

Discovering Periodic high-utility item sets from transactional databases

*Kilaru Gowthami¹, Divvela.Srinivasa Rao²

¹M.Tech.IV Sem.(CSE) Lakireddy Balireddy College of Engineering, Mylavaram

²Sr.Assistant Professor, CSE Department, Lakireddy Balireddy College of Engineering, Mylavaram

Corresponding Author: *Kilaru Gowthami

Abstract: High-utility item set mining is the task of finding high-utility item sets, i.e. sets of things that return a high benefit in a client exchange database. High-utility item sets are helpful, as they give data about profitable set of items purchased by clients to retail store administrators, which can then utilize this data to take strategic marketing decisions. An inherent limitation of customary high-utility item set mining calculations is that they are inappropriate to find repeating client buy conduct, although such conduct is normal all things considered, circumstances (for instance, a client may get a few items consistently, week or month). In this paper, we address this limitation by proposing the task of high-utility item set mining. The objective is to find discover of things that are periodically purchased by clients, create a high benefit. A productive calculation named PHM (Periodic High-utility item set Miner) is proposed to effectively identify all periodic high-utility item sets. Exploratory outcomes demonstrate that the PHM calculation is effective, and can filter a huge amount of non occasional examples to uncover just the desired high-utility item sets.

Keywords: high-utility item set, periodic item set, Average periodicity

Date of Submission: 03-07-2017

Date of acceptance: 24-07-2017

I. Introduction

High-utility item set mining (HUIM) [4, 5, 8–10, 13] is a mainstream information mining task in a great deal of consideration in recent years. It extends the conventional issue of Frequent Item set Mining (FIM) [1]. This last comprises of finding visit item sets. I.e. gatherings of things (item sets) showing up as often as possible in an exchange database [1]. FIM has numerous applications. However, an important limitation of FIM is that it expects that everything can't seem more than once in every exchange and that all things have a similar significance (e.g. weight, unit benefit or value). High-Utility Item set Mining (HUIM) addresses this issue by considering that everything may have non double buy amounts in exchanges also, that everything has a weight (e.g. unit benefit). The objective of HUIM is to find item sets having a high utility (e.g. yielding a high benefit) in an exchange database. In addition, market basket analysis, HUIM has a few other applications, for example, site click stream investigation, and biomedical applications [9, 13]. Mining high-utility item sets is broadly recognized as more difficult than FIM on the grounds that the utility measure utilized as a part of HUIM is not anti- monotonic, i.e. a high utility item set may have supersets or subsets having lower, equivalent or higher utilities [4]. In this manner, methods for reducing the search space in FIM can't be specifically reused in HUIM.

In spite of the fact that few calculations have been proposed for HUIM [4, 5, 8–10, 13], a natural constraint of these calculations is that they are wrong to find repeating client buy conduct, although such behavior is normal in real-life situations. For instance, in a retail location, a few clients may purchase some arrangement of items on around a day by day or week after week premise. Recognizing these buy examples is valuable to better comprehend the conduct of clients and along these lines adjust marketing strategies, for instance by offering particular advancements to cross-advance items, for example, reward or indicates clients who are purchasing a set of items occasionally. In the field of FIM, algorithms have been proposed to find occasional successive examples (PFP) [2, 3, 6, 11, 7, and 12] in an exchange database. Be that as it may, these calculations are insufficient to discover occasional examples that return a high benefit, as they just select examples in light of their recurrence. Consequently, these calculations may locate a huge amount of periodic patterns that produce a low benefit and miss numerous uncommon periodic patterns that return a high benefit.

To address this constraint of past work, this paper proposes the task of high-utility item set mining. The objective is to efficiently find all gatherings of things that are purchased together periodically and produce a high benefit, in a client exchange database. The commitments of this paper are fourfold. In the first place, the idea of periodic patterns utilized as a part of FIM is combined with the idea of HUIs to characterize another kind of examples named high-utility item sets (PHIs), and its properties are studied. Second, novel measures of patterns periodicity named average periodicity and minimum periodicity are introduced with give an adaptable

method for surveying the periodicity of patterns. Third, an efficient calculation named PHM (Periodic High-utility item set Miner) is proposed to effectively find the periodic high-utility item sets. Fourth, a broad exploratory assessment is conveyed to contrast the effectiveness of PHM and the best in class FHM calculation for HUIM. Experimental results demonstrate that the PHM calculation is efficient, and can filter a huge amount of non occasional examples to uncover the desired item sets.

Whatever remains of this paper is composed as takes after. Area 2, 3, 4, 5 and 6 separately presents preliminaries identified with HUIM, related work, the PHM calculation, the trial assessment and the conclusion.

II. Related work

This section reviews related work in high-utility item set mining periodic frequent pattern mining.

2.1 High-utility item set mining

Definition 1 (transaction database). Let I be a chance to be an arrangement of things (images). A exchange database is an arrangement of exchanges $D = \{T_1, T_2, \dots, T_n\}$ with the end goal that for every exchange T_c , $T_c \in I$ and T_c has a special identifier c called its Tid . Each thing $i \in I$ is related with a positive number $p(i)$, called its outside utility (e.g. unit benefit). For every exchange T_c with the end goal that $I \in T_c$, a positive number $q(i, T_c)$ is known as the internal utility of i (e.g. purchase quantity).

Example 1. Consider the database of Fig. 1, which will be utilized as running illustration. This database contains seven exchanges (T_1, T_2, \dots, T_7). Exchange T_3 shows that things a, b, c, d , and e show up in this exchange with an internal utility of separately $1, 5, 1, 3$ and 1 . Fig. 2 demonstrates that the outside utility of these things are separately $5, 2, 1, 2$ and 3 .

