# DATA MINING WITH PATTERNS & E-COMMERCE

## Pranisha Katkam*1, Harshali Patil*2

*1Student, MET Institute Of Computer Science, Bandra(W) Mumbai, India.

*2Associate Professor, MET Institute Of Computer Science, Bandra(W) Mumbai, India.
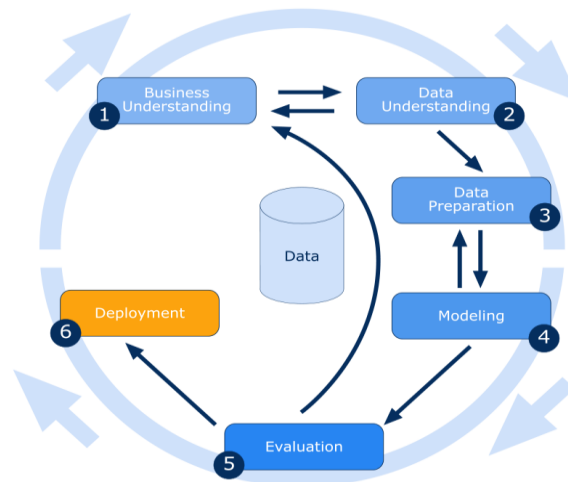
## ABSTRACT

Data mining could be a method of extraction of helpful data and patterns from vast knowledge. it's additionally referred to as information discovery method, information mining from knowledge, information extraction or knowledge /pattern analysis.

**Keywords***: Data Mining, E-commerce, Frequent Patterns, Non-Frequent Pattern, Mining Algorithm**.**

## I.    INTRODUCTION

Data mining could be a logical method that help search through great deal of knowledge or Information so as to search out helpful data. The aim of this approach is to search out patterns that were antecedently unknown. Once these patterns are found, the patterns are used to make certain decisions for development of their businesses.  Three steps concerned are

- Exploration
- Pattern identification
- Deployment



**Fig (1)** data mining process

Data Mining Applications

- Financial Analysis. The banking and finance business depends on high-quality, reliable knowledge.
- Telecommunication business. increasing and growing at a quick pace, particularly with the arrival of the web.
- Intrusion Detection.
- Retail business.

## II.    DATA MINING WITH E-COMMERCE

Having an enormous quantity of information, build some issues for detection of hidden relationships among numerous attributes of information and between many snapshots of information over a time period. These unseen patterns have huge possibility in predictions and personalization in e-commerce.

Ex: Basket analysis, Sales statement, Fraud Detection etc

In business world there's associate degree abundance of obtainable data and an excellent use of it. data should be organized by database tools and data warehouses. Data Mining may be outlined because the art of extracting non-obvious, helpful data from giant databases. This rising field brings a set of powerful techniques that are

relevant for firms to focus their efforts in taking advantage of their knowledge. Data mining tools generate new data for call manufacturers from terribly giant databases. Mechanisms are abstractions, aggregations, summarizations, and characterizations of information. The results of applying refined modeling techniques from the various fields of statistics, artificial intelligence, database management and computer graphics.

**Customer identification**

It may be ascertained that customers drive the revenues of any organization. exploit new customers, delighting and retentive existing customers, and predicting client behavior can improve the provision of merchandise and services and hence the profits. thus the goal of any data processing exercise in e-commerce is to boost processes that contribute to delivering value to the client.

**Web Personalization**

Depending on the goals of the analysis, e-commerce data need to be reworked and aggregated at totally different levels of abstraction. E-commerce knowledge are additional assessed as usage data, content data, structure data, and user data. Usage data consist of details of user sessions and page views. The content data throughout an internet site are unit the gathering of objects and relationships, that are sent to the user. the data compromise mixtures of text material and pictures. Data Mining Applied to Retail E-Commerce Kohavi et al have tried a sensible application of data mining in retail ecommerce data. They share their experience in terms of lessons that they learnt [1]. They classify the important problems in studies, into 2 groups: business and technology. we have a tendency to summarize their findings on the technical problems here. Collection knowledge at the correct level of abstraction is extremely necessary. internet server logs were meant for debugging the server software package from the start. hence they convey little helpful information on customer-related transactions. Approaches together with sessioning the internet logs may yield higher outcomes. Another method would be having the applied server itself log the user connected activities. this can be definitely be richer in linguistics than to the state-less internet logs, and is simpler to keep up compared to state-full internet logs.

Usefulness of UI and users input in e-commerce.

Designing program UI forms must think about the data mining problems in mind. For instance, disabling default values on numerous necessary attributes like Gender, Marital status, Employment standing, etc., can end in richer knowledge collected for demographical analysis. The users ought to be created to enter these values, since it had been found by Kohavi et al that many users left the default values untouched [1].

Certain necessary implementation parameters in retail e-commerce sites just like the automatic time outs of user sessions because of perceived inactivity at the user finish, need to be based not strictly on data processing algorithms, however on the relative importance of the users to the organization. It shouldn't turn out that large clients created to lose their shopping carts due to the time outs that were fastened supported a data mining of the application logs [2].

