

A GA-based Rule Extraction for Bankruptcy Prediction Modeling

Kyung-shik Shin
College of Business Administration
Ewha Womans University
(ksshin@ewha.ac.kr)

.....

Prediction of corporate failure using past financial data is a well-documented topic. Early studies of bankruptcy prediction used statistical techniques such as multiple discriminant analysis, logit and probit. Recently, however, numerous studies have demonstrated that artificial intelligence such as neural networks (NNs) can be an alternative methodology for classification problems to which traditional statistical methods have long been applied. Although numerous theoretical and experimental studies reported the usefulness of neural networks in classification studies, there exists a major drawback in building and using the model. That is, the user can not readily comprehend the final rules that the neural network models acquire.

We propose a genetic algorithms (GAs) approach in this study and illustrate how GAs can be applied to corporate failure prediction modeling. An advantage of GAs approach offers is that it is capable of extracting rules that are easy to understand for users like expert systems. The preliminary results show that rule extraction approach using GAs for bankruptcy prediction modeling is promising.

Key Words: Genetic Algorithms, Rule extraction, Bankruptcy prediction

.....

1. Introduction

Today, Korean financial institutions are paying a heavy price for their indiscriminate practices. Corporate bankruptcies have put several institutions on the brink of insolvency. Many others have been merged with or acquired by other financial institutions. Surviving institutions are rushing to put in place a corporate credit rating system, but are facing difficulties due to lack of data accumulation and scientific credit rating methods.

The present research pertains to a corporate failure prediction modeling which can provide a basis for credit rating system. Prediction of corporate failure using past financial data is a well-documented topic. Early studies of bankruptcy prediction used statistical techniques such as multiple discriminant analysis (Altman, 1968, 1983), logit (Ohlson, 1980) and probit (Zmijewski, 1984). Recently, however, numerous studies have demonstrated that artificial intelligence such as neural networks (NNs) can be an alternative methodology for classification

problems to which traditional statistical method have long been applied (Barniv *et al.*, 1997; Bell, 1997; Boritz and Kennedy, 1995; Chung and Tam, 1992; Etheridge and Sriram, 1997; Fletcher and Goss, 1993; Jo *et al.*, 1997; Odum and Sharda, 1990; Salchenberger *et al.*, 1992; Shin and Han, 1998a; Shin *et al.*, 1998; Tam and Kiang, 1992; Wilson and Sharda, 1994).

Although numerous theoretical and experimental studies reported the usefulness of NNs in classification studies, there are several drawbacks in building and using the model. First, it is an art to find an appropriate NN model which can reflect problem characteristics because there are numerous network architectures, learning methods, and parameters. Second, the user can not readily comprehend the final rules that the neural network models acquire. This characteristic of NNs is often referred to 'Black boxes'.

We propose a genetic algorithms (GAs) approach in this study and illustrate how GAs can be applied to corporate failure prediction modeling. An advantage of this approach offers is that it is capable of extracting rules that are easy to understand for users like expert systems.

The remainder of this paper is organized as follows: The second section provides a brief description of GAs. The third section describes the rule extraction approach using genetic search. The fourth section reports the model development and the results of the experiments. The final section discusses the conclusions and future research issues.

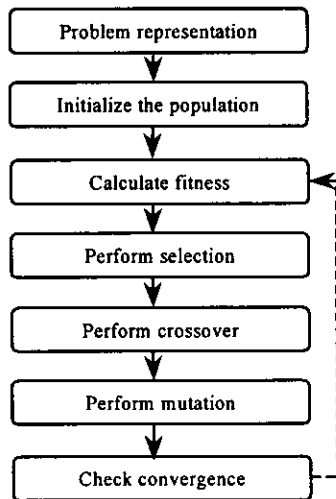
2. Genetic Algorithm Technique

GAs are stochastic search techniques that can search large and complicated spaces on the ideas from natural genetics and evolutionary principle (Davis, 1991; Holland, 1975; Goldberg, 1989). They have been demonstrated to be effective and robust in searching very large spaces in a wide range of applications (Colin, 1994; Han, Jo & Shin, 1997; Koza, 1993; Shin and Han, 1998b). GAs are particularly suitable for multi-parameter optimization problems with an objective function subject to numerous hard and soft constraints. The financial application of GAs is growing with successful applications in trading system (Colin, 1994; Deboeck, 1994), stock selection (Mahfoud and Mani, 1995), portfolio selection (Rutan, 1993), bankruptcy prediction (Kingdom and Feldman, 1995), credit evaluation (Shin and Han, 1998b; Walker *et al.*, 1995) and budget allocation (Packard, 1990).

