

# Online Credit Card Fraudulent Detection Using Data Mining

S.Aravindh<sup>1</sup>, S.Venkatesan<sup>2</sup> and A.Kumaravel<sup>3</sup>

<sup>1&2</sup> *Department of Computer Science and Engineering, Gojan School of Business and Technology,  
Chennai - 600 052, Tamil Nadu, India*

<sup>23</sup> *Department of Computer Science and Engineering & Information Technology,  
Bharath University, Chennai - 600 073, Tamil Nadu, India*

E-mail: aravindhgojan@gmail.com

(Received on 10 July 2012 and accepted on 18 September 2012)

**Abstract** – As e-commerce sales continue to grow, the associated online fraud remains an attractive source of revenue for fraudsters. These fraudulent activities impose a considerable financial loss to merchants, making online fraud detection a necessity. The problem of fraud detection is concerned with not only capturing the fraudulent activities, but also capturing them as quickly as possible. This timeliness is crucial to decrease financial losses. In this research, a profiling method has been proposed for credit card fraud detection. The focus is on fraud cases which cannot be detected at the transaction level. In the proposed method the patterns inherent in the time series of aggregated daily amounts spent on an individual credit card account has been extracted. These patterns have been used to shorten the time between when a fraud occurs and when it is finally detected, which resulted in timelier fraud detection, improved detection rate and less financial loss.

**Keywords:** Fraud Detection, Aggregation, Profile, Credit Card, Time Series

## I. INTRODUCTION

Nowadays fraud detection is a hot topic in the context of electronic payments. This is mostly due to considerable financial losses incurred by payment card companies for fraudulent activities. According to a Cyber Source study conducted in 2010, the percent of payment fraud lost in the United States and Canada was \$3.3 billion in 2009 which is a considerable number [1].

A good fraud detection system should be able to identify the fraudulent activities accurately and also as quickly as possible. Fraud detection approaches can be divided into two main groups: misuse detection and anomaly detection. A misuse detection system is trained on examples of normal and fraudulent transactions. So they can only recognize known frauds. While an anomaly detection system is trained only on normal transactions and they have a potential to detect novel frauds. Difficult access to labeled data and the evolving nature of fraudulent activities, leads to more concentration on anomaly detection techniques. In these techniques the

cardholder's profile is constructed based on his normal spending habits and any inconsistency with regards to this normal profile is considered as a potential fraud. The problem with this approach is the large number of false alarms due to normal changes in cardholder's behavior.

Using anomaly detection techniques for fraud detection involves constructing an efficient profile which considers all aspects of a card holder behavior. Usually a fraudster is not familiar with the spending habits of a card holder, while try to get the most profit from a stolen card. Hence they tend to perform high value transactions, which usually have a different characteristic from the normal card holder transactions.

In this context the transactional profile can reveal the frauds. Many researches consider this kind of fraudulent activities and construct a transactional profile [7], [8], [9] and [10]. But more cautious fraudsters try to follow the normal behaviors of card holder or perform low value transactions in short time intervals. In this case the frequency or volume of transactions is a much better indicator of fraud compared to the characteristics of each individual transaction. For instance, in these frauds the total number or total amount spent on a credit card over a specific time window increases. A few researches consider this type of frauds and construct an aggregated profile. The problem with this approach is the late detection because the system has to wait until the end of the aggregation period before it can make a decision. This problem seems more crucial when the aggregation period is considerable. Also some useful information like the order of data is lost during the aggregation. This order of data is another aspect of a cardholder behavior which can be used to detect some types of frauds.

In this research, we approach the credit card fraud detection problem with an improved aggregated profile. For this purpose the sequence of aggregated daily amounts spent on an individual card holder in a time window has been considered. Then the inherent patterns in these time series

have been extracted to shorten the time between when a fraud occurs and when it is finally detected. Indeed we have taken advantage of the order of data to timelier fraud detection. We demonstrate that the proposed approach leads to improved detection rate and timeliness while it decreases the cost involved in some circumstances.

## II. RELATED WORKS

Misuse detection and anomaly detection are the two main approaches used for credit card fraud detection. The emphasis on misuse detection approaches is usually upon applying classification methods at transaction level. For a recent survey of applying misuse detection techniques in the area of credit card fraud detection see [2], [3], [4] and [5]. In these researches various classification methods like neural networks, decision trees, logistic regression and support vector machine have been used and compared against each other in the area of credit card fraud detection. Also a recent research in [6] various classification methods have been applied on aggregated transactions. This research has demonstrated that aggregated values are a better indicator of frauds in some circumstances.