Table 1: A transaction database

TID	Transaction
T_1	$(a, 1), (c, 1),$
T_2	$(e, 1)$
T_3	$(a, 1), (b, 5), (c, 1), (d, 3), (e, 1)$
T_4	$(b, 4), (c, 3), (d, 3), (e, 1)$
T_5	$(a, 1), (c, 1), (d, 1)$
T_6	$(a, 2), (c, 6), (e, 2)$
T_7	$(b, 2), (c, 2), (e, 1)$

Table 2: External utility values

Item	a	b	c	d	e
Unit profit	5	2	1	2	3

Definition 2 (utility of a item/item set). The utility of a thing i in a exchange T_c is meant as $u(i, T_c)$ and characterized as $p(i) \times q(i, T_c)$. The utility of an item set X (a gathering of things $X \subseteq I$) in an exchange T_c is signified as $u(X, T_c)$ and characterized as $u(X, T_c) = \sum_{i \in X} u(i, T_c)$. The utility of an item set X (in a database) is meant as $u(X)$ and characterized as $u(X) = \sum_{T_c \in g(X)} u(X, T_c)$, where $g(X)$ is the arrangement of exchanges containing X .

Example 2. The utility of thing a in T_6 is $u(a, T_6) = 5 \times 2 = 10$. The utility of the item set $\{a, c\}$ in T_6 is $u(\{a, c\}, T_6) = u(a, T_6) + u(c, T_6) = 5 \times 2 + 1 \times 6 = 16$. The utility of the item set $\{a, c\}$ (in the database) is $u(\{a, c\}) = u(a) + u(c) = u(a, T_1) + u(a, T_3) + u(a, T_5) + u(a, T_6) + u(c, T_1) + u(c, T_3) + u(c, T_5) + u(c, T_6) = 5 + 5 + 5 + 10 + 1 + 1 + 1 + 6 = 34$.

Definition 3 (high-utility item set mining). The issue of high-utility item set mining is to find all high-utility item sets [4, 5, 8–10, 13]. An item set X is a high-utility item set if its utility $u(X)$ is no not as much as a client determined least utility edge minutil given by the client.

Example 3. If $minutil = 30$, the total arrangement of HUIs is $\{a, c\} : 34, \{a, c, e\} : 31, \{b, c, d\} : 34, \{b, c, d, e\} : 40, \{b, c, e\} : 37, \{b, d\} : 30, \{b, d, e\} : 36$, and $\{b, e\} : 31$, where each HUI is clarified with its utility.

In HUIM, the utility measure is not monotonic or anti-monotonic [10, 13], i.e., an item set may have a utility lower, equivalent or higher than the utility of its subsets. A few HUIM calculations go around this issue by overestimating the utility of item sets utilizing the Transaction-Weighted Utilization (TWU) measure [10, 13], which is against monotonic, and characterized as takes after.

Definition 4 (Transaction weighted utilization). The exchange utility (TU) of an exchange T_c is the entirety of the utility of the considerable number of things in T_c . i.e. $TU(T_c) = \sum_{x \in T_c} u(x, T_c)$. The exchange weighted usage (TWU) of an item set X is characterized as the entirety of the exchange utility of exchanges containing X , i.e. $TWU(X) = \sum_{T_c \in g(X)} TU(T_c)$.

Example 4. The TUs of $T_1, T_2, T_3, T_4, T_5, T_6$ and T_7 are individually $6, 3, 25, 20, 8, 22$ and 9 . The TWU of single things a, b, c, d, e are separately $61, 54, 90, 53$ furthermore, 79 . $TWU(\{c, d\}) = TU(T_3) + TU(T_4) + TU(T_5) = 25 + 20 + 8 = 53$.

Theorem 1 (Pruning search space using the TWU). Give X a chance to be an item set, on the off chance that $TWU(X) < \text{minutil}$, then X and its supersets are low utility. [10]

Algorithms, for example, Two-Phase [10], BAHUI [8], PB [5], and UP Growth+ [13] use the above property to prune the search space. They work in two stages. In the principal stage, they distinguish competitor high utility item sets by figuring their TWUs. In the second stage, they filter the database to figure the correct utility of all applicants found in the primary stage to dispose of low utility item sets. Recently, an alternative algorithm called HUI-Miner [9] was proposed to mine HUIs specifically utilizing a single stage. At that point, depth-first search algorithm FHM [4] was proposed, which amplifies HUI-Miner. In FHM, each item set is related with a structure named utility-list [4, 9]. Utility-records permit computing the utility of an item set rapidly by making join operations with utility-arrangements of shorter examples. Utility-records are characterized as takes after.

Definition 5 (Utility-list). Give a chance to be any aggregate request on things from I . The utility-list of an item set X in a database D is an arrangement of tuples with the end goal that there is a tuple $(tid, iutil, rutil)$ for every exchange $Ttid$ containing X . The $iutil$ component of a tuple is the utility of X in $Ttid$. i.e., $u(X, Ttid)$. The $rutil$ component of a tuple is characterized as $P \in Ttid \wedge \forall x \in X \ u(i, Ttid)$.

Definition 5 (Utility-list). Accept that i is the sequential request. The utility-list of $\{a\}$ is $\{(T1, 5, 1), (T3, 5, 20), (T5, 5, 3), (T6, 10, 12)\}$. The utility-list of $\{d\}$ is $\{(T3, 6, 3), (T4, 6, 3), (T5, 2, 0)\}$. The utility-list of $\{a, d\}$ is $\{(T3, 11, 3), (T5, 7, 0)\}$.

To find HUIs, FHM plays out a single database scan to make utility lists of examples containing single things. At that point, longer patterns are acquired by playing out the join operation of utility-lists of shorter patterns. The join operation for single things is executed as takes after. Consider two things x, y such that $x \subseteq y$, and their utility-records $ul(\{x\})$ and $ul(\{y\})$. The utility-list of $\{x, y\}$ is acquired by making a tuple $(ex.tid, ex.iutil + ey.iutil, ey.rutil)$ for each sets of tuples $ex \in ul(\{x\})$ and $ey \in ul(\{y\})$ with the end goal that $ex.tid = ey.tid$. The join operation for two items sets $P \cup \{x\}$ and $P \cup \{y\}$ with the end goal that $x \subseteq y$ is performed as takes after. Let $ul(P)$, $ul(\{x\})$ and $ul(\{y\})$ be the utility-arrangements of P , $\{x\}$ and $\{y\}$. The utility-list of $P \cup \{x, y\}$ is acquired by making a tuple $(ex.tid, ex.iutil + ey.iutil - ep.iutil, ey.rutil)$ for each arrangement of tuples $ex \in ul(\{x\})$, $ey \in ul(\{y\})$, $ep \in ul(P)$ with the end goal that $ex.tid = ey.tid = ep.tid$. Calculating the utility of an item set utilizing its utility-lists and pruning the search space is done as takes after.

