

# Exploring Fraudulent Financial Reporting with GHSOM

Rua-Huan Tsaih<sup>1</sup>, Wan-Ying Lin<sup>2</sup>, and Shin-Ying Huang<sup>1</sup>

<sup>1</sup> Department of Management Information Systems, National Chengchi University, Taipei 11605, Taiwan

<sup>2</sup> Department of Accounting, National Chengchi University, Taipei 11605, Taiwan

**Abstract.** The issue of fraudulent financial reporting has drawn much public as well as academic attention. However, most relevant researches focus on predicting financial distress or bankruptcy. Little emphasis has been placed on exploring the financial reporting fraud itself. This study addresses the challenge of obtaining an enhanced understanding of the financial reporting fraud through the approach with the following four phases: (1) to identify a set of financial and corporate governance indicators that are significantly correlated with fraudulent financial reporting; (2) to use the Growing Hierarchical Self-Organizing Map (GHSOM) to cluster data from listed companies into fraud and non-fraud subsets; (3) to extract knowledge from the fraudulent financial reporting through observing the hierarchical relationship displayed in the trained GHSOM; and (4) to provide justification to the extracted knowledge.

**Keywords:** Financial Reporting Fraud, Growing Hierarchical Self-Organizing Map, Knowledge Extraction.

## 1 Fraudulent Financial Reporting and GHSOM

Fraudulent financial reporting can lead to not only significant investment risks for stockholders, but also financial crises for the capital market. Fraudulent financial reporting were often detected with a very low frequency but with severe impacts [1]. Given the infrequency of synthetic reporting, most auditors cannot develop sufficient experiences and knowledge on fraudulent detection [7]. Furthermore, top management may be involved in providing less fairly represented financial statements. Beasley found that 83% of top managements of the U.S. listed firms, chief executive officer, chief financial officer, sometimes even both, are related to financial statement fraud [3]. Tipgos notes that internal control is designed in a “giver-receiver” model [14]. It means that management implements the internal control and employees are expected to follow it. The internal control mechanism aims to prevent employee frauds, not management frauds. In other words, since managers could bypass the internal control, it created a significant condition of financial statement fraud to lead to bamboozle auditors deliberately [10]. The standard audit procedures are insufficient to detect malfeasance for managers who understand the limit of audit [7].

**Table 1.** Literature summary of fraud detection techniques

Author	Technique	Variable	Sample	Findings
Persons (1995)	Stepwise logistic model	<ul style="list-style-type: none"> <li>➤ 9 financial ratios</li> <li>➤ Z-score</li> </ul>	Matched-pairs design	The study found four significant indicators: financial leverage, capital turnover, asset composition and firm size
Fanning and Cogger (1998)	Self-organizing artificial neural network	<ul style="list-style-type: none"> <li>62 variables</li> <li>➤ Financial ratios</li> <li>➤ Other indicators: corporate governance, capital structure etc.</li> </ul>	Matched-pairs design: 102 fraud and 102 non-fraud	<ul style="list-style-type: none"> <li>➤ Neural network is more effective</li> <li>➤ Financial ratios are over half of 8 significant indicators such as debt to equity, ratios of accounts receivable to sales, trend variables etc.</li> </ul>
Bell and Carcello (2000)	Logistic regression	46 fraud risk factors	77 fraud samples and 305 non-fraud samples	Logistic regression model was significantly more effective than auditors for fraud samples, but for non-fraud samples both made no difference.
Kirkos et al. (2007)	<ul style="list-style-type: none"> <li>● Decision tree</li> <li>● Back-propagation neural network</li> <li>● Bayesian belief network</li> </ul>	<ul style="list-style-type: none"> <li>➤ 27 financial ratios</li> <li>➤ Z-score</li> </ul>	Matched-pairs design: 38 fraud and 38 non-fraud	<ul style="list-style-type: none"> <li>➤ Training dataset: neural network is the most accurate</li> <li>➤ Validation dataset: Bayesian belief network is the most accurate</li> </ul>
Hoogs et al. (2007)	Genetic algorithm	<ul style="list-style-type: none"> <li>➤ 38 financial ratios</li> <li>➤ 9 qualitative indicators</li> </ul>	1 fraud vs. 8 non-fraud design	Integrated pattern had a wider coverage for suspected fraud companies while it remained lower false classification rate for non-fraud ones

There are numerous studies dealing with prediction of financial statement fraud using either logistic regression or neural network [4] [7] [9] [11][15]. Table 1 summarizes the fraud detection techniques.

