

WP 97-18  
October 1997



# Working Paper

Department of Agricultural, Resource, and Managerial Economics  
Cornell University, Ithaca, New York 14853-7801 USA

## **INTRODUCING RECURSIVE PARTITIONING TO AGRICULTURAL CREDIT SCORING**

by

**MICHAEL NOVAK**  
**MANAGER, AGRICULTURAL FINANCE**  
**FEDERAL AGRICULTURAL MORTGAGE CORPORATION**  
**WASHINGTON, DC**

AND

**EDDY L. LADUE**  
**PROFESSOR, AGRICULTURAL FINANCE**  
**CORNELL UNIVERSITY**  
**ITHACA, NY**

It is the policy of Cornell University actively to support equality of educational and employment opportunity. No person shall be denied admission to any educational program or activity or be denied employment on the basis of any legally prohibited discrimination involving, but not limited to, such factors as race, color, creed, religion, national or ethnic origin, sex, age or handicap. The University is committed to the maintenance of affirmative action programs which will assure the continuation of such equality of opportunity.

## INTRODUCING RECURSIVE PARTITIONING TO AGRICULTURAL CREDIT SCORING

The farm financial crisis in the mid 1980's brought increased interest in credit scoring models. Many agricultural lenders and financial advisors have adopted formal credit scoring models to monitor and forecast financial performance (LaDue et al.). Various nonparametric and parametric methods have been utilized to estimate the models, such as: experience-based algorithms (Alcott, Splett et al.); mathematical programming (Hardy and Adrian, and Ziari et al.); logistic regression (Mortensen et al.); probit regression (Lufburrow et al. and Miller et al.); discriminate analysis (Hardy and Weed, Dunn and Frey, and Johnson and Hagan); and linear probability regression (Turvey). There is not unanimous agreement as to the best method for estimating credit scoring models and new methods continue to be researched.

Most recently, the logistic regression has dominated the agricultural credit scoring literature (Miller and LaDue, Turvey and Brown, Novak and LaDue, Splett et al.). The logistic regression succeeded discriminant analysis as the parametric method of choice, primarily based on its more favorable statistical properties (McFadden). Turvey reviews and empirically compares agriculture credit scoring models using four parametric methods with a single data set. He recommends the logistic regression over the probit regression, discriminant analysis and the linear probability regression based on predictive accuracy and ease of use, in addition to the more favorable statistical properties. While the logistic regression improves on some of the statistical properties of the discriminant analysis and linear probability regression, it still possesses numerous statistical problems common to most parametric methods. These problems include: 1) the selection of the explanatory variables; 2 )

identifying the individual variable's relative importance; 3) reduction of the information space's dimensionally; and 4) explicitly incorporating misclassification costs<sup>1</sup>.

Non-agricultural studies have used Recursive Partitioning Algorithm (RPA) to classify financially stress firms. RPA is a computerized, nonparametric classification method that does not impose any a priori distribution assumptions. The essence of RPA is to develop a classification tree, that partitions the observations according to binary splits of characteristic variables. The selection and partitioning process occurs repeatedly until no further selection or division of a characteristic variable is possible, or the process is stopped by some predetermined criteria. Ultimately the observations in the terminal nodes of the classification tree are assigned to classification groups. RPA was originally developed by Friedman. A thorough theoretical exposition of RPA is presented in Breiman, et al. A more practical exposition of the computational aspects of RPA is presented in the CART software documentation (Steinberg and Colla). RPA has not been applied to agricultural creditworthiness classification.

The results of several non-agricultural financial stress classification studies indicate RPA outperforms the other parametric and judgmental models based on predictive accuracy. The problem with these studies is none of them use intertemporal (ex ante) model validation methods. Marais, Patell and Walfson compared RPA with a polytomous probit regression to classify commercial loans for publicly and privately held banking firms. Frydman, Altman and Kao compare RPA with discriminant analysis to classify firms according to their financial stress. In another study, Srinivasan and Kim compare RPA with discriminant analysis, logistic regression, goal programming, and a judgmental model, Analytic

---

<sup>1</sup> Each of these problems will be discussed in detail later in the text.

Hierarchy Process, to evaluate the corporate credit granting process. Each of these studies use within sample observations to estimate predictive ability, but a "true" estimate of RPA's prediction ability should also consider intertemporal (ex ante) predictions. The most practical use of financial stress classification and credit scoring is to focused towards future financial stress or creditworthiness, and not only current financial stress or creditworthiness.

The purpose of this study is to introduce RPA as a method for classifying creditworthy and less creditworthy agricultural borrowers. This study uses the agricultural firm's debt repayment capacity as the indicator of creditworthiness. Due to the relative newness of the method in application to agricultural credit scoring, RPA is presented in detail to provide a thorough background. Secondly, this study compares RPA to the logistic regression. The comparison considers relative misclassification costs, an aspect of borrower classification not fully explored in previous agricultural credit scoring studies. This study also challenges the RPA's superior prediction accuracy, as purported in the financial stress classification literature. The study uses within-sample cross-validation to select the models and intertemporal (ex ante) predictive accuracy, based on the minimization of misclassification costs, to compare the models. The financial stress classification literature only uses within-sample cross-validation to select and compare the models.

The remainder of the paper is divided into five sections. The first section presents the specifics of the RPA. This is followed by a section that discusses the advantages and disadvantages, as well as the differences between, the RPA and logistic regression. The third section describes the data. The fourth and fifth sections present the creditworthiness models and empirical results, respectively. The final section summarizes the paper's results.

## Recursive Partitioning Algorithm

In this section, a hypothetical RPA tree growing process is presented and the terminology is introduced. Following a brief introduction, a more detailed explanation of the optimal splitting rule, assignment of terminal nodes, and selection of appropriate tree complexity is presented.

