

# Combination of Feature Selection and Optimized Fuzzy Apriori Rules: The Case of Credit Scoring

Seyed Sadatrasoul, Mohammad Gholamian and Kamran Shahanaghi  
Faculty of Industrial Engineering, Iran University of Science and Technology, Iran

**Abstract:** Credit scoring is an important topic and banks collect different data from their loan applicants to make appropriate and correct decisions. Rule bases are favourite in credit decision making because of their ability to explicitly distinguish between good and bad applicants. This paper, uses four feature selection approaches as features pre-processing combined with fuzzy apriori. These methods are stepwise regression, Classification And Regression Tree (CART), correlation matrix and Principle Component Analysis (PCA). Particle Swarm is applied to find the best fuzzy apriori rules by searching different support and confidence. Considering Australian and German University of California at Irvine (UCI) and an Iranian bank datasets, different feature selections methods are compared in terms of accuracy, number of rules and number of features. The results are compared using T test; it reveals that fuzzy apriori combined with PCA creates a compact rule base and shows better results than the single fuzzy apriori model and other combined feature selection methods.

**Keywords:** Fuzzy apriori, feature selection, particle swarm, credit scoring.

Received September 7, 2012; accepted June 20, 2013; published online April 23, 2014

## 1. Introduction

In today's competitive economy, credit scoring is widely used in banking industry. Every day, the individual's and company's records of past borrowing and repaying actions are gathered and analyzed by information systems. Banks use this information to determine the individual and company's creditworthiness. The process of lending can be divided into four main phases, including pre-application, application, performance and collection [34]. In this paper, we will address the credit scoring problem which is the critical issue in the application phase. Credit scoring is used to answer one key question-what is the probability of default within a fixed period, usually one year.

There are many methods suggested to classify loan applicants in the credit scoring including statistical and intelligent methods. Logistic regression and linear discriminant analysis are statistical methods that are effectively used in credit scoring [11]. There are also many intelligent methods applied to the problem including neural networks, Support Vector Machines (SVMs), Bayesian networks, case based reasoning and decision trees [17]. Intelligent methods neural network and SVMs are used more frequent and owing to their nonlinear fitness and generalization capabilities, better classify the loan applicants [9, 15, 18]. Some studies have shown the superiority of neural networks, SVM, decision trees and other intelligent methods to statistical methods [7, 16, 26].

Recent studies are focus on hybrid methods with the aim of making synergy with their strengths and

covering their weaknesses. In some methods, both statistical and intelligent methods are used together. Lee *et al.* [20] developed a hybrid neural discriminant method along with BP neural network and discriminant analysis and demonstrated the accuracy dominance of hybrid method rather than individual applications. A two-stage hybrid procedure with artificial neural networks and multivariate adaptive regression is also proposed and provides better results [19]. Tsai and Chen [32] developed hybrid machine learning methods in four main areas and consequently, found the best accuracy in the hybrid application of neural network and logistic regression. There are also studies which show that hybrid meta heuristic methods are used together with intelligent methods. An integration of SVMs, genetic algorithms and F-score is also studied [15]. In the last decade, the use of ensemble methods increased in the area [33, 34, 35, 36, 37]. Neural network ensemble strategies include cross validation, bagging and boosting for financial decision applications which have been studied and shown better accuracy rate and generalization ability [37]. Ensemble learning is an open issue in recent year's studies [22].

However, in practice, many of the above-mentioned methods cannot be used; more because of robustness and transparency needs as well as the regulator's auditing on the credit scoring [31]. Instead, it seems the rule-based method can be successfully used since, the banks can easily interpret the results and explore the rejecting reasons to the applicant and regulatory auditors. Unfortunately, there is actually a little literature in the field of rule based credit scoring. Ben-David [2] has provided a new method for rule pruning

and examined his method on the credit scoring data set. Martens *et al.* [24] used SVM for rule induction in the credit scoring problems. Baesens *et al.* [1] used and evaluated three rule extraction neural network methods including Neurorule, Trepan and Nefclass for rule extraction in three real life data bases including German, Bene 1 and Bene 2 credit databases. They showed nerorule and Trepan yield better classification accuracy compared to the C4.5 method and the logistic regression. Finally, they visualized the extracted rule sets using the decision table. Hoffmann *et al.* [12] introduced a new learning method for fuzzy rule induction based on the evolutionary methods. Malhotra and Malhotra [23] used the Adaptive Neuro Fuzzy Inference Systems (ANFIS) for rule induction and showed that this method works better than discriminant analysis on their credit scoring dataset, gathered from credit unions. They used the back propagation in the learning process of the rules' membership function. Fuzzy rules are more attractable and robust because the rules are expressed in terms of linguistic features, which are usually used by the experts and easier to interpret.