Although several studies have shown the benefits of fraud prediction using regression model or neural network, they are often criticized that they are difficult to deal with high-dimensional data and limited samples. The specific criticism of neural network was a black box of classification process so that auditors were unable to understand the adopted factors and to verify validity of the model.

This study presents the application of Growing Hierarchical Self-Organizing Map (GHSOM) [12] to extract rules and knowledge about financial reporting fraud. Self-Organizing Map (SOM) is designed with the concept of unsupervised learning network to handle high-dimensional data and visualize results. It can also help to discover hidden hierarchies and to extract rules and knowledge from data. Unfortunately, SOM requires predefined number of neural processing units, static architecture of this model and has limited capabilities for the representation of hierarchical relations of the data. The GHSOM will automatically grow additional layers to display refined maps of individual nodes exhibiting high degrees of heterogeneity, thus providing higher levels of classification as needed. In addition, the size of each map was relatively smaller so that users can easily analyze and quickly obtain an overview. In

the practical aspects, several studies [5][6][13] applied GHSOM into information extraction field. Hence, the study applies GHSOM solution to dynamic network structure and hierarchical relationship to help auditors effectively extract rules or features of financial reporting fraud.

## 2 Experimental Design and Results

The research process can be described as Figure 1. In sampling stage, we first use the following sources to identify the occurrence of financial statement fraud. The first source is Summary of Indictments and Sentences for Major Securities Crimes issued by the Securities and Futures Bureau of the Financial Supervisory Commission in Taiwan. The second source is Summary of Group Litigation Cases issued by Securities and Futures Investors Protection Center in Taiwan. The third source is the law and regulations retrieving system of The Judicial Yuan which provides information to verify whether the accused companies committed the financial statement fraud.

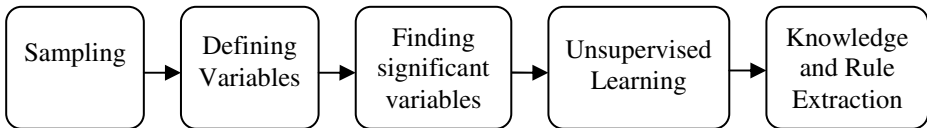


Fig. 1. Research Framework

Table 2. Definition of Fraud code and Color code

Fraud code		Color code		
0	Non-fraud Year	0	□ White	For fraud firms, the years which might be not investigated after detection year were assigned 0. All years of non-fraud firms would be assigned the same color code.
		1	■ Light blue	For fraud firms, the years under investigation and shown non-fraud after the last fraud year and before detection year were assigned 1.
1	Fraud Year	2	■ Deep blue	For fraud firms, the year preceding the first fraud year could be investigated and shown non-fraud.
		3	■ Gray	For fraud firms, other years before the first fraud year with low possibility of investigation were given 3.
		4	■ Red	For fraud firms, the years that appeared fraud in the indictments or judgments were marked 4.

If a company was prosecuted or judged according to the following enactments of Taiwan, it is a fraud firm:

- (1) Securities and Exchange Law : Paragraph 2, Article 20;
- (2) Securities and Exchange Law : Sub-paragraph 5, Paragraph 1, Article 174;
- (3) Business Accounting Law : Article 71;
- (4) Criminal Law : Article 342.

For the fraud firm, indictments and judgments do state the detection year that was investigated by prosecutors’ offices and the fiscal year of financial statements that are fraudulent. For each fraud company, there is a five-year sampling period of financial statements, which covers two years before and after the first fault year. Based on the fraud years and detection year, we further distinguish the difference among non-fraud year financial statements of fraud firms. The detail of fraud code and color code can be explained as Table 2.

We use the matched-pairs concept to create a sample pool of 58 fraud firms and 58 non-fraud firms, all of them are Taiwan publicly traded companies. For each fraud firm, we pick up a non-fraud counterpart that is in the same industry and the total assets are near to those of the fraud firm in the year before the first fraud year. There are totally 117 fraud financial-statement samples and 463 non-fraud financial-statement samples. The imbalanced samples consistent with facts which fraud cases were infrequent relatively.

