# An Association Rule Mining Model for Finding the Interesting Patterns in Stock Market Dataset

Sachin Kamley
Deptt.of Computer Applications
S.A.T.I.
Vidisha, India

Shailesh Jaloree
Deptt.of Maths & Com. Sci.
S.A.T.I.
Vidisha, India

R.S. Thakur
Deptt. of Computer Applications
M.A.N.I.T.
Bhopal, India

## ABSTRACT
In these days, stock market forecasting is one of the most interesting issues, which has gained a more attention due to vast profits. To precisely predict the price of share and making profits has been always challenging task since the longest period of time. This has engrossed the interest and attention of stock brokers, economists and applied researchers. Traditional methods like Fundamental analysis, Technical analysis, and Regression methods are not suitable for this task because these tools and techniques are based on totally different analytical approaches and requiring highly expertise and justification in the area. In this sequence, Association Rule Mining is one of the most interesting research areas for finding the associations, correlations among items in a database. It can discover all useful patterns from stock market dataset. The aim of this research study is to help stock brokers, investors so that they can earn maximum profits for each trading.

## Keywords
Stock Market, Data Mining, Prediction, BSE, Association Rule Mining, Frequent Pattern, SAS 9.2.

## 1. INTRODUCTION

### 1.1 Stock Market
Stock market is a place where company's securities and derivatives are bought and sold at an approved stock price. The stock market refers to the activity generated by the stock exchanges. Stock market trading generates economic stimulus, which one can follow in stock market news [9]. These trends can be analyzed in order to make informed purchases. There are two most important markets are "Primary" and "Secondary" markets.

Primary markets are locations for corporations and government bodies to raise direct financial capital by selling stocks and bonds to investors, usually for specific ventures.

Secondary markets are where stock trading occurs and investors sell their stocks to other investors without the involvement of the issuing companies, monitored by a regulatory body called SEBI (Security and Exchange Board of India) [1]. Stock market behavior is dynamic because stock prices fluctuate every time. It strongly depends upon demand and supply. The prices will high when the demand is high and the prices will low when the share is heavy to sell.

There are two well known Indian stock market indexes Sensex and Nifty. Sensex is one of the oldest market indexes for equities and it includes shares of 30 firms indexed on the BSE (Bombay Stock Exchange of India), which represent about 45% of the index's free-float market capitalization. It was started in 1986 and provides time series data from April 1979, onward. Another index is the S&P (Standard & Poor)

CNX Nifty; it includes 50 shares listed on the NSE, which represent about 62% of its free-float market capitalization. It was started in 1996 and provides time series data from July 1990, onward [1] [20].

### 1.2 Association Rule Mining
Data mining has attracted a great deal of attention in the areas like Medical and Health care, Telecommunication, bioinformatics, financial analysis etc. In each of these application areas, the huge amounts of data available for analysis purpose in recent years. Due to the large size of databases, finding potential information and hidden patterns in data has become increasingly challenging task [12]. In data mining, Association Rule Mining is a most popular research area for discovering interesting patterns and associations in huge amount of databases. It is very interestingness to identify strong rules discovered in databases based on some useful constraints.

Following is the original mathematical definition by Agrawal et al. (1993) [3] is defined as: Let $I = \{ I_1, I_2, \ldots\ldots\ldots\ldots I_n \}$ be a set of items. Let $D$, the database consist a set of transactions, where each transaction $T$ has a unique transaction $Id$ and contains a subset of the items in $I$.

An association rule is defined as $X \rightarrow Y$ where $X = \{x_1, x_2, \ldots\ldots x_n\}$, and $Y = \{y_1, y_2, \ldots\ldots y_n\}$ are sets of items, with $x_i$ and $y_i$ being distinct items for all $i$ and all $j$. The association rule states that if a customer purchase X, items than he or she is also likely to purchase Y items. In general any association rule has the form LHS (Left Hand Side) $\rightarrow$ RHS (Right Hand Side), where LHS and RHS are two sets of items. The set LHS $\bigcup$ RHS is called an itemset, the set of items purchased by the customers [2].

To discover interesting rules from the dataset, two common interest measures are support and confidence.

$support( X \rightarrow Y ) = P( X U Y )$

$confidence (X \rightarrow Y) = P (Y \mid X)$

Support (s) can be measure as the probability $P( X U Y )$ of percentage of transactions in $D$ that contain $X U Y$ i.e., the union of sets $X$ and $Y$ or say, both $X$ and $Y$ occur frequently [12].

Confidence (c) can be measure as the conditional probability $P(Y|X)$ of percentage of transactions in $D$ containing $X$ that also contain $Y$ i.e., shows the confidence of a rule [12].

The aim of support count finding the itemsets which occur infrequently within the data set and hence are irrelevant for final associations. The confidence establishes the intensity of the associations whether it is weak, average or strong. Minimum support (or minsup) is a predefined threshold parameter which defines that the itemsets whose support is less than this threshold are not interesting. Its value varies within the range of 0 to 1 [8].