Among the researches which have been conducted on credit card fraud detection we have concentrated on the ones which apply anomaly detection techniques, the so called behavioral or profile-base techniques.

Typically they have constructed a cardholder profile based on normal training data and then tried to detect fraudulent activities based on the inconsistencies with the normal behavior. Most of these researches have applied data mining techniques like clustering and association rules to construct a transactional profile. For instance, in [7] self organization map has been used to cluster customer transactions. The density of each cluster is the basis of distinction between normal and rare behavior of customers which can be used for detect suspicious activities. Also in [8] DBSCAN, which is a density based clustering algorithm, has been used to create clusters of customer transactions and build a transactional profile. An example of using association rules can be found in [9]. In this research recent transactions of a customer have been dynamically profiled using association rules, to indicate how unusual a new transaction is. The word recent is defined by a sliding window.

In a few researches in this area, the sequence of transactions has been considered for building customer profiles. An example of which can be found in [10]. In this

research a Hidden Markov Model for each customer has been built during the training phase based on a sequence of transaction amounts. When a new transaction arrives, a new sequence is constructed by dropping the first member of the old sequence and appending the new transaction at the end. If the new sequence is not accepted by the trained model, it is considered as fraud. In another research in which combines anomaly and misuse detection techniques, normal and fraudulent sequences of quantized transaction amounts have been formed to capture the cardholder behavior. Then a sequence alignment technique has been used to measure the similarity between a new sequence and the training model. In a different research in for each target cardholder, sequences of daily transaction amounts have been compared against the other cardholders to find the k nearest ones. These similar sequences have been grouped to form the peer group of that cardholder. If the future sequences of that cardholder deviate from its peer group, a fraud alarm is raised. The basis of this research is the assumption that when a group of cardholders are behaving similarly until a specific time, it is very likely that they will continue to have the same behavior for a while.

## III. PROPOSED METHOD

In this research, we have explored the application of transaction sequence for the purpose of timelier credit card fraud detection. The focus in this work is on fraud cases which cannot be detected at transaction level. Indeed, we have proposed an improved aggregated profile which exploits the inherent patterns in time series of transactions. Some extensive modeling on real data reveals strong weekly and monthly periodic structure in cardholder spending behavior. Based on these observations we believe that instead of looking at individual transactions, it makes more sense to look at sequences of transactions. But it is impractical to consider the entire series of cardholder transactions because of the high dimensionality of this data. So we model the time with a sequence of aggregated transactions which can reduce the dimensionality. Also aggregated transactions are more robust to minor shift in cardholder behavior.

To form the time series, the total amount of transactions in each day of year has been calculated. Then the ordered series of these aggregated values form the time series. Like the aforementioned researches [13], [14] which consider 7 days for aggregation, we form 7-day time series. So each time series consists of 7 dimensions each of which corresponds to the total amount of transactions in one day.

As it is mentioned before, based on some observation on real data, there are some periodic structures in transactions, so we expect to find similar trends in yearly 7-day time series. Also since the first year of each year is considered as the starting point of the 7-day period of that year, the time series for each year would be different in terms of days of the week. For example one year may start on Sunday while the next year starts on Friday. This implies that for each year the 7-day time series, of a cardholder that follows a stable weekly trend, should be aligned in terms of days of the week accordingly. Furthermore, a cardholder himself may have some shift in purchasing days. Another pattern is some occasional behavior that can be seen due to holidays and occasions which are repeated in all years like the Christmas holidays. In this research we want to extract these inherent patterns in time series of aggregated transactions, and apply them to detect fraudulent activities more timely and accurately. In fact, by exploiting these patterns we can detect fraud cases before the end of an aggregation period. The details of constructing profiles and fraud detection will be explained in the following sub sections.

**A. Make Profile**

To construct a cardholder profile, his normal transactions is needed as training data. As mentioned earlier a preprocessing step is performed to build time series. Then the inherent patterns in these time series should be extracted to build an efficient profile. In this research two possible patterns are extracted from the training data in two steps. The first possible inherent pattern in a 7-day period could be following the same trend in all years. For extracting this pattern, time series have been clustered using k-means, the most popular clustering algorithm, with Euclidean distance. Since Euclidean distance is used as the similarity measure, the time series which have almost the same trend will be placed in the same cluster. After clustering, if all yearly time series for a specific 7-day period are placed in the same cluster, this period has been labeled as stable-trend period. Then all of the time series that belong to these periods are excluded from the training data and the other ones remain for further analysis in the next stage.