GAs are distinct from many conventional search algorithms in the following ways (Karr, 1995):

- 1) GAs consider many points in the search space simultaneously, not a single point, reducing the chance of converging to local optima;
- 2) GAs work directly with strings of characters representing the parameter set, not the parameters themselves;
- 3) GAs use probabilistic rules to guide their search, not deterministic rules.

GAs perform the search process in four stage: initialization, selection, crossover, and mutation (Davis, 1991; Wong & Tan, 1994). Figure 1 shows the basic steps of GAs.



<Figure 1> Basic steps of genetic algorithms

In the initialization stage, a population of genetic structures, called chromosomes, that are randomly distributed in the solution space, is selected as the starting point of the search. After the initialization stage, each chromosome is evaluated using a user-defined fitness function. The role of the fitness function is to numerically encode the performance of the chromosome. For real-world applications of optimization methods such as GAs, the choice of the fitness function is the most critical step.

The mating convention for reproduction is such that only the high scoring members will preserve and propagate their worthy characteristics from generations to generation and thereby help in continuing the search for an optimal solution. The chromosomes with high performance may be chosen for replication several times whereas poor-performing structures may not be chosen at all. Such a selective process causes the best-performing chromosomes in the population to occupy an increasingly

larger proportion of the population over time.

The crossover forms a new offspring between two randomly selected 'good parents'. The crossover operates by swapping corresponding segments of a string representation of the parents and extends the search for new solution in far-reaching direction. The crossover occurs only with some probability, the crossover rate. There are many different types of crossover that can be performed: the one-point, the two-point, and the uniform type (Syswerda, 1989).

The mutation is a GA mechanism where we randomly choose a member of the population and change one randomly chosen bit in its bit string representation. Although the reproduction and the crossover produce many new strings, they do not introduce any new information into the population at the bit level. If the mutant member is feasible, it replaces the member which was mutated in the population. The presence of mutation ensures that the probability of reaching any point in the search space is never zero.

3. Rule Extraction Using GAs

In building corporate failure predicting model, we use the similar approach that Bauer (1994), and Mahfoud and Mani (1995) suggest in their stock selection applications. We apply GAs to find thresholds (cutoffs) for one or more variables, above or below which a company is considered 'dangerous'. For instance, if the model's structure consists of two variables representing a particular company's quick ratio and a debt ratio, the final rule the GA returns might look like the following:

IF [Debt ratio > 1.50 and Quick ratio < 0.35]
THEN Dangerous

In many cases, the simplistic rule like above the example is insufficient to model relationships among financial variables. Our rule structure contains five conditions using 'AND' relations for this study. The general form of the rule that GAs generate is as follows:

IF [the VAR1 is GREATER THAN OR EQUAL TO (LESS THAN) C1,
AND the VAR2 is GREATER THAN OR EQUAL TO (LESS THAN) C2,
AND,
AND the VAR5 is GREATER THAN OR EQUAL TO (LESS THAN) C5]
THEN Prediction is Dangerous.

If the all of the five conditions are satisfied, then the model will produce 'dangerous' signal on an evaluated company. C1 to C5 denotes the cutoff values which are found through genetic search process. The cutoff values range from 0 to 1, and represent the percentage of the data source's range. This allows the rules to refer to any data source, regardless of the

values it takes on. Above rule structure is summarized in Table 1. In the table, 'which data' means data source the rule refers to. The general rule structure is illustrated in Table 1.

