# Enhanced Data Using Positive Negative Association Mining

Komal Parmar[1], Asst. Prof. Ravi Shukla[2]

*M.E. Computer Engineer, Silver Oak College of Engineering & Technology, Gujarat, India [1]*
*Assistant Professor, Computer/IT Department, Silver Oak College of Engineering & Technology, Gujarat, India [2]*

## ABSTRACT

*Knowledge discovery is the process of analysing data from different perspectives and summarizing it into useful information. Association rule mining is a data mining process used widely in traditional databases to find the positive association rules. Association rules are created by analysing data for frequent patterns and by using the criteria support and confidence to identify the most important relationships. However, there are some other challenging rule mining topics like negative association rule mining. In this research, a rule mining approach has been proposed that provides efficient and secure solution using positive and negative association rule computation on Asynchronous JavaScript and XML (AJAX) data. By using AJAX, we get the faster and secure retrieval of data. Whenever you given the input dataset it will retrieve percentage value of positive negative association rule and attributes of datasets that positively or negatively associated.*

**Keywords:** *Association rule mining, Database, AJAX, Apriori Algorithm, Frequent Pattern-Growth (FPGrowth) Algorithm, Dynamic positive negative mining association algorithm.*

## 1. INTRODUCTION

Due to the event of database technology and systems in recent years, the importance of the data mining has been enhanced in the sectors as well as for business domains like selling, financing and telecommunication. Data processing can be done by using techniques such as data mining. A few examples are, a shop keeper asks all the customers to fill in feedbacks for data processing, hotel management offering the likelihood of online reservation additionally forms information in the event that it requires visitor names, the dates of their stay and their visa number. Association rule mining is one of the main sector data mining techniques that mines the data and finds frequent patterns or associations in large data sets.

The prospective of presented research origins from the association rule mining is that a data mining process is being used widely in databases to find the positive as well as the negative association rules. Association rules are created by analyzing the data from the recurring patterns and by using such criteria and measures; it helps in identifying the most important relationships in the data sets. In this research, a data mining methodology has been derived and proposed that provides an effective, efficient as well as a safe way from the security point of view solution using Positive and Negative association rule computations on (AJAX) Asynchronous JavaScript and XML datasets. By using AJAX, we get the search result in the form of semantic data.

Association rules provides certain guidelines those if/then announcements that find a connection between irrelevant information in a social database or other data storehouse. A case of an association standard would be "If a client purchases twelve eggs, he is 80% prone to likewise buy milk." An association rule is divided into two divisions, (i) antecedent (if) and (ii) consequent (then). An antecedent is an item found in the data. A consequent is an item or an individual data set that is found in along with the antecedent. Association rules are created by analysing and experimenting data for frequent patterns that are found in the datasets and using those matched data and criteria support it is being further processed to identify the important relationships among the data of the entire datasets. Support is considered to an indication on how frequently the data items appear in the entire database. Confidence indicates the

count of the statements i.e. if/then to be found true in the relationships. In data mining, association rules are useful for analysing and predicting customer behaviour.

Moreover, programmers uses the association rules of the data mining to build and execute programs that are capable of learning the machine language. Such programs that are build using machine learning processes is considered to be a type of artificial intelligence (AI) that seeks to build programs with the ability to become more efficient without being explicitly programmed. [1] To identify the most important relations in a dataset using the association rule mining process it makes use of combined analyses of the frequent dataset patterns & statements (if/then) and also support & criteria to identify the relations. Association rule mining is a data mining process that is extensively used in mordern as well as traditional databases to find the association rules, specifically positive association rules. However, there are some other topics which are more challenging and tricky in rule mining methods like negative association rule mining. [2]

Basically all traditional association rule mining algorithms were invented and proposed to find positive associations between data set items. By referring to positive associations we correlate an association between items that exists in transactions i.e. items that are purchased. Moreover, in addition to the positive association, there is another important type of association rule mining which is negative association rule mining process which can provide important data to the company and help in planning future implementation. [3] Interestingly, there are very few researchers that have made their focus and direction on negative association rules due its difficulty in identifying the rules. Association rule mining is a data mining method that allows finding patterns that are frequent and found too often in the transaction in large data sets. The following are the characteristics (accuracy, complete, flexible, reliable, relevant, retrievable and verified) of the patterns that are examined in the mining rules which are preferred for the information to be valuable [4]. An important are where the data mining methods can be applied is in the field of data processing as well as data processing and outsourcing which is provided and beneficial in various industries, for example BPO, Leads and E-commerce. Data Processing Outsource Services range from: Data conversion, Data entry, Word processing, Image processing, Forms processing, Survey processing, Database management, and Script processing. [5]

Data mining on AJAX data stream is another challenging research area. AJAX is extensively used to transmission and storage of the data. Data mining using AJAX is an approach which is considered to be important from the interaction point of view between the client i.e. web browser and server. The structured data in AJAX Web pages cannot be extracted easily due to its asynchronous loading. In this research we have proposed a rule mining method that mines the data using the positive and negative association rules on AJAX data streams which works as a service and also efficiently as well as securely.