Minimum confidence (or minconf) is a user defined threshold indicating that the association rules with confidence less than this threshold are not interesting. Its value varies within the range of 0 to 1 [3].

Maximum Itemset size specifies the maximum size of an itemset. The default value 0 indicates that there is no size limit on the itemset. Reducing the maximum itemset size also reduces the processing time because it saves time or no of iterations. Minimum Itemset Size specifies the minimum size of the itemset. The default value is 0. Reducing Minimum Itemset Size will not reduce the processing time because the algorithm starts with 1-itemset size and increases the size step by step [14].

In this sequence, an optimum value of support and confidence are chosen so that meaningful rules and patterns might be generated. Association rules are usually required to satisfy a user-specified minimum support (minsup) and confidence (minconf) at the same time.

Association rule generation process is usually works in two phases:

1. Firstly, minimum support is applied to find all frequent itemsets and the minimum confidence constraint.

2. Second phase is most important to find rules based on these frequent itemsets and the minimum confidence constraint [23].

Discovering association rules is one of the important data mining problems and there has been increased considerable research in the field of data mining. Determining the correlations among items that occur synchronously in the database and to obtain information which is useful for decision makers. Therefore, the main purpose of implementing the association rule mining algorithm is to find synchronous relationships by analyzing the random data and to use these relationships as a reference during decision making [11] [18].

In this study, ARM algorithm apply on stock data set, finding the relationship among stock variables open price, high price, low price and close price and generating the rules with these variables.

## 2. RELATED WORK

Argiddi et al. (2012) [4] have proposed the FITI (First Intratransaction Then Intertransaction) algorithm which is adaptive to intertransaction association mining. The algorithm works in two phases (1) mining frequent Intratransaction itemsets (2) rule generation. The result was excellent when they compared the performance of FITI with EH-Apriori w.r.t. high running time.

Srisawat (2011) [19] has proposed an application of stock rules in stock market. They applied Association Rule Mining technique for discovering the relationships between individual stocks and they used the transactional dataset consists of 242 trading days from 4 January 2010 to 30 December 2010. The more promising rules were generated from dataset.

Borisov (2011) [6] has proposed Association Rule Mining (ARM) methodology based on detecting fab tool commonality of affected lots. They compared the performance of ARM method with several traditional methods such as ANOVA (Analysis of Variance) and contingency tables using eight actual production cases. The experimental results were more promising.

Paranzape et al. (2013) [16] have proposed a recommender system based on Association Rule Mining. The system generates buy and sell signal time to time and helps to investors for predicting the stock market price.

Nandgopal et al. (2012) [15] have proposed an efficient algorithm called Inter-transaction Association Rule Miner (IAR Miner) for mining inter-transaction itemsets. The proposed algorithm works in two phases (1) First it scans the database for finding the frequent itemsets. For each frequent item found, then IAR Miner converts the original transaction database into a set of domain attributes, called a dataset. (2) Second it enumerates inter-transaction itemsets using an Item set-Dataset tree, called an ID-tree. The IAR Miner can embed effective pruning strategies. It also avoiding costly candidate generation process and repeated support counting.

## 3. PROPOSED METHODOLOGY AND DATA PREPROCESSING

### 3.1 Apriori Algorithm

Apriori algorithm proposed by Agarwal and Srikant 1994. It is the most popular algorithm to find association rules on large scale dataset and makes use of the downward closure property. The algorithm employs level by search or an iterative approach, where K-itemsets are used to explore (K+1) - itemsets. The algorithm terminates when no more frequent K-itemsets can be found [12].

### 3.2 Data Preprocessing

For the purpose of this study, last 6 years TCS Company monthly stock data employed from Bombay Stock Exchange of India site [21]. The data employed in this study contain variables open price, high price, low price and close price of TCS index. The data set encompassed the trading months from 1-1-2006 to 30-12-2011 [21]. The task is to predict the stock prices of TCS Company will be up or down from today to tomorrow by using the past values of company stock. The stock market data sample is given by Table 1.