To understand the tree growing process, a hypothetical tree is illustrated in Figure 1. It is constructed using classification groups  $i$  and  $j$ , and characteristic variables<sup>2</sup>  $A$  and  $B$ . Throughout the paper the classification groups are limited to two. But, in general, classification groups can be greater than two. To start the tree growing process all the observations in the original sample, denoted by  $N$ , are contained in the parent node which constitutes the first subtree, denoted  $T_0$  (not really a tree, but we will call it one anyway).  $T_0$  possess no binary splits and can be referred to as the naive classification tree. All observations in the original sample are assigned to either group,  $j$  or  $i$ , based on an assignment rule. The means to assign  $T_0$  to either group  $i$  or  $j$  depends on the misclassification costs and prior probabilities. When misclassification costs are equal and prior probabilities are equal to the sample proportions of the groups,  $T_0$  is assigned to the group with the greatest proportion of observations, minimizing the number of observations misclassified. When misclassification costs are not equal and prior probabilities are not equal to the sample proportions of the groups,  $T_0$  is assigned to the group that minimizes the observed expected cost of misclassification<sup>3</sup>. Later, it is demonstrated that

---

<sup>2</sup> Characteristic variables are analogous to independent variables in a parametric regression.

<sup>3</sup> Observed expected cost of misclassification is calculated using the misclassification rate of group assignments, sample probabilities, cost of type I errors (classifying a less creditworthy borrower as creditworthy) and cost of type II errors (classifying a creditworthy borrower as less creditworthy). An exact definition of observed expected cost of misclassification will be given later in the text.

minimizing the observed expected cost of misclassification is the same as minimizing the number of observations misclassified, when misclassification costs are equal to each other and prior probabilities are equal to the sample proportions of the groups.

To begin the tree growing process, RPA searches each individual characteristic variable and split value of the characteristic variable in a methodological manner. The computer algorithm then selects a characteristic variable, in this case A, and a split value of characteristic variable A, in this case  $a_1$ , based on the optimal splitting rule. The optimal splitting rule implies that no other characteristic variable and split value can increase the amount of correctly classified observations in the two resulting descendent nodes. In this particular illustration, observations with a value of characteristic variable A less than  $a_1$  will "fall" into the left node and the observations with a value of characteristic variable A greater than  $a_1$  will "fall" into the right node. The resulting subtree, denoted by  $T_1$ , consists of a parent node and, a left and right terminal nodes, denoted by  $t_L$  and  $t_R$ , respectively. The terminal nodes in each subtree that are then assigned to groups,  $i$  or  $j$ , based on the assignment rule of minimizing observed expected cost of misclassification.  $T_0$  and  $T_1$  are the beginning of a sequence of trees that ultimately concludes with  $T_{max}$ . However, in some cases,  $T_1$  may also be  $T_{max}$  depending on the predetermined penalty parameters specified. If  $T_1$  is not  $T_{max}$  then the recursive partitioning algorithm continues.

In this illustration,  $T_1$  is not  $T_{max}$ , so the partitioning process continues. Now, B is the characteristic variable selected and  $b_1$  is the "optimal" split value selected by the computer algorithm. The right node becomes an internal node and the observations within it are partitioned again. Observations with a value of characteristic variable B less than  $b_1$  "fall" into a new left node and observations with a value of characteristic variable B greater than  $b_1$  "fall" into a new right node. The new left and right nodes become terminal nodes in  $T_2$ , while the left node in  $T_1$  still remains a terminal

model in  $T_2$ . All three terminal nodes in  $T_2$  are then assigned to classification groups,  $i$  and  $j$ , based on the assignment rule of minimum observed expected cost of misclassification. Here again,  $T_2$  does not minimize the observed expected cost of misclassification of the original sample, therefore the partitioning process continues. Variable  $A$  is selected again to develop  $T_3$ . As stated above, when the recursive partitioning process is finished, the resulting classification tree is known as  $T_{\max}$ . In this illustration,  $T_3 = T_{\max}$ .  $T_{\max}$  is the tree that minimizes the expected observed cost of misclassification of the original sample. Obviously the development method will over fit the tree, therefore, a method is needed to prune back the tree, such as cross validation, discussed below. Once the classification tree is developed and pruned back, it can be used to classify observations from outside the original sample.

*Optimal Splitting Rule:* In essence, the optimal split rule<sup>4</sup> divides the observations in a parent node into two descendent nodes according to characteristic variable  $A$ 's cut-off value,  $a_1$ , maximizing the amount of correctly classified observations. More formally stated, the optimal splitting rule maximizes the decrease in the sum of the impurities of the two resulting descendent nodes compared with the impurity of the parent node. The sample impurity is represented by the Gini Impurity Index:

$$I(t) = i(t) p(t) = [c_{ij} p(i|t) p(j|t) + c_{ji} p(j|t) p(i|t)] p(t)$$

Where  $p(i|t) = (\pi_i n_i(t) / N_i) / (\sum_{k=i}^j \pi_k n_k(t) / N_k)$  is the conditional probability that an observation in node  $t$  belongs to group  $i$ ,  $p(j|t) = (\pi_j n_j(t) / N_j) / (\sum_{k=i}^j \pi_k n_k(t) / N_k)$  is the conditional probability that an observation in node  $t$  belongs to group  $j$ , and  $p(t) = \sum_{k=i}^j \pi_k n_k(t) / N_k$  is the probability of an

---

<sup>4</sup> The univariate splitting rule implies splitting an axis of one variable at one point. This study is limited to univariate splitting rules, however, CART has the capability to split variables using linear combinations of variables. The resulting classification trees are usually very cumbersome and difficult to interpret when linear combination splitting rules are used.



object falling into node  $t$ .  $N_i$  ( $N_j$ ) is the number of original observations from group  $i$  ( $j$ ).  $\pi_i$  ( $\pi_j$ ) is the prior probability of an observation belonging to group  $i$  ( $j$ ).  $c_{ij}$  ( $c_{ji}$ ) is the cost of misclassifying a group  $i$  ( $j$ ) observation as a group  $j$  ( $i$ ) observation. And,  $n_i(t)$  ( $n_j(t)$ ) is the number of group  $i$  ( $j$ ) observations in node  $t$ .

