# Constructing credit auditing and control & management model with data mining technique

S.C. Chen *, M.Y. Huang

Department of Industrial Engineering & Management, National Chin-Yi, Institute of Technology, Taichung, Taiwan, ROC

## ARTICLE INFO

## ABSTRACT

The 2008 financial tsunami, hitting the globe across all types of industries, causing tides of bankruptcies and severe unemployment, had its epicenter at American subprime in the housing market. In fact, the US subprime storm was just a premonition, while the root cause of the financial tsunami lied in the oversupply of structured credit products. Credit card business, one of the structured credit products, which under an intensively competitive environment, have been released by many banks with high spread, high return, and easy-to-apply appeals to carter to consumers needs. In order to allure the customers, some banks even go to the extent as simplify the credit rating, which in turn has increased credit risk, causing high non-performing ratio, increased debt collection cost, and growing bad debt counts. Accordingly, credit risk auditing plays a vital role in the successful management of credit card business. In response to such needs, the present study aims to conduct analysis and investigation on the current status of the industry with CRISP-DM model. First, customers' demographic data and payment-related statistics were analyzed to identify feature variables, which were then sorted out as demographic data, debt data, payment rating etc. Next, by utilizing artificial neural network of data mining technique, the study tries to predict customer's regular pattern of consumption, payment and/or default and bad debt, and to develop a set of credit granting principle by employing the decision tree technique. Since data mining classification model has a greater power in discriminating credit card granting, it can thus be used to construct accurate credit variable rules and predictive model, to further improve credit checking effect and credit risk control. Using the credit auditing data of a certain bank as a case study, the study intends to verify that the model constructed by the researcher can effectively identify the potential key factors of its credit card granting rule, to minimize the cost loss of Model I and Model II credit business, and eventually enhance the stability and profitability of the bank's credit card business.

© 2010 Elsevier Ltd. All rights reserved.

## 1. Introduction

To date, credit card has become one of the indispensible market phenomena of currency transaction system. This is particularly true of the world with a predominance of commercial interests, where either on the bank end or the consumers end, a finance management idea of "Enjoy first, pay later" is encouraged. Taking the department stores as an example, almost all of them offer interest-free installment at the anniversary day to attract the consumers to conduct advance consumption with their credit. The shadow of the duel-card (credit card and cash card) storm triggered 3 years ago, still lingers in the public mind. According to the latest statistics, until October 2008, total circulation of domestic credit cards amounts to 36.4 millions. And, this is only the new low since the 2005 credit card bad debt storm, because the number of newly issued cards were outnumbered by the cut and suspended cards. Prior to the explosion of duel-card issue, domestic credit card market had experienced its peak period. 2005 set the historic record high of 45.49 million credit cards, which in 2006, dropped to a circulation of 38.32 million due to the impact of credit card bad debt. Since its official emergence in September 2005, total write-off of bad debt by various credit card issuing banks amounted to NT$13.4 billions. According to statistics released by the Financial Supervisory Commission, Executive Yuan in February 2006, sum total of the circulation interests of credit card and the granted loan balance of cash card had climbed to NT$76.49 billions, with 520 thousand overdue card holders, averaging default payment was NT$300 thousands per capita.

Oversupply of structured credit products and over expansion of credit have planted the root cause for the financial tsunami, created the most devastating global-scale financial crisis in nearly 50 years. Such incident makes it imperative that the banking industry should reexamine the way they judge and review the applicants' credits. Excessive issuance and overdue payments of credit cards have caused grave economic problems. Excessive use

* Corresponding author.
E-mail address: scchen@chinyi.ncut.edu.tw (S.C. Chen).

by card holders is not the sole cause for credit card problem. A more serious cause is the simplified credit rating among other reviewing processes that the bank used to win over customers in a competitive environment. This has created tens of billions of bad debt by the excessive consumption of insolvent card holders. Accordingly, financial banks, in conducting credit granting, should adopt a set of standards, stringent reviewing mechanism, and on the basis of revenues, try to make the right selection, to minimize the occurrence of bad debt, and to enhance the management performance of card issuance banks. A majority of previous literature focused on constructing credit card classification or behavior classification model with high accuracy, without taking into account the Model I and Model II errors resulted from misclassification. Here, Model I error refers to the misjudgment of applicants with good credit for high-risk group. Conversely, Model II error happens when applicants with bad credit are misclassified as low-risk group. As shown in the following Table 1:

In this study, the researcher intends to use the CRISP-DM 6-step cycle of improvement procedure to identify the influential factors causing default discrimination control in the reviewing process. Furthermore, by applying artificial neural network (ANN) and rule of decision tree, to cut down misjudged credit reviewing that cause bad debt resulted from credit expansion, and hopefully, to establish a set of relevant rules that can effectively eliminate those errors.