**Table 1. Sample of Stock Dataset**

| Month | Open | High | Low | Close |
|---|---|---|---|---|
| Jan-06 | 1189 | 1258 | 1169 | 1251 |
| Feb-06 | 1255 | 1301 | 1230 | 1295 |
| . | . | . | . | . |
| . | . | . | . | . |
| . | . | . | . | . |
| Jul-11 | 2326 | 2361 | 2254 | 2256 |
| Aug-11 | 2270 | 2285 | 1955 | 2061 |
| Sep-11 | 2086 | 2132 | 1968 | 2028 |
| Oct-11 | 2012 | 2168 | 1948 | 2156 |
| Nov-11 | 2141 | 2163 | 1881 | 1953 |
| Dec 11 | 2212 | 2304 | 1715 | 2190 |

Before applying predictions on the data, data is analyzed and processed for problem; it must be completely inspected, cleaned and normalized. Even the best predictor will also fail on improper data, so data quality and data preparation is crucial. Since predictor can use certain data features, so it would be very important to detect best data pre processing features. Data cleaning routines work to "clean" the data by filling in missing values, smoothing noisy data, and resolving inconsistencies [8] [12]. Stock market stores a huge amount of data throughout the year. Incomplete, noisy, and inconsistent data are very familiar properties of large real-world databases, so majority of time was spent in these experiments on data pre-processing task. The noise exists in the stock dataset can be generally classified into three main categories, which are duplicate records, inconsistent, and incomplete price. Firstly, we clean and normalize the data (removing decimal points), fill missing values in data, removing inconsistency from data, which is shown by Table 1. Stock market data is very difficult to interpret because prices of stock are numerically large and not feasible for computation. There is a need to change the data in symbolic representation and to make use of the variables Open Price, High Price, Low Price and Close Price to carry out this conversion, which is shown in Table 2.

**Table 2. Symbolic Conversion of Stock Data Set**

| S.NO. | Variable | Symbolic Conversion |
|---|---|---|
| 1 | Open | { Open_rise, Open_fall } |
| 2 | High | { High_rise, High_fall } |
| 3 | Low | { Low_rise, Low_fall } |
| 4 | Close | { Close_rise, Close_fall } |

The Table 2 shows symbolic conversion of data values. For this, data values are classified under the below average and the above average. The data values below average comes under the category of falling price and the values which are above average comes under the category of rising price. The Table 3 & Table 4 shows this classification.

**Table 3. Data Classification of Open and High Price**

| | Open Price | High Price |
|---|---|---|
| | lies between(1030-2665) | lies between(1141-2777) |
| Average | **1875** | **1972** |
| | open fall(1030-1866) | high fall(1141-1952) |
| | open rise(1875-2665) | high rise(1985-2777) |

**Table 4. Data Classification of Low and Close Price**

| | Low Price | Close Price |
|---|---|---|
| | lies between(922-2513) | lies between(1045- 2657) |
| Average | **1752** | **1879** |
| | low fall(922-1725) | close fall(1045-1855) |
| | low rise(1785-2513) | close rise(1895-2657) |

The Statistical Analysis Software (SAS) 9.2 is used to generate association rules from stock market data. The preprocessed data is shown in Figure. 1.

**Fig 1: Pre-processed Data of SAS**

After preprocessing step, the next task to apply mining algorithm and to predict stock prices. Here last 6 years data i.e. 1 January 2006 to 30 December 11 are considered for study [21] and applied the Association Rule Mining (ARM) approach to generate association rules for prediction. ARM finds interesting associations and relationships among large set of data items. Different parameters are used in ARM for rule generations are: minimum support count, no. of records (items), CPU time etc. The values of different parameters are used for prediction purpose is shown by Figure 2.

## 4. EXPERIMENTAL RESULTS

This section presented detailed study of experiment work, which is conducted for stock prediction using Association Rule Mining (ARM). The Table 5 shows the frequency of variable (items) i.e. number of times the variable (items) appears in the database.

**Table 5. Item Frequency of Variable**

| Id | Count | Item |
|----|-------|------|
| 1 | 40 | Open_rise |
| 2 | 39 | Low_rise |
| 3 | 39 | High_rise |
| 4 | 39 | Close_rise |
| 5 | 33 | Low_fall |
| 6 | 33 | High_fall |
| 7 | 33 | Close_fall |
| 8 | 32 | Open_fall |