The Gini impurity index is best illustrated when the same special case presented before is considered. That is, when the misclassification costs are equal (i.e.  $c_{ij} = c_{ji} = 1$ ) and the prior probabilities are equal to the sample proportion of the sample groups (i.e.  $\pi_i = N_i/N$  and  $\pi_j = N_j/N$ ). In this case,  $p(i|t)$  and  $p(j|t)$  reduce to  $n_i(t)/N(t)$  and  $n_j(t)/N(t)$ , respectively. Where  $N(t)$  is the number of total observations in node  $t$ . Following through, the Gini Impurity Index for a nodes reduces to  $I(t) = i(t)p(t) = 2 [(n_i(t)/N(t)) (n_j(t)/N(t))]p(t)^5$ . The Gini Impurity Index equals its maximum value when all the observations in the node are equally divided between the groups (i.e.  $p(i|t) = n_i(t)/N(t) = 0.5$ ,  $p(j|t) = n_j(t)/N(t) = 0.5$ ,  $i(t) = 0.5$  and  $I(t) = 0.5p(t)$ ) and equals zero when all the observations in the node belong to the same classification group, (i.e.  $p(i|t) = n_i(t)/N(t) = 1$ ,  $p(j|t) = n_j(t)/N(t) = 0$  and  $I(t) = 0$ ). Put more intuitively, a node containing equal observations from the two classification groups has the most diversity and the Gini Impurity Index will obtain it's maximum value. The maximum value will depend on the number of groups and the probability of an observation falling into node  $t$ . Conversely, a node with observations representing only one group has no diversity, and the Gini Impurity Index measure will obtain it's lowest value, zero.

---

<sup>5</sup> Substituting  $p(j|t) = 1 - p(i|t)$  and  $p(i|t) = 1 - p(j|t)$  in equation 1, and again assuming  $c_{ij} = c_{ji} = 1$ ,  $\pi_i = N_i/N$  and  $\pi_j = N_j/N$ , the Gini Impurity Index can also be expressed as  $I(t) = [1 - p(i|t)^2 - p(j|t)^2]p(t)$ . Also, suppose group  $j$ 's observations have a value of 1 and group  $i$ 's observations have a value of 0, then the measure of impurity,  $i(t)$ , is also the sample variance,  $p(j|t)(1-p(j|t))$ .

The decrease in the impurity between a parent node and two descendent nodes can be denoted by  $\Delta I(t) = I(t) - I(t_L) - I(t_R) = i(t)p(t) - i(t_L)p(t_L) - i(t_R)p(t_R)$ , where  $t$ ,  $t_L$ , and  $t_R$  refers to the parent node, left descendent node, and right descendent node, respectively.  $\Delta I(t)$  is non-negative and its magnitude depends on the characteristic variable and split value selected. Selecting the optimal characteristic variable and corresponding split value for each node  $t$  maximizes  $\Delta I(t)$ , and is consistent with selecting the characteristic variable and corresponding split value that decreases the impurity of the overall classification tree. When the impurity of a tree can not be decreased any further the splitting process is terminated and  $T_{max}$  is obtained.

Assignment of terminal nodes: The terminal nodes for each subtree and  $T_{max}$  are assigned to a group that minimizes the observed expected cost of misclassification <sup>6</sup> for the specific node and correspondingly for the entire tree. The observed expected cost of misclassification of assigning node  $t$  to group  $i$  is defined as  $R_i(t) = c_{ij}p(j,t)$ , where:  $p(j,t) = \pi_j n_j(t)/N_j$  is the probability that an observation is from group  $j$  and falls into node  $t$ . Conversely,  $R_j(t) = c_{ij}p(i,t)$ . The terminal node assignment rule follows a Bayesian Rule,  $R(t) = \min[R_j(t), R_i(t)]$ . That is, the node is assigned to the group with minimum observed expected cost of misclassification. The observed expected cost of misclassification of the entire tree is defined as  $R(T) = c_{ij}\pi_j M_j(T)/N_j + c_{ij}\pi_i M_i(T)/N_i$ , where  $M_i(T)$  ( $M_j(T)$ ) is the number of group  $i$  ( $j$ ) observations misclassified in tree  $T$ <sup>7</sup>.

---

<sup>6</sup> The observed expected cost of misclassification is also known as "resubstitution risk". Where risk is another term for expected cost of misclassification, and resubstitution refers to the fact that risk is being evaluated with the original sample. The selection process classifies the original sample the best, but the resulting model does not necessarily represent the population structure.

<sup>7</sup> By substituting terms, the Gini Impurity Index can also be expressed as  $I(t) = R_j(t)p(j|t) + R_i(t)p(i|t)$ , since  $R_i(t) = c_{ij}p(j|t)p(t)$ ,  $R_j(t) = c_{ij}p(i|t)p(t)$ ,  $p(j|t) = p(j,t)/p(t)$  and  $p(i|t) = p(i,t)/p(t)$ . The

Again for exposition purposes, assume the costs of misclassification are equal (i.e.  $c_{ji}=c_{ij}=1$ ) and the prior probabilities are equal to the sample proportion of each group (i.e.  $\pi_i=N_i/N$  and  $\pi_j=N_j/N$ ). The observed expected cost of misclassification of assigning node  $t$  to group  $i$  reduces to  $R_i(t) = n_j(t)/N$ . Conversely,  $R_j(t) = n_i(t)/N$ . The observed expected cost of misclassification for group  $j$  (i) is equal to the sample proportion of group  $i$  (j) objects falling into node  $t$ . Therefore, the assignment rule becomes particularly simple. The terminal node is assigned to the group with the largest sample proportion of observations in the node. Minimizing the observed expected cost of misclassification in this specific case is the same as assigning the terminal node to the classification group with the majority representation in the node, and observed expected cost of misclassification corresponds with minimizing the overall observed expected cost of misclassification of the tree.

Selecting appropriate tree complexity:  $T_{\max}$  usually over fits the data, therefore, the observed expected cost of misclassification underestimates the "true" risk. In other words,  $T_{\max}$  classifies the original sample the best, but it is not necessarily representative of the population structure. Therefore, a method is needed to select the appropriate tree complexity from the nested sequence of trees. Some of the methods suggested are  $v$ -fold cross-validation, jackknifing, expert judgment, bootstrapping and holdout samples. In this study, cross-validation is used to select that appropriate correct tree complexity and minimize over-fitting of the data. That is, the tree with the smallest cross-validation cost of misclassification is selected and used to predict out-of-sample observations.