Property 1 (Calculating utility of an item set using its utility-list). The utility of an item set is the entirety of $iutil$ values in its utility-list [9].

Theorem 2 (Pruning search space using utility-lists). Give X a chance to be an item set. Give the augmentations of X a chance to be the item sets that can be acquired by attaching a thing y to X to such an extent that $y \subseteq X, \forall i \in X$. On the off chance that the whole of $iutil$ and $rutil$ values in $ul(X)$ is not as much as minutil , X and its expansions are low utility [9].

FHM is extremely productive. However, an imperative constraint of current HUIM calculations is that they are not intended for finding occasional examples.

2.2 Periodic Frequent Pattern Mining

In the field of FIM, Algorithms have been proposed to find periodical continuous designs (PFP) [2, 3, 6, 11, 7, 12] in an exchange database. Finding PFP has applications in numerous spaces, for example, web mining, bioinformatics, and market basket analysis [12]. The idea of PFP is characterized as takes after [12].

Definition 6 (Periods of an item set). May there be a database $D = \{T1, T2, \dots, Tn\}$ containing n transactions and an item set X . The arrangement of exchanges containing X is signified as $g(X) = \{Tg1, Tg2, \dots, Tgk\}$, where $1 \leq g1 < g2 < \dots < gk \leq n$. Two exchanges $Tx \supset X$ and $Ty \supset X$ are said to be back to back as for X if there does not exist an exchange $Tw \in g(X)$ to such an extent that $x < w < y$. The time of two back to back exchanges Tx and Ty in $g(X)$ is characterized as $pe(Tx, Ty) = (y - x)$, that is the quantity of exchanges between Tx and Ty . The times of an item set X is a list of periods characterized as $ps(X) = \{g1 - g0, g2 - g1, g3 - g2, \dots, gk - gk-1, gk+1 - gk\}$, where $g0$ and $gk+1$ are constants characterized as $g0 = 0$ and $gk+1 = n$. accordingly, $ps(X) = S_{1 \leq z \leq k+1} (gz - gz-1)$.

Example 6. For the item set $\{a, c\}$, The list of exchanges containing $\{a, c\}$ is $g(\{a, c\}) = \{T1, T3, T5, T6\}$. Accordingly, the times of this item set are $ps(\{a, c\}) = \{1, 2, 2, 1, \text{ and } 1\}$.

Definition 7 (Periodic Frequent Pattern). The most extreme periodicity of an item set X is characterized as $\text{maxper}(X) = \max(ps(X))$ [12]. An item set X is a periodic frequent pattern (PFP) if $|g(X)| \geq \text{minsup}$ and $\text{maxper}(X) < \text{maxPer}$, where minsup and maxPer are client characterized thresholds [12].

The principal algorithm for mining PFPs is PFP-Tree [12]. It uses a tree-based also, design development approach for finding PFPs. At that point, the MTKPP calculation [2] was proposed for finding the k most successive PFPs in a database, where k is a client indicated parameter. MTKPP uses a vertical structure to keep up data about item sets in the database. A variety of the PF-Tree algorithm named the ITL-Tree was also introduced [3] with reduce the time for mining PFPs by approximating the periodicity of item sets. Another rough calculation for PFP mining was as of recently proposed [7]. Different expansions of the PF-Tree calculation named MIS-PF-tree [6] and MaxCPF [11] were individually proposed to mine PFPs utilizing various minsup limits, and different minsup and minper limits.

An important limitation of traditional algorithm for PFP mining is that they are lacking to discover periodic patterns that return a high benefit, since they just consider the support (recurrence) of examples. Henceforth, they may locate a huge amount of periodic patterns that return a low benefit and miss numerous uncommon periodical designs that return a high benefit.

III. The PHM algorithm

To address the previously mentioned constraint of HUI and PFP mining calculations, this segment presents the idea of occasional high-utility item sets (PHUIs). The in the first place subsection show novel measures to survey the periodicity of HUIs, while the second subsection presents and effective calculation named PHM (Periodic High-Utility Item set Miner) to find PHUIs productively.

3.1 Measuring the periodicity of high-utility patterns

A drawback of the maximum periodicity measure utilized by most PFP calculations is that an item set is naturally disposed of on the off chance that it has a single period of length more than the max threshold. Consequently, this measure might be seen as well strict. To give a more adaptable method for assessing the periodicity of patterns, the idea of average periodicity is presented in the proposed calculation.

Definition 8 (Average periodicity of an item set). The average periodicity of an item set X is characterized as $avgper(X) = \sum_{g \in ps(X)} |g| / |ps(X)|$.

Example 7. The periods of item sets $\{a, c\}$ and $\{e\}$ are separately $ps(\{a, c\}) = \{1, 2, 2, 1, 1\}$ and $ps(\{e\}) = \{2, 1, 1, 2, 1, 0\}$. The average periodicities of these item sets are separately $avgper(\{a, c\}) = 1.4$ and $avgper(\{e\}) = 1.16$.

Lemma 1 (Relationship between average periodicity and support). Let X be an item set showing up in a database D . An option and proportionate way of computing the average periodicity of X is $avgper(X) = |D| / (|g(X)| + 1)$.

