



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

Procedia Computer Science 3 (2011) 987–991

---

---

**Procedia  
Computer  
Science**

---

---

[www.elsevier.com/locate/procedia](http://www.elsevier.com/locate/procedia)

WCIT-2010

## Suspicious activity reporting using dynamic bayesian networks

Saleha Raza<sup>a\*</sup>, Sajjad Haider<sup>a</sup><sup>a</sup> Artificial Intelligence Lab, Institute of Business Administration, Karachi, Pakistan

---

### Abstract

Suspicious activity reporting has been a crucial part of anti-money laundering systems. Financial transactions are considered suspicious when they deviate from the regular behavior of their customers. Money launderers pay special attention to keep their transactions as normal as possible to disguise their illicit nature. This may deceive the classical deviation based statistical methods for finding anomalies. This study presents an approach, called SARDBN (Suspicious Activity Reporting using Dynamic Bayesian Network), that employs a combination of clustering and dynamic Bayesian network (DBN) to identify anomalies in sequence of transactions. SARDBN applies DBN to capture patterns in a customer's monthly transactional sequences as well as to compute an anomaly index called AIRE (Anomaly Index using Rank and Entropy). AIRE measures the degree of anomaly in a transaction and is compared against a pre-defined threshold to mark the transaction as normal or suspicious. The presented approach is tested on a real dataset of more than 8 million banking transactions and has shown promising results.

© 2010 Published by Elsevier Ltd. Open access under [CC BY-NC-ND license](http://creativecommons.org/licenses/by-nc-nd/3.0/).  
Selection and/or peer-review under responsibility of the Guest Editor.

*Keywords:* Anomaly detection; clustering; suspicious financial transactions; suspicious activity reporting; dynamic Bayesian network; anti-money-laundering

---

### 1. Introduction

Identification of suspicious financial transactions is a crucial part of anti-money laundering (AML) systems. Money laundering broadly refers to the process of taking illegally obtained money and inducting it in the cycle of financial system in a manner that disguises its origin and makes it appear to be legitimate. There are various forms of money laundering. One common method is to decompose one large transaction into large number of medium-to-small-sized transactions so that they appear normal and do not come under the radar of any monitoring entity. Certain other methods focus on keeping the illegal money continuously in flow by performing various operations (like investment in a business venture, sales and purchase of instruments or donations to charity organizations, etc). The purpose is to keep this money distant from its actual source and making its tracing extremely difficult.

Suspicious activity reporting (SAR) normally works on the basic principle that any transaction that does not comply with the normal behaviour of a customer or similar group of customers is anomalous and can be considered suspicious. Various anomaly detection techniques have been suggested in the literature to uncover these suspicious financial transactions [1-8]. Most of the AML approaches perform clustering to identify anomalies. Few techniques employ artificial neural networks (ANN) [3], decision trees [7] and support vector machine [8] to identify suspicious transactions. Hidden Markov model (HMM) [9] and Probability suffix tree (PST) [4] based approaches look for sequential anomalies. Sudjianto et al. [10], Patcha et al. [11] and Chandola et al. [12] present comprehensive surveys

---

\*Saleha Raza. Tel.: +92-213-111-677-677  
E-mail address: [saleha.raza@khi.iba.edu.pk](mailto:saleha.raza@khi.iba.edu.pk)

of various anomaly detection methods, used in different application domains, and assess their relative strength or potential weaknesses.