Cross-validation randomly divides all the observations in the original sample into  $V$  groups of approximately equal size. The observations in  $V-1$  groups are used to grow a tree corresponding to

---

relationship between the Gini Impurity Index and observed expected cost of misclassification can be seen in this new Gini Impurity Index expression.

the range of penalty parameters values for which the tree, based on the original sample, was optimal. The observations withheld are then passed through the newly constructed tree and classified. The procedure is repeated  $V$  times. Each time the group withheld is passed through the newly constructed tree and classified. The misclassification risk is summed and averaged for all the  $V$ -fold cross-validation trials to obtain an overall cross-validation expected cost of misclassification estimate. The appropriate tree is typically less complex than  $T_{max}$ .

While statistical resampling schemes, such as cross-validation, are sufficient to alleviate the statistical over-fitting bias, however, they do not account for intertemporal (ex ante) predictions, the basic objective of credit scoring models (Joy and Tofeson). Credit scoring models are not only used to classify borrowers, but to classify borrowers over time and predict future creditworthiness. Previous RPA financial stress classification studies have only evaluated the models based on their cross-validated expected cost of misclassification. In this study, the tree with the minimum cross-validation expected cost of misclassification will be used to predict creditworthy and less creditworthy borrowers in the forthcoming period. The predicted borrower classifications will then be compared to actual borrower classifications in same out-of-sample period. The misclassification risk of the out-of-sample predictions are calculated and the models are ultimately evaluated with regards to the out-of-sample misclassification risk.

### RPA and Logistic Regression Comparison

In this section the advantages and limitations, as well as the differences of the RPA and logistic regression are discussed. The logistic regression is well documented and well received as a qualitative choice model predicated on conventional parametric techniques, while RPA is a new parametric

method; a flexible modeling tool that exploits the power of computers. RPA substitutes a computer algorithm for the traditional method of mathematical analysis to get a numerical answer.

One basic difference between the two methods is the selection of variables. For credit scoring models there is no well-developed theory to guide the selection of financial and economic variables. Most proceed heuristically, selecting variables suggested by expert opinions or based on previous credit scoring models. Credit scoring models developed using RPA do not require the variables to be selected in advance. The computer algorithm selects the variables from a predetermined group of variables. This feature is especially advantageous if there are a large number of variables. In this context, RPA is somewhat analogous to forward stepwise regression, except RPA is not limited by the mathematical tractability or assumptions of conventional statistics, like forward stepwise regression.

In addition to selecting a group of variables, RPA analyzes univariate attributes of individual variables, providing insight and understanding to their predictive structure. The algorithm selects the variable that best classifies the observations and the optimal split value of the variable. When selecting a variable RPA places no limit on the number of times a variable can be utilized. Often the same variable can appear in different parts of the tree. Furthermore, RPA is not significantly influenced by outliers, since splits occur on non-outlier values. Once the split value is selected, the outlier is assigned to a node and the RPA procedure continues.

In contrast, the logistic regression model provides a linear approximation of the individual variables. Little is learned about the univariate attributes of each variable and each variable only appears once in the model. The model can be severely affected by outliers values. Logistic regression does have the advantage that the significance of the group of selected variables can be evaluated from the regression's summary statistics. A limitation of RPA's variable selection method is that once a

variable is selected all the succeeding variables are predicated on the original selected variable, again similar to forward stepwise regression. The tree growing process is intentionally myopic. RPA never looks ahead to see where it is going nor does it try to assess the overall performance of the tree during the splitting process.

In addition to selecting variables and their optimal split value, the CART software also provides competitive and surrogate variables and cut-off values listed in order of importance, for each node in the classification tree. Some variables may not appear in the final classification tree, but still can rank high as a competitive and surrogate variable. This list of competitive and surrogate variables give additional insight to the variable's usefulness. Competitive variables are alternative variables with slightly less ability to reduce impurity in the descendent nodes. Surrogate variables are variables that mimic the selected variables and split values, not only on size and composition of the descendent nodes, but also with respect to which observation lands in the left and right descendent node.

While lacking in univariate attributes, the logistic regression's major advantage is the predicted probabilities of creditworthiness assigned to each borrower. RPA can only classify observations into creditworthy or less creditworthy classes, and can not estimate an overall credit score. The predicted probabilities of creditworthiness provide additional quantitative information regarding a borrower. Furthermore, the predicted probabilities can also be converted to binary creditworthy/less creditworthy classification scheme when a prior probability is specified. Often lenders want to assess the predicted probability of creditworthiness, not only classify borrowers as creditworthy/less creditworthy.

Another difference between the two methods is the way they divide the information space into classification regions. RPA repetitiously partitions the information space as the binary tree is formed. A graphical illustration is presented in Figure 2, it is based on the hypothetical RPA tree in Figure 1.

RPA partitions the information space into four rectangular regions according to characteristic variables, A and B, and their respective optimal split values,  $a_1$  and  $b_1$ . Observations falling in regions 1 and 2 are classified as group i and those falling in region 3 and 4 are classified as group j. In contrast, the logistic regression if implemented as a binary qualitative choice model, partitions the information space into two regions. The logistic regression usually partitions the observations with respect to a prior probability, say c. The example line  $f(Z_m) = c$  divides the information space into two regions.  $Z_m$  is a linear function of variables A and B corresponding to observation m, and  $f(x)$  is the cumulative logistic probability function. The observations are assigned to class i if  $f(Z_m) \geq c$  or group j if  $f(Z_m) < c$ . The difference between the RPA's and logistic regression's information space is the shaded regions.

The two models also differ in the manner in which they incorporate misclassification cost and prior probabilities. RPA uses misclassification costs and prior probabilities to simultaneously determine variable selection, optimal split value and terminal node group assignments. Changes in the misclassification cost and prior probabilities can change the variables selected and the optimal split value, and, in turn, alter the structure of the classification tree. In contrast, the logistic regression is usually estimated without incorporating misclassification cost and prior probabilities. However, after the logistic regression is estimated a prior probability is used to classify borrowers as creditworthy/less creditworthy. Changes in the prior probability value can affect the predictive accuracy of the logistic regression (Mortensen et al.). Maddala (1983 and 1986) argues that prior probabilities should be taken to be the sample rate for the two groups, even with unequal sampling rates.

Despite the differences in the two models, RPA and logistic regression can be integrated to assess borrowers creditworthiness. RPA can select the relevant variables from a large set of variables to be estimated by the logistic regression. Likewise, the logistic regression can be used to estimate the

predicted probabilities and, in turn, the predicted probabilities can be used as a variable in RPA. How RPA utilizes or ranks the predicted probability, as a variable, can provide evidence for or against the logistic regression.

