# Energetic Data Mining Law To Identify Frequenty Item Sets over Huge Database via IFCM

[1] Pragathi Vulpala
[1] Asst. Professor, Department Of Computer Science &Engineering, TKR College Of Engineering And Technology (TKRCET), Medbowli R.N Reddy, Meerpet 500097

*Abstract -* **Association Rule Mining is one of the most important tasks in data mining industry, especially in terms of huge database handling strategies; this association rule mining plays a vital role to deal with the commercial and non-commercial data. The main process of association rule mining is to identify the frequent items from the itemset. In past analysis there are several methodologies available to identify the frequent items over the inputting itemset. Most of the researchers and developers commonly used Apriori' and Frequent Pattern-Tree algorithms to identify the frequent items over the itemset. Generally Apriori investigates multiple number of iterations over the huge databases to identify the frequent items over the inputting itemset, globally Apriori law follows the candidate introduction process for identifying the frequent items. The approach of Frequent Pattern-Tree algorithm is different from Apriori law, it investigates the database twice without' includes the generation of candidates. In this proposed system, a new algorithm is proposed to improve the time constraint and accuracy levels of identifying the frequent items over the large itemset database, called Intelligent Frequency Count Model [IFCM]. This algorithm of IFCM uses the scanning method different from the past two terminologies such as multiple scanning and twice scanning schemes, instead of these two concepts, the IFCM performs scanning with one time as well as it belongs with the candidate generation process. So, this methodology is also termed as Hybrid Intelligent Frequency Count Model [HIFCM]. For all the proposed logic, clearly shows that the IFCM provides better results compare to Apriori and Frequent Pattern-Tree process.**

*Keywords* — **Frequency Count, Itemset, Intelligent Frequency Count Model, IFCM, Hybrid IFCM, HIFCM.**

## 1. INTRODUCTION

As of late there has been an exponential advance in the new-generation and treatment of electronic data as an ever increasing number of undertakings are modernized/ computerized. Any association or commercial and non-commercial organizations has begun to understand that the data accumulated over years is a critical vital advantage and it additionally perceives that there are potential insights discharged in the substantial measure of information. So it expects procedures to remove the most important data from aggregated information [1]. The data/Data/Information mining gives such strategies to discover helpful shrouded data. Data/Information mining is a gathering of strategies for successful robotized disclosure of beforehand obscure, substantial, novel, significant and reasonable example in vast databases [1][2].

It has as of late pulled in significant consideration from database specialists and analysts on account of its pertinence in a few regions, for example, choice help, advertise technique and monetary gauges [3]. The Data/Information mining methods or undertakings can be for the most part named distinct or prescient. Enlightening/Descriptive mining indicates to the technique in which the basic qualities or regular

properties of the information in the information base are portrayed. The clear methods include assignments like grouping, affiliation and consecutive mining [4]. Prescient Data/Information mining assignments are those that perform deduction on input information to land at shrouded learning and make intriguing and valuable gauge [5]. The prescient mining strategies include undertakings like arrangement, relapse and deviation [4].

Key research subjects or difficulties in Data/Information mining are execution, mining strategy, client collaboration and information assorted variety. So the Data/Information mining calculation and systems must be able and versatile well to the span of information base and their execution times [5]. A standout amongst the most mainstream distinct information mining systems is association rule/administer mining [4]. Since its presentation [7], association rule mining has turned out to be one of the center information mining errands and has pulled in amazing enthusiasm among information mining inquires about and specialists [8]. It is a decent strategy for discovering connections between factors in extensive database. For instance,

$$V[x] = Person[y], Buying[x, "Biscuit"] \rightarrow Buying[x, "Cream"] \quad (1)$$

where, x is a variable and Buying[i, j] is a predicate that expresses that individual i buying item j. This decide indicates that a high level of individuals who buy biscuit likewise purchase margarine [9].

*Associative [P➔Q] = ∑[PŨQ]/N   (2)*
*Confident [P➔Q] = ∑[PŨQ]/∑P   (3)*

The greater part of the calculations are made on conventional calculation of association rule mining [7][10]. Association rule mining can be formally characterized as takes after. Let I = {I1, I2…, In} is an arrangement of items. Any subset of I is called itemset. Give D a chance to be an arrangement of exchanges. Every exchange T is an arrangement of items with the end goal that TΩI. Every exchange has a novel identifier. Let P, Q be an arrangement of items, Association run has the shape PαQ, P∞Q = ∑P.Q, where P is a predecessor and Q is the subsequent of the run the show. It utilizes two factual techniques that control the movement of association govern mining are support and certainty [2]. Right off the bat, it decides visit item set in view of least help check. From that point onward, least certainty is utilized to discover association controls between visit items. The help and certainty can be scientifically spoken to as takes after.

