

Survey on Privacy Preservation of Item sets Mining in data mining

Ms. Suman Ahlawat¹, Dr. Anoop Sharma², Dr. Aman Jain³

1. Research Scholar in Singhania University, Jhunjhunu
2. Dean Dept. of CS and IT, Singhania University, Jhunjhunu
3. Professor in Maharishi Arvind Institute of Science and Management, Jaipur

Abstract—Data mining denotes to knowledge mining or extracting from huge amounts of data. Discovering association rules are at the heart of data mining. Mining of suggestion instructions among items in huge database of trades conventions has been accepted such as a main space of database investigation. Extracting meaningful information plays an important role in the mining process. More accurate data can give better result. Solitude protection of gainful items is to be required. We have presented the simple data mining; efficacy mining, unusual item fixed mining and repeated item fixed mining. A brief numerous algorithm overview and methods defined in various research papers has been provided in this paper.

Keywords—Data Mining, Association Rule Mining, Itemset Mining, Utility Mining, Privacy Preservation, k-anonymity, Anonymization, data mining

I. INTRODUCTION

Data mining is the procedure of determining the exhaustive data around the big amount of data which is put away in data warehouses and data sources. Data mining is the data novelty scheme from the big extent of data kept in numerous databases. Here the knowledge belongs to the valuable information which can be used further computation. The simple goal of the data mining is to mine greater-level invisible data as of raw information profusion. Data mining has been recycled in multiple areas of the data. Data mining can be observed as an algorithmic process that proceeds data as per input and produces several designs, for instance instructions of the organization, item sets, and rules of association, or summaries, as output.

Association Rule Mining (ARM) is a well-designed way that recognizes repeated itemsets as of datasets and produces suggestion instructions via supposing that all substances have the similar implication and occurrence of incidence without seeing their convenience [1]. However in a number of actual-world uses like trade advertising, medical diagnosis, client separation, etc., efficacy of itemsets is established on rate, revenue or profits. Efficacy Mining goals to

classify itemsets using maximum benefits through seeing revenue, amount, rate or further user references [2].

Data Mining contains an algorithmic procedure, which proceeds preprocessed input information and abstracts designs. Several methods occur, like association rule mining, organization, clustering, etc. A significant and broadly used data mining procedure is the unearthing of suggestion instructions. Relationship rule excavating aims at determining recurrent itemsets from market carrier data and producing suggestion instructions. Maximum association rule mining algorithms indirectly study the benefits of the itemsets to be the same [3]. A utility is a value attached to an item depending on its evaluation, e.g. if coke has supported 20 and profit of 2%, cookies may have support 10 but with a profit of 20%.

II. ASSOCIATION RULE MINING

There are various methods sufficient these purposes of data mining. Mining Suggestions are one of the methods convoluted in the procedure. These instructions can be successfully used to expose unknown associations, creating outcomes that can deliver a base for estimating and judgment creating. The unique problem addressed by association rule mining was to discover an association among trades of many products as of the study of a huge set of information [4].

Association rule mining (ARM) is the procedure of producing instructions built on the association among the set of items that the consumers buying. Of dawn, data mining detectives have improved upon the superiority of connotation rule mining for occupational growth over assimilating issues such as charge (utility), size of items retailed (weight) and revenue. The instructions quarried lacking seeing efficacy principles (revenue border) will top to a credible harm of gainful instructions.

Current work largesse an Apriori-based isolated rare element set (recurrent itemset) procedure. It exacts the boundary by shortening connections. To address the tasks confronted via remaining work, a solitude conserving FP-

growth (PFP-growth) algorithm, which contains preprocessing phase and mining steps, is planned. Now the preprocessing phase, the database is transmuted to perimeter the length of communications. To apply like a limit, lengthy connections must be split end alternatively reduced. i.e., uncertainly a contract has further items than the bound, it is distributed into various subsets and assurance that every subset is in the limit. To reserve other occurred data in subsets, a graph-placed method is suggested to expose the association of items inside trades and use like association to escort the excruciating procedure. In the mining stage, established on the specified converted database and a user-described inception, recurrent itemsets were revealed. In spite of the possible advantages of transaction splitting, it may bring frequency information loss. Runtime calculation method is used to offset such information loss. In specific, set the loud sustenance of an itemset in the database renovated by contract excruciating, 1st assessment its real provision in the converted database, and formerly other calculate its real provision in the unique database. In calculation, using averaging the descending closure assets (that is, any supersets of an infrequent itemset are infrequent), dynamic reduction method was used.