We allow the model to select 5 variables among 9 alternative financial ratios. We also allow to choosing one variable more than once in the rule structure because the bankruptcy prediction is often highly nonlinear. For example, a 10% increase in sales may result in a good signal, while a 100% increase in that same variable could result in a bad signal. This means that use of multiple cutoff points in extracting knowledge from financial variables is recommended, if necessary. In addition, if we consider the interactions among conditions, this is essential for increasing flexibility of financial modeling.

In setting up the genetic optimization problem, we need the parameters that have to be coded for the problem and an objective or fitness function to evaluate the performance of each string. The parameters that are coded are the cell values of Table 1. As we mentioned above, they are input variables, above or below, and the cutoff values. The varying parameters generate a number of combinations of our general rules.

<Table 1> The general rule structure

Number(j)	Cond1	Cond2	Cond3	Cond4	Cond5	Description
Which data	VAR _{1j}	VAR _{2j}	VAR _{3j}	VAR _{4j}	VAR _{5j}	VAR _{ij} (i=var. number, j = condition number)
Less than / greater than or equal to	L/G _{1k}	L/G _{2k}	L/G _{3k}	L/G _{4k}	L/G _{5k}	L/G _{jk} (k= 1: less than /2: greater than or equal to)
Cutoff values	C ₁	C ₂	C ₃	C ₄	C ₅	Cutoff C _j (j= condition number)

The string encoded for the experiments is as follows:

String {VAR_{1i}, VAR_{2i}, VAR_{3i}, VAR_{4i}, VAR_{5i},
L/G_{1k}, L/G_{2k}, L/G_{3k}, L/G_{4k}, L/G_{5k}, C₁,
C₂, C₃, C₄, C₅ }

The GAs maintain a population of strings which are chosen at random. This initialization allows the GAs to explore the range of all possible solutions, and this tends to favor the most likely solutions. Generally, the population size is determined according to the size of the problem (bigger population for larger problem). The common view is that a larger population takes longer to settle on a solution, but is more likely to find a global optimum because of its more diverse gene pool. We use 100 strings in the population.

The task of defining a fitness function is always application specific. In this study, the objective of the system is to find a rule which would yield the highest hit ratio if rules are fired across the company. We apply the hit ratio of the rule to the fitness function for this study.

The genetic operators such as crossover and mutation which are described in the previous section are used to search for the optimal solutions. Several parameters must be defined for the above operators, and the values of these parameters can greatly influence the performance of the algorithm. The crossover rate ranges 0.5 - 0.7 and the mutation rate ranges 0.06 - 0.12 for our experiment. As a stopping condition, we use 3,000 trials. These processes are done by the genetic algorithms software package EvolverTM 4.0, called from an Excel macro.

4. Experiments and Results

4.1 Data and Variables

The data set contains 528 externally audited mid-sized manufacturing firms which filed for bankruptcy (264 cases) and non-bankruptcy (264 cases) during the period 1995-1997. We apply two stages of input variable selection process. In the first stage, we select 55 variables by factor analysis, independent-samples t-test (between input variable and output variable) and Mann-Whitney U test (for qualitative variables). In the second stage, we select 9 financial variables using the stepwise methods to reduce the dimensionality. The aim of input variable selection approach is to select the input variables satisfying the univariate test first, and then select significant variables by stepwise method for refinement. As we mentioned above, these variables are not the final ones that are used to form a rule, but are provided as the alternative variables for the final selection. Table 2 illustrates the pre-selected variables for this study.

<Table 2> Selected variables

Variables	Name
X1	Value added to total asset
X2	Net income to stockholder's equity
X3	Quick ratio
X4	Liquidity ratio
X5	Current liability to total assets
X6	Retained earnings to total assets
X7	Stockholders' equity to total assets
X8	Financial expenses to sales
X9	Operating income to operating expenses

The data set is split into two subsets, a training set and a validation (holdout) set of 90 and 10 percent of the entire data, respectively. The training data are used for learning rules, and the validation data which have not been used to develop the systems are used to test the results.

4.2 Results

We extract five bankruptcy rules by genetic search process. The rules generated and the corresponding descriptions are illustrated in Table 3 and Table 4.