Generating logs for many million transactions may be an expensive exercise. It may be wise to generate acceptable logs by conducting sampling, as done in statistical quality control. however, sampling might not capture rare events, and in some cases like in advertisement referral based mostly compensations, the information capture could also be obligatory. Techniques that need to be in place try this sampling in an intelligent fashion.

Auditing of knowledge procured for mining, from data warehouses, is obligatory. This is due to the actual fact that data warehouse might need collated data from other diverse systems with a high probability of data being duplicated or lost throughout the ETL operations.

Mining knowledge at the correct level of roughness is important. Otherwise, the results from the data mining exercise might not be correct [2].

## III.     TECHNIQUES AND ALGORITHMS

**PREDICTIVE**

**Classification** is that the most generally applied processing technique that employs a bunch of pre-classified examples to develop a model which is able to classify the population of records at large. Fraud detection and

credit risk applications are significantly well matched to the current sort of analysis. This approach often times employs call tree or neural network-based classification algorithms.

Types of classification models:

- Classification by decision tree induction
- Bayesian Classification
- Neural Networks

**Regression** is analogous to classification, therein it's another dominant sort of supervised learning and is helpful for prognostic analysis. They dissent therein classification is employed for predictions of knowledge with distinct finite categories, whereas regression is employed for predicting continuous numeric information.

As a sort of supervised learning, training/testing information is a vital construct in regression also. regression toward the mean may be a common sort of regression "mining."

What is regression helpful for? Like classification, the potential is limitless. a number of specific examples include: Predicting home costs, as homes tend to be priced on the financial time, as critical being categorical, Trend estimation, within the fitting of trend lines, to applied mathematics info estimation of health connected indicators, like lifetime.

it's simply that the best-fit line isn't linear, and it, instead, takes another kind.This will be noted as curve-fitting, however it's basically no totally different than regression toward the mean and fitting straight lines, although the strategies used for estimation are going to be totally different.

**Frequent pattern mining** could be a idea that is used for a long period of time, while to explain a side information mining that several would argue is that the basis of the term data mining: taking a collection of information and applying applied math ways to search out attention-grabbing and previously-unknown patterns among aforementioned set of information.

we have an inclination for attempting to classify instances or perform instance clustering; we have a tendency to merely need to be told patterns of subsets that emerge at intervals a dataset and across instances, which of them emerge often, that things area unit associated, and that things correlate with others. it is easy to ascertain why the higher than terms become conflated.

Frequent pattern mining is most closely glorious with market basket analysis, that's the identification of subsets of finite superset of product that are unit purchased at some level of every absolute and ratio.

This idea will be generalized on the far side the acquisition of items but, the underlying principle of item subsets remains unchanged

**DESCRIPTIVE**

**Clustering** may be same as identification of comparable categories of objects. By using Clustering techniques, we are able to additionally establish dense and distributed regions in object space and might discover overall distribution pattern and correlations among knowledge attributes. Classification approach may be used for effective suggests that of characteristic teams or categories of object however it becomes expensive therefore cluster may be used as preprocessing approach for attribute set choice and classification.

 Types of Clustering Methods

- Partitioning method
- Grid-based method
- Model-based method

## IV.    DATA MINING WITH FREQUENT PATTERNS

The information obtained with facilitate of data mining tools may be used for resolution advanced issues like detection of fraud identification to boost client shopping for behavior. As most of the users don't seem to be professionally trained to research the patterns of the information, data processing tool in such cases resolve the matter to spot patterns better decision making. The support of associate degree item larger than the required user-defined threshold price, the item is taken into account as frequent item set within the data.

The association rule mining helps in identification of the rule with the interest for decision making and market analysis. the requirement of rule mining becomes vital in each sector. the provision and high spatiality of

information becomes a drag for locating the principles. Therefore, the big databases with numerous techniques for straightforward handling then that extraction of frequent patterns may be done simply. Assigning maximum and minimum occurrence for the item set itself filters out the items with certain threshold value. Hence the frequent items set needs to be preserved for finding optimal items set.

Rule mining is completed by assignment fuzzy value to the data items that have multiple support values that obtains minimum variety of non-frequent items. This approach is accomplished by the application of residual trees that mines the info at every level with given threshold value. It generates items that are non-frequent however whose all subsets are frequent. This approach extracts the minimum variety of related non-frequent item sets for analyzing great amount of data [3].

Different frequent pattern mining algorithms

- Apriori Algorithm
- Rapid Association Rule Mining (RARM)
- Equivalence CLAss Transformation (ECLAT)
- Frequent Pattern (FP) Growth Algorithm