In common, the suggestion instruction is an appearance of the form $X \Rightarrow Y$, where X is predecessor and Y is resultant. Suggestion instruction displays how many times Y has followed in case that X has now followed reliant on the provision and sureness value. Provision: It is the possibility of an item or item sets in the certain transactional database:

$$\text{Provision}(X) = n(X) / n$$

Where n is the complete number of connections in the database and n(X) is the number of connections that encloses the item set X. So, provision ($X \Rightarrow Y$) = provision (XUY). Assurance: It is a provisional possibility, for an suggestion instruction $X \Rightarrow Y$ and definite as per

$$\text{Assurance}(X \Rightarrow Y) = \text{provision}(XUY) / \text{provision}(X)$$

Recurrent itemset: Let A be a set of items, T be the contract database and σ be the user identified minimum support. An itemset X in A (that is X is a subgroup of A) is assumed to be a numerous item fixed in T with deference to σ , if providing $(X)_{\tau} \geq \sigma$. Mining suggestion instructions can be ruined down into the resulting 2 sub-problems:

1. Creating all itemsets that have provision greater than, or equal to, the user identified least provision. i.e., creating all huge itemsets.
2. Creating all the instructions that have least sureness. We can produce the suggestion instruction using more than 1 number of resultant items is produced through the resulting process:

- a. Discover the instruction in which number of consequence = 1.
- b. For the given rules $p(x \rightarrow y)$ and $p(x \rightarrow z)$, the rule $p(x \rightarrow yz)$ is generated by the intersection of both the association rules and get a new rule $p(x \rightarrow yz) = p(xy) / p(x)$.

III. ASSOCIATION RULE MINING APPROACHES FOR ITEMSET MINING

A. Utility Mining

In the data mining association rule mining approaches consider an items utility through transaction set presence. As we know frequent item set mining is used to indicate the frequent items. But we can't say if any item set which have sold frequently will make a profit. Maybe those item sets which are less frequent or rare item set can make more profit than frequent item set. One of the most stimulating tasks of data mining is the highest utility item sets mining efficiently. Identification sets of item with the high utilities is known as Utility Mining. Utility can be dignified in relations of profit, cost or other different user Preferences expression. Such as, a computer system may be more gainful than a telephone in profit terms.

For example- if in a mobile shop, 100 mobile sets of nokia worth rupees -2000/- are sold frequently, but at the same time in another shop a iPhone sold in 60,000/- rarely so its cleared that if any item which sold frequently but with less prices and at the same time another item which sold rarely can make more profit. Utility is amount of an itemset how gainful or beneficial X is. Item set X utility, that is, $u(X)$, which is the abstract of all itemset utilities X in enclosing X all the connections. An itemset X is called an itemset of great utility supposing $u(X)$ greater than or equal to the min_utility , where min_utility is a user definite beginning of minimum utility. High-utility itemset mining objective is to define every that itemsets enclosing utility greater or equal to the user- definite least efficacy beginning.

In the mining of utility based the term utility refers to the user preference quantitative representation, i.e. an itemset utility value is the itemset important measurement in the consumer's perspective. For e.g. if an analyst of sales concludes in few retail research requirements to discover out which itemsets in stores earn revenue of maximum sales for the stores user will describe the any itemset utility as monetary profit that store earns through selling all itemset units. Now note that predictor of trades is not involved in the several connections that itemset hold, then the user is one troubled around the profits created composed concluded each operation comprising the element set. In practice the itemset utility value can be page-rank, profit, popularity, measure of few aesthetic aspect, for example, design or beauty or few other different processes of customer's reference.