## 2. BACKGROUND THEORY

Samet and Taflan (2012) proposed an algorithm named positive negative rule mining (PNRMXS) on XML streaming in database as a service which is based on FP-Growth approach. The processes in algorithm (PNRMXS) take place at two sides, i.e. client side (web browser) and server side (web server). At the client side, some pre-conditioning is done on the pre-defined dataset. At the service provider side, the mining takes place [6]. This algorithm is divided into three steps as shown in the below diagram:
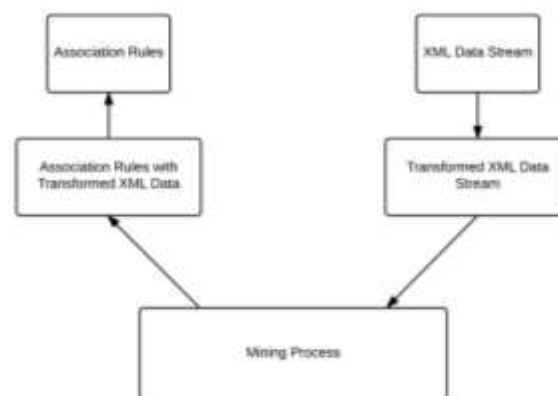


**Figure 2.1** Flowcharts of PNRMXS

One is information change where XML information stream is changed to level record for mining, which utilizes security of balanced mapping which gives protection to the information the changed things are with arbitrary number generator that makes speculating the first thing content troublesome. Second step is mining process, where in creators use milestone windows information handling model. The change that they made to the FP Growth calculation is to filter the information stream just once. Information stream is handled piece by square. Every square contains an altered number of exchanges. The negative guidelines are mined from the current positive principles. Finding substantial and adequately huge number of negative associations is as essential as sparing memory in this methodology. Here development of FP-Tree is like unique FP-Growth calculation.

.
I.    Data Transformation - This is used to hide the data from the service provider by using encryption technique i.e. one-to-one mapping technique. By using this technique, mining algorithms can be applied with 100% accuracy. [7]

II.   Mining Process - It is applied on transformed XML information streams and uses landmark windows processing model that relies on FP-Growth algorithmic rule approach. In this approach the info stream is processed block by block and every block contains constant number of transactions. [8]

III.  Data Re-Transformation - In Mining Operation, the algorithmic rule finds the positive and negative association rules with transformed items and sends them back to the data 2014 International Conference on Circuits, Organization, and Communication and Information Technology Application program (CSCITA) 387 owner. In this step re-transform the items using the single-valued function table generated. [9]

## 3. APRIORI ALGORITHM

Apriori is a well-known association rule mining algorithm. This model, basically, is divided the info mining methodology into two easy steps. In the beginning, the algorithm generates the 1−k giant item sets wherever k is the count of separate things within the transactions. After the candidate item sets area unit generated, the algorithm finds the frequent giant item sets that have support worth bigger than the predefined minimum support value. The next basic step generates the association rules from the large data sets which contains frequent items that fulfil the min confidence constraint [11]. The advantage of the Apriori algorithm is its easy implementation. However, there is a major disadvantage of Apriori algorithm that it requires too many scans over the entire data sets which comprises of the database to find rules which leads to high usage of the system as more memory is required to complete the need of increasing I/O cost. Because of the reiterative scans, it is not suitable for knowledge stream mining during which knowledge ought to be scanned just one occasion. [10]

Association rule generation is usually separate into 2 separate steps:

1. First, minimum support is applied to find all frequent itemsets in an exceedingly info.
2. Second, these frequent itemsets and the minimum confidence constraint are wont to kind rules.