This paper presents a novel approach, called SARDBN (Suspicious Activity Reporting using Dynamic Bayesian Network), that employs a hybrid model of distance based clustering and dynamic Bayesian network to identify anomalies in a sequence of transactions. The core idea behind the approach is that the overall monthly transactions of a customer establish a defined pattern. This pattern closely matches among customers of similar characteristics and can be captured in the form of a probabilistic temporal model. Any deviation from this pattern is marked as suspicious. To improve accuracy of the model and to minimize false alerts, the paper suggests an anomaly scoring metric, termed AIRE (Anomaly Index using Rank and Entropy), that measures the degree of anomaly in each incoming transaction. The transaction is marked as suspicious only when its AIRE value exceeds a certain threshold. Unlike other AML approaches that work on aggregated summaries of customer transactions, like average amount or frequency, SARDBN identifies abnormalities in sequence of transactions. A transaction that apparently looks normal, with respect to the aggregates, may found to be anomalous when considering the particular sequence that it followed. The entropy based scoring mechanism reduces false positives, without compromising on true positives, and ensures that high anomalies are raised only in situations when the model has enough conclusive evidences to support them. The data set used in the experiment has more than 8 million banking transactions and the results show good predictive accuracy for SARDBN.

The rest of the paper is organized as follows. Section 2 presents SARDBN, while Section 3 provides the experimental design and results. Finally, Section 4 concludes the paper and provides directions for future research.

## 2. SARDBN

This section explains SARDBN that is designed to identify anomalies in sequence of transactions. The approach can be divided into three major steps: 1) Clustering, 2) Learning DBN and 3) Anomaly Detection. These steps have been further explained below:

### 2.1. Clustering

As different groups of customers exhibit different transactional pattern, the first step performed by SARDBN is to form clusters of customers that exhibit similar transactional pattern. The similarity in the transactional behaviour is assessed by a customer's average monthly credit and debit amounts, frequency of credit and debit transactions, and delay in two consecutive transactions. The quality of cluster is determined by inter-cluster and intra-cluster distances. The initial clustering of customers may have profound impact on the final performance of DBN – the more cohesive the clusters the better the predictive accuracy of the corresponding DBNs.

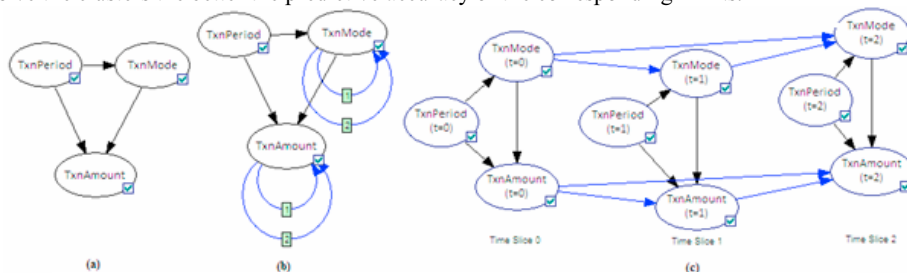


Fig. 1. (a) BN (b) DBN (c) Unrolled DBN, for financial transaction

### 2.2. Learning Dynamic Bayesian Network (DBN)

Once clustering is performed, the next step is to learn the parameters of a DBN for each cluster. The structure of these DBNs can vary in different scenarios and is devised with the aid of subject matter experts. For this study, the three variables under consideration are transaction amount (TxnAmount), mode of transaction (TxnMode) (e.g. cheque withdrawal/deposit, ATM withdrawal, salary transfer, POS payment, bank draft), and period of transaction (TxnPeriod) (i.e. start/middle/end of month). These three variables serve as random variables in a BN and are

connected as depicted in Fig.1<sup>2</sup> (a). When considering sequence of transactions and assuming dependencies within this sequence, the standard BN is converted into a dynamic Bayesian network (DBN) where each transaction belongs to a different time slice. Fig. 1(b) shows a complete DBN of order two, while Fig. 1 (c) gives the unrolled version of this DBN for three time slices. Once the structure of a DBN is finalized, its prior and conditional probabilities are learned from the dataset. For the purpose of this experiment, DBNs of order one and two are discussed in the sequel although the order can be adjusted as per the degree of interrelationship between transactions.