The conventional Association standard mining methodologies consider the utility of the items by its presence in the exchange (transaction) set. The recurrence of itemset is not adequate to mirror (reflect) the real utility of an itemset. For instance, the business administrator may not be occupied with continuous itemsets that don't create significant benefit. Recently, one of the most difficult information mining undertakings has been the mining of high utility itemsets productively [5]. Differentiating evidence of the itemsets using HU (high utilities) is known as UM (Utility Mining). The utility can be measured as far as expense, benefit or different articulations of client's inclinations.

Information mining is the procedure of uncovering nontrivial, previously unknown and conceivably helpful data from huge (large) databases. An important role in multiple data mining challenges, like recurrent design mining, weighted recurrent design mining, and high efficacy design mining show for discerning beneficial designs secreted in a database. Between them, incessant design mining is an important. Mining high efficacy item sets as of databases mentions (refers) to determining the itemsets with great profits. Now, the consequence of itemset efficacy is interestingness, consequence, or usefulness of an item for users. The utility of an item in an exchange (transaction) DB comprises of two perspectives:

- 1) The significance of particular items, which is known as external utility, and
- 2) The significance of items in communications, which is known as internal utility.

Utility mining of an itemset is considered as per the outcome of its outside efficacy and its inside efficacy. An itemset is called as a great efficacy mining itemset in case that its utility is no reduce than a client resolution minimum utility edge; else, it is named a little-utility element groups. Removal high utility element sets from databases is a serious responsibility has an extensive variability of operations, for ex., site instant stream analysis occupational improvement in sequence hypermarkets, irritated publicity in trade sites, online e-trade management, and movable business atmosphere repositioning, and even determining authoritative strategies in biomedical customs.

Table I: Transaction Database

| Transaction Id | X | Y | Z |
|----------------|---|---|---|
| Tr1 | 1 | 1 | 1 |
| Tr2 | 2 | 1 | 0 |
| Tr3 | 3 | 0 | 2 |
| Tr4 | 1 | 2 | 1 |
| Tr5 | 0 | 1 | 0 |
| Tr6 | 5 | 3 | 4 |
| Tr7 | 2 | 2 | 0 |

| | | | |
|------|---|---|---|
| Tr8 | 3 | 1 | 1 |
| Tr9 | 4 | 1 | 1 |
| Tr10 | 2 | 0 | 2 |

Table II: Unit Profit Associated With Items

| Item Name | Profit |
|-----------|--------|
| X | 3 |
| Y | 10 |
| Z | 8 |

B. Frequent Itemset mining

Frequent itemsets [6] are the sets of item that present frequently in the any database transactions. Recurrent element set excavating, basic purpose is to find out every transaction data set item groups. Mining of frequent itemset perform a significant role in the practice and theory of numerous significant tasks of data mining, for example rule of the mining association, emerging pattern, long patterns. It has applied in the telecommunications field, census analysis and analysis of text. Frequent criterion is expressed in itemsets support value terms. The itemset support value is the transaction percentage that include the itemset after that the support value will be compared with predefined threshold value, which was user generated. If support is equal or greater than the minimum threshold value than those values will be further processed for 2k mining of the frequent pattern, those which not succeed the least beginning will be unwanted.

C. Rare Itemset mining

Itemset that do not occur frequently in the database, Or we can say infrequent items in the database. Rare circumstances justify specific considerations since they signify algorithms of data mining main difficulties.

Rare itemsets finding, and in rare suggestion instructions originating order from rare itemsets, may be generally appreciated in medicine and biology. Suppose an expert in biology is involved to find out the cardiovascular diseases (CVD) cause for a particular medical records database. A repeated itemset for instance "{prominent cholesterol level, CVD}" may be validate hypothesis that these two altered items are repeatedly connected, prominent to possible interpretation "people containing a high cholesterol level are at high CVD risk". Another different hand, point that "{vegetarian, CVD}" is a rare itemset might be authenticated that 2 altered itemsets suggestion is relatively extraordinary, important to the conceivable understanding "vegetarian individuals are at a CVD small threat". Moreover, the itemsets {CVD} and {vegetarian} can be both different frequent, while the itemset {CVD, vegetarian} is rare.