Proof. Let $g(X) = \{Tg_1, Tg_2, \dots, Tg_k\}$ be the arrangement of exchanges containing X , to such an extent that $g_1 < g_2 < \dots < g_k$. By definition, $avgper(X) = \sum_{g \in ps(X)} |g| / |ps(X)|$. To demonstrate that the lemma holds, we have to demonstrate that $\sum_{g \in ps(X)} |g| = |D|$, as follows:

(1) We first demonstrate that $\sum_{g \in ps(X)} |g| = |D|$, as follows:

$$\begin{aligned} \sum_{g \in ps(X)} |g| &= (g_1 - g_0) + (g_2 - g_1) + \dots + (g_k - g_{k-1}) + (g_{k+1} - g_k) \\ &= \sum_{g \in ps(X)} |g| = g_0 + (g_1 - g_0) + (g_2 - g_1) + \dots + (g_k - g_{k-1}) + (g_{k+1} - g_k) \\ &= g_{k+1} - g_0 = |D|. \end{aligned}$$

(2) We then demonstrate that $|ps(X)| = |g(X)| + 1$, as follows:

By definition, $ps(X) = \{g_z \mid 1 \leq z \leq k+1, g_z - g_{z-1}\}$. Therefore, the set $ps(X)$ contain $k+1$ component. Since X shows up in k exchanges, $sup(X) = k$, and in this manner $|ps(X)| = |g(X)| + 1$.

Since (1) and (2) holds, the lemma holds.

The above lemma is important as it gives a productive method for ascertaining the average periodicity of item sets in a database D . The term $|D|$ can be ascertained once, and from there on the average periodicity of any item set X can be acquired by just ascertaining $|g(X)| + 1$, and afterward separating $|D|$ by the outcome. This is more effective than ascertaining the normal periodicity utilizing Definition 8. In addition, this lemma is imperative as it demonstrates that there is a connection between the support utilized as a part of FIM and the average periodicity of a pattern.

Although the average periodicity is valuable as it quantifies what is the ordinary period length of an item set, it should not be utilized as the sole measure for assessing the periodicity of an pattern since it doesn't consider whether an item set has periods that differ generally or not. For instance, the item set $\{b, d\}$ has a normal periodicity of 2.33. However, this is misdirecting since this item set as it were shows up in exchange T_3 and T_4 , and its periods $ps(\{T_3, T_4\}) = \{3, 1, 4\}$ change broadly. Naturally, this example should not be a periodic pattern. To abstain from finding designs having periods that differ broadly, our answer is to join the normal periodicity measure with other periodicity measure(s). The accompanying measures are joined with the normal periodicity to accomplish this objective.

To start with, we characterize the base periodicity of an item set as $\text{minper}(X) = \min(\text{ps}(X))$ to abstain from finding item sets having some brief periods. In any case this measure is not reliable since the first and last time of an item set are separately equivalent to 1 or 0 if the item set individually shows up in the first or the last exchange of the database. For instance, the last time of item set $\{e\}$ is 0, since it shows up in the last exchange (T7), and in this manner its base periodicity is 0. Our answer for this issue is to prohibit the first and last periods of an item set from the figuring of the base periodicity. In addition, if the set of periods is empty accordingly of this prohibition, the base periodicity is characterized as ∞ . In whatever remains of this paper, we consider this definition.

Second, we consider the most extreme periodicity of an item set $\text{maxper}(X)$ as characterized in the past segment. The method of reasoning for utilizing this measure in blend with the average periodicity is that it can abstain from finding periodical designs that don't happen for drawn out stretches of time.

As far as figuring cost, a purpose behind picking the base periodicity, most extreme periodicity and normal periodicity as measure is that they can be figured proficiently for an item set X by filtering the list of exchanges $g(X)$ just once. That is, ascertaining these measures doesn't require storing the set of period's $\text{ps}(X)$ in memory. Alternately, other measure, for example, the standard deviation would require to ascertain all times of an item set in advance. Hence, we characterize the idea of periodic high-utility item sets by considering the least periodicity, greatest periodicity and normal periodicity measures.

Definition 9 (Periodic High-Utility Item sets). Let minutil , minAvg , maxAvg , minPer and maxPer be positive numbers, gave by the client. An item set X is an periodic high-utility item set if and just if $\text{minAvg} \leq \text{avgper}(X) \leq \text{maxAvg}$, $\text{minper}(X) \geq \text{minPer}$, $\text{maxper}(X) \leq \text{maxPer}$, and $u(X) \geq \text{minutil}$.

For instance, if $\text{minutil} = 20$, $\text{minPer} = 1$, $\text{maxPer} = 3$, $\text{minAvg} = 1$, and $\text{maxAvg} = 2$, the total arrangement of PHUIs is appeared in table 3.

Table 3: The set of PHUIs in the running example

Itemset	$u(X)$	$ g(X) $	$\text{minper}(X)$	$\text{maxper}(X)$	$\text{avgper}(X)$
$\{b\}$	22	3	1	3	1.75
$\{b, e\}$	31	3	1	3	1.75
$\{b, c, e\}$	37	3	1	3	1.75
$\{b, c\}$	28	3	1	3	1.75
$\{a\}$	25	4	1	2	1.4
$\{a, c\}$	34	4	1	2	1.4
$\{c, e\}$	27	4	1	3	1.4

To build up a efficient algorithm for mining PHUIs, it is critical to plan proficient pruning techniques. To utilize the periodicity measures for pruning the search space, the accompanying hypotheses are introduced.

Lemma 2 (Monotonicity of the average periodicity). Let X and Y are item sets with the end goal that $X \subset Y$. It takes after that $\text{avgper}(Y) \geq \text{avgper}(X)$.

Proof. The average periodicities of X and Y are separately $\text{avgper}(X) = |D| / (|g(X)| + 1)$ and $\text{avgper}(Y) = |D| / (|g(Y)| + 1)$. Since $X \subset Y$, it takes after that $g(Y) \subseteq g(X)$. Thus, $\text{avgper}(Y) \geq \text{avgper}(X)$.

Lemma 3 (Monotonicity of the minimum periodicity). Let X and Y are item sets with the end goal that $X \subset Y$. It takes after that $\text{minper}(Y) \geq \text{minper}(X)$.