After finding the frequency of variable (items), the support for 1- itemset, 2- itemset, 3- itemset, and 4-itemset is calculated. For the purpose of this study 72 items and 288 records are considered. Mining the rules from stock data set is very time consuming process, but proposed algorithm mines the rules very fast and in reasonable time. The summary of frequent itemsets is shown below in Figure 2.
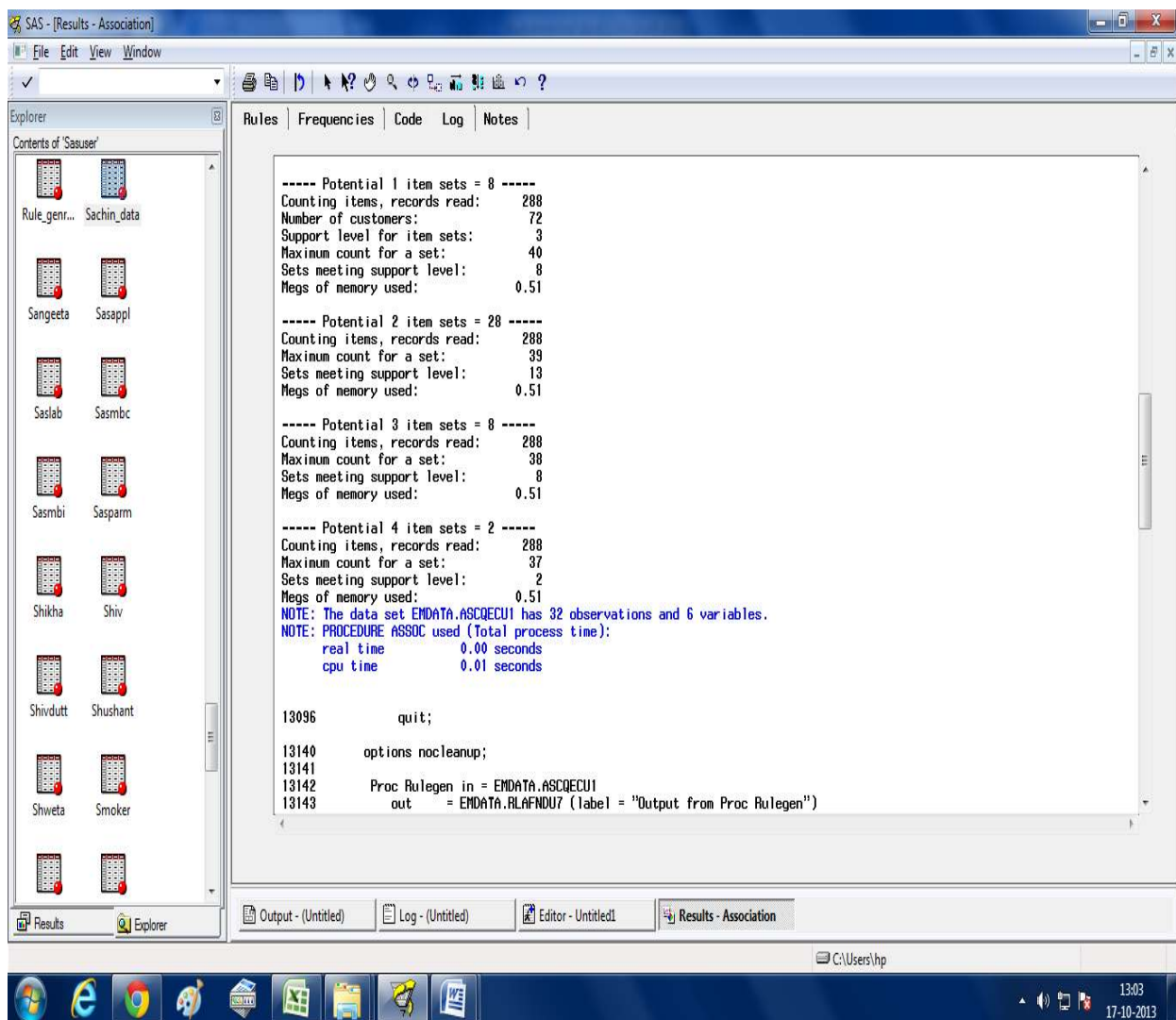


**Fig 2: Summary of Frequent Item Set**

Stock prices are highly correlated with each other i.e. change in one price value may affect other price value. The ARM algorithm easily finding the association between items set (variables) and finds the frequent item set based on support and confidence value [3]. For this purpose different experiments conducted on stock data set. The different support value 2%, 4%, 6%, 8%, 10% & 12% are used for different size rule generation, which are shown through the Tables 6, 7, 8, 9, 10, & 11.

**Table 6. Maximum No. of Rule Generated with Support= 2%**

| Confidence (%) | Max Length 2 | Max Length 3 | Max Length 4 |
|---|---|---|---|
| 10% | 24 | 94 | 150 |
| 20% | 24 | 94 | 150 |
| 40% | 24 | 94 | 149 |
| 60% | 24 | 90 | 135 |
| 80% | 24 | 88 | 133 |
| 100% | 02 | 28 | 53 |

**Table 7. Maximum No. of Rule Generated with Support= 4%**

| Confidence (%) | Max Length 2 | Max Length 3 | Max Length 4 |
|---|---|---|---|
| 10% | 24 | 94 | 150 |
| 20% | 24 | 94 | 150 |
| 40% | 24 | 94 | 149 |
| 60% | 24 | 90 | 135 |
| 80% | 24 | 88 | 133 |
| 100% | 02 | 28 | 53 |

**Table 8. Maximum No. of Rule Generated with Support=6%**

| Confidence (%) | Max Length 2 | Max Length 3 | Max Length 4 |
|---|---|---|---|
| 10% | 24 | 72 | 100 |
| 20% | 24 | 72 | 100 |
| 40% | 24 | 72 | 100 |
| 60% | 24 | 72 | 100 |
| 80% | 24 | 72 | 100 |
| 100% | 02 | 12 | 20 |