## 2. Definition of research model

To cope with the changing environment, many enterprises facing with the surging tide of IT development, expect to benefit from it by gaining some competitive edges. Nevertheless, upon introduction of IT system, they soon find themselves incapable of uncover the wealth of information stored in the huge databank. Thus data mining technique has become a scientific skill to excavate the knowledge and patterns concealed in the diversely complex mountain of data. In defining "data mining", Cabena (1998) explained that data mining is the process of effective accessing and extracting a large volume of information previous unknown, and provide the extracted information to his/her superiors for final decision-making. Berry et al. (1999) pointed out in their study that "data mining" is analyzing and finding meaningful relations or rules from a great amount of data in an automatic or semi-auto manner. Frawley, Piatetsky-Shapiro, and Matheus (1991) instead interpret "data mining" as the process of excavating from databank the non-apparent, implied, unprecedented, and yet may possibly be useful information. Grupe and Owrang (1995) regarded "data mining" as the act of dissecting facts and discover the new relations unknown to experts from the existing information. Hall, Mani, and Barr (1996), however, defined "data mining" as hunting/grabbing knowledge presented in a regular manner or other modes from the sea of data, by combining multiple techniques, such as data visualization, machine learning, statistics, and data warehousing.

In defining the research steps, the study tries to integrate the CRISP-DM 6-step cycle and the DMAIC 5-step data mining process, as shown in Fig. 1:
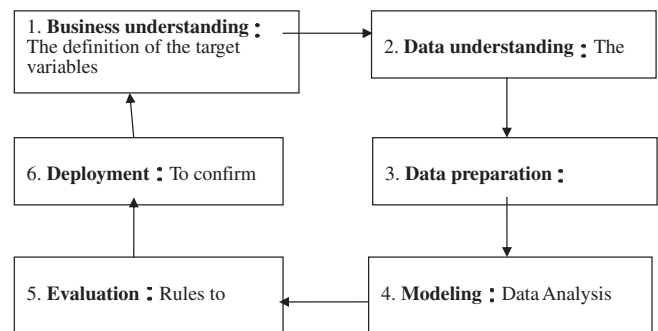
**Table 1**
Type I error and Type II error.

| Actual class | Classified class | |
|---|---|---|
| | Good credit | Bad credit |
| Good credit | Accept | Type I error |
| Bad credit | Type II error | Reject |



**Fig. 1.** CRISP-DM 6 steps in cycling and conform of DMAIC chart.

From Fig. 1, we find:

1. Commercial understanding and defining target variables: In this stage, we can have a clear understanding of the project's goal and need, to achieve the purpose and definition of data mining problems, and perfectly devise an appropriate plan and procedure;
2. Data understanding and defining conditional variables: In this stage, we must have initial collection and understanding of data, evaluate the quality of data, and propose hypotheses for the portions with possibly hidden problems. It is important that we understand the characteristics of data. For instance, in conducting correlation analysis, it is important to explore the positive, negative, or non-correlation between credit risk and consumer's occupation, position, seniority, and income.
3. Data preparation and measurement: In this stage, the raw data construct must be transformed into a final, useful data form, that is, convert them into an information format that can be constructed by instrumental software package. To achieve this purpose, data may be put through several conversions and compilations, depending on the construction need. The process includes table, record, selection of property and modeling, and changing and cleaning of the data fed into the tool kit.
4. Model designing and data analysis: In this stage, we have various models to choose and utilize from. Some of the models, however, may demand strict requirements for specific data forms, e.g., sequential, scattered data, data loss etc. Accordingly, we are often required to move back to the last stage to recheck and recompile the data.
5. Assessment and improvement of credit rules: In this stage, we built one or several models, and analyze those data. Before actual applying the model (s), an overall reviewing is needed to see whether there is any excluded important commercial considerations, in order to ensure that the model (s) meets the initial professional requirements. With this process, we can induce a set of practicable credit improvement rules.
6. Confirmation and Control: To confirm the application effect of the conclusion, to conduct continuous improvement based on the correctness/incorrectness of the responsive model, and to make mutual verification between the acquired knowledge and commercial understandings.