Feature selection increases the learning ability of models by reducing the effect of the curse of dimensionality and over fitting, increasing the model generalization ability and learning speed. Feature selection approaches have two main dimension filter and wrapper. Filter approaches select important features independent from learning algorithms. They rely on different measures of features extractable knowledge which include information, distance and dependency [8]. Wrapper models usually use the accuracy of a selected model to select the subsets. Filters are faster than wrappers, but the latter may find better subsets. There are many studies in credit scoring, which consider feature selection as a preprocessing step. Šušteršič *et al.* [30] used PCA and genetic algorithm as a feature selection step for building neural network model, the result is better for PCA. Chen *et al.* [4] used CART and MARS with SVM and showed that the feature selection improves the accuracy of the SVM. Chen and Li [3] used decision tree,  $F$  score, Rough set and discriminant analysis and SVM; they showed that feature selection almost improves accuracy and Area Under Curve (AUC) in UCI credit datasets. Ping and Yongheng [27] used different feature selection approaches including  $T$  test, correlations, stepwise, Classification And Regression Tree (CART), rough set and MARS combined with CART, SVM and neural networks; the rough set and SVM combination yields to the best results. Many of the mentioned studies emphasize on designing more sophisticated models through feature selection to improve the accuracy of the models.

In this study, a fuzzy association rule is built in a fuzzified credit data set with different feature selection approaches including filter and wrapper approaches.

For designing an optimized fuzzy apriori, Fuzzy Support (FS) and Fuzzy Confidence (FC) must be found. The particle swarm search algorithm is an

appropriate method to find the best of the mentioned pairs. Because of conflict between fewer numbers of rules on the one hand and higher performance in the rule base quality on the other hand, a multi objective fitness function was developed to achieve the best parameters of support and confidence. The first objective of fitness function is to maximize the classification accuracy and the second is to minimize the number of rules. To compare the different credit scoring models, the procedure is done for the whole of feature space and the other four reduced feature spaces.

This study is divided into five major sections: Section 2 describes the basic concepts of hybrid model building. Section 3 introduces the experimental design including datasets, feature selection parameters, Fuzzy apriori model parameters and evaluation criteria. Results and discussions are presented in section 4 and finally, study concluded in section 5.

## 2. Basic Concepts of the Hybrid Model Building

### 2.1. Balancing the Data

Model building on imbalanced data has many problems, including over fitting and having poor rate of learning. In general, the more the balance between good and bad ones, the more accurate the resulting score is [10]. Real credit scoring datasets are often imbalanced because the number of bad applicants is often fewer than the number of good ones. This study uses a real world Iranian credit dataset in addition to, UCI credit datasets; therefore, data balancing must be taken into consideration for Iranian credit dataset.

These techniques are different and include random over sampling, random under sampling, model based and stratified; each of them has its own strengths and weaknesses.

### 2.2. Stepwise Regression

Stepwise regression is one of the methods to find the best combination of features when using regression.

The main procedure is to find one single best feature and add other features iteratively to find the appropriate regression.

### 2.3. Pearson Correlation

Correlation is used to show the correlation of two groups of data. It is used to show how closely the two groups of features are related. Two groups of features are called highly correlated if changes in one feature results in similar change in other features at the same side. There are several correlation coefficients; the most common of them is Pearson. Pearson correlation was developed by Karl Pearson [29]. Pearson correlation can be used for feature selection.

### 2.4. Classification and Regression Tree

Decision tree models represent knowledge in a tree form. Quinlan introduced the first decision tree algorithms [28]. Since, then several decision tree algorithms have been introduced such as ID3, C5, CHAID and CART [28]. CART has been widely used in many areas of science including credit scoring [36]. CART is a nonparametric statistical method and uses a generalization of the binomial variance called the Gini index [21]. CART produces a classification tree when the dependent feature is categorical and a regression tree when it is continuous.

### 2.5. Principal Component Analysis

Principal Component Analysis (PCA) reduces the dimensions of data sets with interrelated features. PCA transforms the dependent features to a new set of independent features that can be used by classifiers directly. The transformation is done in such a way that the most appropriate information is collected using a smaller number of features called Principal components. Each component is combined from a linear function of the variance-covariance matrix of original dependent features. The analysis provides the percentage of variance explained by each Principal component and the correlations between each principal component and the original features [30].

### 2.6. Fuzzy Apriori

Better classification result from apriori could be achieved if support and confidence parameters are adjusted. In this paper, an extension to the method proposed by Hu *et al.* [14]. Genetic algorithm method is replaced with the particle swarm because of its better convergence and ability to reach the better results is less time. This study used particle swarm to find two parameters of the fuzzy association rules which are FS and FC. Also, a multi objective fitness function is used; the first objective is to maximize the classification accuracy and the second objective is to minimize the number of rules. In order to, ensure the validity of results the datasets are randomly partitioned to train and test sets with the ratio of 70/30. Using particle swarm, pairs of (FS, FC) are searched and the one with the best fitness is chosen. This algorithm pseudocode is described below:

Algorithm 1: Fuzzy apriori rule set Algorithm

// main variables and constants  
 $n$  = Number of initial positions.  
 $W_a$  = Accuracy weight.  
 $W_g$  = The number of rules weight.  
 $K$  = Number of fuzzy linguistic variables.  
 Positions pair = Order pair of  $(FS_i, FC_i)$  between zero and one.  
 $T_{max}$  = Maximum number of generations.  
 $G_{best}$  = Global best position.  
 $L_{best}$  = Local best position.  
 $FC$  = Array  $[1*n]$  of fuzzy confidence of fuzzy apriori rule set.