Proof. Since $X \subset Y$, $g(Y) \subseteq g(X)$. If $g(Y) = g(X)$, then X and Y have the same periods, and subsequently $\text{minper}(Y) = \text{minper}(X)$. if $g(Y) \subset g(X)$, then for every exchange $T_x \in g(X) \setminus g(Y)$, the comparing time frames in $\text{ps}(X)$ will be replaced by a bigger period in $\text{ps}(Y)$. Consequently, any period in $\text{ps}(Y)$ can't be littler than a period in $\text{ps}(X)$. Thus, $\text{minper}(Y) \geq \text{minper}(X)$.

Lemma 4 (Monotonicity of the maximum periodicity). Let X and Y are item sets with the end goal that $X \subset Y$. It takes after that $\text{maxper}(Y) \geq \text{maxper}(X)$ [12].

Theorem 3 (Maximum periodicity pruning). Let X be an item set showing up in a database D . X and its supersets are not PHUIs if $\text{maxper}(X) > \text{maxPer}$. Along these lines, if this condition is met, the search space comprising of X what not its supersets can be disposed of.

Proof. By definition, if $\text{maxper}(X) > \text{maxPer}$, X is not a PHUI. By Lemma 4, supersets of X are likewise not PHUIs.

Theorem 4 (Average periodicity pruning). Let X be an item set showing up in a database D . X is not a PHUI and the greater part of its supersets if $\text{avgper}(X) > \text{maxAvg}$, or comparably if $|g(X)| < (|D|/\text{maxAvg}) - 1$. Along these lines, if this condition is met, the inquiry space comprising of X and all its supersets can be disposed of.

Proof. By definition, if $\text{avgper}(X) > \text{maxAvg}$, X is not a PHUI. By Lemma 2, supersets of X are additionally not PHUIs. The pruning condition $\text{avgper}(X) > \text{maxAvg}$ is modified as: $|D| / (|g(X)| + 1) > \text{maxAvg}$. Along these lines, $1 / (|g(X)| + 1) > \text{maxAvg} / |D|$, which can be further reworked as $|g(X)| + 1 < |D| / \text{maxAvg}$, and as $|g(X)| < (|D| / \text{maxAvg}) - 1$.

3.2 The algorithm

The proposed PHM algorithm is a utility-list based calculation, propelled by the FHM calculation [4], where the utility-list of each item set X is clarified with two extra values: $\text{minper}(X)$ and $\text{maxper}(X)$. The primary methodology of PHM (Algorithm 1) takes an exchange database as input, and the minutil , minAvg , maxAvg , minPer and maxPer threshold. The calculation first outputs the database to ascertain $\text{TWU}(\{i\})$, $\text{minper}(\{i\})$, $\text{maxper}(\{i\})$, and $|\text{g}(\{i\})|$ for everything $i \in I$. At that point, the calculation ascertains the esteem $\gamma = (|D|/\text{maxAvg}) - 1$ to be later utilized for pruning item sets utilizing Theorem 4. At that point, the calculation distinguishes the set I^* of all things having a TWU no less than minutil , a most extreme periodicity no more prominent than maxPer , and showing up in no not as much as γ exchanges (different things are overlooked since they can't be a piece of a PHUI by Theorem(1, 3 and 4). The TWU estimations of things are then used to set up an aggregate request on things, which is the request of climbing TWU values (as recommended in [9]). A database scan is then performed. During this database filter, things in exchanges are reordered by the aggregate request, the utility-list of everything $i \in I^*$ is assembled and a structure named EUCS (Estimated Utility Co-Occurrence Structure) is assembled [4]. This last structure is characterized as an arrangement of triples of the frame $(a, b, c) \in I^* \times I^* \times R$. A triple (a, b, c) demonstrates that $\text{TWU}(\{a, b\}) = c$. The EUCS can be actualized as a triangular framework (as appeared in Fig. 1 for the running case), or as a hash map of hash maps where just tuples of the frame (a, b, c) with the end goal that $c \neq 0$ are kept. After the development of the EUCS, the depth-first investigation of item sets begins by calling the recursive strategy Search with the empty item set \emptyset , the arrangement of single things I^* , γ , minutil , minAvg , minPer , maxPer , the EUCS structure, and $|D|$.

Algorithm 1: The PHM algorithm

Input: D: a transaction database,

minutil , minAvg , maxAvg , minPer and maxPer : the thresholds

Output: the set of periodic high-utility item sets

- 1 Scan D once to calculate $\text{TWU}(\{i\})$, $\text{minper}(\{i\})$, $\text{maxper}(\{i\})$, and $|\text{g}(\{i\})|$ for each item $i \in I$;
 - 2 $\gamma \leftarrow (|D|/\text{maxAvg}) - 1$;
 - 3 $I^* \leftarrow$ each item i such that $\text{TWU}(i) \geq \text{minutil}$, $|\text{g}(\{i\})| \geq \gamma$ and $\text{maxper}(\{i\}) \leq \text{maxPer}$;
 - 4 Let σ be the total order of TWU ascending values on I^* ;
 - 5 Scan D to build the utility-list of each item $i \in I^*$ and build the EUCS Structure;
 - 6 Search $(\emptyset, I^*, \gamma, \text{minutil}, \text{minAvg}, \text{minPer}, \text{maxPer}, \text{EUCS}, |D|)$;
-

Item	a	b	c	d
b	25			
c	61	54		
d	33	45	53	
e	47	54	76	45

Fig. 1: The EUCS

Item	a	b	c	d
b	1			
c	4	3		
d	2	2	3	
e	2	3	4	2

Fig. 2: The ESCS

The Search strategy (Algorithm 2) takes as input, an item set P, augmentations of P having the shape Pz implying that Pz was already acquired by affixing thing z to P, γ , minutil , minAvg , minPer , maxPer , the EUCS, and $|D|$. The search method plays out a circle on every expansion Px of P. In this circle, the normal periodicity of Px is gotten by separating $|D|$ by the quantity of components in the utility-list of Px in addition to one (by Lemma 1). At that point, if the normal periodicity of Px is in the $[\text{minAvg}, \text{maxAvg}]$ interval, the aggregate of the util estimations of the utility-list of Px is no not as much as minutil (cf. Property 1), the minimum/maximum periodicity of Px is no less/not more than $\text{minPer}/\text{maxPer}$ as indicated by the qualities put away in its utility-list, then Px is a PHUI and it is yield. At that point, if the aggregate of iutil and rutil values in