**TABLE :** COMPARISON OF DIFFERENT ALGORITHMS[4]

|   |   | **Apriori** | **RARM** | **ECLAT** | **FP-Growth** |
|---|---|---|---|---|---|
| 1 | **Technique** | Breadth first search & Apriori property | Depth first search on SOTrieIT to generate 1-Itemset & 2-Itemset. | Depth first Search & Intersection of transaction ids to generate candidate itemset | Divide and conquer |
| 2 | **Time** | Execution time is considerable as time is consumed in scanning database for each candidate item set generation | Less execution time as compared to Apriori Algorithm and FP Growth algorithm | Execution time is less then apriori algorithm | Less time as compared to Apriori algorithm |
| 3 | **Drawback** | Too many Candidate Itemset.Too many passes over database. Requires large memory space. | Difficult to use in Interactive system mining, Difficult to use in incremental Mining | It requires the virtual memory to perform the transformation | FP-Tree is expensive to build ,Consumes more memory. |
| 4 | **Advantage** | Use large itemset property.Easy to Implement. | No candidate generation. Speeds up the process for generating candidate 1-Itemset & 2-Itemset. | No need to scan database each time a candidate Itemset is generated as support count information will be obtained from previous Itemset. | Database is scanned only two times. No candidate generation |
| 5 | **Data Format** | Horizontal | Horizontal | Vertical | Horizontal |
| 6 | **Storage** | Array | Tree | Array | Tree |

| | Structure | | | | |
|---|---|---|---|---|---|

Algorithm for Non-frequent Pattern

The discovery of non-frequent items from the set of things that square measure but the required threshold value have gained a great deal of interest in current days. For getting the occasional itemset from the weighted transactional data set a weighted support live has been thought-about for data pruning. The min-support threshold and max-support threshold are acquired from the weighted transactional database. The min-support threshold means the minimum value of the item existing within the database. The max-support threshold means the utmost value of prevalence of the item within the transaction. The minimum support value is outlined because the determined frequency wherever every transaction item support is the chosen weighted maximum or minimum operate. the same weighted dataset is that the union of all equivalent transactional set related to every weighted transaction. The rule initiates with pruning of the items with the utmost threshold price. The mining procedure is comparable to the FP-growth rule that performs projection on itemsets. The tree construction and recursive continues for mining process.

The equivalent dataset is generated for insertion to the FP-tree with the weighted allotted to every data item. However, to cut back the increasing time quality the FP-tree pruning with weighted threshold square measure applied for pruning the information set that discards the items at starting of the tree construction. because the tree-construction is completed. it is initiated for getting the itemsets and corresponding minimum patterns.

In association rule mining, every itemset contains a bound prevalence frequency that might be termed because the weight of the itemset. the burden might be positive, negative or null. Mining of such weighted transactional datasets for locating frequent patterns square measure known as weighted itemsets. within the state of art of the sporadic itemset mining algorithms, the power of taking the tiny frequent itemset into thought is negligible.

The non-frequent mining of the item sets could be a FP-Growth formula finding the non-frequent items from the given set of frequent things. The formula performs the mining method in 2 steps: 1st by scheming the non-frequent things by pruning the transactional DB given the user-defined support values. The mining method is analogous to the FP-Tree approach wherever the burden is related to every item within the transactions. Apart from sorting the items supported their support value, the formula types the items with their associated weight issue. associate item in FP-Tree is pruned if it seems solely in these tree methods from the foundation to a leaf node characterised by a weighted support value bigger than a planned threshold value. The pruning method continues till all the nodes square measure encountered, finally leading to fascinating non-frequent item sets [3].

## V.　CONCLUSION

Data mining has importance regarding finding the patterns, improving the experience of the client in e-commerce sites as well as to use different algorithms and techniques for frequent and non-frequent patterns. Comparison of Apriori, FP-Tree, RARM algorithms are shown and used for finding non-frequent patterns, some drawbacks and advantages are also listed. As more algorithms and techniques are developed and implemented we can use them to find much more advanced and rare patterns using data mining in e-commerce as well as databases.

## VI.　REFERENCES

[1] R. Kohavi, "Lessons and Challenges from Mining Retail E-Commerce Data, " 2004.

[2] Hamid Rastegari, Mohd Noor Md. Sap, " DATA MINING AND E-COMMERCE: METHODS,APPLICATIONS, AND CHALLENGES" 2008.

[3] Karamjit Kaur, Rajeev Bedi, R.C.Gangwar, " Comparative Analysis Of Non-Frequent Pattern Mining Approach"2015.

[4] Shamila Nasreen, Muhammad Awais Azam, Khurram Shehzad, Usman Naeem,

[5] Mustansar Ali Ghazanfar ,"Frequent Pattern Mining Algorithms for Finding Associated Frequent Patterns for Data Streams: A Survey" 2014.

[6] IdhebaMohamad Ali O. Swesi, Azuraliza Abu Bakar, AnisSuhailis Abdul Kadir, "Mining Positive and Negative Association Rules from Interesting Frequent and Infrequent Itemsets", 9th International Conference on Fuzzy Systems and Knowledge Discovery, 2012.

[7]     WeiminOuyang, "Mining Positive and Negative Fuzzy Association Rules with Multiple Minimum Supports", International Conference on Systems and Informatics, 2012.

[8]     R. Prabamanieswari, "A Combined Approach for Mining Fuzzy Frequent Itemset", International Journal of Computer Applications (0975 – 8887), 2013.

[9]     K.Suriya Prabha and R.Lawrance, "Mining Fuzzy Frequent itemset using Compact Frequent Pattern (CFP)tree Algorithm", International Conference on Computing and Control Engineering (ICCCE 2012), 12 & 13 April, 2012.