## 3. Defining target variables and input of variables

Data used in this study were collected from the consumers' basic credit checking data and post-checking consumption records of a certain bank. According to the credit card databank of the case bank, a total of 266,083 credit cards were then in circulation, among them, 123,592 were valid. Number of card issued in that month was 307, with an average monthly account of NT$ 2398,

**Table 2**
Standards of credit score table.

| List | Score content | | | | |
|------|------|------|------|------|------|
| Annual income | ≥1200000 | ≥800000 | ≥400000 | ≥250000 | <250000 |
| | 35–40 point | 30–35 point | 25–30 point | 15–20 point | 10 point |
| Deposits go throw | Total deposits ≥500000 | ≥250000 | ≥100000 | Others | |
| | 15 point | 12 point | | 8 point | 4 point |
| Credit and house | Residential set mortage | No Residential set mortage | Residential set other mortage | | |
| | 10–12 point | 8–10 point | 6–8 point | | |
| | Housing family house dormitory | Others | | | |
| | 4–6 point | 3 point | | | |
| Job | Manager above level | Manager below level | Privately operated enterprise official personnel | Others | |
| | 13–14 point | 10–13 point | 7–10 point | 4 point | |
| Work experience | Serving more than 10 years | Serving more than 6 years | Serving more than 2 years | Others | |
| | 8–10 point | 6–8 point | 4–6 point | 3 point | |
| Family status | Married and had child | Married and no child | No married and others | | |
| | 4 point | 3 point | 2 point | | |
| The case of credit card holders and payment records | Used credit card 5 years and payment records normal | | Used credit card 3 years and payment records normal | | |
| | 5 point | | 4 point | | |
| | Used credit card 1 year and no payment record | | Others | | |
| | 3 point | | 2 point | | |

and 1.06% overdue rate. Manpower credit checking data collected for the study was from December 2008, consumption records were from the same target samples collected 3 months later. A total of 310 individuals of credit checking data and consumption records were provided by the case bank, of which, 267 were normal transactions, 43 are bad debts. Based on the bank's credit rating score (refer to Table 2), the target and input variables are defined as shown in Table 3:

Based on the variables listed in Score Table 3 that may influence credit risk, we will predict the types of loan borrowers, that is, good or bad customers. Detailed descriptions of each variable are given as follows:

1. Annual income (X1): This is a major factor in assessing borrower's solvency. Loan borrowers with an annual income of NT\$ 1 million have a far better solvency than those with only NT\$300 thousands.
2. Record of transactions (X2): The amount of bank deposit is a major indicator to determine loan borrower's solvency.
3. Credit granting and real estate status (X3): Applicants with private-owned house as mortgage usually get high credit score. Bank granting credit under collateral has a better guarantee, and can auction the mortgaged house off to balance any bad debt if there is one. Applicants without a house as collateral can only be allowed a relatively limited credit line.

4. Occupation (X4): Occupational position has a direct impact on the income. In general, those with a position higher than manager will get fairly high score. Public servants, who enjoy a stable job with less chance of dismissal, are viewed by credit grantor as good clients. In contrast, professional military personnel, although enjoy stable income and job condition, are usually categorized as high bad debt group due to their difficulty in finding any civilian job after retirement.
5. Seniority (X5): Seniority is considered as a stable factor. An applicant having only $3 \sim 4$ months of working experience is usually still on a probation period. An applicant employed less than one year has a high separation rate since he/she is still in the adaptive period. Once in bad debt, the most frequently cited reason for insolvency by those loan borrowers is none other than "out-of-work."
6. Family status (X6): Marital status is one variable that may provide a valuable reference for credit rating. Usually, applicants with children have higher degree of self-control over monetary and material desires, compared to childless single, thus tend to get higher credit scores.
7. Credit card holding and record of payment (X7): The reasons for multi-card holding may be favoritism or gift giving. Whatever the reason, the bank has to be on guard against multi-card holders, especially when each of them is in debt, to see if they are used for daily expenditures.
8. Bad debt or not (Z1): For target variable, bad debt or not, we take into account the risk of loss. High credit line imposes high risk for the bank, while low credit line indicates low risk for the bank. Usually, credit line is determined in reference to the loan borrower's monthly income and occupation.

## 4. Data measurement and model analysis

First confirm the analysis sample and variables, select the input variables and target variables. Next, convert the data into the experimental model. To avoid influencing the over-fitness and accuracy of data, we designate the learning data and verifying data. Raw data are randomly segmented with the nodes of Clementine 10.1 into three groups: 60% training data, 20% learning data, and

**Table 3**
The target and input variables.

| Input variable | Variable | List |
|------|------|------|
| | X1 | Annual income |
| | X2 | Deposits go throw |
| | X3 | Credit and house |
| | X4 | Job |
| | X5 | Work experience |
| | X6 | Family status |
| | X7 | The case of credit card holders and payment records |
| Target variable | Z1 | Non-performing loans or not |

20% verifying data, for use in the model. The experimental model is shown in Fig. 2:

Regarding the experimental model, the study uses ANNs to conduct analysis, because its internal nonlinear structure has been widely applied in various fields (Jane & Zhan, 1998; Smith & Gupta, 2002; Widrow, Rumelhart, & Lehr, 1994), especially in business management-related applications and model construction (Zhang, 2004).

