

# Hybrid mining approach in the design of credit scoring models

Nan-Chen Hsieh\*

*Department of Information Management, National Taipei College of Nursing, No. 365, Min-Ten Road 11257, Taipei, Taiwan, ROC*

---

## Abstract

Unrepresentative data samples are likely to reduce the utility of data classifiers in practical application. This study presents a hybrid mining approach in the design of an effective credit scoring model, based on clustering and neural network techniques. We used clustering techniques to preprocess the input samples with the objective of indicating unrepresentative samples into isolated and inconsistent clusters, and used neural networks to construct the credit scoring model. The clustering stage involved a class-wise classification process. A self-organizing map clustering algorithm was used to automatically determine the number of clusters and the starting points of each cluster. Then, the K-means clustering algorithm was used to generate clusters of samples belonging to new classes and eliminate the unrepresentative samples from each class. In the neural network stage, samples with new class labels were used in the design of the credit scoring model. The proposed method demonstrates by two real world credit data sets that the hybrid mining approach can be used to build effective credit scoring models.

© 2005 Elsevier Ltd. All rights reserved.

*Keywords:* Data mining; Credit scoring model; Clustering; Class-wise classification; Neural network

---

## 1. Introduction

In response to the recent growth of the credit industry and the management of large loan portfolios, the industry is actively developing credit and behavioral scoring models. Credit scoring models help to decide whether to grant credit to new applicants using customer's characteristics such as age, income and marital status (Chen & Huang, 2003). Behavioral scoring models help to analyze purchasing behavior of existing customers (Setiono, Thong, & Yap, 1998). This study utilizes a hybrid mining approach in the design of credit scoring models to support credit approval decisions.

A simple parametric statistical model, linear discriminant analysis, was the first model employed for credit scoring. However, analysts have questioned the appropriateness of linear discriminant analysis for credit scoring because of the categorical nature of the credit data and the fact that the covariance matrices of the good and bad credit classes are not likely to be equal. Researchers are now investigating more sophisticated models to overcome some

of the deficiencies of the linear discriminant analysis model. One of the efforts is leading to the investigation of non-parametric statistical methods, that is, neural networks for scoring applications.

Swales and Yoon (1992) used neural networks to differentiate stocks. They found that the neural networks performed significantly better than the linear multiple discriminant models. Tam and Kiang (1992) compared the neural network approach with a linear classifier model, logistic regression model, neural network model, and decision tree model to predicate bank failures. They demonstrated that neural networks are more accurate, adaptive, and robust than to other methods. Zhang et al. (1999) similarly showed that neural networks are significantly better than logistics regression models in bankruptcy predication. Desai, Crook, and Overstreet. (1996) performed an experiment that the performance of discriminant analysis is comparable to the performance of back-propagation neural networks in classifying loan applicants into good and bad credit. They pointed out that more customized architectures might be necessary for building effective generic models to classify consumer loan applications in the credit union environment. West (2000) investigated the credit scoring accuracy of five neural

---

\* Tel./fax: +886 2 2822 7101 2200.

E-mail address: nchsieh@mail1.ntcn.edu.tw

network models, and reported that a hybrid architecture of neural network models should be considered for credit scoring applicants. That is, non-parametric and hybrid design architecture will be very useful in developing effective credit scoring models. The meaning of an effective credit scoring model is two-fold, relating to accuracy and easy interpretation of classified results.

Generally, hybrid design architectures for creating more accurate classifiers that uses neural networks guided by clustering algorithms (Gopalakrishnan, Sridhar, & Krishnamurthy, 1995; Sung, 1998) or genetic algorithms (Kim & Street, 2004) were proposed. The other studies focus on integrating the multivariate analysis and neural network to increase clustering accuracy. Punj and Steward (1983) proposed a two-stage method that combines Ward’s minimum variance method and the K-means method. Balakrishnan et al. (1996) integrated unsupervised FSCL neural networks with the K-means method. Kuo, Ho, and Hu (2002) proposed a two-stage method, which uses the self-organizing map to determine the number of clusters and then employs the K-means algorithm to classify samples into clusters. This study addresses the benefit by investigating a simple but effective hybrid utility of clustering and neural network techniques in the design of a credit scoring model.