### Data

The data for this study were collected from New York State dairy farms in a program jointly sponsored by Cornell Cooperative Extension and the Department of Agricultural, Resource and Managerial Economics at the New York State College of Agriculture and Life Sciences, Cornell University. Seventy farms have been Dairy Farm Business Management (DFBS) cooperators from 1985 through 1993. Data for these seventy farms are analyzed in this study. Such a data set is critical in studying the dynamic effects of farm creditworthiness<sup>8</sup>. The farms represent a segment of New York State dairy farms which value consistent contribution of financial and management information. The financial information collected includes the essential components for deriving a complete set of

---

<sup>8</sup> Two types of estimation biases that typically plague credit evaluation models are choice bias and selection bias. Choice bias occurs when the researcher first observes the dependent variable and then draws the sample based on that knowledge. This process of sample selection typically causes an "oversampling" of financial distress firms. To overcome choice bias, this study selects the sample first and then calculates the dependent variable.

The other type of bias plaguing credit evaluation models is selection bias. Selection bias is a function of the nonrandomness of the data and can asymptotically bias the model's parameters and probabilities (Heckman). There are typically two ways selection bias can affect credit evaluation models. First, financially distressed borrowers are less likely to keep accurate records, therefore, these borrower are not included in the sample (Zmijewski et al.). And secondly, through the attrition rate of borrowers, because panel data are usually employed. In this study, there were probably borrowers that participated in the DFBS program during the earlier years of sample period, but exited the industry or stopped submitting records to the data base before the end of the sample period. In analyzing financial distress models, Zmijweski et al. found selection bias causes no significant changes in the overall classification and prediction rate. Given Zmijweski's results the study does not correct for selection bias and proceeds to estimate the credit evaluation models with the data presented.



sixteen financial ratios and measures as recommended by the FFSC. Table 1 exhibits all sixteen mean values of the financial ratios and measures for the seventy farms over the sample period<sup>9</sup>. Additional farm productivity, cost management and profitability statistics for these farms are summarized in Smith, Knoblauch, and Putnam (1993).

### Creditworthiness Measures

A key available component of this data set was the planned/scheduled principal and interest payment on total debt. It reflects the borrower's expectations of debt payments for the up-coming year. Having this values facilitates the calculation of the coverage ratio<sup>10</sup>. The coverage ratio estimates whether the borrower generated enough income to meet all expected payments and is used as an indicator of creditworthiness in this study. The coverage ratio as an indicator of creditworthiness, based on actual financial statements, has been introduced to credit scoring models as an alternative to loan classification and loan default models<sup>11</sup> (Novak and LaDue (1994), and Khoju and Barry). This

---

<sup>9</sup> Some of the borrowers reported zero liabilities, therefore, their current ratio and coverage ratio could not be calculated. To retain these borrowers in the sample and avoid values of infinity, the current ratios were given a value of 7, indicating strong liquidity, and the coverage ratio value was bound to -4 to 15 interval. The bounded interval of the coverage ratio indicates both extremes of debt repayment capacity.

<sup>10</sup> If not specified otherwise, the coverage ratio refers to the term debt and capital lease coverage ratio as defined by the FFSC.

<sup>11</sup> Historically, agricultural credit evaluation models have been predicated on predicting bank examiners' or credit reviewers' loan classification schemes (Johnson and Hagan; Dunn and Frey; Hardy and Weed; Lufburrow et al.; Hardy and Adrian; Hardy et al., Turvey and Brown, and Oltman). These studies have assessed the ability of statistical, mathematical or judgmental methods to replicate expert judgment. However, these models present some problems when credit evaluation is concerned. It is difficult to determine whether the error is due to the model or bank examiners' or credit reviewers' loan

indicator of creditworthiness is aligned with cash-flow or performance-based lending, as opposed to the more traditional collateral-based lending, and have been facilitated by improvements in farm records and computerized loan analysis (Novak and LaDue (1996).

The coverage ratio, a quantitative indicator of creditworthiness, needs to be converted to a binary variable in order to assist the lender in making a credit decision, to grant or deny a credit request. In this study, an a priori, cut-off level of one is used. A coverage ratio greater (less) than one identifies a borrower as creditworthy (less creditworthy)<sup>12</sup>. That is, the borrower has (not) generated enough income to meet all expected debt obligations.

Furthermore, the two-year and three-year average coverage ratio was found to provide a more stable, extended indicator of creditworthiness (Novak and LaDue, 1997). The annual, two-year average and three-year average measures of creditworthiness, using an a priori cut-off value of one, for the seventy farms are classified in Table 2. The number of borrowers considered creditworthy

---

classification. These problems are not limited to agricultural credit scoring models (Maris et al., and Dietrich and Kaplan).

Some agricultural credit scoring studies have used default (Miller and LaDue, and Mortensen et al.). Default is inherently a more objective measure. However, lenders and borrowers can influence default classification. Lenders can influence default classifications by decisions to forebear, restructure, or grant additional credit to repay a delinquent loan. Borrowers can influence or delay default by selling assets, depleting credit reserves, seeking off-farm employment, and other similar activities. Secondly, default is based on a single lender's criteria. Borrowers with split credit can be current with one lender and delinquent or in arrears with another lender. Thirdly, the severity of some types of default make it less than adequate. Because of these ambiguities surrounding default, and an alternative cash-flow measure of creditworthiness is used.

<sup>12</sup> The terminology less creditworthy is used instead of not creditworthy, because it is recognized that the farms in the data sample have been in operation over a nine year period and most of them have utilized some form of debt over this period. However, the sample includes Farm Service Agency, Farm Credit and various private banks borrowers. The various lending institutions represent varying degrees of creditworthiness amongst the borrowers in the sample. Creditworthiness to one lender may be less creditworthy to another. The data set can be viewed as a compilation of lender's portfolios.

decreases over-time, indicating some of the borrowers in the sample are becoming less creditworthy overtime, given the economic environment. Identifying the borrowers with diminishing debt repayment ability prior to any serious financial problems exemplifies the usefulness of the creditworthiness indicator and should be of value to lenders when evaluating a borrowers credit risk before granting a loan or evaluating their portfolio<sup>13</sup>.