The next example is occupied from the pharma covigilance field, i.e., a pharmacology dedicated detection field, survey and adverse drug effects study. Deliver an opposing drug things database, rare itemset mining allows a official connecting drugs method using opposing effects, that is, finding cases where a drug had fatal or undesired effects on patients. In this technique, a repeated association as “{drug} ∪ {A}”, where “{A}” is an itemset describing a desirable effect kind, means that this suggestion describes an predictable and acting right way for a drug. Overdivergence, a rare itemset for instance “{drug} ∪ {B}” may be inferred as the point that “{B}” describes an unusual technique of drug acting, probably leading to an undesirable effect.

So that this search can be fulfilled by identifying rare item set in the database. so in these type of condition rare item set convert more appropriate than regular item set.

In this item we current an example of occasional and non-current item-set removal. Input data is made up of a database of transaction, and every transaction is recognized through an ID and is made up of a set of items. In the actual world, transactions can be observed through a customer as a basket bought until the set period of time (day, week, month, etc.). Every basket is made up of a set of items that are bought consecutively. In Table 1 we signify an intellectual database, which is represented by D, where the letter of the alphabet is examined an item. Looking at the transaction database like that obtainable in Table 1, our aim is to discover 2 types of set of items, also known as item-set. The 1st type is made up of those item-sets that are present in most of the 2nd contract, and the 2nd type is made up of those item sets that are not in any contract and are made up of the maximum items equivalent to the biggest cardinality contracts. The number of items set in the database is known as item-set provision. Our case is equal to supreme provision 3 [7].

Table III. Transaction database

| ID | Transaction |
|-----|-------------|
| Tr1 | {a,b,c,d} |
| Tr2 | {b,d} |
| Tr3 | {a,b,c,e} |
| Tr4 | {c,d,e} |
| Tr5 | {a,b,c} |

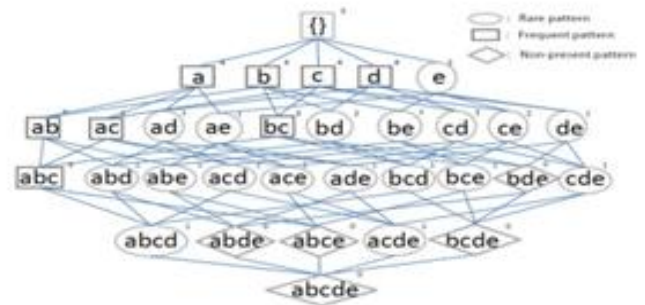


Figure 1: Lattice representing a hierarchically ordered space of item-sets and their frequencies

The set of all item-sets that can be produced as of the contract database is accessible in Figure 1 using a diagram of the subset lattice for 5 items by the related occurrences in the database. In the lattice every level is made up of item-sets consuming the equal length. The highest component in the lattice is the blank set.

IV. LITERATURE SURVEY

Luca Cagliero and Paolo Garza [8] proposed a paper in which the discovering the rare issue and weighted itemsets was handled. i.e., the IWI (infrequent weighted itemset) mining problem. 2 new quality processes are suggested to the drive IWI mining process. Furthermore, two different algorithms that achieve IWI and Minimal IWI mining efficiencies, driven through proposed measures, were presented.

Younghee Kim et al. [9] proposed an efficient algorithm named weighted Support Frequent itemsets (WSFI) was proposed which normalized weight mine over the streams of data, along with that a original tree structure as well suggested which is known as the WSFP-Tree (weighted support FP-tree), that stores compacted serious information around repeated itemsets. The suggested WSFP Tree is a protracted FP-tree built data structure. It is an extended prefix-tree structure to store compressed, critical knowledge about the frequent patterns. The estimation demonstrates that the WSFI-mine outperforms the DSM-FI and THUI-Mine in mining frequent itemsets over the data streams.