$FC_{best}$  = The best fuzzy confidence of fuzzy apriori rule set.  
 $FS$  = Array  $[1*n]$  of fuzzy support of fuzzy apriori rule set.  
 $FS_{best}$  = The best fuzzy support of fuzzy apriori rule set.  
 $frules$  = Fuzzy apriori rule set.  
 $F$  = Rule set fitness function.  
 Position = The position of each particle which stands for an appropriate rules set.

```
//algorithm
Perform fuzzy partitioning (k=3)
Position = Generate n initial positions (n)
For (i=1 to  $T_{max}$ )
{
    [FS, FC] = Calculating FSFC (position)
    For (i=1 to n)
    {
        Item set = Create frequent fuzzy item sets ( $FS_i$ )
        Primary fuzzy rules = Create fuzzy rules (Item set,  $FC_i, FS_i$ )
        Reduced fuzzy rules = Reduce redundant rules (Primary fuzzy rules)
        Rules weights = Use adaptive rules to adjust fuzzy rule's weights (Reduced fuzzy rules)
        Final fuzzy rules = fined final fuzzy rules (Reduced fuzzy rules, Rules weights)
         $F = w_a * Accuracy + w_g * Number\ of\ rules$ 
        //Compute the fitness
    }
    For (i=1 to n)
    {
         $G_{best}$  = Find the  $G_{best}$  (F)
         $L_{best}$  = Find the  $L_{best}$  (F)
        Update Velocity and position ()
        Position = Maintain the best position ()
    }
}
Frules = Create Fuzzy apriori rules ( $FS_{best}, FC_{best}, G_{best}, L_{best}$ )
Return [Accuracy, Number of rules] = find accuracy (frules)
```

Figure 1 shows the overall steps of the decision making process. Balancing the imbalanced data and feature selection steps are illustrated as the pre-processing steps.

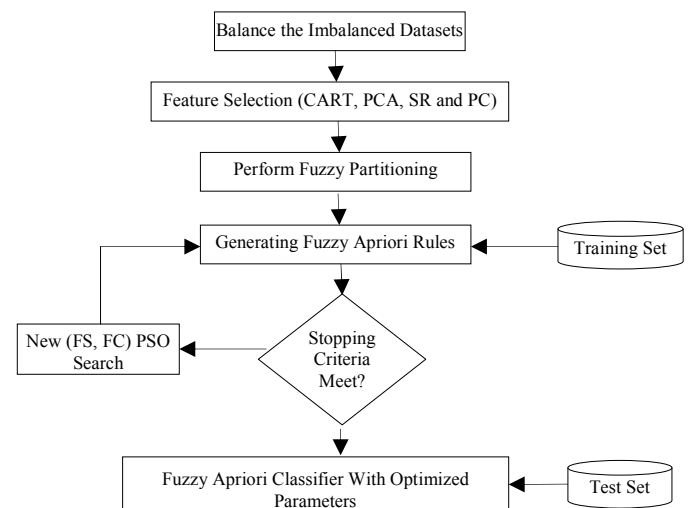


Figure 1. The overall steps of the hybrid decision making process.

### 3. Experimental Design

#### 3.1. The Data Sets

To evaluate the performance of different feature selection methods, three datasets are used. Australian and German credit data sets from University of California at Irvine (UCI) machine learning repository are applied. These datasets can be found at <http://archive.ics.uci.edu/ml/datasets.html>. An Iranian commercial bank dataset is also used to evaluate the proposed algorithm. Table 1 shows the characteristics of the two UCI credit datasets.

Table 1. UCI credit datasets description.

Dataset Name	Data Size	Input Features		
		Total	Continuous	Categorical
German	1000	24	7	17
Australian	690	14	6	8

Australian credit dataset has been successfully used for credit scoring and evaluation systems in many previous works, especially using intelligent methods [12, 13, 16, 25, 33, 35, 37]. The dataset includes 15 characteristics among them; eight characteristics are categorical and six characteristics are continuous. The dataset includes 690 instances of loan applicants; instances are labeled as classes 1 (worthy, 307 instances) and 0 (unworthy, 384 instances).

German dataset is also used in many works. For each applicant, the dataset include 24 input features describe the credit history, account balances, loan purpose, loan amount, employment status and personal information. This data set only consists of numeric attributes. This dataset includes 1000 instances of loan applicants; the data instances are labelled as classes 1 (worthy, 700 instances) and 2 (unworthy, 300 instances).