The Euclidean distance is very sensitive to small distortions in the time axis. If two time series are identical, but one is different slightly along the time axis, then the Euclidean distance may consider them to be very different from each other. But as it was mentioned before, the second possible inherent pattern in a period could be following the

same trend by permuting the time axis as we can see in Fig.1. So in order to find the similarity between such sequences, the time axis should be best aligned before calculating the Euclidean distance.

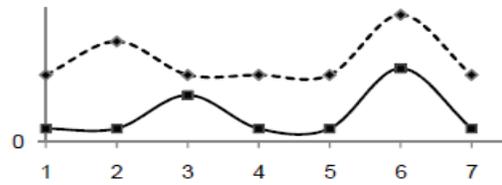


Fig.1 An example of permuted - trend time series

The remaining time series from the first stage have been clustered using this new distance, we call it permuted distance. For this purpose the k-means algorithm should be changed. Briefly k-means algorithm selects k initial points as cluster centers. Then each point is assigned to the closest center using a distance measure. When all points have been assigned, the new centers are recalculated by averaging cluster members. These steps are repeated until the centers no longer move. Usually the Euclidean distance is used as distance measure in the k-means algorithm. This should be modified for the permuted pattern. To find the distance for permuted time series, any permutation of the time axis for the first one is considered, and the Euclidean distance between all of them with the second one is calculated. Then the minimum value is selected as the distance between the two time series. Also the current averaging method for finding new centers may not produce the real average of the time series in our case, thus resulting in incorrect k-means clustering results. Figure 2 is the result of usual averaging method of the two time series while we expect the result which is shown in Fig. 3. So the time series should be aligned in time axis before calculating the average time series.



Fig. 2 Usual averaging of the two time series

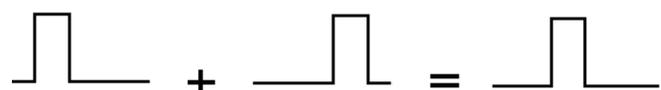


Fig. 3 Desired averaging of the two time series

The remaining time series from the first stage are clustered with this new version of k-means. As a result the time series which are almost the same after alignment in time axis are been placed in a same cluster. We labeled the 7-day periods for which all yearly time series placed in the same cluster as permuted-trend. Moreover there are some yearly occasions in which the cardholder behavior is almost the same for all years like Christmas holidays. So we can improve our profile by identifying these days in permuted-trend periods. For this purpose for all time series of these periods the best alignment for the permuted distance is found. If one day is not permuted in the best alignments, it is flagged as a stable day.

After these two stages, the remaining periods are labeled as unpredictable-trend. So at the end of the training phase we have a time series for each 7-day period of year with the specification about which groups it belongs to and which days are stable days for the second group.

### **B. Fraud Detection**

After the training phase, fraudulent activities can be detected based on the degree of deviation from the cardholder profile. For this purpose when a new transaction arrives they are accumulated to build the current period time series. Based on the type of current period in profile which can be stable-trend, permuted-trend and unpredictable-trend, the fraud detection is performed online, at the end of each day or at the end of period respectively. For the stable-trend periods, since the cardholder behavior in corresponding days are almost the same, the fraud detection can be done online. While the transactions of a day are accumulated, it is compared against the corresponding values in the profile.

Whenever this value exceeds with a ratio of  $\theta_1$  from the corresponding amount in the cardholder's profile, it indicates a fraud. For the permuted trend periods, at the end of each day, the similarity between the current time series with the corresponding one in the profile is computed. Since in the middle of a period the current time series is smaller than the corresponding one in the profile, we should consider all of the subsets of profile time series with the same length as the current time series. Then the minimum permuted distance between them is considered. If this value exceeds from a threshold  $\theta_2$ , it indicates a fraud. Considering all subsets of profile time series, the days which are flagged as stable-days should remain immovable. So at the end of each day we can say that there is some fraud among the days and we don't have to wait until the end of period. One important point is that for this group while we make the time series, whenever a

fraud case has been identified in a day, we should replace this day with the corresponding value from the profile in order to prevent the fraud value from affecting the decision for the next days of the period.