**Table 9. Maximum No. of Rule Generated with Support=8%**

| Confidence (%) | Max Length 2 | Max Length 3 | Max Length 4 |
|---|---|---|---|
| 10% | 24 | 72 | 100 |
| 20% | 24 | 72 | 100 |
| 40% | 24 | 72 | 100 |
| 60% | 24 | 72 | 100 |
| 80% | 24 | 72 | 100 |
| 100% | 02 | 12 | 20 |

**Table 10. Maximum No. of Rule Generated with Support=10%**

| Confidence (%) | Max Length 2 | Max Length 3 | Max Length 4 |
|---|---|---|---|
| 10% | 24 | 72 | 100 |
| 20% | 24 | 72 | 100 |
| 40% | 24 | 72 | 100 |
| 60% | 24 | 72 | 100 |
| 80% | 24 | 72 | 100 |
| 100% | 02 | 12 | 20 |

**Table 11. Maximum No. of Rule Generated with Support=12%**

| Confidence (%) | Max Length 2 | Max Length 3 | Max Length 4 |
|---|---|---|---|
| 10% | 24 | 72 | 100 |
| 20% | 24 | 72 | 100 |
| 40% | 24 | 72 | 100 |
| 60% | 24 | 72 | 100 |
| 80% | 24 | 72 | 100 |
| 100% | 02 | 12 | 20 |

The above results show the interestingness measure when confidence (%) increases than number of rules decreases and the rules are generated with different combinations like length-2 {open, close}, length-3 {open, high, close}, length-4 {open, high, low, close} etc. the graph is shown by Figure 3.
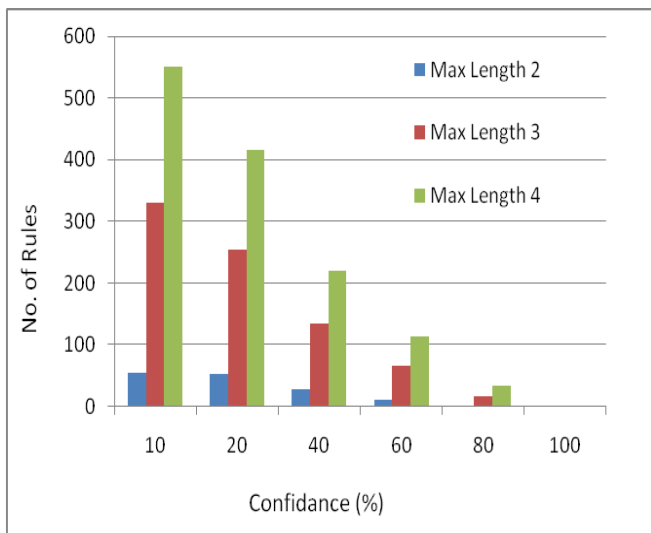


**Fig 3: Total No. of Rules Generated with Different Values of Confidence (%)**

.