The real world database of a major Iranian bank is also used for the experiments. The initial dataset include 1109 corporate applicants 60 financial and non financial features in the period from 2009 to 2012. First, a data cleaning process is done on the data. In general, data cleaning include removing redundant, outliers, data and missing values. There were a few missing values for some corporate, some of them lack financial data and others lack the result of their loans; in fact, they were in the process of debt repay. So, 387 corporate are excluded. From 722 remained corporate, 652 companies are credit worthy and 70 was unworthy.

Once the data cleaning process was completed, the categorical features include the type of industry; type of company and type of book are converted to numerical features using dummy features. The results and descriptions of the changes done are shown in Appendix 1. Using dummy features number of features increased to 60.

The main dataset has nearly a 90/10 class distribution. To avoid over fitting, the G/B odd's ratio of 3:1 which is proposed in other major credit scoring studies is used in this paper [37, 38]. Furthermore, since over-sampling generally gives better performance than under sampling, Random minority Over Sampling (ROS) is used to over sample the applicants labelled as bad. Finally, the total number of data, number of goods and number of bads are reported in Table 2.

Table 2. Iran credit dataset description.

Status	Data Size	Good / All (%)	Inputs Features		
			Total	Continuous	Categorical
Before Cleaning	1109	NA	51	43	8
After Cleaning	722	90.3	60	39	21
Balanced Dataset	869	75.02	60	39	21

#### 3.2. Feature Selection

The parameter's setting designed for different feature selection approaches are described in the following:

- *Pearson Correlation*: Pearson correlation at 90% level of importance is used to find significant features in order to distinguish credit worthy customers from non worthy ones.
- *PCA*: For PCA factors which their Eigen values are greater than one are considered.
- *Stepwise*: 0.05 is used as a probability of  $F$  to model entry and 0.1 is used as a limit of removal for considering the co linearity problem.
- *CART*: Gini measure of impurity is used and maximum tree depth is set to five levels.

#### 3.3. Hybrid Fuzzy Apriori with PSO

The parameters for fitness function and particle swarm are defined as ( $W_{AC}=200$ ,  $W_G=1$ ,  $N_p=100$ ,  $C_1=1$ ,  $C_2=1$ ,  $r_1=1$ ,  $r_2=1$ ,  $T_{max}=100$ ) and the last runs are done using the tuned parameters. It's suggested that the learning rates should be specified as  $0 < \eta_r < \eta_i < 1$  for example,  $\eta_r=0.0001$  and  $\eta_i=0.1$  [14]. Because in the rule bases number of rules and accuracy rate is important together, the results analysis contains both of them. The proposed algorithm is run at least five times with 20 particles for each of the models separately. In each run a rule base discovered and evaluated using the fitness function which is defined.

#### 3.4. Evaluation Criteria

70 percent of the data set is used to train the hybrid model and other 30 percent is used to test the results. The model which has the highest prediction accuracy and lower number of rules is selected as the most suitable model.

### 4. Results and Discussion

Table 3 shows classification accuracy, number of rules and feature extracted for different experiments. The best test set classification accuracy, the lowest number of rules and features are bolded for each data set.

Table 3. Classification accuracy, number of rules and number of features for different datasets and feature selection methods.

Method	German Credit			Australian Credit			Iranian Credit		
	Accuracy	Number of Rules	Number of Features	Accuracy	Number of Rules	Number of Features	Accuracy	Number of Rules	Number of Features
All+FA	73.51	4	24	<b>84.83</b>	7	14	74.23	4	59
PC+FA	73.51	2	16	<b>84.83</b>	3	13	74.23	<b>1</b>	14
SR+FA	73.51	3	10	<b>84.83</b>	3	6	73.46	2	<b>3</b>
CART+FA	71.19	2	8	<b>84.83</b>	3	8	74.23	1	8
PCA+FA	<b>73.51</b>	<b>1</b>	<b>5</b>	77.73	<b>2</b>	<b>5</b>	73.46	1	5

Paired sample *T* test with significant different at the 5% level of importance is used to compare feature selection methods. By evaluating means of results the feature selection methods are also ranked. Table 4 shows the ranking results for the number of rules and number of features separately. The accuracy measure is disregarded in Table 4, because the baseline model and feature selection methods have not significant difference. The aim of this paper is to find out the best method, so the method with the best performance is selected as the first rank and the other method's rank is extracted by considering its *T* statistic value. In other words, greater difference of *T* statistic between two pairs of feature selection methods includes the best performer, distinguishes the next method's rank. This process is done iteratively until all the feature selection methods ranks are extracted.

Therefore, for the number of rules, PCA is ranked first and then the CART and PC together placed the second. For the number of features extracted, PCA is placed first and SR is placed next. Other methods are ranked below the before mentioned methods. It can be seen that PCA is the better feature selection method which provides lower number of rules and number of features at an acceptable accuracy rate. On the other hand, stepwise and CART compete for the second place. Stepwise extract features well but CART shows lower number of rules.