Finally, for the unpredictable-trend periods at the end of 7- day period, we compute the distance between the current time series and the corresponding one in cardholder profile and if it exceeds from a threshold  $\theta_3$ , it indicates fraud. For this group, at the end of the period we have a label which tells us there are some frauds in this period.

The best value for the parameters  $\theta_1$ ,  $\theta_2$  and  $\theta_3$  is obtained by examining the performance of the system over various values for them, and choosing the one with the best average result on all of the profiles using a tuning set.

Clearly, the proposed method improves the timeliness of fraud detection which is proved to be most effective in the stable-trend periods and the permuted-trend, consecutively. However, the mentioned method does not improve the timeliness of the fraud detection for the unpredictable-trend periods.

## **IV. EXPERIMENTAL RESULTS**

The performance of the proposed scheme has been compared with the performance of the aggregational part of the offline system proposed in [14]. In that research, the aggregated profile is constructed based on the weekly behavior of cardholders and the fraud detection is performed at the end of each week. We expect that our proposed method can increase the detection rate and improve the timeliness of that method. Also the aggregational profile proposed in [13] will be compared against our proposed method.

In [13] the model of aggregation consists of a set of descriptors for quantifying time series of cardholder behavior. These time series are built using all of the k-day periods of normal transactions. 1, 3 and 7 days periods are used for evaluation, among them we choose the 7- day one for comparison, which conforms to our approach.

### **A. Dataset**

To evaluate our work we have developed an application to generate synthetic data containing genuine and fraudulent transactions. The profile driven method has been used for generating data like the one applied in [9]. We believe that our dataset can give us a good approximation for evaluation of the proposed method because we use real scenarios to create the data. As it was mentioned before, based on some

observation on real data, there are some periodic structure in credit card transaction data and also some occasional events.

Also there are various weekly and seasonal patterns in cardholder behaviors. These real scenarios have been applied in data generation to justify the results. Also normal distribution, which is the most common observed probability distribution in many natural processes, has been used to create number and amount of transactions.

Five attributes for each transaction have been considered including year, month, week of month, day of week and amount. The first four attributes indicate the time sequence of data and the last one is a good descriptor to quantify the time series. We have created four different profiles to generate different kinds of cardholder behaviors. In the first one the cardholder has almost similar periodic behavior. In the second one the cardholder has similar behavior with some shift in the time axis. In the third profile the cardholder has an unpredictable behavior. Finally in the fourth cardholder has a mixture of different behaviors in different times. Transactions for three years are created for each cardholder as training data. Then a mixture of genuine and fraudulent transactions of one year is generated for test data. Fraudsters usually follow two different scenarios to avoid detection: high value transactions with long gaps or small value ones with short gaps. The first scenario can be detected by a transactional profile and the second one can be detected by an aggregational profile. Because we want to evaluate an aggregational profile, fraudulent activities are created based on the second scenario.

For each profile three datasets are created. The first one which contains normal transactions is a training set. The second and third ones contain a mixture of normal and fraudulent transactions. The second one is a tuning set which is used for obtaining the best values for the system parameters and the last one is a test set used for evaluating the proposed method. Table I shows the number of transactions in each dataset of the four profiles.

TABLE I CHARACTERISTICS DATASET

	<i>Profile 1</i>	<i>Profile 2</i>	<i>Profile 3</i>	<i>Profile 4</i>
Training Set	3314	5299	6709	3951
Tuning Set	1186	2129	2950	1854
Test Set	1206	2114	2968	1809

## B. Performance Measures

The transactions which are flagged by a fraud detection system include the fraudulent and normal transactions which are classified correctly (TP, TN) and the fraudulent and normal transaction flagged erroneously (FN, FP). A good fraud detection system should lead to maximum number of TP and TN and minimum number of FP and FN.

Several performance measures have been applied for fraud detection systems. The appropriate one should take into account the specific issues in fraud detection systems. In a recent research [18] the appropriate performance measures for plastic card fraud detection systems have been proposed. We have applied two measures which are proposed in that research and widely applied in recent fraud detection researches: timeliness ratio and loss function. The first one measures the speed of fraud detection and defined as the proportion of FN to F, the second one measures the cost involved. In this measure different cost consider to different error because the FNs are more serious than the FPs. We use the function used in [6] which is as (1).

$$L(s) = \frac{TP + FP + 100 * FN}{N + 100 * F} \quad (1)$$

Smaller values for these two measures indicate a better performance. Also we use a standard measure, TP%, which is the percent of TPs to all of the fraudulent transactions. Clearly higher values indicate a better performance.