Numerous inquires about have been done in creating productive technique for finding incessant examples [3][6][8][9][10] in the wake of presenting Apriori by Agrawal et al. [7]. Among those techniques, the Apriori [7] and FP-Tree [14] are the most prevalent strategies for finding regular items. FP-tree utilizes distinctive calculation system and structure to discover visit items with less figuring time than Apriori. This system acquaints the TR-FCTM strategy with find noteworthy regular items with less figuring time than FP-tree. In this approach a new terminology Intelligent Frequency Count Model [IFCM] is introduced, which provides better support with the combination of Apriori and FP-Tree methods, so it is called as Hybrid-IFCM.

## II. PROPOSED METHODOLOGY

The proposed approach depends on the perception that databases and its transactions' often-times hold indistinguishable transactions/exchanges. The proposed logic applies Intelligent Frequency Count Model [IFCM], which provides more accuracy and reduce the manipulation time. The method comprises of finding these transactions/exchanges and to swap them with

single transactions/exchange. This can be scientifically characterized as, if an arrangement of indistinguishable transactions/exchanges [Ty1, Ty2,… , Tym] are in a database D at that point it is supplanted by a solitary new transactions/exchange as (TyM, ym) where TyM is the indistinguishable itemset and m is the aggregate number of specific indistinguishable itemset in the database [18]. This diminishes the aggregate transactions/exchanges in database as not exactly or equivalent to [2X-1] transactions/exchanges where, I speaks to add up to the quantity of various things in the present market. So this incredibly lessens registering time of finding incessant itemset. Give it X a chance to be any of an itemset in database D. It expresses that an itemset X of transactions/exchange T is a subset of X[Y⊆I] and an arrangement of such transactions/exchanges shape the database D. So in database D each transactions/exchange itemset X will be a component of [2X-1], where 2X is a control set of I. Power set of I contain every one of the subsets of I that may be as transactions/exchanges itemset in the transactions/exchange database D aside from ϕ. Consequently our calculation utilize one table that is name is Recurrence/Frequency Count Table.

It has two fields, for example, itemset and Recurrence/Frequency check esteem. This table makes passages of Recurrence/Frequency check of each itemset that are identified in transactions/exchange database. The Recurrence/Frequency check of each itemset is the tally of the presence of such itemset in value-based database D. This table is made also, might be kept in memory till the regular itemset are not found [19]. The arrangement of Recurrence/Frequency table include is given Table.1.

*Table-1 Intelligent Frequency Count Model Structure*

| S.No. | Itemset(y) | Intelligent Frequency Count [IFC] |
|---|---|---|
| 1 | I(1) | I(1).FC |
| 2 | I(2) | I(2).FC |
| 3 | I(3) | I(3).FC |
| . | - | - |
| . | - | - |
| . | - | - |
| 2^(x-1) | I(x-1) | I(x-1).FC |

For instance, Let I(i) = {I(1), I(2), I(3)} be the arrangement of things and the distinctive kinds of itemset that can be made from I are {I(1)}, {I(2)}, {I(3)},… , {I(1), I(2), I(3)}. The Recurrence/Frequency check table for above said things with beginning quality is given in Table-2.

*Table-2 Initial Value Assignment of Intelligent Frequency Count Table*

| S.No. | Itemset(y) | Intelligent Frequency Count [IFC] |
|---|---|---|
| 1 | I(1) | 0 |
| 2 | I(1)(2) | 0 |
| 3 | I(1)(2)(3) | 0 |
| 4 | I(1)(2)(3) (4) | 0 |
| 5 | I(1)(2)(3) (4)(5) | 0 |
| 6 | I(1)(2)(3) (4)(5)(6) | 0 |
| 7 | I(1)(2)(3) (4)(5) (6)(7) | 0 |
| 8 | I(1)(2)(3) (4)(5) (6)(7)(8) | 0 |
| 9 | I(1)(2)(3) (4)(5) (6)(7)(8)(9) | 0 |
| 10 | I(1)(2)(3) (4)(5) (6)(7)(8)(9)(10) | 0 |