Table 4. Ranking of different feature selection methods.

Performance Measure	Rank
Number of Rules	PA> CA, PC> SR> FA
Number of Features	PA> SR> CA> FA> PC

Note: PC for Pearson; SR for stepwise; CA for CART; PA for PCA; FA for baseline model (Fuzzy apriori)

Table 5 shows the original, extracted number of and percent of features for different feature selection methods and datasets. It can be seen that PCA extracts the most valuable features and therefore, it has the best average extraction rate; 78.33% of features are extracted on average by PCA. Stepwise regression and CART are placed second and the PC is the worst player.

Table 5. Number of selected features for different feature selection methods and datasets.

Method	German Credit		Australian credit		Iranian credit		Average Extraction Rate
	Number of Features	Feature Extraction Rate	Number of Features	Feature Extraction Rate	Number of Features	Feature Extraction Rate	
Fuzzy Apriori	24	NA	14	NA	59	NA	NA
PC+Apriori	16	33.33%	13	7.14%	14	76.27%	38.92%
SR+Apriori	10	58.33%	6	57.14%	3	94.92%	70.13%
CART+Apriori	8	66.67%	8	42.86%	8	86.44%	65.32%
PCA+Apriori	5	79.17%	5	64.29%	5	91.53%	78.33%

Figure 2 Shows the Accuracy rate versus feature extraction rate for different datasets. The results show that PCA extracts the appropriate features better than other methods at an acceptable accuracy rate, although it is the worst player in the accuracy of Australian credit dataset. On the other hand, PC extracts important features poorly. CART and SR are placed at the middle range; they perform well in all three datasets.

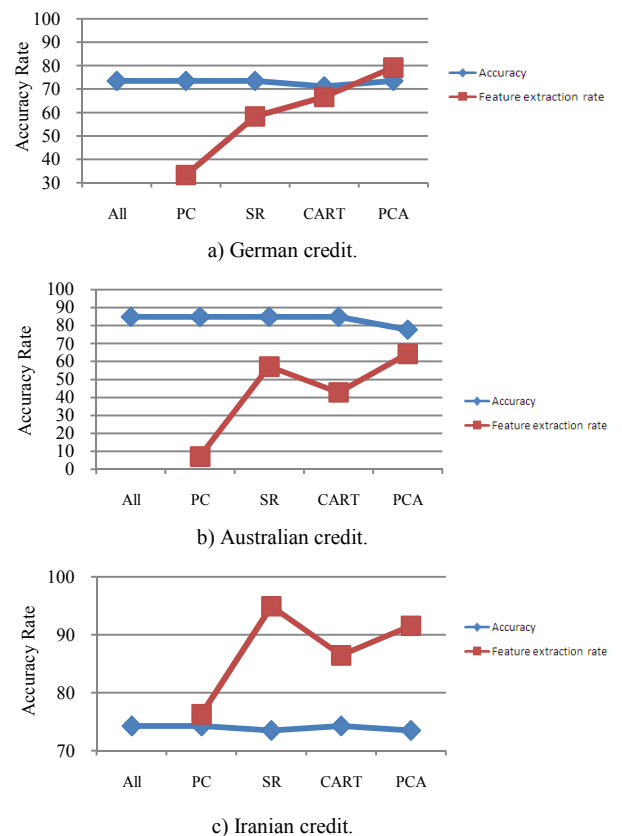


Figure 2. Accuracy rate versus feature extraction rate.

Figure 3 shows the Accuracy rate versus rule extraction rate for different datasets. The results show that although, PCA extracts better rules than other methods, it has the worst accuracy rate in Australian credit dataset. On the other hand, PC extracts important rules better and simultaneously all reported accuracy rates are at the best level. SR is the worst player, in Iranian credit dataset; it shows the lower accuracy with more rules. CART is placed at the middle range. It performs well in all three datasets.