<Table 3> The rules generated

Rule number		Cond 1	Cond 2	Cond 3	Cond 4	Cond 5
Rule 1	Variable code > / < code Cutoffs	2 1 0.426	4 1 0.847	5 1 0.520	7 1 0.595	8 1 0.665
Rule 2	Variable code > / < code Cutoffs	2 1 0.520	2 1 0.595	3 1 0.697	7 1 0.590	8 1 0.503
Rule 3	Variable code > / < code Cutoffs	2 1 0.426	4 1 0.560	6 2 0.082	7 1 0.590	8 1 0.520
Rule 4	Variable code > / < code Cutoffs	2 1 0.560	3 1 0.697	6 2 0.130	7 1 0.577	8 1 0.515
Rule 5	Variable code > / < code Cutoffs	2 1 0.560	3 1 0.697	6 2 0.082	7 1 0.590	8 1 0.520

<Table 4> The description of rules

Rule number	Description
Rule 1	IF Net income to stockholder's equity is less than 0.426* AND Liquidity ratio is less than 0.847 AND Current liability to total assets is less than 0.520 AND Stockholders' equity to total assets is less than 0.595 AND Financial expenses to sales is less than 0.665, THEN Dangerous.
Rule 2	IF Net income to stockholder's equity is less than 0.520 AND Quick ratio is less than 0.697 AND Stockholders' equity to total assets is less than 0.590 AND Financial expenses to sales is less than 0.503, THEN Dangerous.
Rule 3	IF Net income to stockholder's equity is less than 0.426 AND Liquidity ratio is less than 0.560 AND Retained earnings to total assets is greater than or equal to 0.082 AND Stockholders' equity to total assets is less than 0.590 AND Financial expenses to sales is less than 0.590, THEN Dangerous.
Rule 4	IF Net income to stockholder's equity is less than 0.560 AND Quick ratio is less than 0.697 AND Retained earnings to total assets greater than or equal to 0.130 AND Stockholders' equity to total assets is less than 0.577 AND Financial expenses to sales is less than 0.515, THEN Dangerous.
Rule 5	IF Net income to stockholder's equity is less than 0.560 AND Quick ratio is less than 0.697 AND Retained earnings to total assets is greater than or equal to 0.082 AND Stockholders' equity to total assets is less than 0.590 AND Financial expenses to sales is less than 0.520, THEN Dangerous.

* Represent the percentage of the data source's range.

The goal in optimization is ideally to find the best solution to a problem. Since GAs try to find out the optimal or near optimal combination of above searching parameters, the final solution is one. However, in most real-world problems, one does not usually know the best possible solution. Therefore, a more realistic objective is to find alternatively good solutions. We generate multiple rules by choosing multiple strings in the converged population. Since the fitness function of GAs measures the quality of a particular solution, we select the strings with high level of fitness values. So the derived rules are alternatively good rules which show high level of hit ratio although there are minor differences in simulated performance.

The hit ratios calculated from simulation results are summarized in Table 5. In table 5, hit ratio(A) denotes the rate of correct classification if the rule is fired, while hit ratio(B) represents overall classification accuracy of the set.

<Table 5> The performance of derived rules (%)

Rules	Train (476 cases)		Validation (52 cases)		
	Hit ratio(A)	Hit ratio(B)	Hit ratio(A)	Hit ratio(B)	# of cases fired
Rule 1	79.0	78.8	84.6	80.0	30
Rule 2	80.7	80.0	76.9	75.0	28
Rule 3	82.6	80.0	84.6	82.1	28
Rule 4	81.6	79.6	78.9	77.8	27
Rule 5	80.2	80.0	78.9	77.8	27
Average	80.8	79.7	80.8	78.5	28

The average hit ratio if the rules are fired is 80.8% of training and validation sets, respectively. This means if the financial variables

of a company are within the feature ranges of derived rules, the probability of bankruptcy is about 80% of cases.