2.3. Anomaly Detection

After learning DBN for each cluster, the model is ready to be used for anomaly detection. During this phase, each incoming transaction, along with its last n transactions, is passed through the DBN of its respective cluster and its amount and mode are predicted. These predictions are ranked on the basis of their posterior probabilities. The one with the highest probability is assigned a rank of one, the second highest as two and so on. An anomaly index, AIRE (Anomaly Index using Rank and Entropy) has been developed that make use of these ranks and the corresponding entropy:

$$AIRE_i = (r - 1 + e * (.5 + .5k - r)) / (k - 1) \tag{1}$$

Where r and e are rank and entropy of predicted variable and k represents total number of states. The total anomaly score of a transaction is the weighted sum of AIRE value of each predicted attribute.

$$AIRE = \sum_i (w_i * AIRE_i) \tag{2}$$

Where i represents each random variable predicted from DBN and w<sub>i</sub> and AIRE<sub>i</sub> are the corresponding weight (assigned subjectively) and anomaly score.

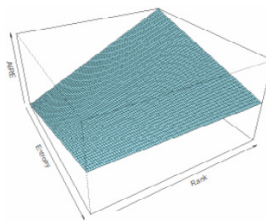


Fig. 2 – The surface showing AIRE as function of rank and entropy

The computed AIRE value is compared against a given threshold to mark the transaction as normal or suspicious. Lower rank represents better compliance of transaction with the past behaviour and hence a low anomaly. While higher rank implies poor compliance and thus high anomaly. Entropy value shows conclusiveness of the model to report anomaly with lower entropy representing better conclusiveness and vice versa. As entropy increases, probability distribution loses its decisiveness; making AIRE converging to the middle value, that is, 0.5. Fig.2 shows the shape of AIRE as a function of rank and entropy while Fig.3 depicts its application using a sequence of three transactions along with the corresponding instantiated DBN.

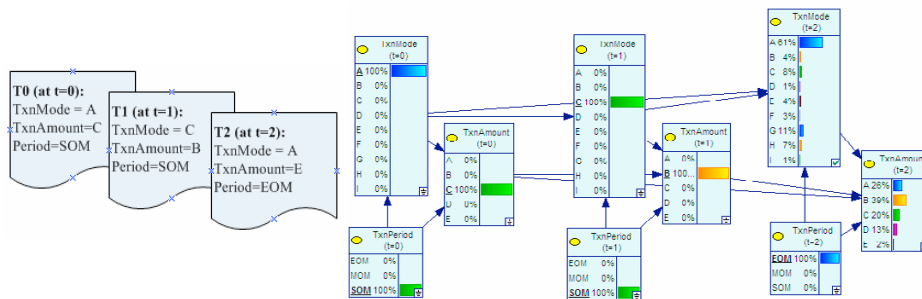


Fig. 3 - Instantiated DBN for a sequence of three transactions

<sup>2</sup> Figs. 1 and 3 are drawn using GeNIe that can be downloaded from <http://genie.sis.pitt.edu>

Table 1 summarizes the anomaly detection phase of SARDBN.

Table 1- SARDBN - Anomaly detection

---

Given  $Tx_n$ , C, DBN, V, W, DBN, n,  $\tau$ :  
 where  $Tx_n$  = Incoming transaction  
 C = customer of  $Tx_n$ , DBN = DBN of the cluster the customer C belongs to.  
 V = Set of variables to be predicted, W = Set of weights for each variable to be predicted  
 n = Order of DBN,  $\tau$  = Threshold