This proposed work utilizes exchange blending and IFCM to locate the noteworthy continuous things. In the first place it consolidates the comparable exchanges in the database and stores the combined exchanges in the principle memory. Later it peruses the blended exchanges one by one from principle memory and refresh the IFCM correspondingly. To locate the continuous itemset for any limit esteem it filters the IFCM not the database. IFCM has passages of Recurrence/Frequency tally of all itemset however not the aggregate help tally of that itemset. The Recurrence/Frequency check of each itemset is the tally of the immediate presence of such itemset in value-based database D. The aggregate tally of specific itemset X is ascertained by looking at whether it a subset of all its greater itemset in the IFCM. On the off chance that aggregate Recurrence/Frequency check of specific itemset is more prominent than or equivalent to ST (Support Threshold/Edge) at that point the itemset is incorporated into the FI (Frequent Itemset).

_____
Algorithm: Intelligent Frequency Count Model
_____
Inputs: Database Db and ST
Output: Frequent Itemset (FI)

Step-1: Begin
Step-2: Construct ST;

Step-3: Manipulate the raw database D record by record until it reaches the end of record
Step-4: Update the ST
Step-5: Count the number of different items involved in the data base;
Step-6: }
Step-7: Construct the FCT;
Step-8: Manipulate ST and update corresponding entry in FCT;
Step-9: FI = { $\phi$ };
Step-10: For ( i=1; i<2I; i++)
Step-11: {
Step-12: $TCX_i$ = 0;
Step-13: For (j=i+1; j< 2I; j++)
Step-14: {
Step-15: If $X_i \subseteq X_j$;
Step-16: {
Step-17: $TCX_i = TCX_i + FCT.FC (j)$;
Step-18: }
Step-19: }
Step-20: if ( $TCx_i >= ST$)
Step-21: {
Step-22: FI = FI U $X_i$;
Step-23: }
Step-24: }
Step-25: End For

_____

## III. LITERATURE SURVEY

In the year of 2006, the author "G.K. Gupta", proposed a paper titled "Introduction to Data Mining with Case Studies" in that he summarized some important details such as: the field of information digging gives procedures to robotized disclosure of important data from the gathered information of modernized operations of undertakings. This framework offers a reasonable and thorough prologue to the two information mining hypothesis and practice. It is composed principally as an investigation for the students of software engineering, administration, PC applications, and data innovation. The proposed framework guarantees that the students comprehend real information mining systems regardless of whether they don't have a solid scientific foundation. The procedures incorporate information pre-preparing, affiliation run mining, managed characterization, group investigation, web information mining, internet searcher inquiry mining, information warehousing and OLAP. To upgrade the comprehension of the ideas presented, and to indicate how the strategies depicted in the framework are utilized as a part of training, every section is trailed by

maybe a couple contextual analyses that have been distributed in academic diaries. Most contextual analyses manage genuine business issues (for instance, advertising, online business, CRM). Concentrate the contextual analyses gives the peruser a more prominent knowledge into the information mining methods.

In the year of 2012, the authors "Saravanan Suba and Christopher T", proposed a paper titled "A Study on Milestones of Association Rule Mining Algorithms in Large Databases" in that they summarized some important details such as: information mining helps in doing mechanized extraction and producing prescient data from huge measure of information. The affiliation govern mining is one of the essential zone of research in Data mining. The Association govern mining recognizes the helpful affiliations or relationship among huge arrangement of information things. In this paper, we give the imperative ideas of Association administer mining and existing calculations and their viability and disadvantages. The references gave in this paper secured the fundamental hypothetical issues and controlling the specialist in a fascinating exploration bearing that presently can't seem to be found.