While the second step is straight forward, the first step wants additional attention. Finding all frequent itemsets in a database is tough since it involves looking all doable itemsets (item combinations). The set of possible itemsets is the power set over I and has size $2n − 1$ (excluding the empty set that isn't a legitimate itemset). Although the size of the ability set grows exponentially within the variety of things n in I, efficient search is doable victimization the downward-closure property of support (also referred to as anti-monotonicity) that guarantees that for a frequent itemset, all its subsets are additionally frequent and therefore for an infrequent itemset, all its supersets must additionally be infrequent . Exploiting this property, efficient algorithms (e.g., Apriori and Eclat) can notice all frequent itemsets. [12]

## 4. FP GROWTH ALGORITHM

FP-Growth is another surely understood association principle mining calculation. The Mining Process of the FP-Growth calculation was partitioned into two stages is as per the following. To start with the FP tree is built as clarified beneath. With a specific end goal to discover the bolster estimation of everything, the information set is initially filtered. The rare things are disposed of, on the grounds that they don't have significance in the Mining Process. The continuous things are sorted in diminishing bolster esteem. After that, the information set is filtered afresh to develop FP-Tree. The main exchange is perused and the hubs are made. Likewise, the bolster estimation of everything is set to 1. On the off

chance that the exchanges don't contain basic prefixes, the procedure is preceded with making new hubs. Something else, if there are exchanges with normal things, their ways is covered somewhat or completely. As a result of covering ways, bolster estimations of basic things are expanded by 1 and bolster estimations of the others are set to 1. The procedure is preceded until the sum total of what exchanges have been mapped. Finally, visit thing sets which have higher bolster esteem than the client indicated bolster limit are mined without competitor s et era. [13]

### 4.1 FP - Growth Algorithm: Example

The FP Growth algorithm requires a simple two step procedure to build an association rule:

1. Build a compact data structure which is known as the FP - tree
2. Extract frequent datasets from the FP - tree

Let us consider the following transaction table

| Id | List of item Ids |
|------|------------------|
| T100 | {I1, I2, I3} |
| T200 | {I2, I4} |
| T300 | {I2, I3} |
| T400 | {I1, I2, I4} |
| T500 | {I1, I3} |
| T600 | {I2, I3} |
| T700 | {I1, I3} |
| T800 | {I1, I2, I3, I5} |
| T900 | {I1, I2, I3} |

**Table 4.1** Transaction table

Now we will a FP tree of the above transactional database. Item sets are considered in order of their descending value of support count. The final generated FP tree according to the above transactional database table is as shown below:



**Figure 4.1** Transaction dataset

The FP – tree construction is completed by following the above mentioned steps. Now, according the generated FP – tree the frequently generated patterns are as mentioned below:

| | Frequent Pattern Generated |
|---|---|
| I5 | {I2, I5:2}, {I1, I5:2}, {I2, I1, I5:2} |
| I4 | {I2, I4:2} |
| I3 | {I2, I3:4}, {I1, I3:4}, {I2, I1, I3:2} |
| I1 | {I2, I1, I4} |

**Table 4.2** Transaction dataset Pattern

**4.2 Comparison of Algorithms**

| Algorithm | Advantages | Disadvantages |
|---|---|---|
| Apriori | – Easy to implement<br>– New pruning technology<br>– Avoids wastage of counting candidate which are infrequent | – Too many scans on database high CPU usage |
| FP Growth | – Only two passes on database<br>– Faster and better than the Apriori algorithm<br>– Computation cost decreased<br>– FP Tree construction | – FP Tree is difficult to use in an interactive mining system<br>– FP Tree is not suitable for incremental mining |

**Table 4.3** Comparison of Apriori

## 5. BASE PAPER APPROACH

Following are the steps that are being summarized in this research paper as per the base paper approach

*Step 1* - Scan TDB to find frequent itemset if the item set does not satisfy the following condition.
   a.   (support<min_sup&&weight<min weight)
   b.   (support*MAXW<min_sup)

*Step 2* - Ascending order sorting of items in weight which forms the weight_order and header of FP tree

*Step 3* - Scan the TDB and global FP tree algorithm using weight_order is built

*Step 4* - Mine global FP tree for weighted frequent itemset mining in bottom up manner forming conditional database using condition (a) and other as (support*MINW<min_sup)

*Step 5* - When all the items in the global header table have been mined, wfim is completed

## 6. BASE PAPER APPROACH

The following are the steps that are being performed by refining the older approach and creating an all new approach as per proposed in the research paper

*Step 1 -* Collect the dataset and take it as an input TDB (Transaction dataset)

*Step 2 -* State min_supp (predicted minimum support) and scan each TDB to learn frequency for each T (transaction) in TDB

***Step 3 -*** For each of the transaction T in TDB, create a FP tree

***Step 4 -*** Call Re-order function

     a.   Compare the weight of the node with its parent node and associate positive negative sign.
     b.   Calculate $^-y = b\_0 + b\_1\,(X)$ , defines associations of particular node with other nodes

             where $b\_0 = min\_degree$

             $b\_1 = $ current node degree
             $X = $ No. of same nodes at that level

     a.   Calculate prediction error $= y - {}^-y$

             where, $y = $ all the nodes at the same level

     b.   Calculate SSTO $= $ sum of square of degrees of all nodes at same level
     c.   Calculate Co-relation $r^2 = $ (SSTO-PE)/SSTO  and find the value of r