ANN is a data processing system that tries to simulate the structure and/or functional aspects of biological neutral network with large amount of simple, linked artificial neurons. Artificial neuron is a simple simulated version of biological neuron. It accesses information from external environment or other artificial neurons which, through simple computation, output the results to the external environment or other ANN. Artificial neuron is also called processing element. Its model is shown in the following Fig. 3:

The output of each processing unit becomes the input of many processing units. In general, the equation of output value and input value is expressed with the function of the weight-sum of the input value, as shown in the following formula:

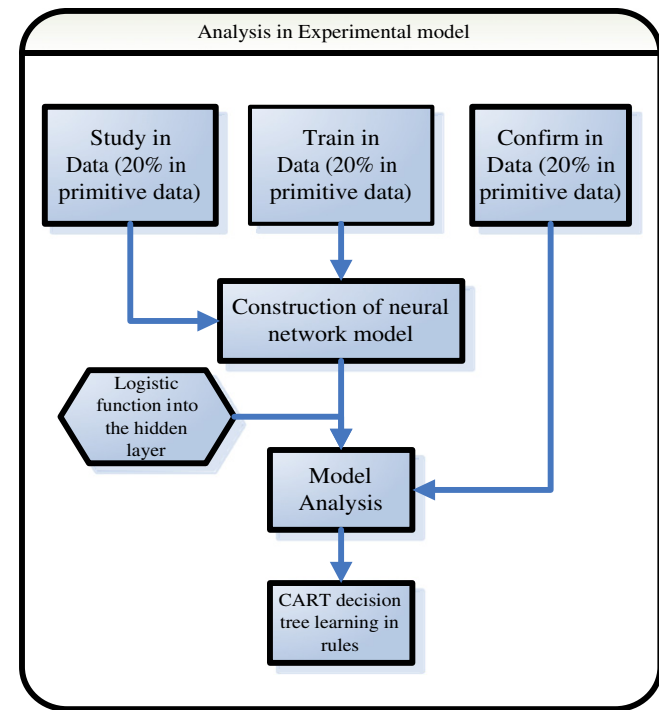$$Y_j = f\left(\sum_i W_{ij}X_i - \theta_j\right), \tag{1}$$



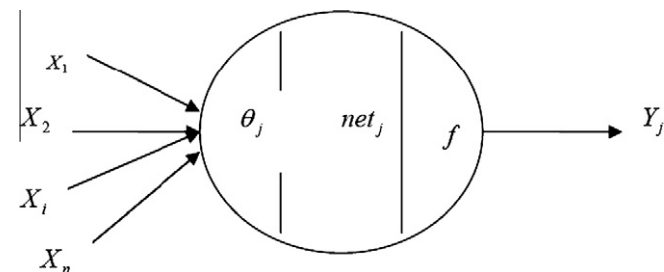Fig. 2. Analysis in experimental model.



Fig. 3. Processing unit in neural network.

Wherein $Y_j$ = output signal of simulated biological neural model; $f$ = transforming signal of simulated biological neural model; $W_{ij}$ = strength of neural node of simulated biological neural model, also called the weighted value; $X_i$ = input signal of simulated biological neural model; $\theta_j$ = threshold value of simulated biological neural model. Connection is the signal between the processing units, also called transmission path. At each connection, there is a weighted value $W_{ij}$ for each value, indicating the influencing strength of $i$ processing unit has on $j$ processing unit. A neural network is composed of many artificial neurons and connections, which also consists of various network models, of which, the back-propagation network (BPN) is the most popularly applied one. One BPN contains many layers, with each layer further contains several processing units. The input processing unit is used for inputting information from the external environment, while the output processing unit is used for outputting information to the external environment. An additional important processing layer, the hidden layer, is responsible for providing neural interaction between various neural networks, and the processing capability of the problems in the internal structure.

During the computation of the entire ANN, each processing unit will react to the input with different strength, either to strengthen or to inhibit the input. While in the ANN, the strength of neural nodes may be regulated via network training, to generate new neural nodes, to suspend non-functional neural nodes, or to regulate the strength of the existing neural nodes (Gluck & Myers, 2001).