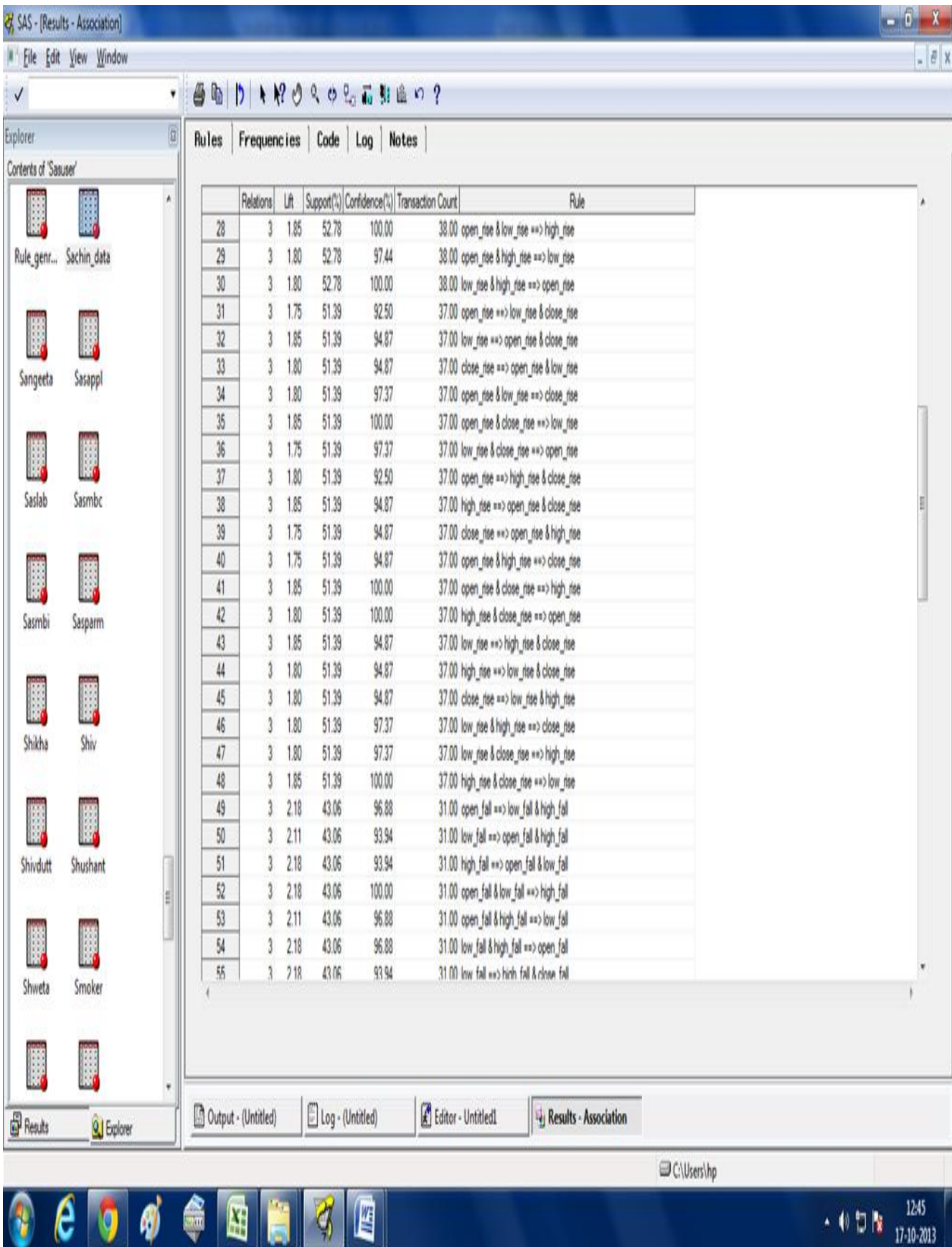| | Relations | Lift | Support(%) | Confidence(%) | Transaction Count | Rule |
|---|---|---|---|---|---|---|
| 28 | 3 | 1.85 | 52.78 | 100.00 | 38.00 | open_rise & low_rise ==> high_rise |
| 29 | 3 | 1.80 | 52.78 | 97.44 | 38.00 | open_rise & high_rise ==> low_rise |
| 30 | 3 | 1.80 | 52.78 | 100.00 | 38.00 | low_rise & high_rise ==> open_rise |
| 31 | 3 | 1.75 | 51.39 | 92.50 | 37.00 | open_rise ==> low_rise & close_rise |
| 32 | 3 | 1.85 | 51.39 | 94.87 | 37.00 | low_rise ==> open_rise & close_rise |
| 33 | 3 | 1.80 | 51.39 | 94.87 | 37.00 | close_rise ==> open_rise & low_rise |
| 34 | 3 | 1.80 | 51.39 | 97.37 | 37.00 | open_rise & low_rise ==> close_rise |
| 35 | 3 | 1.85 | 51.39 | 100.00 | 37.00 | open_rise & close_rise ==> low_rise |
| 36 | 3 | 1.75 | 51.39 | 97.37 | 37.00 | low_rise & close_rise ==> open_rise |
| 37 | 3 | 1.80 | 51.39 | 92.50 | 37.00 | open_rise ==> high_rise & close_rise |
| 38 | 3 | 1.85 | 51.39 | 94.87 | 37.00 | high_rise ==> open_rise & close_rise |
| 39 | 3 | 1.75 | 51.39 | 94.87 | 37.00 | close_rise ==> open_rise & high_rise |
| 40 | 3 | 1.75 | 51.39 | 94.87 | 37.00 | open_rise & high_rise ==> close_rise |
| 41 | 3 | 1.85 | 51.39 | 100.00 | 37.00 | open_rise & close_rise ==> high_rise |
| 42 | 3 | 1.80 | 51.39 | 100.00 | 37.00 | high_rise & close_rise ==> open_rise |
| 43 | 3 | 1.85 | 51.39 | 94.87 | 37.00 | low_rise ==> high_rise & close_rise |
| 44 | 3 | 1.80 | 51.39 | 94.87 | 37.00 | high_rise ==> low_rise & close_rise |
| 45 | 3 | 1.80 | 51.39 | 94.87 | 37.00 | close_rise ==> low_rise & high_rise |
| 46 | 3 | 1.80 | 51.39 | 97.37 | 37.00 | low_rise & high_rise ==> close_rise |
| 47 | 3 | 1.80 | 51.39 | 97.37 | 37.00 | low_rise & close_rise ==> high_rise |
| 48 | 3 | 1.85 | 51.39 | 100.00 | 37.00 | high_rise & close_rise ==> low_rise |
| 49 | 3 | 2.18 | 43.06 | 96.88 | 31.00 | open_fall ==> low_fall & high_fall |
| 50 | 3 | 2.11 | 43.06 | 93.94 | 31.00 | low_fall ==> open_fall & high_fall |
| 51 | 3 | 2.18 | 43.06 | 93.94 | 31.00 | high_fall ==> open_fall & low_fall |
| 52 | 3 | 2.18 | 43.06 | 100.00 | 31.00 | open_fall & low_fall ==> high_fall |
| 53 | 3 | 2.11 | 43.06 | 96.88 | 31.00 | open_fall & high_fall ==> low_fall |
| 54 | 3 | 2.18 | 43.06 | 96.88 | 31.00 | low_fall & high_fall ==> open_fall |
| 55 | 3 | 2.18 | 43.06 | 93.94 | 31.00 | low_fall ==> high_fall & close_fall |