***Step 5 -*** The value of r defines the positive and negative association probability

## 7. PROPSED MODEL

In this proposed work first of all we collect traditional databases and then we performed to find Association Rules using an improved algorithm which is based on FP Growth algorithm with Horizontal Tree Approach. The preprocessing of data and FP tree creation with reordering will improve search of matching association. Once that process is finished Positive and Negative Association will be set by level in the new FP Tree, which will be displayed as the result.



**Figure 7.1** Flow Diagram of Proposed Algorithm

**Figure 7.2** Novel Approach Flow Diagrams

## 8. RESEARCH GAP

Faster retrieval for fuzzy information stored datasets is necessary. Horizontal tree approach does not ensure proper implementation of the algorithm and therefore does have the efficiency and impression for data retrieval. To improve accuracy datasets with different class and association must be experimented. Different tree approaches with branching probability should be examined. Co-relations between frequent itemsets and different support count needs more processing requirements.

## 9. IMPLEMENTATION

The Positive Negative association standards & its rules will be applied to the dataset to remove the relationship among the datasets. Likewise, the re-organizing of the FP calculation will hold the outcomes all the more effectively and rapidly.

In the proposed work of the implementation, we have devised and implemented a new algorithm to generate both the association rules i.e. Negative and Positive association rules. There are very few papers to discuss and discover negative association rules. Our algorithm generates differs from others in the sense of that it generates the association rules from a different candidate set.

We conducted our experiments and used the following datasets which were available on UCI Machine Learning Repository (German Credit Dataset) to study the behavior of the algorithms that are being compared and used in the research. It has a number data items present in the dataset which can be used to determine the results of the above proposed algorithm according to flow of the novel approach. The dataset stores the following information i.e. purpose of selling the product, credit amount, age and personal. We have used large sized dataset in this paper in order to determine the efficiency and amount of time the algorithm takes to mine the data using the association rules. We have randomly selected the dataset for only experimental purpose and to determine the strength of the proposed algorithm. [14]

We have also used two different types dataset to determine the results as per shown in the following sections