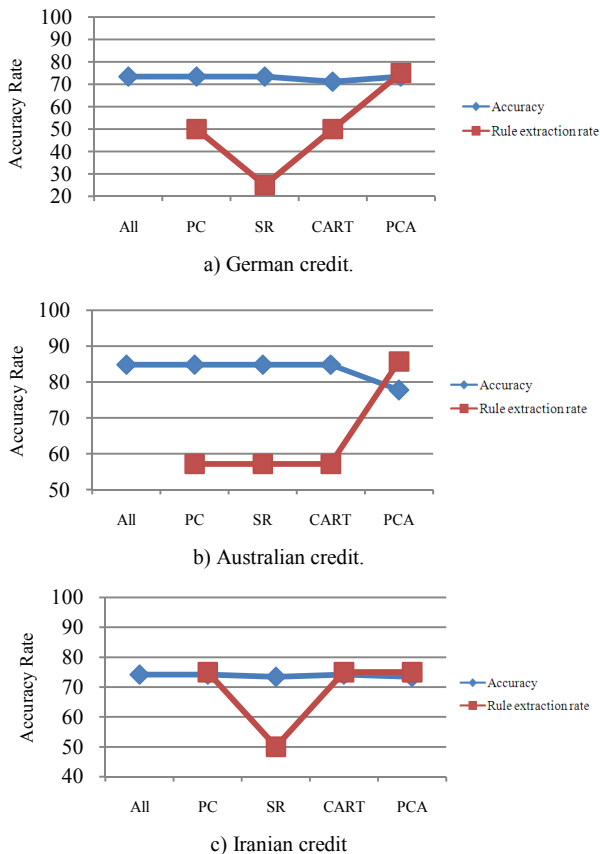


Figure 3. Accuracy rate versus rule extraction rate.

## 5. Conclusions

Predicting the credit score of loan applicants is a critical issue. Different data mining techniques have been used to predict loan applicants score. Rule based methods are one of the most favourite methods introduced better accuracy and lower number of rules enhance their chance of usability. This study is the first study which uses fuzzy apriori and feature selection techniques in credit scoring. Fuzzy apriori, results in rules with linguistic features which can be used by bank's loan officer better, this is a main issue in small and mid-sized banks, which mainly haven't a credit scoring system. Feature selection is used to filter out irrelevant features and consequently, it reduces the model training time and costs; it can also improve the model performance. This paper, used and compared four different feature selection methods over three credit datasets. They are stepwise, Pearson correlation, PCA and CART. Accuracy, number of rules and

number of extracted features is used to evaluate and compare the results.

Feature selection on fuzzy apriori, shows better results in the number of rules and number of features used. On average, PCA is the best method and CART placed the second in the number of rule's measure. On the other hand, stepwise place second and outperforms CART in the feature extraction measure.

The potential future work in the area could be considered in enhancing the algorithm performance using the multiple minimum supports to find frequent item sets in order to, increase the rule bases quality. In the area of applications, findings can be used in other related business domains include bankruptcy prediction, customer churns, stock prices up and down predictions. Furthermore, other feature selection methods can be compared with PCA as the best performer on apriori.

## References

- [1] Baesens B., Setiono R., Mues C., and Vanthienen J., "Using Neural Network Rule Extraction and Decision Tables for Credit-Risk Evaluation," *Management Science*, vol. 49, no. 3, pp. 312-329, 2003.
- [2] Ben-David A., "Rule Effectiveness in Rule-Based Systems: A Credit Scoring Case Study," *Expert Systems with Applications*, vol. 34, no. 4, pp. 2783-2788, 2008.
- [3] Chen L. and Li C., "Combination of Feature Selection Approaches with SVM in Credit Scoring," *Expert Systems with Applications*, vol. 37, no. 7, pp. 4902-4909, 2010.
- [4] Chen W., Ma C., and Ma L., "Mining the Customer Credit using Hybrid Support Vector Machine Technique," *Expert Systems with Applications*, vol. 36, no. 4, pp. 7611-7616, 2009.
- [5] Chi W. and Chch H., "A Hybrid Approach to Integrate Genetic Algorithm into Dual Scoring Model in Enhancing the Performance of Credit Scoring Model," *Expert Systems with Applications*, vol. 39, no. 3, pp. 2650-2661, 2011.
- [6] Chuang J. and Chen J., "The Building of Credit Scoring System on the Residential Mortgage Finance," *the International Journal of Forecasting*, vol. 15, no. 2, pp. 65-90, 2006.
- [7] Crook N., Edelman B., and Thomas C., "Recent Developments in Consumer Credit Risk Assessment," *the European Journal of Operational Research*, vol. 183, no. 3, pp. 1447-1465, 2007.
- [8] Dash M. and Liu H., "Feature Selection for Classification," *Intelligent data analysis*, vol. 1, no. 4, pp. 131-156, 1997.
- [9] Desai V., Conway D., Crook J., and Rj G., "Credit-Scoring Models in the Credit-Union Environment using Neural Networks and Genetic Algorithms," *the IMA Journal of Management Mathematics*, vol. 8, no. 4, pp. 323-346, 1997.