### Development of the Creditworthiness Model

In this section the annual, two-year average and three-year average multiperiod credit scoring models are discussed. The RPA and logistic regression methods are used to estimate these models. The annual models are developed using 1985, 1986, 1987, 1988, and 1989 characteristic values to classify 1986, 1987, 1988, 1989 and 1990 creditworthy and less creditworthy borrowers, respectively. In other words, each model uses one period lagged characteristic values to classify creditworthy and less creditworthy borrowers. The annual models are evaluated using 1990, 1991 and 1992 characteristic values to predict 1991, 1992 and 1993 creditworthy and less creditworthy borrowers, respectfully. Lastly, the predicted creditworthy classifications are compared to the actual classifications for 1991, 1992 and 1993 to determine the intertemporal efficacy of the models.

Similarly, the two year average RPA and logistic regression models uses 1985-1986 and 1987-88 averages of the characteristic values to classify creditworthy borrowers in the average periods 1987-

---

<sup>13</sup> Granted other factors also influence credit risk such as collateral offered, reputation in the community, education, borrower's experience, borrower's personal attributes, and borrower's management ability. Many of the other factors listed have to be evaluated in conjunction with the model, but outside the model, by the loan officer. Creditworthiness models are generally designed to assist rather than replace the loan officer in lending decisions.

88 and 1989-90, respectfully. The models are evaluated using 1989-90 average characteristics values to predict 1991-92 average creditworthy and less creditworthy borrowers. The three-year average model uses 1985-86-87 average characteristic variables to classify 1988-89-90 average creditworthy and less creditworthy borrowers. The three-year average models are evaluated using 1988-89-90 average characteristic values to predict 1991-92-93 average creditworthy and less creditworthy borrowers. In both models, the two-year and three-year average models, the predicted creditworthy and less creditworthy classifications are compared to actual classifications to determine the intertemporal efficiency of the models.

RPA does not require the characteristic variables to be selected. Therefore, all sixteen FFSC recommended ratios and measures, and lagged classification variables are included in the population set. Many of the variables represent similar financial concepts, but are still included in the population set, allowing RPA to select the most appropriate variables. In addition, the predicted probability of creditworthiness from the logistic regression model was included as a possible characteristic variable.

The logistic regression model requires the characteristic or explanatory variables to be selected in advance. As a result, this study follows previous studies and specifies a parsimonious credit scoring model, where a borrower's creditworthiness is a function of solvency, liquidity and lagged debt repayment capacity (Miller and LaDue, Miller et al. and Novak and LaDue, 1997). The specific variables used in the model are debt/asset ratio, current ratio and a binary, lagged dependent variable<sup>14</sup>.

---

<sup>14</sup> Two other logistic regression models, a stepwise and "eight variable" model (the later, was presented in Novak and LaDue, 1994) were also estimated for annual, two-year and three-year average periods. The results are not reported, because the parameters did not always have the expected sign and the within sample and out-of-sample prediction rates were lower than RPA's and paramoninous (three variable) logit model's prediction rates for all the comparable time periods.

The utilization of both estimation methods require the specification of an appropriate prior probability. With RPA, the specified prior probability is essential in the development of the tree and the variables selected. With the logistic regression the prior probability is not needed to develop a model, but is necessary to classify the observations. In this study, the prior probability is determined by proportion of creditworthiness borrowers in the total sample. The values are 0.852, 0.896 and 0.905 for the annual, two-year average and three-year average periods, respectfully. The prior probabilities demonstrate that the percentage of creditworthy borrowers in the sample data set increase as the average period lengthens.

In addition to prior probabilities, misclassification cost also need to be specified. Previous agricultural credit scoring models either ignore misclassification costs or assume they are equal. However, it is reasonable to assume that the misclassification costs may not be the same for all types of decisions. The costs of granting, or renewing, a loan to a less creditworthy borrower is typically not the same as denying, or not renewing, a loan to a creditworthy borrower. This study does not estimate the cost of these misclassifications, but demonstrates the classification sensitivity of these costs. The relative cost of Type I and Type II misclassification errors are varied accordingly from 1:1, 2:1, 3:1, 4:1 and 5:1, with the relatively higher misclassification cost put on Type I error<sup>15</sup>. While the less creditworthy measure used in this model may not be as serious as an actual loan losses or bankruptcy of a borrower. There is still a higher cost associated to loan servicing, payment collection, and loan analysis for less creditworthy borrowers.

---

<sup>15</sup> Type I error is a less creditworthy borrower classified as a creditworthy borrower and a Type II error is a creditworthy borrower classified as a less creditworthy borrower.

### Comparison of RPA and Logit Model Results

Figure 3 presents the classification tree generated from the RPA for the annual time period when the misclassification cost of a type I error is three times greater than that a type II error. The model is simple, comprised of the coverage ratio lagged one period. Borrowers with a coverage ratio greater than 1.50 a year prior are classified as creditworthy and borrowers with a coverage ratio less than 1.50 a year prior are classified as less creditworthy. In other words, to insure all payments will be made by the borrower in the next year the current coverage ratio needs to be greater than 1.50.

In the same figure, below the classification tree, five surrogate variables are listed. These variables were selected on their ability to mimic the selected variable, the coverage ratio, and split value, 1.50. Repayment ability measured by the repayment margin and binary, lagged dependent variable are included in the list and appear to be good surrogate variables. Another variable selected as a surrogate variable is the borrower's predicted probability of creditworthiness from the logistic regression. The selection of predicted probability adds some additional validity to its use as a credit score. Also noteworthy is that the split value of the predicted probability is very similar to the prior probability of the annual sample period.

A list of competitor variables are also listed below the same figure. The repayment margin was listed as the first competitor variable. The first competitor variable implies that if the coverage ratio was restricted or eliminated from the sample the repayment margin would have been selected as the primary variable in the classification tree.

In figure 4 the two-year average classification tree is presented, again for a 3:1 relative misclassification costs, with the higher misclassification cost attribute to a type I error. In this classification tree the repayment margin was selected as the primary characteristic variable and the

coverage ratio was selected as the competitor and surrogate variable. Similar to the annual model, the binary lagged dependent variable was also selected as surrogate and competitor variables and the predicted probability was also selected as a surrogate variable. The other variables selected were net farm income, return on equity and operating expense ratio.

