# An Analysis of Data Mining Applications for Fraud Detection in Securities Market

S.Dhanalakshmi[1], C.Subramanian[2]
[1]Research Scholar, Department of Computer Applications, Dr. M.G.R. Educational and Research Institute University, Chennai
[2]Assistant professor, Department of Computer Applications, Dr. M.G.R. Educational and Research Institute University, Chennai
Email: dhanu212@gmail.com , cspratha@gmail.com

Abstract-In recent securities fraud broadly refers to deceptive practices in connection with the offering for sale of securities. There are many challenges involved in developing data mining applications for fraud detection in securities market, including: massive datasets, accuracy, privacy, performance measures and complexity. The impacts on the market and the training of regulators are other issues that need to be addressed. In this paper we present the results of a Comprehensive systematic literature review on data mining techniques for detecting fraudulent activities and market manipulation in securities market. We identify the best practices that are based on data mining methods for detecting known fraudulent patterns and discovering new predatory strategies. Furthermore, we highlight the challenges faced in the development and implementation of data mining systems for detecting market manipulation in securities market and we provide recommendation for future research works accordingly.

## I. INTRODUCTION

Securities fraud, also known as stock fraud and investment fraud, is a deceptive practice in the stock or commodities markets that induces investors to make purchase or sale decisions on the basis of false information, frequently resulting in losses, in violation of securities laws.[1] Offers of risky investment opportunities to unsophisticated investors who are unable to evaluate risk adequately and cannot afford loss of capital is a central problem. Securities fraud can also include outright theft from investors (embezzlement by stockbrokers), stock manipulation, misstatements on a public company's financial reports, and lying to corporate auditors. The term encompasses a wide range of other actions, including insider trading, front running and other illegal acts on the trading floor of a stock or commodity exchange.

## II. REVIEW OF LITERATURE

### A. Canadian Securities Administrators (CSA) 2010 report

The Canadian Securities Administrators (CSA or we) are examining the mutual fund fee structure in Canada in order to see whether there are investor protection or fairness issues, and to determine whether any regulatory responses are needed to address any issues we find. This paper is intended to be a platform to begin a discussion on the current mutual fund fee structure in Canada. This discussion paper is the first step in the CSA's public consultations about this project. It

• provides an overview of the roles of the market participants in the mutual fund industry (mutual fund manufacturers and advisors who distribute the funds)
• provides an overview of the current mutual fund fee structure
• identifies some investor protection and fairness issues we think arise from the current fee structure
• provides an overview of global regulatory reforms
• describes some regulatory options the CSA could potentially consider, either alone or in combination. Some of the options would impact mutual funds or mutual fund manufacturers directly, and others would impact those who sell the product.

## III. CANADIAN MUTUAL FUND INDUSTRY PARTICIPANTS

The participants in the Canadian mutual fund industry include the mutual fund manufacturers who produce and promote mutual fund products and advisors who distribute those products to investors.

1. The mutual fund manufacturers

There are currently 103 mutual fund manufacturers in Canada. They fall into the following four categories:

**i.** Canadian banks/deposit-takers
The fund management arms of 7 Canadian chartered banks together with the Movement Desjardins in Québec currently account for 43% of mutual fund assets under management. These manufacturers largely distribute their mutual funds through their branch networks, full-service and discount brokerage networks. Most of them also distribute a separate series of securities of their mutual funds, known as the *Advisor* series, through third party advisors.

**ii.** Life insurers
While Canadian life insurance companies primarily produce and promote segregated fund products, they are also involved in manufacturing mutual funds. These manufacturers currently represent 4.6% of mutual fund assets under management.

**iii.** Independents
Independent mutual fund manufacturers are those that are not a subsidiary of one of the large deposit-taker institutions. These independents manage the largest share of industry assets and currently represent 49.4% of mutual fund assets under management.