The preliminary results above demonstrate that GAs are effective methods for extracting rules for the bankruptcy prediction. Their success is due to their ability to learn nonlinear relationships among the input variables. A drawback of this approach is that the model produces predictions only when the rules are fired, while NNs make predictions on every case except when explicitly restricted. The average number of cases that are fired by a specific rule is 28 among 52 cases (53.8%). This problem, however, can be reduced by integrating multiple rules derived. We have many ways to integrate these rules. For example, if one of the five rules makes 'Danger' signal, the model may produce 'Danger' signal to the users.

5. Concluding Remarks

We applied GAs to extract rules that can predict corporate failure. This paper is just a first attempt to explore the potential of genetic-based systems to handle bankruptcy prediction problems systematically. The results show that rule extraction approach using GAs for bankruptcy prediction modeling is promising.

This paper, however, has several limitations. First, although we derived multiple rules using traditional GAs, it is necessary to extend the GAs through use of a niching method (Mahfoud and Mani, 1995). Unlike the traditional GAs, which makes the population eventually converge around a single point in the solution

space, the GA that uses a niching method converges about multiple solutions or niches.

Second, the current rule structure is quite limited. As a next research step, this structure will be considerably extended by incorporating additional features. It is likely that more informative features will possibly lead to improved results, although we should consider the efficiency problem. Further improvements may be obtained by incorporating qualitative factors and quantitative ones. We plan to include qualitative variables in extracting the prediction rules.

References

- Altman, E., "Financial ratios, discriminant analysis and the prediction of corporate bankruptcy," *The Journal of Finance*, vol.23 (1968), 589-609.
- _____, *Corporate financial distress - A complete guide to predicting, avoiding and dealing with bankruptcy*, John Wiley, New York, 1983.
- Barniv, R., Agarwal, A. and Leach, R., "Predicting the outcome following bankruptcy filing: A three-state classification using neural networks," *Intelligent Systems in Accounting, Finance and Management*, vol.6 (1997), 177-194.
- Bauer, R. J., *Genetic Algorithms and Investment Strategies*, John Wiley & Sons, 1994.
- Bell, T., "Neural nets or the logit model? A comparison of each model's ability to predict commercial bank failures," *Intelligent Systems in Accounting, Finance and Management*, vol.6 (1997), 249-264.
- Boritz, J. and Kennedy, D., "Effectiveness of neural networks types for prediction of business failure," *Expert Systems with Applications*, vol.9 (1995), 503-512.
- Chung, H. and Tam, K., "A Comparative Analysis of Inductive Learning Algorithm". *Intelligent Systems in Accounting, Finance and Management*, vol.2 (1992), 3-18.
- Colin, A. M., "Genetic algorithms for financial modeling," In Deboeck, G.J. (Eds.), *Trading On The Edge*, John Wiley, New York, 1994, 148-173.
- Davis, L., *Handbook of Genetic Algorithms*, Van Nostrand Reinhold, New York, 1991.
- Deboeck, G. J., "Using GAs to optimize a trading system," In Deboeck, G.J (Eds.), *Trading On The Edge*, John Wiley, New York, 1994, 174-188.
- Etheridge, H. and Sriram, R., "A comparison of the relative costs of financial distress models: Artificial neural networks, logit and multivariate discriminant analysis," *Intelligent Systems in Accounting, Finance and Management*, vol.6 (1997), 235-248.
- Fletcher, D. and Goss, E., "Forecasting with neural networks: An application using bankruptcy data," *Information and Management*, vol.24, no.3 (1993), 159-167.
- Goldberg, D. E., *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison-Wesley, 1989.
- Han, I., Jo, H. and Shin, K. S., "The hybrid systems for credit rating," *Journal of the Korean Operations Research and Management Science Society*, vol.22, no.3 (1997), 163-173.
- Holland, J. H., *Adaptation in Natural and Artificial Systems*, Ann Arbor, The University of Michigan Press, 1975.