Every sample data has a code composed of industry abbreviation, four-digit stock code, two-digit year abbreviation, one-digit fraud code, and one-digit color code. Take a fraud firm in electron industry for example as Table 3:

**Table 3.** An example of sample encoding.

Industry abbrev	Stock code	Year abbrev.	Fraud code	Color code	Sample code
		00	0	0	E82950000
		99	0	1	E82959901
E	8295	98	1	4	E82959814
		97	0	2	E82959702
		96	0	3	E82959603

We employed one nominal dependent variable-FRAUD, which is dichotomous and expressed as 1 and 0 according to whether the year financial statement is fraud or non-fraud. Initially, we use 25 independent variables from financial and corporate governance dimensions to be indicators of financial reporting fraud. The financial ratios were employed to measure profitability, liquidity, operating ability, financial structure and cash flow ability of a firm. Moreover, corporate governance variables and Z-score were utilized to examine probability of financial distress.

The measurement for profitability indicators we adopt are: Gross profit margin (GPM), Operating profit ratio (OPR), Return on assets (ROA), [8], Growth rate of net sales (GRONS), Growth rate of net income (GRONI).

The liquidity ability indicators we adopt are: Current ratio (CR), Quick ratio (QR).

The operating efficiency indicators we adopt are: Accounts receivable turnover (ART), Total asset turnover (TAT), Growth rate of accounts receivable (GROAR), Growth rate of inventory (GROI), Growth rate of Accounts receivable to gross sales (GRARTGS), Growth rate of Inventory to gross sales (GRITGS), Accounts receivable to total assets (ARTTA), Inventory to total assets (ITTA).

The financial structure indicators we adopt are: Debt ratio (DR), Long-term funds to fixed assets (LFTFA).

The cash flow ability indicators we adopt are: Cash flow ratio (CFR), Cash flow adequacy ratio (CFAR), Cash flow reinvestment ratio (CFRR).

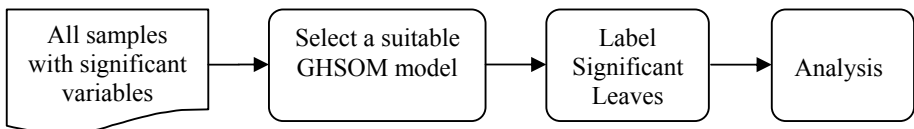
The corporate governance indicators we adopt are: Stock Pledge ratio (SPR), Deviation between control rights and cash flow rights (DBCRCFR), Deviation between ratio of controlled board seats and cash flow rights (DBCBCFR), Z-score.

Then we adopt the CANDISC to do tolerance test for multi-collinearity and significance test to derive significant variables. The result of multi-collinearity test suggested that the GRITGS variable should be excluded since its tolerance was extremely lower than other independent variables. Then, we use the F-value to determine the significance of each of remaining 24 independent variables. The result indicated that the following eight variables had statistically significant effects: ROA, CR, QR, DR, CFR, CFAR, SPR and Z-Score. The result of structure coefficient (i.e. discriminant loadings) shown that ROA had the greatest effect on the function, followed by Z-Score and CFAR has the smallest effect.

These eight significant variables examine a company from different dimensions:

- (1) Profitability: ROA can be used to assess a firm's ability to generate profits by the use of its own assets. [11] indicated that lower profit may give management an incentive to overstate revenues or understate expenses.
- (2) Liquidity: CR and QR can be used to measure a firm's liquidity which means its short-term ability to pay a debt. QR excludes inventory and prepaid expenses whose ability to realize is lower than cash or accounts receivable.
- (3) Financial structure: DR can be employed to inspect a firm's financial structure. [11] found that fraud firms have higher financial leverage than non-fraud firms.
- (4) Cash flow ability: CFR and CFAR can be used to test a company's ability to paying debts and other disbursement such as capital expenditures, inventory additions and cash dividends using cash flows from operating activities.
- (5) Stock pledge ratio: SPR can be utilized to measure the financial pressure on leverage degree of directors and supervisors by pledging their stocks to obtain funds.
- (6) Financial condition: Z-score can be used to measure a company's financial situation to determine the relationship between financial distress and fraud.