## C. Optimization of Parameters

The proposed method has 3 parameters,  $\theta_1$ ,  $\theta_2$  and  $\theta_3$ . In choosing a value for these parameters, there is a tradeoff between TP% and FP%. In this work we choose the best value for each of them using TP/FP(%). The best value for each parameter is obtained by examining the performance of the system over various values of them using the tuning sets and choosing the one with the best average result on all of the profiles. As a result the values 1.4, 0.7 and 0.2 have been obtained experimentally for  $\theta_1$ ,  $\theta_2$  and  $\theta_3$  respectively.

## D. Validation Results

First we study the performance of our aggregation method to the one proposed in [14]. As we can see in Fig. 4 TP% of our proposed method is better than the one proposed in [14]. Also Fig. 5 and 6 indicate that the cost and timeliness of our proposed method is better too. It can be clearly seen from these figures that when a cardholder follows an almost stable trend in the corresponding times of the years, the case which

has occurred in the first profile, the performance of the system increases significantly. It is due to the fact that the fraudulent activities can be detected in real time. As a result, more frauds can be detected by the system, in a timelier manner and with less cost. In the second test case which indicates a cardholder with the permuted behavior, the performance of the system is slightly better, because the fraudulent activities can be detected at the end of each day. But if the cardholder has an unpredictable behavior, which is simulated in the third case, the performance of our method is almost the same as the one proposed in [14] because there is no pattern in the cardholder behavior which can be used for timelier detection and the fraudulent activities can be detected at the end of 7-day periods.

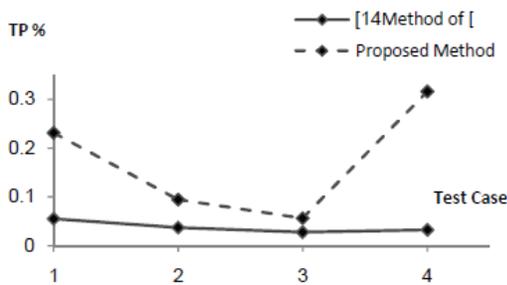


Fig. 4 TP% of four test cases for aggregational part of method proposed in [14] against our method

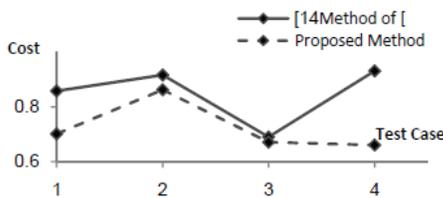


Fig. 5 Cost of four test cases for aggregational part of method proposed in [14] against our method

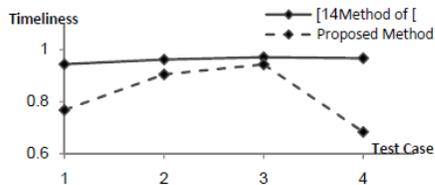


Fig. 6 Timeliness of four test cases for aggregational part of method proposed in [14] against our method

Next the performance of the proposed method is compared against the aggregation method proposed in [13]. In that research the procedure for detecting fraudulent activities is run at the end of each day, considering 7 days before the current day. It can be seen from Fig. 7, 8 and 9 that almost

the same results are obtained as the previous experiment. One of the underlying reasons for this improved result may be considering seasonal behavior in the proposed method. Also the same reasons as discussed for the previous experiment apply to this experiment as well.

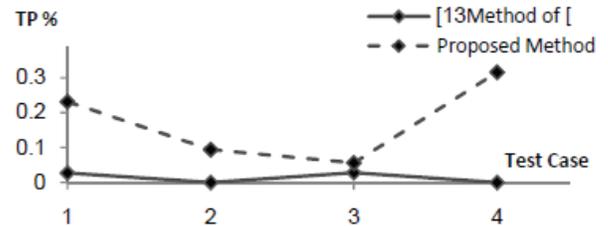


Fig. 7 TP% of four test cases for aggregational part of method proposed in [13] against our method

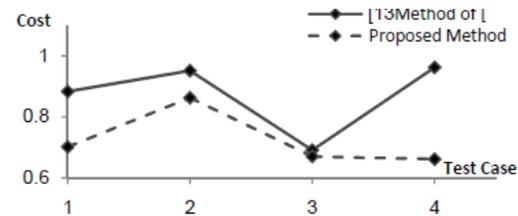


Fig. 8 Cost of four test cases for aggregational part of method proposed in [13] against our method