1. Amazon Commerce Review Set
2. Whole Sale Customers

*Step 1 -* Get input of dataset from UCI machine learning repository



**Figure 9.1** Repository of Dataset

**Step 2 -** Input data in the form of MYSQL file



**Figure 9.2** Input Dataset

**Step 3 -** Retrieval of dataset using AJAX



**Figure 9.3** Retrieval of Data

***Step 4 -*** Re-ordering of rule mining in code analysis
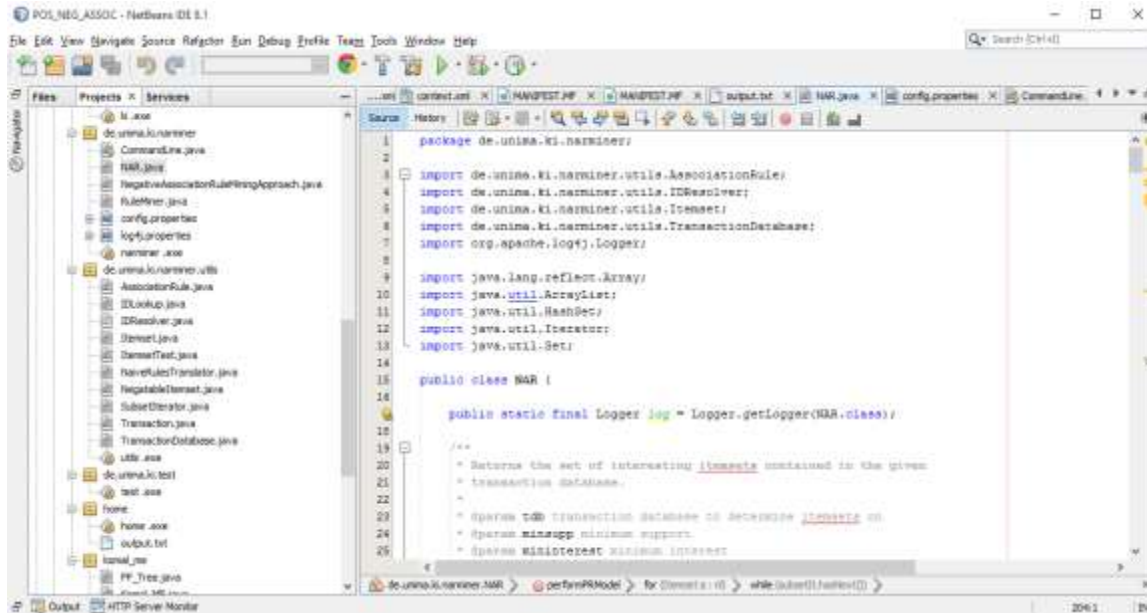