Integrated Intelligent Research (IIR)

International Journal of Data Mining Techniques and Applications
Volume: 03 Issue: 01 June 2014, Page No. 9- 15
ISSN: 2278-2419

**iv**. Unions and Associations

The remaining 3% of mutual fund industry assets are managed by unions and associations. Mutual funds produced and promoted by these manufacturers are generally organized for specific target groups (e.g. teachers, physicians) and generally only members of those groups can buy them.

2. Current mutual fund fees

a. Sales charges

Most Canadian mutual fund manufacturers sell funds under several different purchase options. The options relate generally to the method in which the sales charges are paid. The mutual fund manufacturers set the rate of sales charges that may be payable under the various purchase options.
The different purchase options are:

• Front-end sales charge

• Under this option, investors pay a sales commission directly to the advisor at the time they buy securities of the mutual fund.

• This is often referred to as a "front-end load". The advisor's sales commission is deducted from the total amount paid by the investor, which means only the remaining amount is invested in the fund.

• While the sales commission set by the mutual fund manufacturer may  the purchase amount, investors. typically negotiate a lower sales commission with their advisor. Over the last few years, we understand that Canadian advisors

3. Deferred sales charge (DSC)

• Under this option, investors pay a sales charge at the time they redeem from the mutual fund, rather than at the time of purchase. This is often called a "back-end load". This allows the entire amount paid by the investor to be invested in the mutual fund at the time of purchase.

• The rate of the DSC payable by investors when they redeem declines the longer they hold the investment and becomes nil after a specified holding period. This is known as the "redemption schedule". The DSC paid by an investor is typically around 6% in Depending on the mutual fund manufacturer's DSC policy, the amount of the DSC an investor pays on a redemption can be based either on the original purchase price of the mutual fund securities or their current market value when they are redeemed.

• Investors can avoid DSCs by holding their mutual investment until the end of the redemption schedule or redeeming no more to another within the same fund family without a charge.

• While the investor does not directly pay a sales commission to the advisor at the time of purchase, the advisor typically receives a commission from the mutual fund manufacturer equivalent to 5% of the amount purchased. The mutual fund manufacturer will generally borrow the money necessary to pay these advisor commissions and therefore will incur financing costs. These costs are recouped by the mutual fund manufacturer through ongoing management fees charged to the fund

• DSCs paid by investors who redeem before the end of the redemption schedule are not paid to the advisor or the mutual fund, but rather to the mutual fund manufacturer or third party financing services provider that paid the advisor's sales commission at the time of purchase.

4. Low-load sales charge

Many mutual fund manufacturers offer a low-load sales charge option, which works like the DSC option described above, but on a shorter redemption schedule. The commission paid by a mutual fund manufacturer to the advisor at the time the investor purchases securities of a fund on a low-load basis typically purchase amount. Funds sold on a no-load basis do not offer any sales commission to advisors (either one paid by the investor or the mutual fund manufacturer), nor do they charge a fee at the time the investor redeems Mutual funds purchased on a no-load basis in Canada are generally bought directly from the mutual fund manufacturer or an affiliate, either of which must be a registered dealer firm.The FBI focuses its financial crimes investigations on such criminal activities as corporate fraud, securities and commodities fraud, health care fraud, financial institution fraud, mortgage fraud, insurance fraud, mass marketing fraud, and money laundering .These are the identified priority crime problem areas of the Financial Crimes Section (FCS) of the FBI.Mission**:** The mission of the FCS is to oversee the investigation of financial fraud and to facilitate the forfeiture of assets from those engaging in federal crimes .

5. White Collar Crime (WCC) National Priorities

Based upon FBI field office threat strategies and directives established by the President, the Attorney General, the Director, and the Criminal Investigative Division (CID), the following national priorities for the WCC Program (WCCP) have been established: public corruption, corporate fraud/securities fraud to include Ponzi schemes, health care fraud, FIF (to include bank failures and mortgage fraud),insurance fraud, money laundering, and mass marketing fraud .Although public corruption is a national priority within the WCCP.

## IV.    CORPORATE FRAUD

As the lead agency investigating corporate fraud, the FBI has focused its effort son cases which involve accounting schemes, self-dealing by corporate executives, and obstruction of justice. The majority of corporate fraud cases pursued by the FBI involve accounting schemes designed to deceive investors, auditors, and analysts about the true financial condition of a corporation or business entity.