- [10] Finlay S., "Multiple Classifier Architectures and their Application to Credit Risk Assessment," *the European Journal of Operational Research*, vol. 210, no. 2, pp. 368-378, 2011.
- [11] Harrell E. and Lee L., *A Comparison of the Discrimination of Discriminant Analysis and Logistic Regression under Multivariate Normality*, Elsevier science publishers, New York, USA, 1985.
- [12] Hoffmann F., Baesens B., Mues C., Van Gestel T., and Vanthienen J., "Inferring Descriptive and Approximate Fuzzy Rules for Credit Scoring using Evolutionary Algorithms," *the European Journal of Operational Research*, vol. 177, no. 1, pp. 540-555, 2007.
- [13] Hsieh C., "Hybrid Mining Approach in the Design of Credit Scoring Models," *Expert Systems with Applications*, vol. 28, no. 4, pp. 655-665, 2005.
- [14] Hu C., Chen S., and Tzeng H., "Finding Fuzzy Classification Rules using Data Mining Techniques," *Pattern Recognition Letters*, vol. 24, no. 1-3, pp. 509-519, 2003.
- [15] Huang L., Chen C., and Wang J., "Credit Scoring with a Data Mining Approach Based on Support Vector Machines," *Expert Systems with Applications*, vol. 33, no. 4, pp. 847-856, 2007.
- [16] Huang Z., Chen H., Hsu J., Chen H. and Wu S., "Credit Rating Analysis With Support Vector Machines and Neural Networks: A Market Comparative Study," *Decision Support Systems*, vol. 37, no. 4, pp. 543-558, 2004.
- [17] Lahsasna A., Aïnon N., and Wah Y., "Credit Scoring Models using Soft Computing Methods: A Survey," *the International Arab Journal of Information Technology*, vol. 7, no. 2, pp. 115-123, 2010.
- [18] Lee C., "Application of Support Vector Machines to Corporate Credit Rating Prediction," *Expert Systems with Applications*, vol. 33, no. 1, pp. 67-74, 2007.
- [19] Lee S. and Chen F., "A Two-Stage Hybrid Credit Scoring Model using Artificial Neural Networks and Multivariate Adaptive Regression Splines," *Expert Systems with Applications*, vol. 28, no. 4, pp. 743-752, 2005.
- [20] Lee S., Chiu C., Lu J., and Chen F., "Credit Scoring using The Hybrid Neural Discriminant Technique," *Expert Systems with Applications*, vol. 23, no. 3, pp. 245-254, 2002.
- [21] Loh Y., "Classification and Regression Trees," *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, vol. 1, no. 1, pp. 14-23, 2011.
- [22] Louzada F., Anacleto O., Candolo C., and Mazucheli J., "Poly-Bagging Predictors for Classification Modelling for Credit Scoring," *Expert Systems with Applications: An International Journal*, vol. 38, no. 10, pp. 12717-12720, 2011.
- [23] Malhotra R. and Malhotra K., "Differentiating Between Good Credits and Bad Credits using Neuro-Fuzzy Systems," *the European Journal of Operational Research*, vol. 136, no. 1, pp. 190-211, 2002.
- [24] Martens D., Baesens B., Van T., and Vanthienen J., "Comprehensible Credit Scoring Models using Rule Extraction from Support Vector Machines," *the European Journal of Operational Research*, vol. 183, no. 3, pp. 1466-1476, 2007.
- [25] Nanni L. and Lumini A., "An Experimental Comparison of Ensemble of Classifiers for Bankruptcy Prediction and Credit Scoring," *Expert Systems with Applications*, vol. 36, no. 2, pp. 3028-3033, 2009.
- [26] Ong S., Huang J., and Tzeng H., "Building Credit Scoring Models using Genetic Programming," *Expert Systems with Applications*, vol. 29, no. 1, pp. 41-47, 2005.
- [27] Ping Y. and Yongheng L., "Neighborhood Rough Set and SVM Based Hybrid Credit Scoring Classifier," *Expert Systems with Applications*, vol. 38, no. 9, pp. 11300-11304, 2011.
- [28] Michie D., *Expert Systems in the Micro-electroni*, Edinburgh University Press, UK, 1979.
- [29] Rodgers L. and Nicewander A., "Thirteen Ways to Look at the Correlation Coefficient," *American Statistician*, vol. 42, no.1, pp. 59-66, 1988.
- [30] Šušteršič M., Mramor D., and Zupan J., "Consumer Credit Scoring Models with Limited Data," *Expert Systems with Applications*, vol. 36, no. 3, pp. 4736-4744, 2009.
- [31] Thomas C., *Consumer Credit Models: Pricing, Profit and Portfolios*, Oxford University Press, USA, 2009.
- [32] Tsai F. and Chen L., "Credit Rating by Hybrid Machine Learning Techniques," *Applied Soft Computing*, vol. 10, no. 2, pp. 374-380, 2010.
- [33] Tsai F. and Wu W., "Using Neural Network Ensembles for Bankruptcy Prediction and Credit Scoring," *Expert Systems with Applications*, vol. 34, no. 4, pp. 2639-2649, 2008.
- [34] Van T. and Baesens B., *Credit Risk Management: Basic Concepts: Financial Risk Components, Rating Analysis, Models, Economic and Regulatory Capital*, Oxford University Press, USA, 2009.
- [35] Wang G., Hao J., Ma J., and Jiang H., "A Comparative Assessment of Ensemble Learning for Credit Scoring," *Expert Systems with Applications*, vol. 38, no. 1, pp. 223-230, 2010.
- [36] West D., "Neural Network Credit Scoring Models," *Computers and Operations Research*, vol. 27, no. 11-12, pp. 1131-1152, 2000.
- [37] West D., Dellana S., and Qian J., "Neural Network Ensemble Strategies for Financial Decision Applications," *Computers and Operations Research*, vol. 32, no. 10, pp. 2543-2559, 2005.