**Figure 9.4** Re-ordering
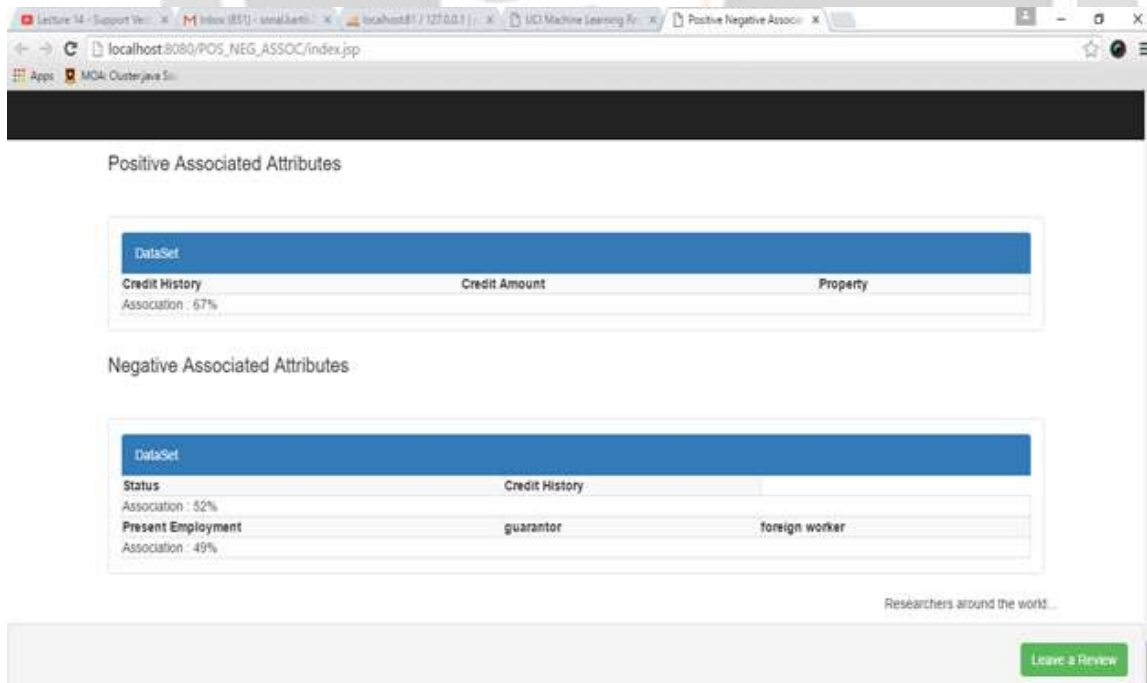
***Step 5 -*** Output of association rule



**Figure 9.5** Output
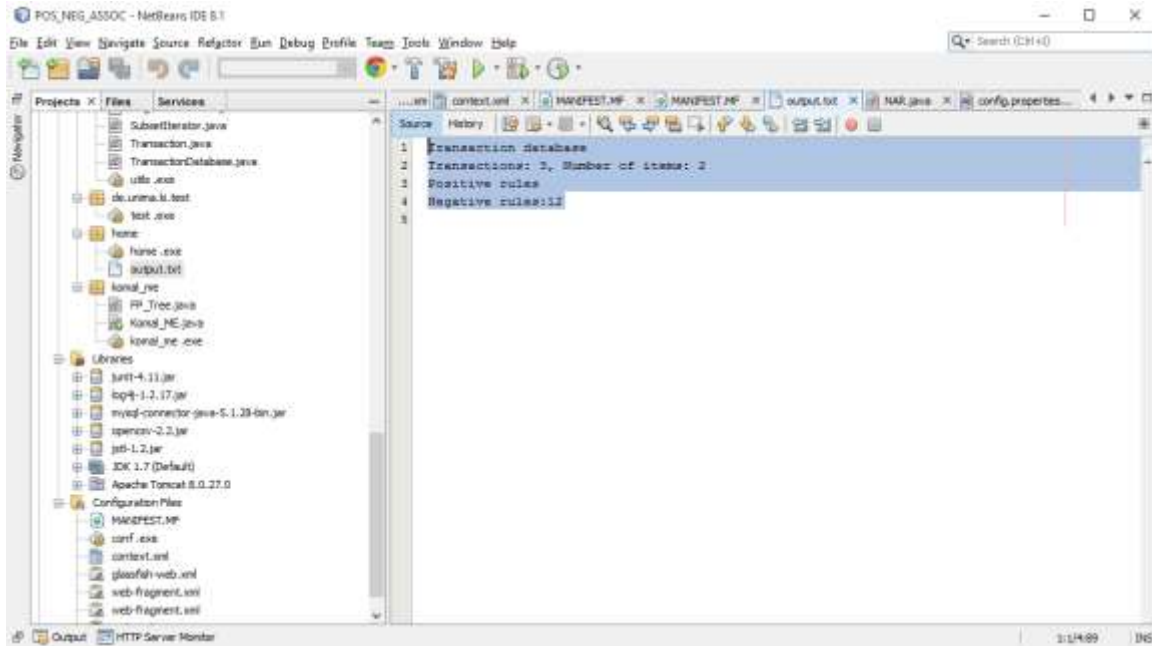
***Step 6 -*** Output information of dataset



**Figure 9.6** Output of dataset

***Step 7 -*** First set of minimum support and minimum confidence analysis used in the algorithm in code analysis
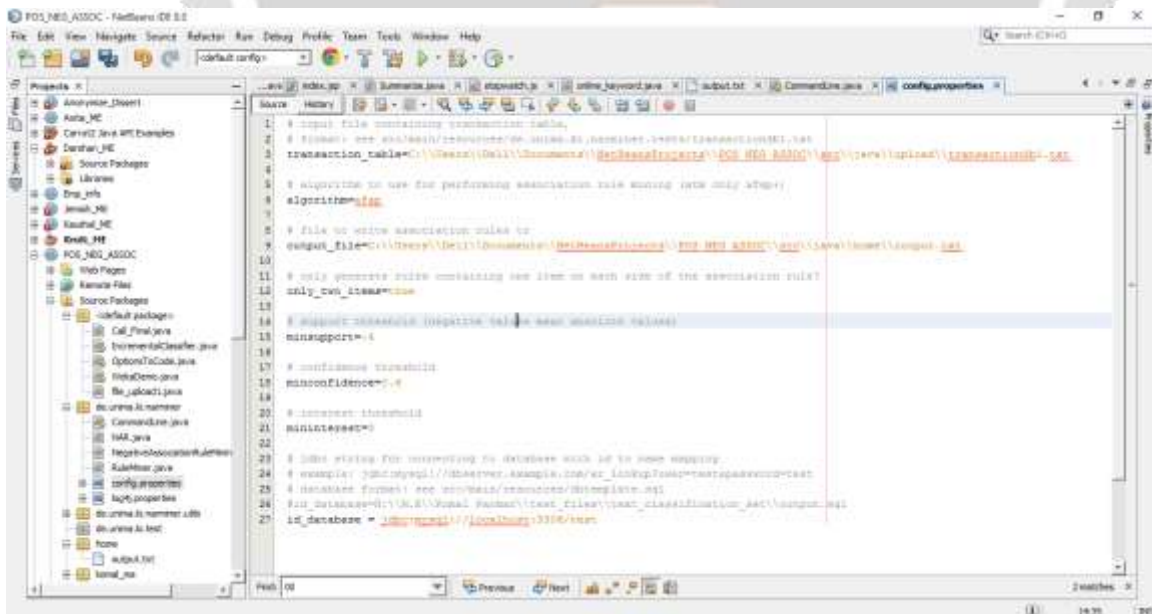


