [14], [15], [20] for providing us the code of GEN_ITEMSETS algorithm for our experiments.

REFERENCES

- [1] Agrawal R., Imielinski T., Swami N, "Mining association rules between sets of items in large databases", in ACM SIGMOD, 1993, pp. 207-216.
- [2] Mai T., Vo B., Nguyen L.T.T. A lattice-based approach for mining high utility association rules. Information Sciences, 399, 2017, pp.81-97.
- [3] Yun U., Lee G., Yoon E, "Efficient high utility pattern mining for establishing manufacturing plans with sliding window control", IEEE Transactions on Industrial Electronics, 64(9), 2017, pp.7239 – 7249.
- [4] Mai T., Nguyen L.T.T, "An Efficient Approach for Mining Closed High Utility Itemsets and Generators. Journal of Information and Telecommunication", 1(3), 2017, pp.193-207.
- [5] Bundit, M., Nunnapp, B., Arnon, R., Athasit, S., Putchong, U., "Parallel association rule mining based on FI-Growth algorithm", in ICPDS'07, 2007, pp. 1-8.
- [6] Lakhil, L., and Stumme, G., "Efficient mining of association rules based on formal concept analysis", in FCA'05, 2005, pp. 180-195.
- [7] Grahne G., Zhu J., "Fast algorithms for frequent itemset mining using FP-Trees", IEEE Transactions on Knowledge and Data Engineering, 17(10), 2005, pp.1347-1362.
- [8] Zaki, M.J. and Hsiao, C.J., "Efficient algorithms for mining closed itemsets and their lattice structure", IEEE Transactions on Knowledge and Data Engineering, 17(4), 2005, pp.462-478.
- [9] Zaki, M.J., "Mining non-redundant association rules". Data Mining and Knowledge Discovery, 9(3), 2004, pp.223-248.
- [10] Sahoo, J., Das, A. K., Goswami, A., "An effective association rule mining scheme using a new generic basis". Knowledge and Information Systems, 43(1), 2015, pp.127-156.
- [11] Negrevergne, B., Termier, A., Mhauth, J., Uno, T., "Discovering Closed Frequent Itemsets on Multicore: Parallelizing Computations and Optimizing Memory Accesses", International Conference on High Performance Computing and Simulation (HPCS), 2010, pp. 521-528.
- [12] Le T., Vo B., "The Lattice-based approaches for mining association rules: a review". WIREs Data Mining and Knowledge Discovery, 6(4), 2016, pp.140-151.
- [13] Vo B., Le B., "Mining traditional association rules using frequent itemsets lattice", in CIE'09, 2009, pp. 1401-1406.
- [14] Tran N.A., Tran C.T., Le H.B., "Structures of association rule set" in ACIIDS'12, 2012, pp. 361-370.
- [15] Truong C.T., Tran N.A., "structure of set of association rules based on concept lattice", in ACIIDS'10, 2010, pp. 217-227.
- [16] Deng Z.H., "DiffNodesets: An efficient structure for fast mining frequent itemsets", Applied Soft Computing, 41, 2016, pp.214-223.
- [17] Deng Z.H., Lv S.L., "PrePost+: An efficient N-lists-based algorithm for mining frequent itemsets via Children-Parent Equivalence pruning", Expert Systems with Applications, 42(13), 2015, pp.5424-5432.
- [18] Vo B., Le T., Coenen F., Hong T.P., "Mining frequent itemsets using the N-list and subsume concepts", International Journal of Machine Learning and Cybernetics, 7(2), 2016, pp.253-265.
- [19] Goethals, B., and Zaki, M., "FIMI '03 Workshop on Frequent Itemset Mining Implementations", 2003, <http://www.cs.rpi.edu/~zaki/PaperDir/FIMI03.pdf>.
- [20] Tran N.A., Duong V.H., Tran C.T., Le H.B., "Efficient algorithms for mining frequent itemsets with constraint", in KSE'11, 2011, pp. 19-25.
- [21] Szathmary, L., Valtchev, P., Napoli, A., Godin, R., "Efficient vertical mining of frequent closures and generators", in Advances in Intelligent Data Analysis VIII, Springer Berlin Heidelberg, 2009, pp. 393-404.
- [22] Anh N. T., Tin C.T., and Bac L.H., "An approach for mining concurrently closed itemsets and generators", in ICCSAMA'13, 2013, pp.355-366.
- [23] Vo B., Hong T.P., Le B., "DBV-Miner: A dynamic Bit-Vector approach for fast mining closed frequent itemsets", Expert Systems with Applications, 39(8), 2012, pp.7196-7206.
- [24] Vo B., Le B., "interestingness measures for association rules: Combination between lattice and hash tables", Expert Systems with Applications, 38(9), 2011, pp.11630-11640.
- [25] Vo B., Le T., Hong T.P., Le B., "An effective approach for maintenance of pre-large-based frequent-itemset lattice in incremental mining", Applied Intelligence, 41(3), 2014, pp.759-775.
- [26] Agrawal R., Shafer J.C., "Parallel mining of association rules". IEEE Transactions on Knowledge and Data Engineering, 8(6), 1996, pp.962-969.
- [27] Han E., Karypis G., and Kumar V., "Scalable parallel data mining for association rules", in ACM SIGMOD'97, 1997, pp. 277-288.

- [28] Zaïane, O.R., El-Hajj, M., and Lu, P., "Fast parallel association rule mining without candidacy generation", in ICDM'01, 2001, pp. 665-668.
- [29] Pasquier N., Taouil R., Bastide Y., Stumme G., and Lakhal L., "Generating a condensed representation for association rules," J. of Intelligent Information Systems, vol. 24, no. 1, 2005, pp. 29-60.
- [30] Ai, D., Pan, H., Li, X., Gao, Y., & He, D., "Association rule mining algorithms on high-dimensional datasets". Artificial Life and Robotics, 23(3), 2018, pp.420-427.
- [31] Fournier-Viger P., Lin J.C.W., Vo B., Truong T.C., Zhang J., Le H.B. "A survey of itemset mining". WIREs Data Mining and Knowledge Discovery, 7(4), 2017, e1207.
- [32] <http://fimi.ua.ac.be/data>.
- [33] http://coron.loria.fr/site/downloads_datasets.php.