G.C.Lan et al. [10] proposed a novel pattern type, known Rare Utility Itemsets, which consider not only individual profits and quantities but also usual current periods and items branches in a multidatabase atmosphere. An original method of mining called as the 2-Phase Algorithm for Mining Rare Efficacy Itemsets in various Databases (TP-RUI-MD) was suggested to efficiently see rare efficacy itemsets. The 2-Phase Algorithm for Mining Rare Efficacy Itemsets in Various Databases algorithm is planned to discover rare-utility itemsets environment. The 1st

one is that we suggested a original itemset type called rare-utility itemset in a multi-database environment.

Hua. Fu. Li et al. [11] Proposed two effective one pass algorithm, which known as MHUI-TID and MHUI-BIT, for mining high utility itemsets from information streams inside of the exchange sliding window. These two distinctive successful thing learning representation and an amplified lexicographical tree-based rundown information structure is created to expand the mining high utility thing sets proficiency.

David j. haglin et al. [12] suggested minimal infrequent itemsets (MINIT) discovery process which was the 1st algorithm created particularly for classifying minimal infrequent itemset (MIIs). The computational period compulsory on the four dataset recommends a connection amongst the amount of MIIs and the volume of calculation necessary. The insignificant occasional itemset problematic is NP-complete.

J. Hu et al. [13] classify high utility item groupings. In transaction to the traditional suggestion instruction and repeated item mining procedures, the goal of the algorithm is to discover data sections, definite using the few items (instructions) sets, which satisfy several conditions to an actual assessment to crack it via specific partition trees, called as high profit partition trees and considered the various splitting schemes performance.

H. Yao et al. [14] suggested the efficacy problem built mining is to find the itemsets that are important agreeing to their efficacy values. In this paper a priori assets and unreliable restraint assets are not valid to the efficacy based itemset mining issue. As an outcome, mathematical itemset utility value properties were analyzed.

V.S. Tseng et al. [15] suggest a novel method, specifically Temporal High Utility Itemsets (*THUI*)-Mine, for the mining of temporal great efficacy itemsets as of data streams excellently and efficiently. For our best information, from data streams. Novel *THUI*-Mine influence is that it can capably temporal great efficacy itemsets classify over creating like that the presentation. Hence, the determining process every window can be succeeded capably using restricted memory space, fewer applicant itemsets and time of CPU I/O. This meets the critical needs on efficiency of time and space for mining data streams.

Liu et al. [16] proposed two different stage algorithm for high utility itemsets discover. In 1st phase, a model relates “transaction-weighted descendant closure assets” to advance the applicant documentation on the search space. In another stage, one additional database scan is the high utility item sets identify performed.

V. COMPARATIVE ANALYSIS OF PRIVACY PRESERVING TECHNIQUES

| Method | Advantage | Disadvantage | Approach |
|---------------|--|---|--|
| k-anonymity | It reduces the granularity of data representation. This granularity is condensed adequately that any specified record maps on partially k further records in the information. | The method is disposed to several types of attacks especially after background information is accessible to the attacker. The adversary can use an association between one or more identifier attributes with the sensitive attribute in order to narrow down possible values of the sensitive field more. | k-anonymous method |
| Randomization | Data is altered using totaling noise to the unique data. Credentials of data openly is not probable. The novel record values cannot be simply estimated after the inaccurate data. It is relatively simple, and does not require knowledge of the distribution of other records in the data | The method on its own is weak and does not offer complete reliability, hence it is used in combination with other algorithms. The quality of data is disturbed and the procedure is irreversible. Reconstruction leads to the leakage of Privacy, which relates to the possible risks | Additive Perturbation Perturbation by random projection technique |
| Encryption | The method groups the data | It involves complex | Integer dividing |

| | | | |
|-------------|---|--|----------------------|
| | into various classes and the encryption is based on the key values generated within each class. Since the key is not a constant private or public key, the method provides a greater amount of protection. | mathematical computations. | uilt encryption |
| Cryptograph | Isolated events can mutually calculate any function of their inputs, lacking illumination any further data. It covers all data apart from for the selected yield of the function | There may exist Ruined events, who select their inputs freely of the truthful events' inputs. This asset is critical in a closed mart | Unconscious transfer |

[4] Farah Hanna AL-Zawaidah and Yosef Hasan Jbara, "An Improved Algorithm for Mining Association Rules in Large Databases", World of Computer Science and Information Technology Journal (WCSIT) ISSN: 2221-0741 Vol. 1, No. 7, 311-316, 2011.