**Fig 4: Association Rule Generation Window in SAS 9.2**

The some of association rules generated from stock market data is shown by Figure 4. To find interesting rules confidence value is fixed on above 90% and considered only 10 rules, which are shown by Table 12.

**Table 12. Association Rules with Confidence > 90%**

| Id | Rules | Confidence (%) |
|----|-------|----------------|
| 1 | Open_rise & low_rise=>high_rise | 100 |
| 2 | Open_rise & high_rise=>low_rise | 100 |
| 3 | Low_rise & high_rise=>open_rise | 100 |
| 4 | Open_rise=>Low_rise & Close_rise | 97.44 |
| 5 | Close_rise & low_rise =>open_rise | 97.37 |
| 6 | Open_rise & Low_rise=>Close_rise | 97.37 |
| 7 | Open_rise & Close_rise=>Low_rise | 94.87 |
| 8 | Low_rise & Close_rise=>Open_rise | 94.87 |
| 9 | Open_rise=>High_rise & Close_rise | 94.87 |
| 10 | Close_rise=>Open_rise & High_rise | 92.5 |

The extracted rules (shown in Table 12) can be analyzed in a following manner. When open price of shares rises and also low price rises than high price of shares also rises. (According Association Rule Id-1) with 100% confidence that means investing on particular month is beneficial. In the same way when low price and high price of shares are rises than open price also rises. (According Association Rule Id-3) with 100% confidence that means market will open at high values from last month and investing on particular month would be beneficial. Finally we can conclude that predicting the stock market trends the association rules are very useful. These rules help us to know the market conditions and guide us when to invest in the market. The code window which is generated during ARM process is shown in Figure 5.
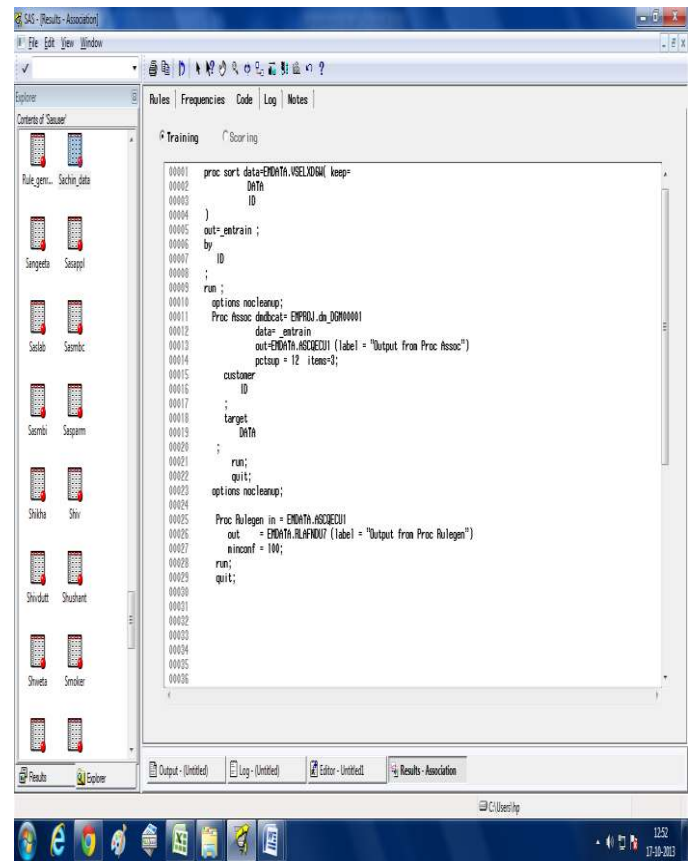


**Fig 5: Association Rule Code Window in SAS 9.2**

## 5. CONCLUSIONS

In this work we applied Association Rule Mining (ARM) approach. This method helps to learn investment plan for stock brokers and investors and understanding the market conditions. ARM has wide range of applicability in stock market research, but problem is that sometimes it does not compute results in reasonable time i.e. when number of records grows rapidly, than number of rules also grows. Some of them are redundant and infrequent. Stock market nature is highly susceptible. To increase performance of prediction there is need of method which can read multiple records simultaneously at a time and give better results. Neural Network (NN) methods have great ability to take multiple variables in relationship and can train multiple records simultaneously. In future more variables like (fundamental and technical) would be considered for analysis and financial performance of share market would be compared with other stock exchanges like China Stock Exchange (CSE), Brazil Stock Exchange (BSE).