The procedure of GHSOM experiment can be mentioned as Figure 2. We use SOM toolbox and GHSOM toolbox in the platform of Matlab R2007a to conduct the GHSOM experiment. The trial and error would be performed in different breadth, depth and normalization method to get a suitable GHSOM model for analysis.



**Fig. 2.** Procedure of GHSOM experiment

At the beginning, we used GHSOM and significant variables for total samples (i.e. fraud samples and non-fraud samples) to get several hierarchical structures under different parameters including breadth, depth and normalization method. Next, we compared these GHSOM models and selected a suitable model base on the depth and map size. Finally, we chose some significant leaves with better explanatory power to label and analyze.

The growing process of GHSOM is primarily dominated by two parameters,  $\tau_1$  and  $\tau_2$ . While  $\tau_1$  controls the size of individual SOM,  $\tau_2$  determines the minimum data granularity and the global termination criterion. Therefore, a resultant hierarchy is not necessarily balanced and individual SOM could have different number of units and configurations.

(Where is on one hand?)On the other hand, each map will be named according to its upper layer, current layer which it was located on and its order in the same upper map. For example, a map named “L1m2-L2m1 “ indicates that it was from the second map of layer 1 and its current location is in the first map of layer 2.

The study attempts to present data hierarchy and to extract knowledge from clusters. An applicable model was also capable of comparing among clusters and analyzed within a cluster. Hence, the criteria of a GHSOM model can be defined as:

- (1) The depth of a model should be greater than two layers.
- (2) The breadth of individual map should consist of two firms at least. It meant that a map would be expected to have over ten samples.
- (3) New maps shouldn't extremely cluster in a minority of the parent.

The study performed canonical discriminant analysis to do tolerance test for multi-collinearity and significance test for selecting significant variables. In addition, the analysis also presented prediction rate of discriminant function.

The result of multi-collinearity test suggested that one variable, namely GRITGS, were excluded since their tolerance were extremely lower than other independent variables. The detail was described as Table 4. As a result, we acquired 24 independent variables as input to the Canonical Discriminant Analysis after deleting the variable.

**Table 4.** Tolerance Test of Variable

Deleted Variables	Variance within groups	Tolerance
GRITGS	89043100.972	.001

The study verified whether or not the discriminant function could show the significant difference by means of Wilks'  $\Lambda$  statistic. The corresponding P-value of Wilks'  $\Lambda$  value which was less than level of significance( $\alpha=0.05$ ) proved that a significant effect of the discriminant function. The result can be shown as Table 5.

**Table 5.** Wilks'  $\Lambda$  Statistic

Function	Wilks' Lambda value	$\frac{\lambda_1^2}{\lambda_1^2 + 1}$	D.F.	Significance
1	.766	151.095	24	.000

We made use of F-value to determine the significance of each independent variable. The result of corresponding p-value indicates that eight variables have statistically significant effects, including ROA, CR, QR, DR, CFR, CFAR, SPR and Z-Score.

Next, we utilized structure coefficient (i.e. discriminant loadings) to compare the discriminant power of individual variable. In brief, it is used to estimate the relative importance of each variable to the discriminant function based on the absolute value of structure coefficient. The result shows that ROA had a greatest effect on the function, the secondary is variable Z-Score and variable CFAR is the smallest. The study would analyze the variation of significant variables across a number of years.

The result shown in Table 6 lists the consistency between significance and relative importance from the rank of influence. The significant variables, such as ROA and Z-score, also provided stronger discrimination for the function. With regard to direction of influence, the result pointed out variable DR and SPR among all significant variables appeared negative correlation. It also indicates that a company whose most of significant indicators became bigger may tend to health.