**Seyed Sadatrasoul** is a PhD student in industrial engineering and systems management at Iran University of Science and Technology Tehran. He received his Bs degree in industrial engineering from IUST, in 2006 and obtained MS degree in information technology management from

Tarbiat modares university, Tehran, in 2009. Presently he is the assistant of faculty member of IT Group in School of Industrial Engineering and is actively engaged in conducting academic, research and development programs in the field of data and process mining. He has contributed more than 20 research papers to many national and international journals and conferences. He has also published two books by reputed publishers. His research interests include data mining and its applications with operation research, e-commerce and credit allocation in financial institutes.



**Mohammad Gholamian** is an assistant professor in School of Industrial Engineering at the Iran University of Science and Technology, Tehran. He received his MS degree in industrial engineering from Isfahan University of Technology, Isfahan in 1998 and

obtained PhD degree in Industrial Engineering from Amirkabir University of Technology, Tehran in 2005 for the work in the field of Hybrid Intelligent Decision Making Systems. Presently he is a faculty member of IT Group in School of Industrial Engineering and is actively engaged in conducting academic, research and development programs in the field of Industrial Engineering and Information Technology. He has contributed more than 120 research papers to many national and international journals and conferences. Besides this, he has published four books by reputed publishers. His research interests include data mining, soft computing, decision theory and e-business models.



**Kamran Shahanaghi** is an assistant professor in School of Industrial Engineering at the Iran University of Science and Technology, Tehran. He received his MS degree in Industrial Engineering from IUST in 1986 and obtained PhD degree in 2000.

Presently, he is a faculty member of optimization Group in School of Industrial Engineering and is actively engaged in conducting academic, research and development programs in the field of Industrial Engineering and optimization. He has contributed more than 140 research papers to many national and international journals and conferences. His research interests include operation research and uncertainty.

## Appendix 1

Features included in Iran’s credit dataset, and their types are sorted alphabetically and shown in Table 6.

Table 6. List of features in Iran commercial bank credit dataset.

Feature	Type
Accounts Receivable	Continuous
Accumulated Gains or Losses	Continuous
Active in Internal Market	Categorical
Audit Report	Categorical
Average Exports Over the Past Three Years	Continuous
Capital	Continuous
Company Background (Number of Years)	Continuous
Current Account Weighted Average	Continuous
Current Accounts Creditor Turn Over	Continuous
Current Assets	Continuous
Current Liabilities	Continuous
Current Period Assets	Continuous
Current Period Sales	Continuous
Current Period Shareholder Equity	Continuous
Experience With Bank(Number Of Years In 5 Categories)	Continuous
Export Price Index	Continuous
Financial Costs	Continuous
Gross Profit	Continuous
Inflation Rate	Continuous
Inventory Cash	Continuous
Last Three Years Average Imports	Continuous
Long-Term Financial Liabilities	Continuous
Mangers History	Continuous
Net Profit	Continuous
Non-Current Assets	Continuous
Non-Current Liabilities	Continuous
Number of Countries That The Company Export to	Continuous
Other Accounts Receivable	Continuous
Prior Period Assets	Continuous
Prior Period Sales	Continuous
Prior Period Shareholder Equity	Continuous
Sale	Continuous
Seasonal Factors	Categorical
Shareholder Equity	Continuous
Short-Term Financial Liabilities	Continuous
Stock	Continuous
Target Market Risk (From 1 To 5)	Continuous
Tehran Stock Exchange Index	Continuous
Three Prior Year Foreign Exchange Rate	Continuous
Top Mangers History	Categorical
Total Assets	Continuous
Total Liabilities	Continuous
Two-Prior Period Assets	Continuous
Two-Prior Period Sales	Continuous
Two-Prior Period Shareholder Equity	Continuous
Type of Book: Accredited Auditor (=1,Other=0)	Categorical
Type of Book: Audit Organization (=1,Other=0)	Categorical
Type of Book: Tax Declaration(=1,Other=0)	Categorical
Type of Company: Cooperative (=1, Other =0)	Categorical
Type of Company: Limited And Others (=1, Other =0)	Categorical
Type of Company: PJS (=1, Other =0)	Categorical
Type of Company: Stock Exchange (=1, Other =0)	Categorical
Type of Company: Stock Exchange(LLP) (=1, Other =0)	Categorical
Type of Industry: Agricultural (=1, Other =0)	Categorical
Type of Industry: Chemical (=1, Other =0)	Categorical
Type of Industry: Industry And Mine (=1, Other =0)	Categorical
Type of Industry: Infrastructure and Service (=1, Other =0)	Categorical
Type of Industry: Oil and Petrochemical (=1, Other =0)	Categorical
Year of Financial Ratio	Categorical
Basel: Creditworthy (=1, Other =0)	Categorical