ANN has a random feature, with different algorithmic rules generating model of different accuracies. We can choose from several models the one with the optimal overall performance, or use the final prediction acquired from all models. Using six algorithmic rules: high-speed, dynamics, plurality, revision cancellation, expansion selection, and RBFN (Radial Basis Function Network), Clementine can generate models with different accuracies. However, in selecting algorithmic rules, we will have to consider the balance between timing and accuracy. In the study, we select cancellation of revision algorithm, to remove those weak neurons and abort unnecessary inputs. Cancellation of revision begins with one or two hidden layers. The learning method is similar to the high-speed algorithm, that is, to conduct significance analysis on the hidden neurons, and then eliminate the weakest hidden neurons from the network. Such algorithmic operation of learning and elimination will be continuously repeated until some improvements are made to the network. Through repeated learning, the data may become over-learning, and such over-learned sample data tend to lose capability of generalization. To prevent over-learning from occurring, we randomly select only part of them for network learning. If, in case, all such data pass through network learning, other data are then used for intersectional appropriation, in order to evaluate the current performance of network. Based on the accuracy of model testing group, we then decide the time for stopping the learning process. Moreover, since the ANN was initially set according to the randomly specific weight, we can thus reconstruct the dynamic direction of the network by setting the random seed. Under the default condition, Clementine is capable

Table 4
Neural network model analysis.

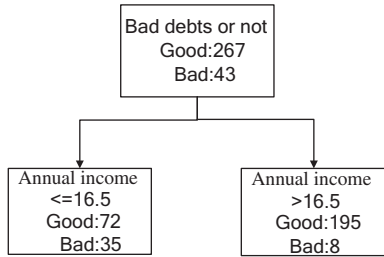| Input variables | Importance |
|---|---|
| Credit and house | 0.318188 |
| Years and service | 0.307392 |
| Annual income | 0.29905 |
| Family status | 0.14107 |
| Credit card holders and contributions | 0.091413 |
| Job | 0.0836199 |
| Deposits between the cases | 0.0123309 |

**Fig. 4.** The first differences of decision tree.

of determining when the network has reached the ultimate learning status, and then set the time for suspension of learning. Analysis of the case ANN model is shown in Table 4:

As shown in Fig. 4, the (estimated) accuracy of the intersectional appropriation of ANN prediction is about 97.7%, representing the accuracy of the data passing through the optimal ANN. Accuracy is the absolute percentage of the calculated error of each data column to the predicted value, by subtracting the average percentage of the entire data column from 1. In our case network, there are 7 neurons on the input layer, 1 on the hidden layer, and 2 on the output layer. The input columns are arranged in order of their significance levels. Significance level ranges from 0.0 ∼ 1.0, with 0.1 indicating no significance at all, and 1.0 represents the greatest significance. In general conditions, a value of >0.35 is very rare. For all variables of the case network, the most significant one is "credit granting and real estate status", followed by "seniority," "annual income," "family conditions," "credit card holding status," "payment record," "occupation," and "account transactions." The "accuracy analysis" nodes of the evaluation data fitness are deployed at the rear of data stream (bad debt or not). On the aspect of reliability, the average accuracy of the model has reached up to 96%, confirming an excellent fitness between the predicted value and the actual value.

Lastly, all output data are exported in the SPSS format to facilitate model verification. To conduct cross-sectional verification, evaluation of the learning results is used to evaluate the predicted accuracy of the verified data with the same procedure. The physical verification model is produced with exactly the same procedure as that of ANN. Of which, the average accuracy of all verification models reaches up to 0.972, a better score than the 0.96 of learning data. As shown in the verification data, the model formed with the results of accuracy analysis has greater predictive power, which can be used to determine the general models.

## 5. Improving credit rules

On the aspect of improving credit granting rules, although ANN has analysis indicating the relative significance of the output variables, it is incapable of acquiring any information in function forms, it, accordingly, lacks the power of making any specific explanation like those of linear regression analysis and logistic regression analysis. In light of this, the study intends to conduct regular extraction by using decision tree to track the learning results of ANN. In other words, we use the predicted results of ANN as the output variables, and select the input variables to conduct decision tree analysis.

Decision tree is a tree-like structure used for classification. It adopts the ramification framework of the tree to generate rules applicable to all classification problems. It is widely applied to a variety of decision-making problem-solving tools. The main feature of a decision tree is the classification the input variables according to a certain set of rules or method, and demonstrating them in a tree branch manner, to indicate the discrimination of

classifications caused by input variables. By following the analyzing rules of the decision tree, we can quarry out the factors having a significant impact on the results, by conducting hierarchical classification of the decision-making data.