## V.    SECURITIES AND COMMODITIES FRAUD

The continued uncertainty and volatility of today's financial markets could be measured by the Dow Jones Industrial Average movement from 12,681 on July 22, 2011, to10, 655 on October 3, 2011. As a result of such tumultuous markets,

Integrated Intelligent Research (IIR)

International Journal of Data Mining Techniques and Applications
Volume: 03 Issue: 01 June 2014, Page No. 9- 15
ISSN: 2278-2419

the FBI witnessed a steady rise in securities and commodities frauds as investors sought alternative investment opportunities.

Investment Fraud: These schemes, sometimes referred to as High YieldInvestment Fraud, involve the illegal sale or purported sale of financial instruments. Financial instruments are defined broadly as any contract that gives rise to a financial asset of one entity and a financial liability or equity instrument to another entity.

The following are definitions of the most common investment fraud scheme variations:

Ponzi Schemes - A Ponzi scheme is an investment fraud that involves thepayment of purported returns to existing investors from funds contributed by new investors .Ponzi schemes often share common characteristics such as offering overly consistent returns ,unregistered investments, high returns with little or no risk, or secretive or complex strategies.

Affinity Frauds - Perpetrators of affinity frauds take advantage of the ten of people to trust others with whom they share similarities such as religion or ethnic identity to gain their trust and money.

Pyramid Schemes - In pyramid schemes, as in Ponzi schemes, money collected from new participants is paid to earlier participants. In pyramid schemes, however, participants receive commissions for recruiting new participants into the scheme.

A. Alberta Securities Commission (ASC) 2010 report

The Alberta Securities Commission (ASC) is pleased to publish its 23rd annual Corporate Finance Disclosure Report (Report). The Report continues to be an important vehicle to share with market participants our observations on public disclosure provided by Alberta reporting issuers (RIs).

B. Liquidity and Capital Resources

In last year's report we noted several areas relating to liquidity and capital resources where RI disclosures could be improved. These areas included financial statement, MD&A and prospectus disclosures regarding financial covenants, risks, and working capital.

C. Overview of the Research work

Analysis the securitites Market using Data Mining:
Data mining is the process of discovering patterns in large data sets involving methods at the database systems. The overall goal of the data mining process is to extract information from a data set and transform it into an understandable structure for further use. Aside from the raw analysis step, it involves database and data management aspects, data preprocessing, model and inference considerations, interestingness metrics, complexity considerations, post-processing of discovered structures, visualization, and online updating.

D. Relational Databases

A relational database is a collection of tables, each of which is assigned a unique name. Each table consists of a set of attributes (*columns* or *fields*) and usually stores a large set of tuples (*records* or *rows*). Each tuple in a relational table represents an object identified by a unique *key* and described by a set of attribute values.

Data mining involves six common classes of tasks:
• Anomalydetection(Outlier/change/deviation detection) – The identification of unusual data records, that might be interesting or data errors that require further investigation.
• Association rule learning (Dependency modeling) – Searches for relationships between variables. For example a supermarket might gather data on customer purchasing habits. Using association rule learning, the supermarket can determine which products are frequently bought together and use this information for marketing purposes. This is sometimes referred to as market basket analysis.
• Clustering – is the task of discovering groups and structures in the data that are in some way or another "similar", without using known structures in the data.
• Classification – is the task of generalizing known structure to apply to new data. For example, an e-mail program might attempt to classify an e-mail as "legitimate" or as "spam".
• Regression – attempts to find a function which models the data with the least error.
• Summarization – providing a more compact representation of the data set, including visualization and report generation.