the utility-list of P_x are no not as much as $minutil$, the quantity of components in the utility-list of P_x is no not as much as γ , and $maxper(P_x)$ is no more than $maxPer$, it implies that augmentations of P_x ought to be investigated (by Theorem 1, 3 and 4). This is performed by combining P_x with all augmentations P_y of P to such an extent that $y \supset x$ to shape expansions of the frame P_{xy} containing $|P_x| + 1$ things. The utility-list of P_{xy} is then developed by calling the Construct system (calculation 3), to join the utility-arrangements of P , P_x and P_y . This last method is mainly the same as in HUI-Miner [9], with the special case that periods are computed amid utility-list development to acquire $maxPer(P_{xy})$ and $minPer(P_{xy})$ (not appeared). At that point, a recursive call to the Search method with P_{xy} is done to ascertain its utility and investigate its extension(s). The Search strategy begins from single things; recursively investigates the search space of item sets by affixing single things, and just prunes the search space utilizing Theorem 1, 3 and 4. Consequently, it can be easily observed that this technique is right and finish to finding all PHUIs.

Algorithm 2: The Search procedure

Input: P : an itemset, $ExtensionsOfP$: a set of extensions of P , γ , $minutil$, $minAvg$, $minPer$, $maxPer$, the EUCS structure, $|D|$

Output: the set of periodic high-utility itemsets

```

1 foreach itemset  $P_x \in ExtensionsOfP$  do
2    $avgperP_x \leftarrow |D| / (|P_x.utilitylist| + 1)$ ;
3   if  $SUM(P_{xy}.utilitylist.iutils) \geq minutil \wedge$ 
       $minAvg \leq avgperP_x \leq maxAvg \wedge P_x.utilitylist.minp \geq$ 
       $minPer \wedge P_x.utilitylist.maxp \leq maxPer$  then output  $P_x$ ;
4   if  $SUM(P_x.utilitylist.iutils) + SUM(P_x.utilitylist.rutils)$ 
       $avgperP_x \geq \gamma$  and  $P_x.utilitylist.maxp \leq maxPer$  then  $\geq minutil \wedge$ 
5      $ExtensionsOfP_x \leftarrow \emptyset$ ;
6     foreach itemset  $P_y \in ExtensionsOfP$  such that  $y \supset x$  do
7       if  $\exists (x, y, c) \in EUCS$  such that  $c \geq minutil$  then
8          $P_{xy} \leftarrow P_x \cup P_y$ ;
9          $P_{xy}.utilitylist \leftarrow Construct(P, P_x, P_y)$ ;
10         $ExtensionsOfP_x \leftarrow ExtensionsOfP_x \cup \{P_{xy}\}$ ;
11      end
12    end
13    Search( $P_x, ExtensionsOfP_x, \gamma, minutil, minAvg, minPer, maxPer,$ 
          EUCS,  $|D|$ );
14  end
15 end

```

Besides, in the usage of PHM, two extra enhancements are included, which are quickly described next.

Optimization 1. Estimated Average Periodicity Pruning (EAPP). The PHM calculations creates a structure called EUCS to store the TWU of all sets of things happening in the database, and this structure is utilized to prune any item set P_{xy} containing a couple of things $\{x, y\}$ having a TWU lower than $minutil$ (Line 7 of the search methodology). The procedure EAPP is a novel methodology that utilizes a similar thought however prunes item sets utilizing the normal periodicity rather than the utility. During the second database check, a novel structure called ESCS (Estimated Support Co-occurrence Structure) is made to store $|g(\{x, y\})|$ for each match of things $\{x, y\}$ (as appeared in Figure 2). At that point, Line 7 of the search strategy is altered to prune item set P_{xy} if $|g(\{x, y\})|$ is not as much as γ by Theorem 4.

Optimization 2. Abandoning List Construction early (ALC). Another methodology introduced in PHM is with quit developing the utility-list of an item set if a particular condition is met, showing that the item set can't be a PHUI. By Theorem 4, an item set $P \cup x$ can't be a PHUI, if it appears in less than $\gamma = (|D|/\maxAvg) - 1$ exchanges.

Algorithm 3: The Construct procedure

```

Input: P: an itemset, P x: the extension of P with an item x, P y: the
        Extension of P with an item y
Output: the utility-list of P xy
1  UtilityListOfP xy  $\leftarrow \emptyset$ ;
2  foreach tuple ex  $\in$  P x.utilitylist do
3      if  $\exists$  ey  $\in$  P y.utilitylist and ex.tid = ey.tid then
4          if P.utilitylist  $\neq \emptyset$  then
5              Search element e  $\in$  P.utilitylist such that e.tid = ex.tid;
6              exy  $\leftarrow$  (ex.tid, ex.iutil + ey.iutil - e.iutil, ey.rutil);
7          end
8          else
9              exy  $\leftarrow$  (ex.tid, ex.iutil + ey.iutil, ey.rutil);
10         end
11         periodexy  $\leftarrow$  calculate period (exy.tid, UtilityListOfPxy);
12         Update MinPer MaxPer (UtilityListOfP xy, periodexy);
13         UtilityListOfP xy  $\leftarrow$  UtilityListOfP xy  $\cup$  {exy};
14     end
15 end
16 return UtilityListPxy;
    
```

The procedure ALC comprises of adjusting the Construct system (Algorithm 3) as takes after. The main change is to instate a variable max with the value γ in Line 1. The second adjustment is to the accompanying lines, where the utility-list of Pxy is developed by checking if each tuple in the utility-arrangements of P x shows up in the utility-list of P y (Line 3). For each tuple not showing up in P y, the variable max is decremented by 1. In the event that max is littler than γ , the development of the utility-list of P xy can be stopped since $|g(P \cup x)|$ won't be higher than γ . Along these lines P xy is not a PHUI by Theorem 4, and its augmentations can also be overlooked.