Being the easiest and most popularly accepted approach to express knowledge, decision tree has been adopted widely in the data mining field. It's been especially regarded as the most effective approach in solving classification problems, such as credit card granting, direct-effect marketing response, predicting customer loss etc. A decision tree can be exhibited in chart, graphics, or rules. In the case of rules, they are easy-to-interpret and understandable enough as to process series or classified variables. And by utilizing the maximum information gain, it can be used to select segmented variables, and to display the relative significance of variables. Additionally, it can process large data set in a very efficient manner. Besides, since there is no correlation between the sizes of the tree and the databank, it has the advantage of smaller but flexible computing capacity. When there are a great number of variables input the model, the decision tree can still be constructed. Related rules of the CART decision tree proposed in the study are shown in the following Table 5:

Fig. 5 concisely shows the contents and node points of the entire data, and under what conditions ramification of data is formed. From the results, we can produce various subsets. In the case study, we first identify the variables appear on the decision tree. After "annual income" appears at the initial ramification, "seniority" appears at the node 1. And at node 6, "credit card holding status and payment record" appears as another ramification variable. This means that in predicting ANN, "annual income" is the most influential variable in determining bad debt. Following "annual income," "seniority" and "credit card holding status and payment record" are the next two most influential variables.

At the 1st ramification, those with an annul income score of ≦16.5, the bad debt probability increases from the overall 13.871% to 32.710%. On the other hand, those with annual income score of >16.5, the bad debt probability is reduced to 3.941%. For all those with an annual income score of 16.5, equivalent to NT$300 thousands, classification 0 and 1 represent normal client and bad debt client respectively. The 1st ramification is shown in Fig. 4:

At the 2nd ramification, those with annual income score of ≦16.5, and with a seniority score of >9.5, the bad debt probability increases from 32.710% to 100%. Those with an annual income score of >16.5, and with credit card holding and payment record

**Table 5**
The Rules indicate of CART decision tree.

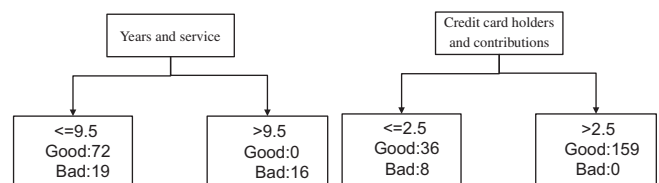| Annual income | | | |
|---|---|---|---|
| < = 16.5 | | > 16.5 | |
| Years and service | | Credit card holders and contributions | |
| < = 9.5 | >9.5 | < = 2.5 | >2.5 |
| Credit and house | | Annual income | |
| < = 5.5 | >5.5 | < = 22.5 | >22.5 |
| | | Family status | |
| | | < = 2.5 | > 2.5 |



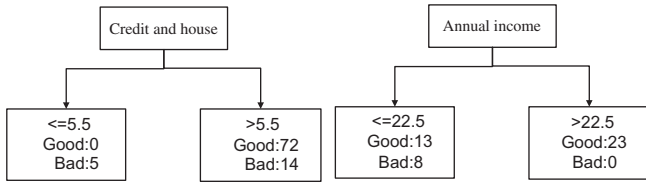**Fig. 5.** The second differences of decision tree.

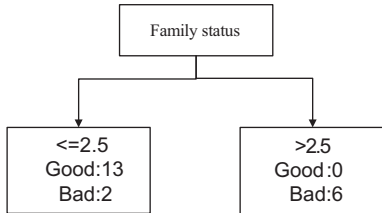Fig. 6. The third differences of decision tree.



Fig. 7. The fourth differences of decision tree.

score of >2.5, have a 0% bad debt probability. Those who meet these requirements have reached 51.29% of the total samples. From the 2nd ramification, we found two important rules: (1) Those with an annual income of ≤NT$300 thousands, and a seniority of more than 10 years, have 100% bad debt probability, occupying 5.161% of the total samples; and (2) Those with an annual income of >NT$300 thousands, and holding a credit card for more than one year, and having no delay payment for nearly a year, have 0% bad debt probability, occupying 51.29% of the total samples. Details of the 2nd ramification are shown in Fig. 5:

The 3rd ramification is described with node 3 and node 4: those with an annual income score of ≤16.5, seniority score of ≤9.5, with credit granting and real estate status as a ramifying point, scoring >5.5 and ≤5.5, have a bad debt probability of 16.279% and 100% respectively. Another ramification node starts from node 7. After incorporating annual income as a ramification variable, we found that those with an annual income score of >22.5, have 0% bad debt probability; those with a score of ≤22.5 have 38.095% bad debt probability. We identify three important rules in the 3rd ramifica-tion: (1) Those with an annual income of ≤NT$300 thousands, with a seniority of less than 10 years, and credit granting and real estate status of rented house or governmental dorm, have 100% bad debt probability, amount 1.613% valid sample of all samples; (2) Those with an annual income of ≤NT$300 thousands, with a seniority of less than 10 years, and credit granting and real estate status of pri-vate-owned house, have 16.279% bad debt probability; and (3) Those with an annual income of >NT$325 thousands, without his-tory of holding credit card, but having a record of delay payment, have 0% bad debt probability, valid sample accounting for 7.419% of the total samples; those with an annual income of <NT$325 thousands, have 38.095% bad debt probability. The 3rd ramifica-tion is shown in Fig. 6:

The 4th ramification starts from the point of high bad debt probability (node 8), then retrieve to the node of high bad debt ra-tio, and end at the starting point of the 3rd ramification. In this bracket, when family status scores >2.5, the bad debt probability is 100%. A score of ≤2.5 indicates a bad debt probability of between $100 \sim 13.333\%$. As shown in the 4th ramification, those with an an-nual income between 300 to 325 thousands, and having no history of holding a credit card or delay payment, and with a family status of married, have a bad debt probability of100%, accounting for 1.935% of valid sample compared with total samples. The 4th ram-ification is shown in Fig. 7:

## 6. Confirmation and control

The study focuses on the issue of classification. In assessing the effectiveness of a classification system, the conventional approach is to build a confusion matrix in advance, which was most fre-quently used model by previous researches in measuring the per-formance of a classification system. The higher the overall correctness of a model indicates a higher overall accuracy of the mode. Confusion matrix is effective in determining any misjudg-ments existing in the predicted results, and can further identifying those classifications that are easily misjudged.

After classifying with experimental model, some individuals are correctly classified, while some others are probably being misclas-sified. As shown in Table 4, credit ratings classified as bad debt "Yes" and "No" may have the following four combinations: (1) Acceptance; (2) Refusal; (3) Model I error; and (4) Model II error.

1. Sensitivity: One element used to make correct prediction of a good client:

$$\text{Sensitivity} = \frac{\text{a predicted normal client is actually a normal client}}{(\text{a predicted normal client is actually a normal client} + \text{a predicted bad debt client is actually a normal client})}$$
$$= \frac{1}{1+2}.$$

2. Specificity = One element used to make correct prediction of a bad client:

$$\text{Specificity} = \frac{\text{a predicted bad debt client is actually a bad debt client}}{(\text{a predicted bad debt client is actually a bad debt client} + \text{a predicted normal client is actually a bad debt client})}$$
$$= \frac{4}{3+4}.$$

3. . Total correct prediction ratio=

$$\frac{(\text{a predicted normal client is actually a normal client} + \text{a predicted bad debt client is actually a bad debt client})}{(\text{a predicted normal client is actually a normal client} + \text{a predicted bad debt client is actually a bad debt client} + \text{a predicted bad debt client is actually a normal client} + \text{a predicted normal is actually a bad debt client})}$$
$$= \frac{1+4}{1+2+3+4}.$$

4. Model I error (misjudging ratio of good client) = 1− sensitivity (correct prediction ratio of good client).
5. Model II error (misjudging ratio of bad client) = 1− specificity (correct prediction ratio of bad debt client).

**Table 6**
Neural network confusion matrix.

|  | Good | Bad | All |
|---|---|---|---|
| Determine good | 259 | 14 | 273 |
| Determine bad | 8 | 29 | 37 |
| All | 267 | 43 | 310 |

**Table 7**
Neural network combine CART decision tree.

|  | Good | Bad | All |
|---|---|---|---|
| Determine good | 258 | 2 | 260 |
| Determine bad | 9 | 41 | 50 |
| All | 267 | 43 | 310 |

**Table 8**
Error value in experimental model.

| Experimental model | Type I error | Type II error | All correct rate |
|---|---|---|---|
| Neural network | 3% | 32.6% | 92.9% |
| Neural network combine CART decision tree | 3.4% | 4.7% | 96.5% |

In terms of statistics testing, the loss caused by Model I error (convict an innocent man as guilty) is far severer than that of Model II (provisional release a guilty man). Facing with the mounting pressure of NPL (non-performing loan) created by a series of card delinquency, the negative impact of misclassifying a bad client as good client (a typical Model II error) is greater than misclassifying a good client as a bad client (a typical Model II error). Calculation of the confusion matrix is as follows:

The greater the Model II error represents greater ratio of misjudging a bad client as a good client, thus creates greater bad debt loss for the bank, thus a greater failure ratio of the model's predictions. Greater Model I error, that is, misjudging a good client as a bad client, will reduce the bank's chance of profits. This approach was frequently used by previous researchers in measuring the performance of classification: the higher the overall correctness is, the higher the accuracy of the model.