In figure 5 the classification tree for the three-year average period is presented. Again as a comparison to the previous two trees, a 3:1 relative misclassification cost ratio is used. The repayment margin was selected as the primary variable characteristic and the coverage ratio was selected as the surrogate and competitor variables. In this average time period, the binary lagged dependent variable or predicted probability were not selected as either competitor or surrogate variables. The selected competitor and surrogate variables were operating expense ratio, net farm income, rate of return on assets, operating profit margin ratio and interest expense ratio.

All the ratios and measures selected as surrogate or competitive variables in the two-year and three-year time periods represent a borrower's repayment capacity, financial efficiency or profitability. A borrower's solvency and liquidity position does not appear as useful in classifying two-year and three-year average indicators of borrowers creditworthiness.

The results are consistent with expectations. The indicator of creditworthiness is repayment capacity. The repayment capacity is predicated on operating profits and losses, hence profitability and financial efficiency. Similarly, the repayment capacity is not predicated on total debt, but scheduled debt which reduces the efficacy of the solvency measures and does not consider the liquidation of assets or inventory changes which reduces the efficacy of the liquidity measures.

As stated the selected characteristics variables as competitor and surrogate variables appear reasonable, but their exclusion from the actual classification tree may at first appear to be a

concern. The naive model is selected when relative misclassification costs are low and the other classification trees only have a maximum of two splits. However, this is consistent with other studies. Frydman et al. found the naive model also did best in classifying their data when misclassification costs were assumed equal, and found that the cross-validation classification trees had considerable less splits than the non-cross-validation classification trees. Their largest cross-validation classification tree had a maximum of three splits. As a result, for exposition purposes the non-cross-validation trees were presented. These trees are aesthetically more appealing. They are not pruned, have considerably more splits and classify more observations, but of course have less generalizability.

The estimated logistic regression models are presented in Table 3. All the parameters, for each of the models have the expected sign. In the annual model the debt/asset ratio and the binary lagged dependent parameters are significant at the 95% level. In the two-year average model the binary lagged dependent variable is significant at the 99% level.

The within- and out-of-sample misclassification rates of the RPA and logistic regression models are presented in Table 4 and 5, respectively. Historically, agricultural credit scoring models have been evaluated on their misclassification rates. Using annual data, the within sample misclassification rates indicates that the RPA model classifies the observation better than the logistic regression for relative misclassification cost of 1:1 and 2:1. Moreover, the RPA model is also the naive model. The computer algorithm concluded that if relative misclassification cost were equal or occurring at a 2:1 ratio then a lenders should classify all borrowers as creditworthy. When the relative misclassification cost ratios are assumed to be higher (i.e. 3:1, 4:1 and 5:1) the logistic regression does better at classifying the borrowers. Note that the logistic regression model does not change with the



misclassification cost scheme as does the RPA model, because the logistic model does not inherently consider misclassification costs. A comparison of the out-of-sample misclassification rates, using annual observations, indicates that the RPA model does better at classifying borrowers in 1991 for all relative misclassification cost scenarios, and the logistic regression model do better at classifying borrowers in 1992 and 1993 for all the relative misclassification cost scenarios.

When the two-year average data are used to develop the models, the misclassification results indicate that RPA does better at classifying the within-sample borrowers for all relative misclassification cost scenarios, but the logistic regression does better when classifying out-of-sample borrowers. Again, as in the annual data, when relative costs are equal RPA concludes that all borrowers should be classified as creditworthy in order to minimize the cost of misclassification. When the three-year average data are used to develop the models, the misclassification results indicate that the RPA model does best at classifying both the within- and out-of-sample borrowers for all relative classification cost scenarios.

An alternative method to evaluate the models, is to compare minimum expected cost of misclassification instead of overall misclassification. The evaluation results change since minimum expected cost of misclassification takes into account the reality of unequal costs of type I and II errors. In some cases, as the relative misclassification costs increase, the overall misclassification rate also increases, in order to minimize the expected cost of misclassification.

In Table 6 the expected cost of misclassification for each model and relative misclassification cost is presented. The RPA model does best at minimizing the expected misclassification cost for the all within-sample time periods for all relative misclassification costs scenarios. This is not surprising. The objective of RPA is to minimize the expected cost of misclassification, while the objective of the

logistic regression is to maximize the likelihood function for the specific data set. This is also where nonagricultural financial stress studies have concluded that RPA is a better model than other models. If this study was to conclude here, it would also conclude RPA is a better method of classification. However, this study continues by comparing intertemporal, out-of-sample observations.

Using the annual time period data, RPA model performs best in 1991 for all relative misclassification costs scenarios, and in 1992 and 1993, but only when the misclassification costs are equal. Recall that the annual RPA model was also the naïve model when the misclassification costs are equal. It is interesting to note that previous agricultural credit scoring research have typically assumed equal misclassification costs, but did not always compare the estimated model with the naïve model. In this case, the naïve model out performs the logistic regression model. However, the assumption that misclassification costs are equal is not very realistic in credit screening models.

Using the same data, the logistic regression model does best at classifying observations from 1992 and 1993, for the remainder of the relative misclassification cost scenarios. Likewise, the logistic regression does better at minimizing the expected cost of misclassification when the two-year average out-of-sample observations for each relative misclassification scenarios, except when misclassification costs are equal then RPA, represented by the naïve model, does better. Lastly, RPA does better at minimizing the expected cost of misclassification than the logistic regression model when predicting three-year average, out-of-sample observation for each of the relative misclassification costs scenarios.

Assuming the minimization of the expected cost of misclassification is the appropriate method for evaluating these two methods, we can not conclude a superior model using this data set. A different data set may have different results, and would warrant exploration.

## Conclusion

This study introduces RPA to agricultural credit scoring and because of its relative newness, provides a detailed exposition on how RPA classifies observations. The study also demonstrates RPA's advantages and disadvantages in relation to the logistic regression. RPA attributes include not requiring variable pre-selection, analyzing the univariate attributes of individual variables, not being affected by outliers, providing a competitive and surrogate variable summary, and explicitly incorporating misclassification costs. On the other hand, the logistic regression possesses some desirable attributes not contained by RPA, such as overall summary statistics and an individual quantitative credit score for each observation.