E. Object-Relational Databases

Object-relational databases are constructed based on an object-relational data model. This model extends the relational model by providing a rich data type for handling complex objects and object orientation.
These patterns and trends can be collected and defined as a *data mining model*. Mining models can be applied to specific scenarios, such as:
• Forecasting: Estimating sales, predicting server loads or server downtime
• Risk and probability: Choosing the best customers for targeted mailings, determining the probable break-even point for risk scenarios, assigning probabilities to diagnoses or other outcomes
• Recommendations: Determining which products are likely to be sold together, generating recommendations
• Finding sequences: Analyzing customer selections in a shopping cart, predicting next likely events
• Grouping: Separating customers or events into cluster of related items, analyzing and predicting affinities
Temporal Databases, Sequence Databases, Time-Series Databases and Spatial Databases
• A temporal database typically stores relational data that include time-related attributes. These attributes may involve several timestamps, each having different semantics.
• A sequence database stores sequences of ordered events, with or without a concrete notion of time. Examples include

Integrated Intelligent Research (IIR)

International Journal of Data Mining Techniques and Applications
Volume: 03 Issue: 01 June 2014, Page No. 9- 15
ISSN: 2278-2419

customer shopping sequences, Web click streams, and biological sequences.

• A time-series database stores sequences of values or events obtained over repeated measurements of time (e.g., hourly, daily, weekly). Examples include data collected from the stock exchange, inventory control, and the observation of natural phenomena (like temperature and wind).

• Spatial databases contain spatial-related information. Examples include geographic (map) databases, very large-scale integration (VLSI) or computed-aided design databases, and medical and satellite image databases.

Data mining works

While large-scale information technology has been evolving separate transaction and analytical systems, data mining provides the link between the two. Generally, any of four types of relationships are sought:

• Classes: Stored data is used to locate data in predetermined groups. For example, a restaurant chain could mine customer purchase data to determine when customers visit and what they typically order. This information could be used to increase traffic by having daily specials.

• Clusters: Data items are grouped according to logical relationships or consumer preferences. For example, data can be mined to identify market segments or consumer affinities.

• Associations: Data can be mined to identify associations. The beer-diaper example is an example of associative mining.

• Sequential patterns: Data is mined to anticipate behavior patterns and trends. For example, an outdoor equipment retailer could predict the likelihood of a backpack being purchased based on a consumer's purchase of sleeping bags and hiking shoes.

F. Data mining consists of five major elements:

• Extract, transform, and load transaction data onto the data warehouse system.

• Store and manage the data in a multidimensional database system.

• Provide data access to business analysts and information technology professionals.

• Analyze the data by application software.

• Present the data in a useful format, such as a graph or table.

Different levels of analysis are available:

• Artificial neural networks: Non-linear predictive models that learn through training and resemble biological neural networks in structure.

• Genetic algorithms: Optimization techniques that use processes such as genetic combination, mutation, and natural selection in a design based on the concepts of natural evolution.

• Decision trees: Tree-shaped structures that represent sets of decisions. These decisions generate rules for the classification of a dataset. Specific decision tree methods include Classification and Regression Trees (CART) and Chi Square Automatic Interaction Detection (CHAID) . CART and CHAID are decision tree techniques used for classification of a dataset. They provide a set of rules that you can apply to a new (unclassified) dataset to predict which records will have a given outcome. CART segments a dataset by creating 2-way

splits while CHAID segments using chi square tests to create multi-way splits. CART typically requires less data preparation than CHAID.

• Nearest neighbor method: A technique that classifies each record in a dataset based on a combination of the classes of the $k$ record(s) most similar to it in a historical dataset (where $k$ 1). Sometimes called the $k$-nearest neighbor technique.

• Rule induction: The extraction of useful if-then rules from data based on statistical significance.

• Data visualization: The visual interpretation of complex relationships in multidimensional data. Graphics tools are used to illustrate data relationships.

G. Cluster analysis

Cluster analysis or clustering is the task of grouping a set of objects in such a way that objects in the same group (called a cluster) are more similar (in some sense or another) to each other than to those in other groups (clusters).

a. Centroid based clustering

In centroid-based clustering, clusters are represented by a central vector, which may not necessarily be a member of the data set.When the number of clusters is fixed to k, $k$-means clustering gives a formal definition as an optimization problem: find the $k$ cluster centers and assign the objects to the nearest cluster center, such that the squared distances from the cluster are minimized.

b. Connectivity based clustering (hierarchical clustering)

Connectivity based clustering, also known as *hierarchical clustering*, is based on the core idea of objects being more related to nearby objects than to objects farther away. These algorithms connect "objects" to form "clusters" based on their distance.

c. Distribution-based clustering

The clustering model most closely related to statistics is based on distribution models. Clusters can then easily be defined as objects belonging most likely to the same distribution.