[5] Mengchi Liu, Junfeng Qu, "mining high utility itemsets without candidate generation", Cikm '12 Proceedings Of The 21st Acm International Conference On Information And Knowledge Management, Pages 55-64.

[6] Endu Duneja and A.K. Sachan, "A Survey on Frequent Itemset Mining with Association Rules", International Journal of Computer Applications (0975 – 8887) Volume 46– No.23, May 2012.

[7] Mehdi Adda , Lei Wu , Sharon White , Yi Feng, "PATTERN DETECTION WITH RARE ITEM-SET MINING", International Journal on Soft Computing, Artificial Intelligence and Applications (IJSCAD), Vol.1, No.1, August 2012.

[8] Luca Cagliero and Paolo Garza, "Infrequent Weighted Itemset Mining Using Frequent Pattern Growth", in IEEE Transactions on Knowledge and Data Engineering 26(4):903-915 · April 2014

[9] Younghee Kim, Wonyoung Kim and Ungmo Kim, "Mining Frequent Itemsets with Normalized Weight in Continuous Data Streams", Journal of Information Processing Systems 6(1):79-90 · March 2010

[10] G.C.Lan, T.P.Hong and V.S. Tseng, "A Novel Algorithm for Mining Rare-Utility Itemsets in a Multi Database Environment", The 26th Workshop on Combinatorial Mathematics and Computation Theory, 2009, pp. 293-297.

[11] H.F.Li, H.Y. Huang , Y.Cheng Chen, y. Liu, S.Lee, "Fast and memory efficient mining of high utility itemsets in data streams", in :Eighth International Conference of Data Mining 2008.

[12] D. J. Haglin and A.M. Manning, "On minimal infrequent itemset mining, in DMIN, 2007, pp. 141-147.

[13] Jianying Hu, Aleksandra Mojsilovic, "High-utility pattern mining: A method for discovery of high-utility item sets", Pattern Recognition 40 (2007) 3317–3324.

[14] Hong Yao, Howard J. Hamilton, Liqiang Geng, "A Unified Framework for Utility Based Measures for Mining itemsets", In Proc. of the ACM Intel. Conf. on Utility-Based Data Mining Workshop (UBDM), pp. 28–37, 2006.

[15] Vincent S. Tseng, Chun-Jung Chu, Tyne Liang, "Efficient Mining of Temporal High Utility Itemsets from Data streams", Proceedings of Second International Workshop on Utility-Based Data Mining, August 20, 2006.

[16] Ying Liu, Wei-keng Liao and Alok N. Choudhary, "A Two-Phase Algorithm for Fast Discovery of High Utility Itemsets", in Lecture Notes in Computer Science 3518:689-695 · May 2005.

V. CONCLUSION

PPDM (Privacy preserving data mining) is a novel time of study in data mining. Its eventual objective is to progress effective algorithms that agree one to abstract applicable information after huge data amounts, though check difficult data as of disclosure or inference.

Utility mining discovers each itemsets whose utility values are equal or higher than a user identified threshold in a transaction database. But, the itemset utility value does not justify the "descendent closure assets". i.e., a greater efficacy itemset subset may not be a greater efficacy itemset. The utility mining task is in limiting the applicant set scope and simplifying the efficacy computing calculation. Therefore, consideration was rewarded on privacy preserving utility mining (PPUM) and suggested limited algorithms for it.

REFERENCES

[1] Abhijit Raorane, R.V.Kulkarni, "Data Mining Techniques: A Source For Consumer Behavior Analysis", International Journal of Database Management Systems, Vol.3, No.3, Aug. 2011, pp.45-56.

[2] Sudip Bhattacharya1 , Deepty Dubey, "High Utility Itemset Mining", International Journal of Emerging Technology and Advanced Engineering, ISSN 2250-2459, Volume 2, Issue 8, August 2012.

[3] H. Yao and H. J. Hamilton, "Mining itemset utilities from transaction databases," Data and Knowledge Engineering, vol. 59, pp. 603-626 2006.