## REFERENCES

[1] Abdoh Tabrizi.H, and Jouhare H. "The Investigation of Efficiency of stock price index of T.S.E", *Journal of Financial Research;* Vol.13, PP. 11-12, 1996.

[2] Agrawal R., and Srikant R. "Fast Algorithms for Mining Association Rules in Large Databases", *In Proc. 20th VLDB,* PP. 478-499, Sept. 1994.

[3] Agrawal R., Imeielinski T., and Swami A. "Mining Association Rules between Sets of Items in Large Databases", *Proceedings of the ACM SIGMOD*

*Conference on Management of Data",* Washington, D.C., PP. 207-216 ,1993.

[4] Argiddi Rajesh V, and Apte Sulabha S "Fragment Based Approach to Forecast Association Rules from Indian IT Stock Transaction Data", *International Journal of Computer Science and Information Technologies(IJCSIT),* Vol. 3, Issue (2), PP. 3493-3497, 2012.

[5] Bayardo R.J. "Effcient Mining Long Patterns from Databases", *SIGMOD,* PP. 85-93, 1998.

[6] Borisov Alexander "Rule Induction for Identifing multilayer Tool Commonalities" , *IEEE Transactions on Semiconducter Manufacturing ,* Vol. 24, PP. 197-201, May 2011.

[7] Brown and Jennings, on Technical Analysis "*The Review of Financial Studies",* Vol. 2, Issue (4), PP. 527-551, 1989.

[8] Connolly, and Begg T. C. "Database Systems: A Practical Approach to Design, Implementation and Management", *Addison-Wesley Longman,* ISBN 0201342871, 1998.

[9] Das Ambika Prasad "Security analysis and portfolio Management", *I.K. International Publication,* 3$^{rd}$ Edition, New Delhi (India), 2008.

[10] Fama Eugene F. "Random Walks in Stock Market Prices", *Financial Analysts journal,* Vol. 51, Issue (1), PP. 75-80, January -Febaruary, 1995.

[11] Hajizadeh E., Ardakani H., and Shahrabi J. "Application of Data Mining Techniques in Stock Markets: A Survey", *Journal of Economics and International Finance,* Vol. 2, Issue (7), PP. 109-118, July 2010.

[12] Han J., and Kamber M. "Data Mining: Concepts and Techniques", *Morgan Kaufmann,* 2nd edition, San Francisco, CA, 2006.

[13] Khan Aurangzeb, and Khan Khairullah "Frequent Patterns Mining of Stock Data Using Hybrid Clustering Association Algorithm", *University Technology PETRONAS,* 2009.

[14] Liu B., Hsu W., and M.A. Y. "Mining Association Rules with Multiple Minimum Supports", *International Conference on Knowledge Discovery and Data Mining,* PP. 337-341, 1999.

[15] Nandagopal S., Arunachalam V.P., and Karthik S. "A Novel Approach for Mining Inter-Transaction Itemsets", *European Scientific Journal ,* Vol. 8, PP. 14-22, 2012.

[16] Paranjape Preeti, and Deshpande Umesh "A Stock Market Portfolio Recoommender System Based on Association Rule Mining", *Journal of Applied Soft Computing,* Vol. 13, Issue (2), PP. 1055-1063, Febrarury 2013.

[17] Saeedmanesh M., izadi T., and Ahvar E. "A Hybrid Data Mining Technique for Stock Exchange Prediction", *International Multi Conference,* China, 2002.

[18] Srikant R., Vu Q., and Agrawal R. "Mining Association Rules with Item Constraints", *In Proc. 1997 International Conference Knowledge Discovery and Data Mining (KDD '97'),* Newport Beach, CA, PP. 67-73, Aug. 1997.

[19] Srisawat "An Application of Association Rule Mining Based on Stock Market", *3$^{rd}$ International Conference on Data Mining and Intelligent Information Technology Applications (ICMIA) ,* 24-26 Oct., 2011.

[20] Stock Market Information Available at "http://www.sebi.gov.in".

[21] Stock Market Dataset Available "http://www.bseindia.com".

[22] Tan P., Kumar V., and Shrivastava J. "Selecting the Right Interesting Measure for Association Patterns", *Inf. Syst.,* Vol. 29, Issue (4), PP. 293-331, 2004.

[23] Toivonen H. "Sampling Large Databases for Association Rules", *In Proc. 22 nd VLDB,* PP. 134-145, Sept. 1996.