**Figure 9.7** First set of input of dataset in code analysis

**Step 8 -** Output of the algorithm in code analysis using the first set of data
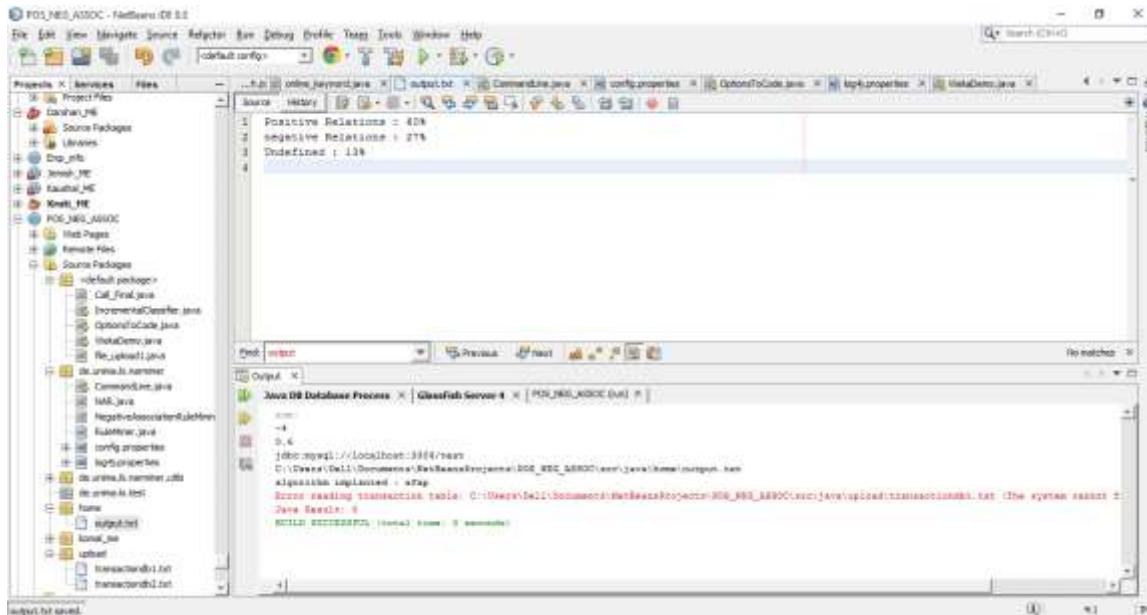


**Figure 9.8** Result of first set in code analysis

**Step 9 -** Second set of minimum support and minimum confidence analysis used in the algorithm in code analysis
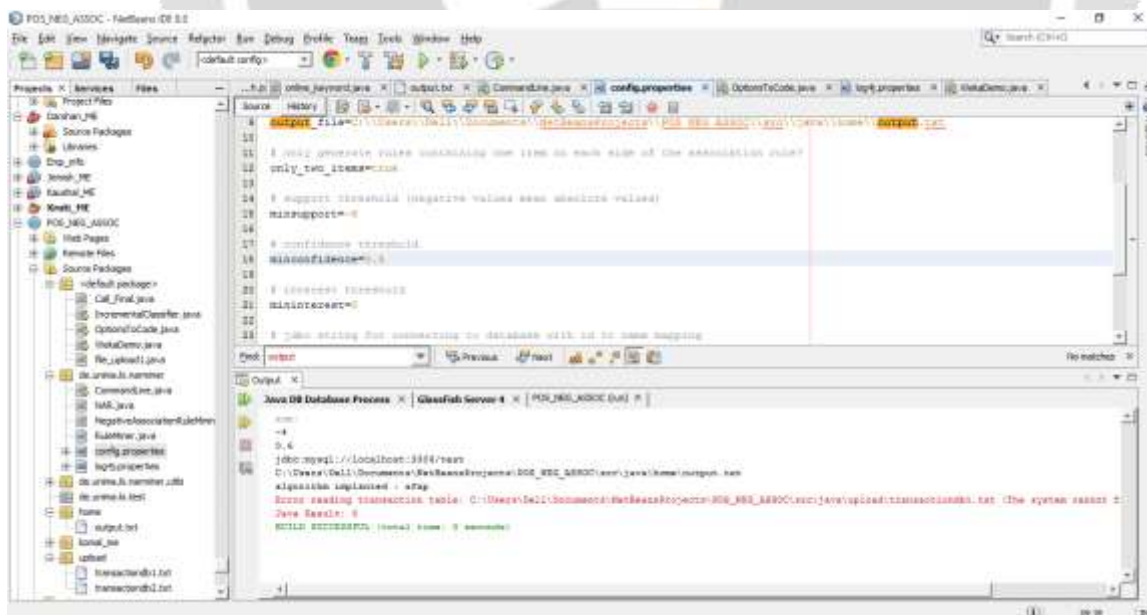