**Table 6.** Significance and Relative Weights of Independent Variable

Variable	Structure Coefficient	F value	Significance	Rank of Influence	Direction of Influence
GPM	0.14	3.51	0.061		
OPR	-0.03	0.16	0.688		
ROA	0.77	105.82	<b>0.000***</b>	1	+
GRONS	0.06	0.63	0.427		
GRONI	-0.02	0.05	0.822		
CR	0.34	20.59	<b>0.000***</b>	5	+
QR	0.28	13.42	<b>0.000***</b>	7	+
ART	0.09	1.58	0.210		
TAT	0.19	6.38	0.012		
GROAR	0.03	0.12	0.731		
GROI	0.07	0.90	0.344		
GRARTGS	0.00	0.00	0.997		
ARTTA	0.11	2.25	0.134		
ITTA	0.12	2.37	0.125		
DR	-0.42	30.46	<b>0.000***</b>	4	-
LFTFA	0.02	0.09	0.764		
CFR	0.33	19.21	<b>0.000***</b>	6	+
CFAR	0.24	9.89	0.002***	8	+
CFRR	0.19	6.41	0.012		
SPR	-0.47	38.85	<b>0.000***</b>	3	-
SMLSR	-0.19	6.18	0.013		
DBCRCFR	0.02	0.04	0.835		
DBCBSFR	-0.05	0.41	0.524		
Z-score	0.64	72.74	<b>0.000***</b>	2	+

Note. \*\*\* p < 0.01.

**Table 7.** Prediction rate of discriminant function

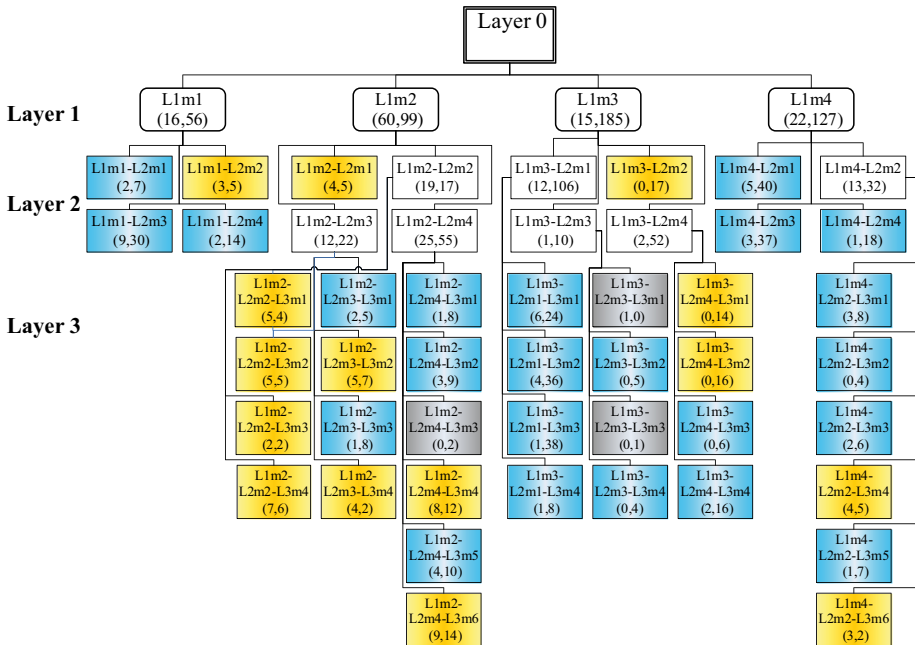
		Class	Predict	
			0	1
<b>Original</b>	No.	0	390	77
		1	44	69
	%	0	83.5	16.5
		1	38.9	61.1

As a whole, the canonical discriminant function incorrectly classify 21.9% of samples, that is to say, prediction power of the function achieved 79.1%. The result was listed in Table 7. In the study dealing with financial reporting fraud prediction, it is less costly to classify a non-fraud firm as a potential candidate for fraud than a potentially fraudulent firm as non-fraud. Based on the principle, type I error should be lower to reduce misclassification costs. Type I error is 16.5% and type II error is 38.9%.

H<sub>0</sub>: The firm did not commit financial reporting fraud

H<sub>1</sub>: The firm committed financial reporting fraud

The selected GHSOM model develops a tree with three layers and 41 leaves as Figure 3. The tree creates four maps including L1m1, L1m2, L1m3 and L1m4 in the first layer. By means of comparing ratio of fraud to non-fraud sample among them, L1m2 having the highest fraud ratio indicated over half of fraud sample could be classified into its clusters in the following layers. By contrast, L1m3 with the lowest fraud ratio was expected to probably produce some pure non-fraud clusters. The detail can be listed in Table 8.