IV. Experimental Study

We played out a test study to survey the execution of PHM. The experiment was performed on a PC with a 6th era 64 bit Core i5 processor running Windows 10, and furnished with 12 GB of free RAM. We looked at the execution of the proposed PHM calculation with the state-of-the-art FHM calculation for mining HUIs. All memory estimations were done utilizing the Java API. The experiment was carried on four genuine datasets normally utilized as a part of the HUIM literature: retail, mushroom, chain store and food mart. These datasets have differed attributes and speaks to the fundamental sorts of information commonly experienced, all things considered, situations (dense, sparse and long transactions). Let $|I|$, $|D|$ and A represents to the quantity of exchanges, particular things and normal exchange length of a dataset. Retail is a sparse dataset with various things ($|I| = 16,470$, $|D| = 88,162$, $A = 10, 30$). Mushroom is a dense dataset with long exchanges ($|I| = 119$, $|D| = 8,124$, $A = 23$). Chain store is a dataset that contains a huge number of exchanges ($|I| = 461$, $|D| = 1,112,949$, $A = 7.23$). Food mart is a sparse dataset ($|I| = 1,559$, $|D| = 4,141$, $A = 4.4$). The chain store and food mart datasets are real-life client exchange databases containing real outside and inner utility qualities. The retail and mushroom datasets contains engineered utility qualities, created randomly [9, 13]. The source code of all calculations and datasets can be downloaded from <http://goo.gl/Y6eBdz>.

In the test, PHM was run on each dataset with settled minper and minAvg values, while shifting the minutil limit and the estimations of the maxAvg and maxper parameters. In these examinations, the qualities for the periodicity thresholds have been discovered exactly for each dataset (as they are dataset particular), and were demonstrated the trade-off between the quantity of periodic patterns found and the execution time. Take note of

those outcomes for shifting the minper and minAvg qualities are not demonstrated on the grounds that these parameters have less impact on the patterns found than alternate parameters. From that point, the documentation PHM V-W-X-Y represents the PHM calculation with minper = V, maxper = W, minAvg = X, and maxAVG = Y.

Fig. 3 analyzes the execution times of PHM for different parameter values furthermore, FHM. Fig. 4 thinks about the quantity of PHUIs found by PHM for different parameter values, and the quantity of HUIs found by FHM.

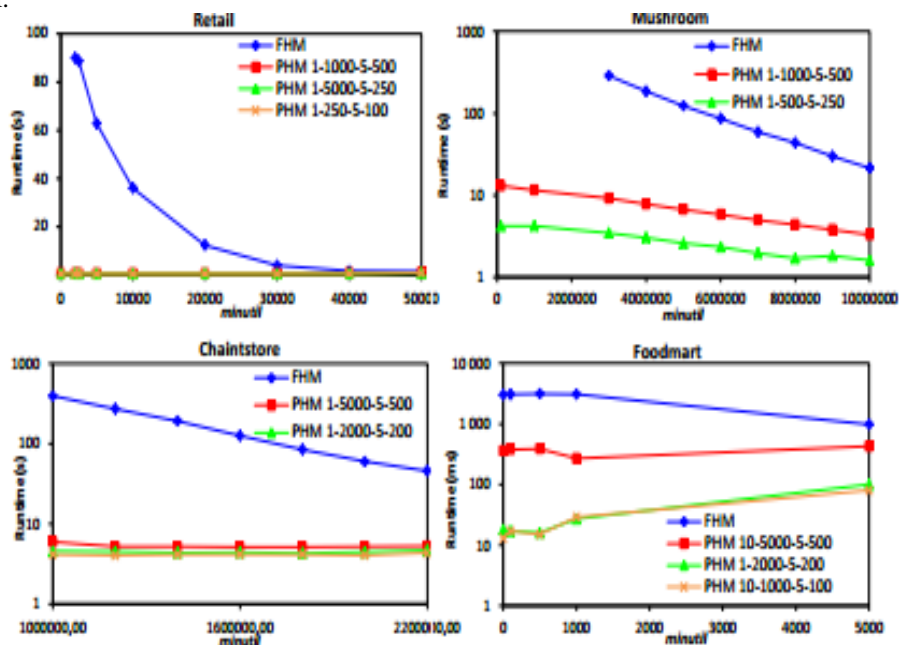


Fig. 3: Execution times

It can first be watched that mining PHUIs utilizing PHM can be significantly speedier than mining HUIs. The explanation behind the great execution of PHM is that it prunes an extensive piece of the search space utilizing its outlined pruning systems in view of the maximum and average periodicity measures. For all datasets, it can be found that a huge amount of HUIs are non periodic, and in this