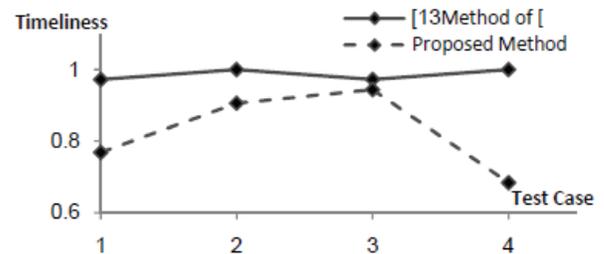


Fig. 9 Cost of four test cases for aggregational part of method proposed in [13] against our method

## V. CONCLUSION

In this paper we have addressed the general problem of credit card fraud detection using anomaly detection techniques, by exploiting the sequence of transactions in constructing cardholders' profiles. We have investigated how this affects detection performance. The focus is on fraud cases which cannot be detected at the transaction level. A new method for constructing an aggregated profile is proposed.

To this end the pattern of aggregated daily purchases of cardholders are extracted from the training data. Due to the seasonal behavior of cardholders these patterns are time dependent. Then these extracted patterns have been

used for more accurate fraud detection in a timelier manner. Experimental results show that the proposed method can improve the fraud detection in the situations where cardholders follow some purchasing patterns in corresponding times of the years.

#### REFERENCES

- [1] CyberSource; "11th Annual Online Fraud Report"; 2010. <http://forms.cybersource.com/forms/FraudReport2010NACYBSwwwQ109> last accessed on 2010/09/10.
- [2] R. Brause, L. T., and M. Hepp, "Neural data mining for credit card fraud detection," *11th IEEE International Conference on Machine Learning and Cybernetics*, Vol 7, 2008, pp.3630-3634.
- [3] R. Chen, S. Luol, X. Liang, and V.C. Lee, "Personalized approach based on SVM and ANN for detecting credit card fraud", *International Conference on Neural Networks and Brain*, 2005, pp. 810-815.
- [4] A. Shen, R. Tong, and Y. Deng, "Application of classification models on credit card fraud detection," *International Conference on Service Systems and Service Management*, June 2007, pp. 1-4.
- [5] M.F. Gadi, X. Wang, and A.P. Lago, "Comparison with parametric optimization in credit card fraud detection," *Seventh International Conference on Machine Learning and Applications*, 2008, pp. 279-285.
- [6] C. Whitrow, D.J. Hand, P. Juszczak, D. Weston, and N.M. Adams, "Transaction aggregation as a strategy for credit card fraud detection," *Data Mining and Knowledge Discovery*, Vol.18, No. 1, 2009, pp.30-55.
- [7] J. Quah and M. Sriganesh, "Real-time credit card fraud detection using computational intelligence," *Expert Systems with Applications*, Vol. 35, No. 4, 2008, pp. 1721-1732.
- [8] S. Panigrahi, A. Kundu, S. Sural, and A. Majumdar, "Credit card fraud detection: A fusion approach using Dempster-Shafer theory and Bayesian learning," *Information Fusion*, Vol. 10, No. 4, 2009, p. 9.
- [9] J. Xu, A.H. Sung, and Q. Liu, "Behaviour mining for fraud detection," *Journal of Research and Practice in Information Technology*, Vol. 39, No. 1, 2007, pp. 3-18.
- [10] A. Srivastava, A. Kundu, S. Sural, and A. Majumdar, "Credit card fraud detection using Hidden Markov Model," *IEEE Transactions on Dependable and Secure Computing*, Vol. 5, No. 1, 2008, pp. 37-48.
- [11] A. Kundu, S. Sural, and A. Majumdar, "Two-stage credit card fraud detection using sequence alignment," *Information Systems Security*, Springer Berlin / Heidelberg, 2006, pp. 260-275.
- [12] D.J. Weston, D.J. Hand, N.M. Adams, C. Whitrow, and P. Juszczak, "Plastic card fraud detection using peer group analysis," *Advances in Data Analysis and Classification*, Vol. 2, No. 1, 2008, pp. 45-62.
- [13] M. Krivko, "A hybrid model for plastic card fraud detection systems," *Expert Systems with Applications*, vol 37, no 8, 2010, pp 6070-6076.
- [14] L.Seyedhossein, M.R. Hashemi, "A hybrid profiling method to detect heterogeneous credit card frauds", *7th International ISC Conference on Information Security and Cryptology*, 2010, pp 25-32.