**Fig. 3.** GHSOM Tree with three layers

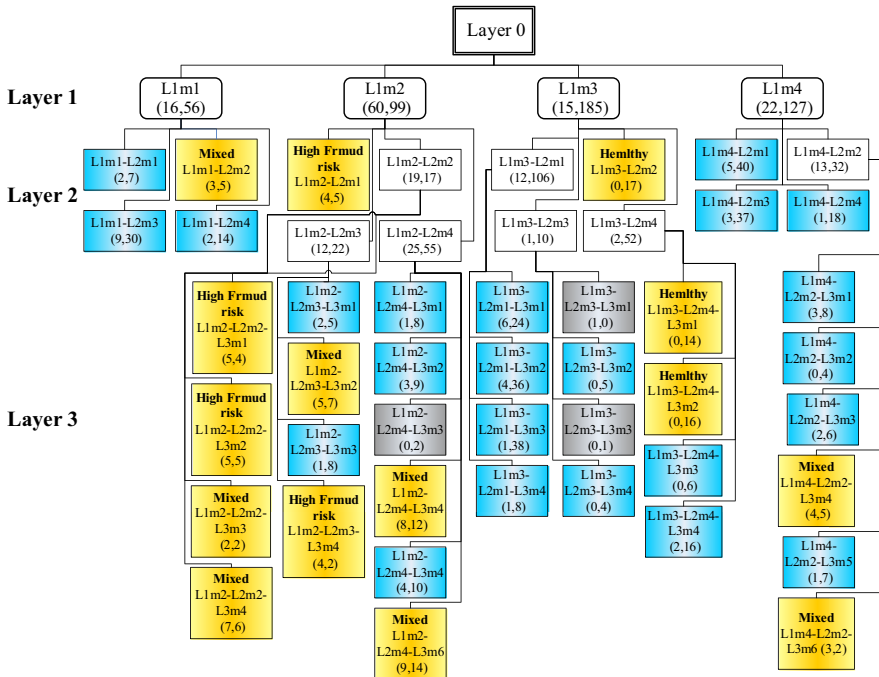


**Table 8.** Ratio of fraud to non-fraud in Layer1

Layer1	Number		Fraud ratio (%)
	Fraud	Non-fraud	
L1m1	16	56	<b>28.57</b>
L1m2	60	99	<b>60.61</b>
L1m3	15	185	<b>8.11</b>
L1m4	22	127	<b>17.32</b>

This GHSOM model generates 41 leaves, in which L1m2-L2m4-L3m3, L1m3-L2m3-L3m1 and L1m3-L2m3-L3m3 leaves have less than three sample data. These three leaves are excluded from the further discussion. The overall ratio of number of fraud sample data to number of non-fraud sample data is 113:467; this ratio information is adopted as the norm for picking up the leaves whose ratio is quite deviated from this value. Figure 3 shows 15 pick-up leaves in yellow, deleted leaves in gray, and other leaves in blue. Amongst 15 pick-up leaves, there are three leaves whose number of non-fraud sample data is much larger than of fraud ones.

The study would classify 15 significant leaves into high fraud risk group, mixed group and healthy group according to their characteristics are shown in Figure 4. In GHSOM model, all of high fraud risk groups were generated from L1m2 with the highest fraud ratio while all of healthy groups were produced from L1m3 with the lowest fraud ratio. That is to say, L1m2 and L1m3 had superior discrimination power.



**Fig. 4.** GHSOM Tree with Label

### 1. High fraud risk group

The group had totally four leaves whose non-fraud samples all came from fraud firms. The significant features were the worst profitability, liquidity, financial structure and Z-score. The results showed techniques which fraud samples cooked the books probably were not skillful and thorough so financial adversity could be detected from the preceding dimensions. On the other hand, fraud samples have manipulated the related accounts so their financial situation may be more terrible.

### 2. Mixed group

The group consisted of fraud samples and non-fraud ones which belonged to fraud firms or non-fraud firms. In terms of comparison among leaves, some leaves had few but unique features which were primarily in the worst cash flow ability and higher stock pledge level.

### 3. Healthy group

The group with three leaves had excellent characteristics in financial structure, stock pledge level and Z-score. Because few fraud firms had good financial state and lower stock pledge level, they probably represented financial statements correctly.