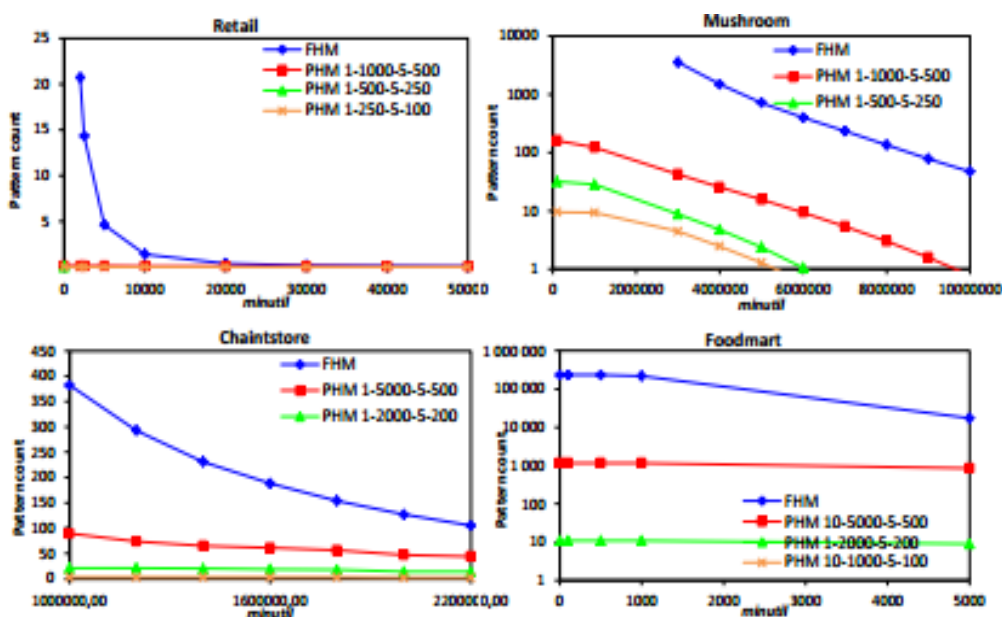


Fig. 4: Number of patterns found

Way pruning non periodic patterns prompts an enormous execution change. For instance, for the least minutil, maxPer and maxAvg values on these datasets, PHM is individually up to 214, 127, 100 and 230 times speedier than FHM. When all is said in done, the progressively the periodicity limits are prohibitive, the more the gap between the runtime of FHM and PHM increments.

A second observation is that the quantity of PHUIs can be significantly less than the quantity of HUIs (see Fig. 4). For instance, on retail, 20,714 HUIs are found for minutil = 2, 000. However, just 110 HUIs are PHUIs for PHM 1-1000-5-500, and 7 for PHM 1-250-5-150. A portion of the patterns found are quite interesting as they contain a few things. For instance, it is found that things with item ids 32, 48 and 39 are occasionally purchased with a normal periodicity of 16.32, a minimum periodicity of 1, and a maximum periodicity of 170. Huge reduction in the quantity of examples is observed on alternate datasets. These results demonstrate that the proposed PHM calculation is valuable as it can filter huge amount of non periodic HUIs experienced in real datasets, and can run quicker.

Memory utilization was additionally looked at, although detailed outcomes are not appeared as a figure because of space confinement. It was watched that PHM can go through to 10 times less memory than FHM depending upon how parameters are set. For example, on chain store and minutil = 1,000,000, FHM and PHM 1-5000-5-500 respectively consumes 1,631 MB and 159 MB of memory.

V. Conclusion

This paper investigated the issue of mining periodic high-utility item sets (PHUIs). A proficient calculation named PHM (Periodic High-utility item set Miner) was proposed to proficiently find PHUIs utilizing novel least and normal periodicity measures. A broad test think about with real datasets has demonstrated that PHM can be more than two requests of size quicker than FHM, also, find more than two requests of greatness less examples by sifting non periodic HUIs. Source code of PHM, FHM and datasets can be downloaded from <http://goo.gl/Y6eBdz>. For future work, we will consider outlining elective calculations to mine PHUIs.

References

- [1]. Agrawal, R., Srikant, R.: Fast algorithms for mining association rules in large databases. In: Proc. Int. Conf. Very Large Databases, pp. 487–499, (1994)
- [2]. Amphawan, K., Lenca, P., Surarerks, and A.: Mining top-k periodic-frequent pattern from transactional databases without support threshold. In: Proc. 3rd Intern. Conf. on Advances in Information Technology, pp. 18–29 (2009)
- [3]. Amphawan, K., Surarerks, A., Lenca, and P.: Mining periodic-frequent item sets with approximate periodicity using interval transaction-ids list tree. In: Proc. 2010 Third Intern. Conf. on Knowledge Discovery and Data Mining, pp. 245–248 (2010)
- [4]. Fournier-Viger, P., Wu, C.-W., Zida, S., Tseng, V. S.: FHM: Faster high-utility item set mining using estimated utility co-occurrence pruning. In: Proc. 21st Intern. Symp. On Methodologies for Intel. Syst., pp. 83–92 (2014)
- [5]. Lan, G. C., Hong, T. P., and Tseng, V. S.: An efficient projection-based indexing approach for mining high utility item sets. Knowl. And Inform. Syst. 38(1), 85–107 (2014)
- [6]. Kiran, R. U., Reddy, P. K.: Mining Rare Periodic-Frequent Patterns Using Multiple Minimum Supports. In: Proc. 15th Intern. Conf. on Management of Data (2009)
- [7]. Uday, U. R., Kitsuregawa, M., Reddy, P. K.: Efficient Discovery of Periodic-Frequent Patterns in Very Large Databases. Journal of Systems and Software, 112, 110–121 (2015)
- [8]. Song, W., Liu, Y., Li, J.: BAHUI: Fast and memory efficient mining of high utility item sets based on bitmap. Intern. Journal of Data Warehousing and Mining. 10(1), 1–15 (2014)
- [9]. Liu, M., Qu, J.: Mining high utility item sets without candidate generation. In: Proc. 22nd ACM Intern. Conf. Info. And Knowl. Management, pp. 55–64 (2012)
- [10]. Liu, Y., and Liao, W., Choudhary, A.: A two-phase algorithm for fast discovery of high utility item sets. In: Proc. 9th Pacific-Asia Conf. on Knowl. Discovery and Data Mining, pp. 689–695 (2005)
- [11]. Surana, A., Kiran, R. U., Reddy, P. K.: An efficient approach to mine periodic frequent patterns in transactional databases. In: Proc. 2011 Quality Issues, Measures of Interestingness and Evaluation of Data Mining Models Workshop, pp. 254–266 (2012)
- [12]. Tanbeer, S. K., Ahmed, C. F., Jeong, B. S., and Lee, Y. K.: Discovering periodic frequent patterns in transactional databases. In: Proc. 13th Pacific-Asia Conference on Knowledge Discovery and Data Mining, pp. 242–253 (2009)
- [13]. Tseng, V. S., Shie, B.-E., Wu, C.-W., Yu, P. S.: Efficient algorithms for mining high utility item sets from transactional databases. IEEE Trans. Knowl. Data Eng. 25(8), 1772–1786 (2013)

IOSR Journal of Computer Engineering (IOSR-JCE) is UGC approved Journal with Sl. No. 5019, Journal no. 49102.

Kilaru Gowthami. "Discovering Periodic high-utility item sets from transactional databases." IOSR Journal of Computer Engineering (IOSR-JCE) 19.4 (2017): 33-42.