Distribution-based clustering is a semantically strong method, as it not only provides you with clusters, but also produces complex models for the clusters that can also capture correlation and dependence of attributes.

d. Density-based clustering

In density-based clustering, clusters are defined as areas of higher density than the remainder of the data set. Objects in these sparse areas - that are required to separate clusters - are usually considered to be noise and border points.

e. Algorithms for mining data streams

Frequent Pattern Mining

(1) Mining Multiple Datasets

In many situations, such as in a data warehouse, the user usually has a view of multiple datasets collected from different data sources or from different time points.

(2) Mining Data Streams

Integrated Intelligent Research (IIR)

International Journal of Data Mining Techniques and Applications
Volume: 03 Issue: 01 June 2014, Page No. 9- 15
ISSN: 2278-2419

In recent years, database and data mining communities focus on a new model of data processing, where data arrives in the form of continuous streams. Because it is not feasible to store all data, it is quite challenging to perform the traditional data mining operations in a streaming environment. Our current and proposed research focuses on many of the challenges associated with mining streaming data. Particularly, we are driven by 1) the need to have low memory requirements, as stream mining is often carried out in small and hand held devices that do not have much memory, 2) some variants of the mining that may be desirable in a streaming environment, and 3) the need for having high accuracy and confidence in the results generated.

(3) New Data Structures for Maintaining Frequent Itemsets When mining frequent itemsets on data streams, or
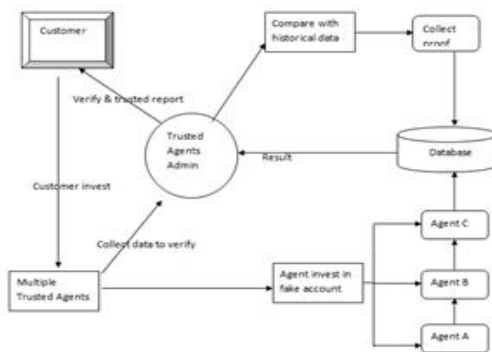
## VI. METHODOLOGY

Architecture Diagram



Figure-1

Online transaction processing or OLTP is a class of information systems that facilitate and manage transaction-oriented applications, typically for data entry and retrieval transaction processing.

Customer can buy the shares through the Multiple Trusted Agents. The buying details is stored in separate databases. Multiple Agents is sending the datas to the Admin Agent. In between time the Multiple Agents can misuse the amounts. i.e., Agent invest in fake account. Then the fake investment is stored in separate databases. So Admin agent can check original collected proof which is sent by customers and the fake investment details.Then fake investment details are shown in admin agent.

### A. Trusted agent Detector

Detector is the assignment of a label to a given input value. An example of pattern recognition misclassification, which attempts to assign each input value to one of a given set of *classes* (for example, determine whether a given email is "spam" or "non-spam"). However, pattern recognition is a more general problem that encompasses other types of output as well.

caching the results of mining operators in CMT queries, there are several common operations.

Efficient and Exact K-Means Clustering on Very Large Datasets-K-means requires several passes on the entire dataset, which can make it very expensive for large disk-resident datasets. In view of this, a lot of work has been done on various approximate versions of k-means, which require only one or a small number of passes on the entire dataset.Efficient and Effective Decision Tree Construction on Streaming Data Decision tree construction is a well studied problem in data mining. Recently, there has been much interest in mining streaming data. Domingos and Hulten have proposed a one-pass algorithm for decision tree construction. Their work uses Hoeffding inequality to achieve a probabilistic bound on the accuracy of the tree constructed.