## 3 Conclusions

The research attempted to present a hierarchical structure from data and to extract knowledge related to financial reporting fraud through Growing Hierarchical Self-Organizing Map (GHSOM) as well as a set of financial and corporate governance indicators. Our financial sample excludes financial industry and includes 113 frauds and 467 non-frauds firm-year observations. First of all, discriminant analysis and GHSOM was applied to obtain eight significant variables and to develop a suitable hierarchical model. Next, 15 leaves with explanation power was selected to perform analysis among and within clusters.

The research result appeared distinctions among high fraud risk group, mixed group and healthy group were described as below: In terms of consistency, all of leaves had smaller coefficient of variation in liquidity, financial structure and stock pledge level while the discrimination among them was in profitability, cash flow ability and Z-score.

To sum up, many fraud samples in 15 leaves belonged to Iron& Steel and Building Material& Construction industry and committed shenanigans in 1998 and 1999 which East Asian Financial Crisis occurred seriously. Because of bear market, the operation of fraud firms deteriorated sharply and could not create cash flow. Under the pressure of capital, they borrowed short-term loans to meet operating demand so financial structure became worse. The vicious circle motivated evil top management to misappropriate corporate money for keeping stock price rising. In the meantime, fraudulent financial reporting could conceal the embezzlement to make investors and banks trust them. Moreover, frequent schemes included overstating revenues through fictitious sales, obtaining money in nominal accounts such as temporary payment or prepayment for purchases, recording loans from related party into accounts receivable, falsified some accounts like accounts receivable etc.

## References

1. Association of Certified Fraud Examiners. Report to the nation on occupational fraud & abuse [Electronic Version], <http://www.acfe.com/documents/2006-rttn.pdf>
2. Altman, E.I.: Financial ratios, discriminant analysis and the prediction of corporate bankruptcy. *The Journal of Finance* 23(4), 589–609 (1968)
3. Beasley, M.S., Carcello, J.V., Hermanson, D.R.: Fraudulent financial reporting: 1987-1997 an analysis of U.S. public companies (1999)
4. Bell, T.B., Carcello, J.V.: A Decision Aid for Assessing the Likelihood of Fraudulent Financial Reporting. *Auditing* 19(1), 169–184 (2000)
5. Dittenbach, M., Merkl, D., Rauber, A.: The Growing Hierarchical Self-Organizing Map. In: *The Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural Networks- IJCNN 2000* (2000)
6. Dittenbach, M., Rauber, A., Merkl, D.: Uncovering hierarchical structure in data using the growing hierarchical self-organizing map. *Neurocomputing* 48(1-4), 199–216 (2002)
7. Fanning, K.M., Cogger, K.O.: Neural network detection of management fraud using published financial data. *International Journal of Intelligent Systems in Accounting, Finance & Management* 7(1), 21–41 (1998)
8. Hoogs, B., Kiehl, T., Lacombe, C., Senturk, D.: A genetic algorithm approach to detecting temporal patterns indicative of financial statement fraud. *Intelligent Systems in Accounting Finance and Management* 15(1/2), 41–56 (2007)
9. Kirkos, E., Spathis, C., Manolopoulos, Y.: Data Mining techniques for the detection of fraudulent financial statements. *Expert Systems with Applications* 32(4), 995–1003 (2007)
10. Loebbecke, J.K., Eining, M.M., Willingham, J.J.: Auditors' experience with material irregularities: frequency, nature, and detectability. *Auditing* 9(1), 1–28 (1989)
11. Persons, O.S.: Using financial statement data to identify factors associated with fraudulent financial reporting. *Journal of Applied Business Research* 11(3), 38–46 (1995)
12. Rauber, A., Merkl, D., Dittenbach, M.: The Growing Hierarchical Self-Organizing Map: Exploratory Analysis of High-Dimensional Data. *IEEE Transactions on Neural Networks* 13(6), 1331–1341 (2002)
13. Shih, J.-Y., Chang, Y.-J., Chen, W.-H.: Using GHSOM to construct legal maps for Taiwan's securities and futures markets. *Expert Systems With Applications* 34(2), 850–858 (2008)
14. Tippos, M.A.: Why management fraud is unstoppable. *CPA Journal* 72(12), 34–41 (2002)
15. Virdhagriswaran, S., Dakin, G.: Camouflaged fraud detection in domains with complex relationships. In: *The Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining* (2006)