- Jo, H., Han, I. And Lee, H., "Bankruptcy prediction using case-based reasoning, neural networks, and discriminant analysis," *Expert Systems With Applications*, vol.13, no.2 (1997), 97- 108.
- Karr, C., "Adaptive control of an exothermic chemical reaction system using fuzzy logic and genetic algorithms," In Medsker, L. R. (Eds.), *Hybrid Intelligent Systems*, Kluwer Academic Publishers, 1995.
- Kingdom, J. and Feldman, K., *Genetic algorithms for bankruptcy prediction*, Search Space Research Report, No.01-95, Search Space Ltd. London, 1995.
- Koza, J., *Genetic programming*, The MIT Press, 1993.
- Mahfoud, S. and Mani, G., "Genetic algorithms for predicting individual stock performance," *Proceedings of the 3rd International Conference on Artificial Intelligence Applications on Wall Street*, 1995, 174-181.
- Odom, M. and Sharda, R., "A neural networks model for bankruptcy prediction," *Proceedings of the IEEE International Conference on Neural Network*, vol.2, 1990, 163-168.
- Ohlson, J., "Financial ratios and the probabilistic prediction of bankruptcy," *Journal of Accounting Research*, vol.18, no.1 (1980), 109-131.
- Packard, N., "A genetic learning algorithm for the analysis of complex data," *Complex Systems*, vol.4 (1990), 543-572.
- Rutan, E., "Experiments with optimal stock screens," *Proceedings of the 3rd International Conference on Artificial Intelligence Applications on Wall Street*, 1993, 269-273.
- Salchenberger, L., Cinar, E. and Lash, N., "Neural networks: A new tool for predicting thrift failures," *Decision Sciences*, vol.23 (1992), 899-916.
- Shin, K. S. and Han, I. (a), "Bankruptcy Prediction Modeling Using Multiple Neural Networks Models," *Proceedings of Korea Management Science Institute Conference*, 1998.
- _____ (b), "Using Genetic Algorithm to Support Case-Based Reasoning: Application to Corporate Bond Rating Integration," *Proceedings of Second Asia Pacific Decision Sciences Institute (DSI) Conference*, Taipei, 1998.
- Shin, K. S., Shin, T. S. and Han, I., "Corporate Credit Rating System Using Bankruptcy Probability Matrix", *Proceedings of IV International Meeting on Artificial Intelligence and Emerging Technologies in Accounting, Finance and Taxation*, Spain, 1989.
- Syswerda, G., "Uniform crossover in genetic algorithms," In Schaffer, J.D. (Eds.), *Proceedings of 3rd Int'l Conference of Genetic Algorithms*, San Maeto, Morgan Kaufmann, 1989.
- Tam, K. and Kiang, M., "Managerial applications of neural networks: the case of bank failure predictions," *Management Science*, vol.38, no.7 (1992), 926-947.
- Walker, R., Haasdijk, E. and Gerrets, M., "Credit evaluation using a genetic algorithm," In Coonatilake, S. and Treleaven, P. (Eds.), *Intelligent Systems for Finance and Business*, John Wiley, 1995, 39-59.
- Wilson, R. and Sharda, R., "Bankruptcy prediction using neural networks," *Decision Support Systems*, vol.11, no.5(1994), 545-557.

Wong, F. and Tan, C., "Hybrid neural, genetic and fuzzy systems," In Deboeck, G.J. (Eds.), *Trading On The Edge*, John Wiley, New York, 1994, 245-247.

Zmijewski, M. E., "Methodological Issues Related to the Estimated of Financial Distress Prediction Models". *Journal of Accounting Research*, vol.22, no.1 (1984), 59-82.

국문요약

유전자 알고리즘을 활용한 부실예측모형의 구축

신경식*

기업부실예측은 과거로부터 많은 연구가 이루어진 분야로, 주로 통계기법에 의한 분류예측문제로 다루어져 왔다. 최근에는 인공신경망, 의사결정나무 등 비선형성을 반영할 수 있는 인공지능 기법을 적용한 연구가 많이 수행되고 있다.

본 연구에서는 최적화에 주로 활용하는 인공지능 기법인 유전자 알고리즘을 규칙추출을 통한 기업부실예측 모형의 개발에 적용하고, 활용가능성을 검증하였다.

* 이화여자대학교 경영대학 교수