Identify Fake AgentAgent administrator will veify all the data which is going out from the network and compare with original dataset and if there is any mismatch in the data with original dataset then it is considered as fake transcation and it will easily track the fake agent from the multiple agent groups. Finally, the fake agent will be marked in black list.

we annotated papers based on the data mining approach (i.e. supervised, unsupervised and semisupervised learning), data mining method (eg. classification, outlier detection, social network analysis), data mining challenges, and issues.

## VII. RESULT

### A. High Frequency Trading

High Frequency Trading (HFT) is an evolution of securities market that adopts sophisticated algorithms to automaticallfay analyze and react to market data in milliseconds. It is estimated that HFT accounts for 35% of the stock market trades in Canada and 70% of the stock trades in USA [18].

HFT strategies can be divided into five categories:

1) Passive Rebate Arbitrage: providing liquidity and receiving incentives from exchanges similar to market makers that is one of the major HFT strategies
2) Latency Arbitrage: making profit through buying and selling the same security between domiciles (i.e. Inter-listed Arbitrage) or between domestic marketplaces (i.e. Intra-listed Arbitrage)
3) Information Arbitrage: making profit through buying and selling the same security at a higher price by taking advantage of "mispricing between the various forms of a tradable index"
4) Statistical Arbitrage: applying statistical methods to detect trading patterns for *relative value trading*
5) Market Structure Trading: making profit through trading opportunities that are a consequence of the new market structure such as flash orders, locked markets, and dark pools While HFT supports liquidity and contributes to price formation in the market, it might have negative impacts in adverse market conditions. Regulators have been considering trading obligations and supervision on HFT especially after the May 6 flash crash. Both growth and impact of HFT in stock market have brought great interests of industry and thus

Integrated Intelligent Research (IIR)

International Journal of Data Mining Techniques and Applications
Volume: 03 Issue: 01 June 2014, Page No. 9- 15
ISSN: 2278-2419

requires regulators to establish an environment to support fair and orderly trading market. Unlike traditional trading, HFT is not subject to significant trading obligations and there is very little public information regarding fraudulent patterns and activities of HFT systems. Data mining techniques can be employed to identify fraudulent activities and predatory strategies in HFT.

### B. Massive Data

The datasets in securities market are huge. There are over 2700 securities listed in NAZDAQ and Super Montage (NASDAQ's trading platform) facilitates more than 5000 transactions per second. Similarly, the number of transactions show a significant growth from 250,000 to 5 million orders per day a few years ago, to 700,000 to 40 million orders each day3. The rate of growth in the amount of data is rapidly increasing due to changes in trading strategies by both buy-side and sell-side firms. Thus, there is extra pressure on industry to accommodate faster trading systems and on regulatory organizations to adapt to new strategies.
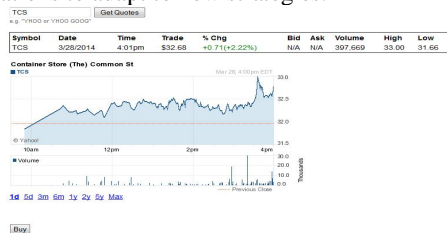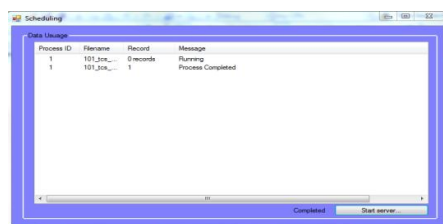


Figure-2

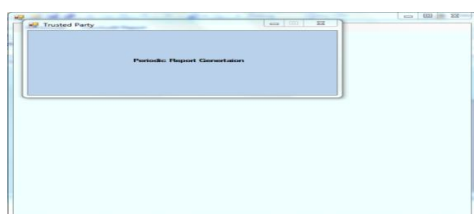

Figure-3



Figure-4



Figure-5



Figure-6

Company name is given by customer for searching the share details. Then it will show all details of shares like date, trade, volume, high, low and chart.If customer needs to buy the share of company then click the Buy button. Then customer can buy the number of shares.When the customer buy the share, that details stored in Agent side. Afterthat the agent can change the amount as duplicate before sending to the company. when the server is started, the share details stored in database.When the company side is report is generated, the company can detect the fraud in share values which is marked by red and green colour. The Fraud shares is filled by red. Green colour represented the non fraud share values.