In the year of 2006, the authors "Ya-Han Hu and Yen-Liang Chen", proposed a paper titled "Mining Association Rules with Multiple Minimum Supports: A New Mining Algorithm and A Support Tuning Mechanism" in that they summarized some important details such as: mining affiliation rules with various least backings is a vital speculation of the affiliation manage mining issue, which was as of late proposed by Liu et al. Rather than setting a solitary help limit for all things, they enable clients to indicate numerous base backings to mirror the natures of the things, and an Apriori-based calculation, named MSapriori, is produced to mine all regular itemsets. In this paper, we consider a similar issue yet with two extra changes. To start with, we propose a FP-tree-like structure, MIS-tree, to store the pivotal data about continuous examples. As needs be, a productive MIS-tree-based calculation, called the CFP-development calculation, is created for mining all continuous itemsets. Second, since every thing can have its own base help, it is exceptionally troublesome for clients to set the suitable limits for all things at once. By and by, clients need to tune things' backings and run the mining calculation over and over until the point that a tasteful end is come to. To accelerate this tedious tuning process, a productive calculation which can keep up the MIS-tree structure without rescanning database is proposed. Tests on both

manufactured and genuine datasets demonstrate that our calculations are considerably more effective and versatile than the past calculation.

## IV. EXPERIMENTAL RESULTS

To contemplate the execution of this proposed calculation, it has been completed a few trials. The intel core™ i5-2450m CPU @2.5GHZ, 4.0GB RAM, 64 bit windows 7 working framework and NetBeans IDE 8.0.2 were utilized to lead the test. The manufactured Data set of 100,1000, 2000 and 4000 with 5 things were made to contrast this proposed IFCM and Apriori and FPtree as far as time required to discover critical regular itemsets from given datasets. The main trial thinks about the execution time devoured of Apriori, FP-tree and proposed calculation by applying the above said four gatherings of datasets. The Table.8 indicates execution time to find noteworthy regular things created by Apriori, FP-tree and IFCM for four gatherings of informational collections with 20% least bolster edge.

The following table illustrates the performance of IFCM compare to other algorithms in terms of execution time scenario.

*Table-3: Execution Time Analysis of Different Algorithms with Different Datasets*

| Iterations | Apriori | FP-Tree | IFCM |
|---|---|---|---|
| 10 | 4 | 8 | 2 |
| 100 | 12 | 18 | 8 |
| 1000 | 22 | 19 | 12 |
| 10000 | 58 | 78 | 49 |
| 100000 | 86 | 150 | 80 |
| 1000000 | 155 | 200 | 112 |
| 10000000 | 359 | 246 | 225 |

It is watched that the execution time is diminished straightly from Apriori to FP-tree and FP-tree to IFCM and the contrast expands increasingly as the quantity of exchanges increments. The Fig.1 exhibits the execution of Apriori, FP-tree what's more, IFCM as indicated by the execution time of every calculation for given four gatherings of exchanges. It is effectively watched that the IFCM outflanks than Apriori and FP-tree.
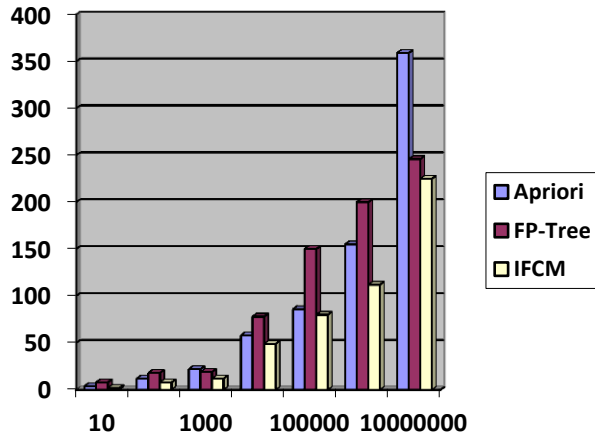
*Fig.1 Execution Time Analysis of Different Algorithms with Different Datasets*

## V. CONCLUSION

Frequent Pattern Mining calculations are critical to discover intriguing examples in huge information base. This IFCM is planned utilizing exchange combining, finding direct Frequency/recurrence check and aggregate Frequency/recurrence mean an itemsets. The trial comes about demonstrate that the proposed philosophy devours less time for exchange filtering and applicant age while it is contrasted and Apriori and FP-tree. So the IFCM outflanks than Apriori and FP-tree for level versatility of exchanges for artificially produced datasets. The vertical versatility and constant informational index ought to be connected in future to enhance its productivity further.

## REFERENCES

[1] G.K. Gupta, "Introduction to Data Mining with Case Studies", Prentice-Hall of India Pvt. Limited, 2006.