In the experimental model, we combine ANN with the CART decision tree, to verify whether such combination can bring better prediction results. Therefore, we induce the concept of confusion matrix, and Model I & II errors into the verified experiment results, to interpret the concrete effect of the strengthened experimental model with various values.

1. ANN confusion matrix: predictions of ANN model on bad debt data of the study case, and the values of Model I error and Model II error are shown in Table 6.
2. ANN combining with CART decision confusion matrix: predictions of ANN and CART decision tree on bad debt data of the study case, and the values of Model I error and Model II error are shown in Table 7.
3. After computing the confusion matrix, we compared the results of ANN and ANN combined with CART decision. Findings show that the prediction results of ANN model has a better discriminating power (3%) in determining Model I error (misjudging good as bad); and a relatively higher error (32.6%) in determining Model II error (misjudging bad as good), reaching an overall

correct prediction rate of 92.9%. The prediction results of ANN combined with CART decision has a better discriminating power (3.4%) in determining Model I error (misjudging good as bad); and a fairly good effect (4.7%) in determining Model II error (misjudging bad as good), reaching an overall correct prediction rate of 96.5%. Related data are shown in Table 8.

## 7. Conclusion

With its effective prediction of the probability of applicant's future default, the credit scoring model not only can greatly upgrade the handling efficiency of consumer finance, but can also effectively minimize the credit risk encountered the bank. The bank's profits and loss depends on the quality of its credit granting system. However, the current credit granting policy popularly adopted by most banks relies most on a scoring method as the guideline, which has been proven weak in constructing accurate credit granting decision-making, and worse even, it has neglected the trend of environmental changes. By combining the ANN and the rules of CART decision tree, the study has effectively reduced the errors of Model I and Model II to 3.4% and 4.7% respectively, and brought up the forecast accuracy of the entire credit granting to as high as 96.5%. As shown in the experimental analysis model, there are four factors that have a significant impact on the applicant's default probability: (1) Annual income, with an annual income of NT$300 thousands as the baseline; (2) Seniority, with 10-year work history as the baseline; (3) Credit granting and real estate status, with ownership/non-ownership of private house as the baseline; and (4) Family condition, with married/single as the baseline. In making credit checking, the bank should be particularly cautious against the item (s) that has any probability of default. Applicants with more than two possible default items are the high bad debt rate cohort, and should be denied request for card granting. The above rules summed up in the study may provide the frontline auditors a valuable reference to credit banking auditing business, so that their credit scoring sheet may become a more accurate decision-making tool.

## Acknowledgements

## References

Cabena, P. (1998). *Discovering Data Mining: From Concept to Implementation*. Upper Saddle River, New York: Prentice Hall.
Frawley, W. J., Piatetsky-Shapiro, G., & Matheus, C. J. (1991). Knowledge discovery in databases: an overview. In *Knowledge Discovery in Databases* (pp. 213–228). Cambridge, MA: AAAI/MIT.
Gluck, M. A., & Myers, C. E. (2001). *Gateway to Memory: An Introduction to Neural Network Modeling of the Hippocampus and Learning*. Cambridge, Mass: MIT Press.
Grupe, F. H., & Owrang, M. M. (1995). Database mining discovering new knowledge and cooperative advantage. *Information Systems Management*.
Hall, J., Mani, G., & Barr, D. (1996). *Applying computational intelligence to the investment process. Proceedings of CIFER-96: Computational Intelligence in Financial Engineering, New York, March 1996*. Piscataway, N.J: IEEE Press.
Jane, Zhen-Fu, & Zhan, Feng-lung (1998). The production portfolio strategy of the decision analysis- in a semiconductor plant, for example. *Technology Management Journal, 3*(1), 137–156.
Michael, J.A. Berry and Gordon Linoff. (1999). Mastering data mining. The Art & Science of Customer Relation Management.
Smith, Nina, & Gupta, Nabanita-Datta (2002). Children and career interruptions: the family gap in Denmark. *Economica, 69*, 609–629.
Widrow, B., Rumelhart, D., & Lehr, M. A. (1994). Neural networks: applications in industry business and science. *Communications of the ACM, 37*(3), 93–105.
Zhang, G. P. (2004). *Neural Networks in Business Forecasting*. Hershey, Pa: Idea Group Publishing.