## VIII. CONCLUSION

The significant growth of the capital market due to increasing interests in investing in securities market requires regulatory organizations to expand their efforts to ensure a fair and orderly market for the participants. Data mining methods are effective in detecting different types of fraud in securities market. In this paper we reviewed the most significant data mining methods that are applied to monitoring securities market and detecting market manipulation. We highlighted numerous challenges that are involved in developing data mining methods for detection and prevention of fraudulent activities in securities market. Some of the challenges in designing and developing data mining methods include massive datasets, different sources and forms of data and using appropriate performance measures to evaluate the method. We also provided directions for new research in this field.

### REFERENCES

[1] Han, J. and Kamber M. (2006), Data Mining Concepts and Techniques, Morgan Kaufmann, pp.4-27.
[2] Bose, N.K. and Liang, P. (1996). Neural Networks Fundamental With Graphs Algorithms And Applications. Mcgraw-Hill: New York, NY.
[3] Delmater R., and Handcock M. (2001), 'Data Mining Explained: A Manager's Guide to Customer- Centric Business Intelligence', Digital Press, Boston.
[4] J. Neville, D. Jensen, J. Komoroske, K. Palmer, H. Goldberg. "Using Relational Knowledge Discovery to Prevent Securities Fraud", 11th ACM SIGKDD International Conference on Knowledge Discovery in Data mining series, pp. 449-458, 2005
[5] Guilford, J. P. (1967). The nature of human intelligence. New York: McGraw-Hill.
[6] Hertz, J.K.( 1991). Introduction to the Theory Of Neural Computation. Addison –Wesley: New York, NY.
[7] Phua, V. Lee, K. Smith, and R. Gayler. A comprehensive survey of data mining-based fraud
[8] detection research. Artificial Intelligence Review, 2005
[9] Caruana, R. & Niculescu-Mizil, A. (2004). Data Mining in Metric Space: An Empirical Analysis of Supervised Learning Performance Criteria. Proc. of SIGKDD04, 69-78

Integrated Intelligent Research (IIR)

International Journal of Data Mining Techniques and Applications
Volume: 03 Issue: 01 June 2014, Page No. 9- 15
ISSN: 2278-2419

[10] T. Fawcett, F. Provost. "Activity monitoring: Noticing Interesting Changes in Behavior", In proceedings of SIGKDD99, pp. 53-62, 1999

[11] S. Viaene, R. Derrig, G. Dedene. "A Case Study of Applying Boosting Naive Bayes to Claim Fraud Diagnosis", IEEE Transactions on Knowledge and Data Engineering, vol. 16, pp. 612-620, 2004

[12] W. Lee, D. Xiang. "Information-theoretic Measures for Anomaly Detection", In proceedings of 2001 IEEE Symposium on Security and Privacy, pp. 130-143, 2001

[13] K. Yamanishi, J. Takeuchi, G. Williams, P. Milne. "On- Line Unsupervised Outlier Detection Using Finite Mixtures with Discounting Learning Algorithms", Journal of Data Mining and Knowledge Discovery, vol. 8, pp. 275-300, 2004

[14] P. Burge, J. Shawe-Taylor. "An Unsupervised Neural Network Approach to Profiling the Behavior of Mobile Phone Users for Use in Fraud Detection", Journal of Parallel and Distributed Computing, vol. 61, pp. 915-925, 2001

[15] R. Bolton, D. Hand. "Unsupervised Profiling Methods for Fraud Detection", In proceeding of Credit Scoring and Credit Control VII, 2001

[16] P. Chapman, J. Clinton, R. Kerber, T. Khabaza, T. Reinartz, C. Shearer, and R. Wirth. "CRISP-DM 1.0 Step-by-step data mining guide", The CRISP-DM consortium, 2000