As shown in Fig. 1, the proposed hybrid scoring model has two processing stages. In the clustering stage, samples are grouped into homogeneous clusters. In particular, unrepresentative samples are indicated as ‘isolated’ and ‘inconsistent’ clusters. Herein, isolated clusters are thinly populated clusters; inconsistent clusters are those clusters with inconsistent class values. To improve the performance

of the credit scoring model, isolated clusters should be eliminated. Since inconsistent clusters may be important and eliminating them may cause the loss of valuable information, investigating inconsistent clusters in detail will be very helpful in understanding customers. In the neural network stage, the neural network used samples with new class labels to create models for predicate consumer loans.

For a better understanding of our study, Section 2 begins with an overview of credit scoring models in general. Section 3 offers a hybrid credit scoring model and considers why it should perform better than other credit scoring models. Following this discussion, Section 4 empirically tests the hybrid credit scoring model using two real world credit data sets. Finally, Section 5 discusses the findings of the experiment and offers observations about practical applications and directions for future research.

## 2. Description the analysis methodology

### 2.1. Credit and behavioral scoring models

Credit and behavioral scoring models (Thomas, 2000) are one of the most successful applications of statistical and operational research modelling in finance and banking, and the number of scoring analysts in the industry is constantly increasing. The main objective of both credit and behavioral scoring models is to classify samples into homogeneous groups (Lancher, Coates, Shanker, & Fant, 1995). Hence, scoring problems are related to classification by statistical analysis (Hand, 1981; Hsieh, 2004; Johnson & Wichern,

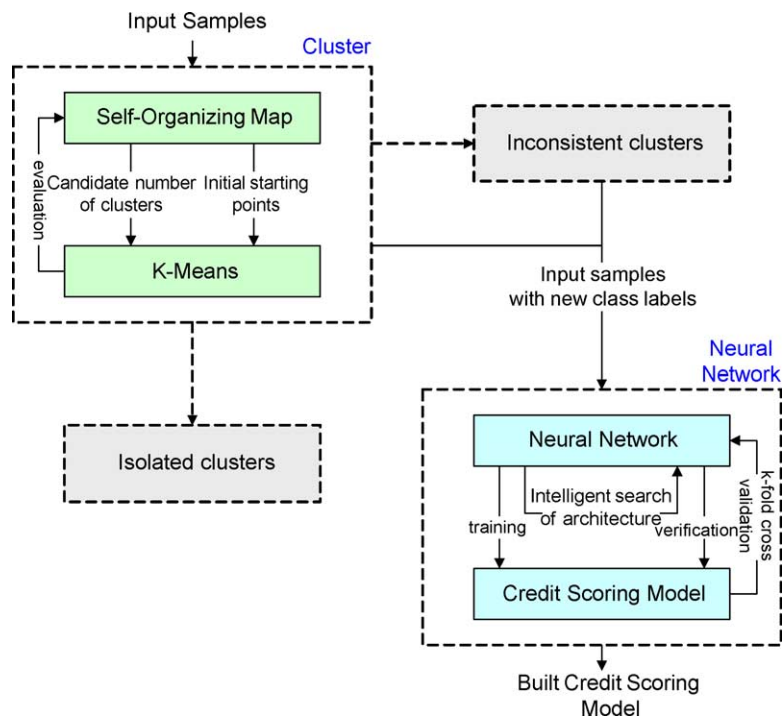


Fig. 1. A hybrid mining credit scoring system.

متن کامل مقاله

دریافت فوری ←

**ISI**Articles

مرجع مقالات تخصصی ایران

- ✓ امکان دانلود نسخه تمام متن مقالات انگلیسی
- ✓ امکان دانلود نسخه ترجمه شده مقالات
- ✓ پذیرش سفارش ترجمه تخصصی
- ✓ امکان جستجو در آرشیو جامعی از صدها موضوع و هزاران مقاله
- ✓ امکان دانلود رایگان ۲ صفحه اول هر مقاله
- ✓ امکان پرداخت اینترنتی با کلیه کارت های عضو شتاب
- ✓ دانلود فوری مقاله پس از پرداخت آنلاین
- ✓ پشتیبانی کامل خرید با بهره مندی از سیستم هوشمند رهگیری سفارشات