1. Load last n transactions ( $Tx_0 - Tx_{n-1}$ ) of customer C.
2. For  $i = 0$  to  $n-1$ 
  - a. Instantiate random variables on timeslice  $t_i$  with the attributes of transaction  $Tx_i$
3. Perform inference in the DBN.
4. For each  $v_i \in V$ 
  - a. Sort the states of  $v_i$  in descending order of their posterior probabilities.
  - b. Let  $r$  = Rank of state that matches with the actual value of  $Tx_n$   
 $p$  = Probability of state that matches with the actual value of  $Tx_n$   
 $k$  = Total number of states of  $v_i$
  - c. Calculate normalized entropy  $e_i$  of  $v_i$  as:  $e = - \sum_k p \log_2(p) / \log_2(k)$
  - d. Calculate AIRE <sub>$i$</sub>  of  $v_i$  as:  $AIRE_i = (r - 1 + e^{*(.5 + .5k - r)}) / (k - 1)$
5. Calculate total anomaly AIRE of transaction  $Tx_n$  as:  $AIRE = \sum_i w_i * A_i$
6. If  $AIRE > \tau$ , Mark  $Tx_n$  as anomalous Else, Mark  $Tx_n$  as normal

---

### 3. Experimentation and Results

SARDBN is tested on a dataset containing banking transactions of non-corporate customers. The dataset contains around 100 thousands customers incurring 8.2 million transactions over the period of a year. The following steps are performed to identify anomalies using SARDBN:

*Data preprocessing:* Data preprocessing is performed to clean the dataset and discard some of the transactions and account types that have either been inactive throughout the year or are unlikely to be involved in money laundering.

*Data segmentation:* After preprocessing, the whole dataset is divided into two parts: Training and Testing. Training part comprises of transaction from January till October and is used for clustering and parameter learning of the DBN. Test data includes transactions during November and December and is used to analyze the model for prediction and anomaly detection.

*Clustering:* Training data is clustered on the basis of the average monthly credit/debit amount and the average frequency of transactions for each customer. Fuzzy c-means algorithm [13] is used to form four clusters of customers.

*Amount Discretization:* Amounts range from very small to very large values and follow different variances in each cluster. For each cluster, amounts are discretized using unsupervised k-bins discretization method [14].

*Model Learning:* For experimental purposes, DBN of order one and two are tested in this paper. Model learning includes extracting prior and conditional probabilities from training dataset.

*Testing:* For the DBN of order two, all sequences of three consecutive transactions of the same account are loaded from test dataset and are passed through the DBN to extract their rank and entropy. These two parameters are then used to calculate AIRE for each transaction. The AIRE value is matched against a given threshold to mark the transaction as normal or suspicious. The process is repeated for the DBN of order one, with sequences of two consecutive transactions. Table 2 summarizes accuracy of SARDBN to predict TxnAmount and TxnMode. The readings are taken from a total of 120000 transactions with 30000 transactions belonging to each cluster. The results show good accuracy of SARDBN in predicting the mode and amount of an incoming transaction. The transactions that deviate extensively from the learned model have higher AIRE values and are considered anomalous. Table 3 gives number of anomalous transactions for different threshold values while Figure 4 plots a graph of number of anomalies against different threshold values. As the threshold increases, number of anomalies decreases and vice versa. The exact value of threshold may vary from situation to situation and can be decided by the subject matter expert.

Table 1 - Accuracy of SARDBN to predict TxnMode and TxnAmount

Cluster	Correctly Predicted TxnMode		Correctly Predicted TxnAmount	
	1 <sup>st</sup> Order DBN	2 <sup>nd</sup> Order DBN	1 <sup>st</sup> Order DBN	2 <sup>nd</sup> Order DBN
1	62 %	71%	73%	79%
2	89 %	90%	94%	95%
3	64 %	51%	67%	69%
4	81 %	84%	79%	80%