**Figure 9.9** Second set of input of dataset in code analysis

***Step 10 -*** Output of the algorithm in code analysis using the second set of data
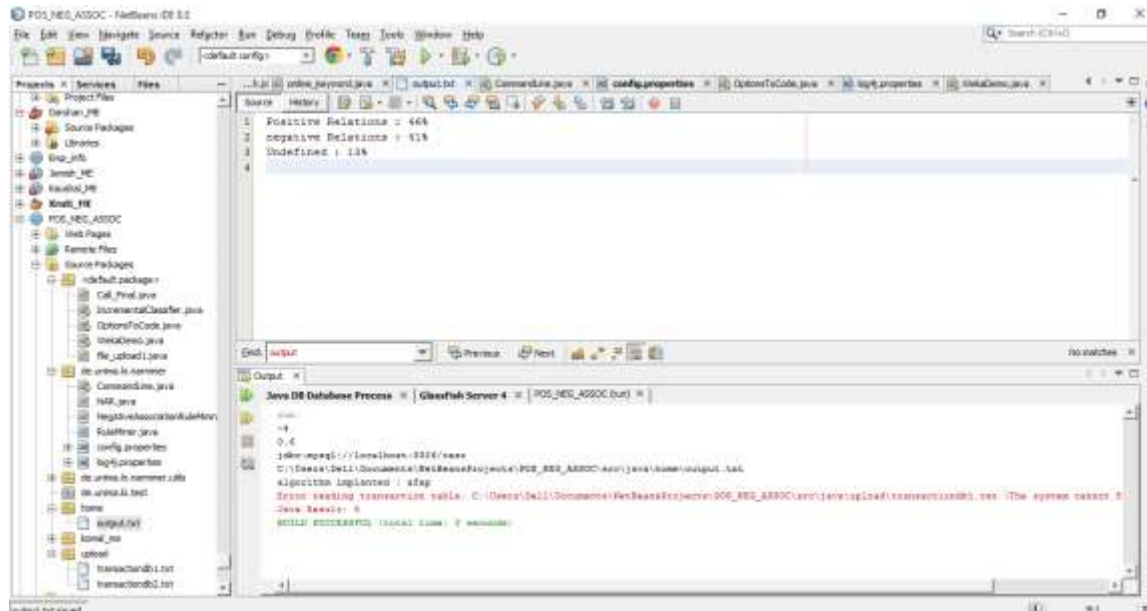


**Figure 9.10** Result of second set in code analysis

## 10. CONCLUSION

As referred in the base paper they using AJAX data stream on which datamining processing consumes more time so they have not transformed AJAX data stream ,which does not provide secure data transactions and they have focus more on the Positive Rule and frequent itemsets which not provides accurate results for both Positive and Negative Rules. So we have applied transformation AJAX stream to provide data transactions data and introducing new algorithm and tries to provide secure transactions and accurate results for both Positive Negative Rules. We will also try to decrease consuming space n time .which will make mining processing more efficient.

## 11. REFERENCES

[1]     Pascal Hitzler, Markus Krötzsch, Sebastian Rudolph, "Foundations of Semantic Web Technologies," Textbooks in Computing, Chapman and Hall/CRC Press, 2009.

[2.]     Victoria Nebot, Rafael Berlanga "Finding association rules in semantic web data", Knowledge-Based Systems, vol. 25, no. 1, pp. 55-62, 2012.

[3.]     E. Willighagen (2013), RRDF - support for the resource description framework [Online], Available: http://cran.r-project.org/web/packages/rrdf/rrdf.pdf.

[4.]     W3C (2008), SPARQL query language for RDF [Online], Available: http://www.w3.org/TR/rdf-sparql-query/.

[5.]     Treese, G.W. and Stewart, L.C., Designing Systems for Internet Commerce, Addison Wesley, 1998.

[6.]     Booch, G, Rumbaugh, J., and Jacobson, I, UML User's Guide, Addison Wesley, 1999.

[7.]     Berthold Daum, Udo Merten, "system Architecture with XML", Morgan Kaufmann publishers, 2003.

[8.]     J. D. Conley and R. P. Whitehurst, "Automatic and transparent denormalization support, wherein denormalization is achieved through appending of fields to base relations of a normalized database. Patent, US5369761".

[9.]     Corbitt, B.J., "Trust and e-commerce a study of e-commerce perceptions", Electronic Commerce Research & Application, Vol .2 No.3, pp.203 -15(2003).

[10.]    Song, I-Y. And LeVan-Schultz K., "Data Warehouse Design for E-Commerce Environments", Lecture Notes in Computer Science, Vol. 1727, Springer, pp. 374-388.

[11.]    Mudra Doshi, Bidhisha Roy, "Enhanced Data Processing Using Positive Negative Association Mining on Ajax Data", (2014 IEEE).

[12.]    Varsha Kavi and Divyesh Joshi, "A Survey Based on Enhancing Data Processing of Positive and Negative Association Rule Mining", (International Journal of Computer Science and Engineering, March 2014).

[13.]    Idheba Mohamad Ali, Azuraliza Abu Bakar ,Anis Suhailis Abdul Kadir, "Mining Positive and Negative Association rules from Interesting Frequent and Infrequent Itemsets", ( IEEE ,2012).

[14.]    Xushan Peng, Yanyan Wu, "Research and Application of Algorithm for Mining Positive and Negative Association Rules", (IEEE 2011).