[2] Saravanan Suba and Christopher T, "A Study on Milestones of Association Rule Mining Algorithms in Large Databases", International Journal of Computer Applications, Vol. 47, No. 3, pp. 12-19, 2012.

[3] Ya-Han Hu and Yen-Liang Chen, "Mining Association Rules with Multiple Minimum Supports: A New Mining Algorithm and A Support Tuning Mechanism", Decision Support Systems, Vol. 42, No. 1, pp. 1-24, 2006.

[4] S. Shankar and T. Purusothaman, "Utility Sentient Frequent Itemset Mining and Association Rule Mining: A Literature Survey and Comparative Study", International Journal of Soft Computing Applications, No. 4, pp. 81-95, 2009.

[5] N.P. Gopalan and B. Sivaselvan, "Data Mining Techniques and Trends", PHI Learning, 2009.

[6] Ashoka Savasere, Edward Omiecinski and Shamkant B. Navathe, "An Efficient Algorithm for Mining Association Rules in Large Databases", Proceedings of the 21st International Conference on Very Large Data Bases, pp. 432-444, 1995.

[7] R. Agrawal, T. Imielinski and A. Swami, "Mining Association Rules Between Sets of Items in Large Databases", Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 207-216, 1993.

[8] M.J. Zaki and C.J. Hsiao, "CHARM: An Efficient Algorithm for Closed Association Rule Mining", Technical Report 99-10, Department of Computer Science, Rensselaer Polytechnic Institute, 1999.

[9] Yin-Ling Cheung and A.W.-C. Fu, "Mining Frequent Itemsets without Support Threshold: with and without Item Constraints", IEEE Transactions on Knowledge and Data Engineering, Vol. 16, No. 9, pp. 1052-1069, 2004.

[10] R. Agrawal and R. Srikant, "Fast Algorithms for Mining Association Rules", Proceedings of 20th International Conference on Very Large Data Bases, pp. 487-499, 1994.

[11] Jong Soo Park, Ming-Syan Chen and Philip S. Yu, "Using a Hash- Based Method with Transaction Trimming and Database Scan Reduction for Mining Association Rules", IEEE Transactions on Knowledge and Data Engineering, Vol. 9, No. 5, pp. 813-825, 1997.

[12] H. Toivonen, "Sampling Large Databases for Association Rules", Proceedings of the 22th International Conference on Very Large Data Bases, pp. 134-145, 1996.

[13] Jong Soo Park, Ming-Syan Chen and Philip S. Yu, "An Effective Hash Based Algorithm for Mining Association Rules", Proceedings of the 1995 ACM

SIGMOD International Conference on Management of Data, pp. 175-186, 1995.

[14] Jiawei Han, Jian Pei and Yiwen Yin, "Mining Frequent Patterns without Candidate Generation", Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 1-12, 2000.

[15] Mohammed Al-Maolegi1, Bassam Arkok, "An Improved Apriori Algorithm for Association Rules", International Journal on Natural Language Computing, Vol. 3, No. 1, pp. 21-29, 2014.

[16] R. Srikant and R. Agrawal, "Mining Generalized Association Rules", Future Generation Computer Systems, Vol. 13, No. 2-3, pp. 161-180, 1997.

[17] S. Sunil Kumar, S. Shyam Karanth, K.C. Akshay, Ananth Prabhu and M. Bharathraj Kumar, "Improved Apriori Algorithm based on Bottom Up Approach using Probability and Matrix", International Journal of Computer Science Issues, Vol. 9, No. 2, pp. 232-246, 2012.

[18] Z. Souleymane, P.F. Viger, J.C.-W. Lin, C.-W. Wu and V.S. Tseng, "EFIM: A Highly Efficient Algorithm for High-Utility Itemset Mining", Proceedings of 14th Mexican International Conference on Artificial Intelligence, pp. 530-546, 2015.

[19] R. Ahirwal, N.K. Kori and Y.K. Jain, "Improved Data Mining Approach to Find Frequent Itemset using Support Count Table", International Journal of Emerging Trends & Technology in Computer Science, Vol. 1, No. 2, pp. 195-201, 2012.

Pragathi Vulpala received his B.Tech degree in CSE in Sree visvesvaraya college of engineering and technology (svits), Mahabubangar in the year 2007 and M.Tech received in CSE in Sri indu college of engineering, hyd in the year 2015 her interest areas in data mining.