Table 3 – Number of Anomalies

Threshold	Number of Anomalies
0.6	917
0.65	338
0.7	143
0.75	65
0.8	27
0.85	4

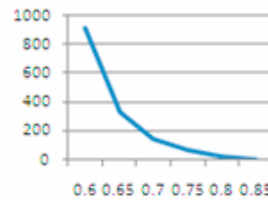


Fig. 4. Anomalies vs Threshold

#### 4. Conclusion

Identification of suspicious financial transactions to unhide money-laundering activities has always been a complex problem. This complexity can be attributed to the vagueness in the criteria of a transaction being suspicious, the consciousness of money-launderers to keep their moves unobserved, and the difficulty in validating the results obtained. This study employed a hybrid model of distance based clustering and dynamic Bayesian networks, called SARDBN, to identify anomalies in sequence of transactions. An anomaly scoring metric called AIRE was also presented that quantifies the degree of anomaly in each incoming transaction. SARDBN has been tested on a huge dataset of financial transactions and is shown to have promising results. However, a thorough comparison of SARDBN with some of the existing AML approaches can shed more light on the pros and cons of each method and is an area of future research.

#### References

- [1] T. Zhu, "An Outlier Detection Model Based on Cross Datasets Comparison for Financial Surveillance," *Proceedings of the 2006 IEEE Asia-Pacific Conference on Services Computing*, IEEE Computer Society, 2006, pp. 601-604.
- [2] T. Jun, "A Peer Dataset Comparison Outlier Detection Model Applied to Financial Surveillance," *Proceedings of the 18th International Conference on Pattern Recognition - Volume 04*, IEEE Computer Society, 2006, pp. 900-903.
- [3] Lin-Tao Lv Na Ji Jiu-Long Zhang, "A RBF neural network model for anti-money laundering," *Wavelet Analysis and Pattern Recognition, 2008. ICWAPR '08. International Conference on*.
- [4] X. Liu, P. Zhang, and D. Zeng, "Sequence Matching for Suspicious Activity Detection in Anti-Money Laundering," *Intelligence and Security Informatics*, 2010, pp. 50-61.
- [5] M.F. Jaing, S.S. Tseng, and C.M. Su, "Two-phase clustering process for outliers detection," *Pattern Recogn. Lett.*, vol. 22, 2001, pp. 691-700.
- [6] X. Wang and G. Dong, "Research on Money Laundering Detection Based on Improved Minimum Spanning Tree Clustering and Its Application," *Proceedings of the 2009 Second International Symposium on Knowledge Acquisition and Modeling - Volume 02*, IEEE Computer Society, 2009, pp. 62-64.
- [7] Su-Nan Wang and Jian-Gang Yang, "A Money Laundering Risk Evaluation Method Based on Decision Tree," *Machine Learning and Cybernetics, 2007 International Conference on*, 2007.
- [8] Jun Tang and Jian Yin, "Developing an intelligent data discriminating system of anti-money laundering based on SVM," *2005 International Conference on Machine Learning and Cybernetics*, Guangzhou, China: 2005, pp. 3453-3457 Vol. 6.
- [9] Y. Li, D. Duan, G. Hu, and Z. Lu, "Discovering Hidden Group in Financial Transaction Network Using Hidden Markov Model and Genetic Algorithm," *Proceedings of the 2009 Sixth International Conference on Fuzzy Systems and Knowledge Discovery - Volume 05*, IEEE Computer Society, 2009, pp. 253-258.
- [10] "Statistical Methods for Fighting Financial Crimes, Agus Sudjianto, Sheela Nair, Ming Yuan, Aijun Zhang, Daniel Kern, Fernando Cela-Díaz. *Technometrics*. February 1, 2010, 52(1): 5-19. doi:10.1198/TECH.2010.07032.."
- [11] A. Patcha and J. Park, "An overview of anomaly detection techniques: Existing solutions and latest technological trends," *Computer Networks*, vol. 51, Aug. 2007, pp. 3448-3470.
- [12] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Comput. Surv.*, vol. 41, 2009, pp. 1-58.
- [13] J.C. Bezdek, *Pattern Recognition with Fuzzy Objective Function Algorithms*, Kluwer Academic Publishers, 1981.
- [14] P. Tan, M. Steinbach, and V. Kumar, *Introduction to data mining*, Pearson Addison Wesley, 2006.