More significantly, the study corroborates the results of the non-agricultural credit classification studies. RPA performs superior to the logistic regression when classifying within sample observations using cross-validation selection methods and minimizing expected costs of misclassification. However, this study takes the validation process one step further and assumes intertemporal (out-of-sample) minimization of expected cost of misclassification is more appropriate for credit scoring model validation. Under this evaluation method the same results are not maintained. In some cases RPA outperforms the logistic regression and in other cases the logistic regression outperforms the RPA model. These findings suggest that the cross-validation method may not be effective enough to surmount overfitting the sample data which limits RPA's intertemporal prediction ability.

In addition, this study considers misclassification costs. Previously, agricultural credit scoring research has only evaluated models on the number of misclassified observations and not considered minimizing expected costs of misclassification. Granted this study only considers

relative misclassification costs, not actual costs, however, the results do indicated that these costs affect the models' performance and, in the case of RPA, the development of the model. Future agricultural credit scoring research should consider minimizing expected costs of misclassification, instead of minimizing misclassification observations, for evaluating models. Similarly, addition effort should be made towards calculating actual misclassification costs. In many cases, it is unrealistic to assume misclassification costs are equal.

In summary, the study has taken strides in introducing RPA, along with misclassification costs, into agricultural credit scoring, however, the results, with regards to RPA's superior performance, are as convincing as the non-agricultural financial stress literature empirical results, when intertemporal model validation is considered. Then again, the need to rigorously prove a new classification method's absolute superiority over existing methods may not always be necessary. Classification methods are always continuing to be refined and improved. From a rigorous theoretical standpoint, as well as its likely appeal to practitioners, RPA does present several very attractive features, and can be employed in conjunction with other existing methods, but does not necessarily have to replace them.

## References

Betubiza, E. and D.J. Leatham. "A Review of Agricultural Credit Assessment Research and Annotated Bibliography." Texas Experiment Station, Texas A&M University System, College Station, Texas, June 1990.

Breiman, L., J.H. Friedman, R.A. Olshen, and C.J. Stone. Classification and Regression Trees. Belmont CA: Wadsworth International Group 1984.

Dietrich, J.R. and R.S. Kaplan. "Empirical Analysis of the Commercial Loan Classification Decision." *The Accounting Review*. 57(1982):18-38.

Dunn, D.J. and T.L. Frey. "Discriminant Analysis of Loans for Cash Grain Farms." *Agr. Fin. Rev.* 36(1976):60-66.

Farm Financial Standard Council. *Financial Guidelines for Agricultural Producers: Recommendations of the Farm Financial Standards Council (Revised)* 1995.

Friedman, J.H. "A Recursive Partitioning Decision Rule for Nonparametric Classification." *IEEE Transactions on Computers*, April (1977): 404-09.

Frydman, H., E.I. Altman, and D. Kao. "Introducing Recursive Partitioning for Financial Classification: The Case of Financial Distress." *The Journal of Finance* 40 (1985):269-91.

Hardy, W.E. Jr. and J.L. Adrian, Jr. "A Linear Programming Alternative to Discriminant Analysis in Credit Scoring." *Agribusiness* 1 (1985):285-292.

Hardy, W.E., Jr., S.R. Spurlock, D.R. Parrish and L.A. Benoist. "An Analysis of Factors that Affect the Quality of Federal Land Bank Loan." *Southern Journal of Agricultural Economics*. 19 (1987):175-182.

Hardy, W.E. and J.B. Weed. "Objective Evaluation for Agricultural Lending." *Southern Journal of Agricultural Economics*. 12(1980):159-64.

Johnson, R.B. and A.R. Hagan. "Agricultural Loan Evaluation with Discriminant Analysis." *Southern Journal of Agricultural Economics*. 5(1973):57-62.

Khoju, M.R. and P.J. Barry. "Business Performance Based Credit Scoring Models: A New Approach to Credit Evaluation." *Proceedings North Central Region Project NC-207 "Regulatory Efficiency and Management Issues Affecting Rural Financial Markets"* Federal Reserve Bank of Chicago, Chicago Illinois, October 4-5, 1993.

LaDue, Eddy L. Warren F. Lee, Steven D. Hanson, and David Kohl. "Credit Evaluation Procedures at Agricultural Banks in the Northeast and Eastern Cornbelt." *A.E. Res.* 92-3, Cornell University, Dept. of Agricultural Economics. February 1992.

- Lufburrow, J., P.J. Barry and B.L. Dixon. "Credit Scoring for Farm Loan Pricing." *Agr. Fin. Rev.* 44 (1984):8-14.
- Marais, M.L., J.M. Patell, and M.A. Walfson. "The Experimental Design of Classification Models: An Application of Recursive Partitioning and Bootstrapping to Commercial Bank Loan Classifications." *Journal of Accounting Research Supplement.* 22 (1984):87-114.
- Maddala, G.S. Limited-Dependent and Qualitative Variables in Econometrics. Cambridge University Press. 1983
- Madalla, G.S. "Econometric Issues in the Empirical Analysis of Thrift Institutions' Insolvency and Failure." Federal Home Loan Bank Board, Invited Research Working Paper 56, October 1986.
- McFadden, D. "A Comment on Discriminate Analysis versus LOGIT Analysis." *Annals of Economics and Social Measurement* 5(1976):511-23.
- Miller, L.H., P. Barry, C. DeVuyst, D.A. Lins and B.J. Sherrick. "Farmer Mac Credit Risk and Capital Adequacy." *Agr. Fin. Rev.* 54 (1994):66-79.
- Miller, L.H. and E.L. LaDue. "Credit Assessment Models for Farm Borrowers:A Logit Analysis." *Agr. Fin. Rev.* 49(1989):22-36.
- Mortensen, T.D., L. Watt, and F.L. Leistriz. "Predicting Probability of Loan Default." *Agr. Fin. Rev.* 48(1988):60-76.
- Novak, M.P. and E.L. LaDue. "An Analysis of Multiperiod Agricultural Credit Evaluation Models for New York Dairy Farms." *Agr. Fin. Rev.* 54(1994):47-57.
- Novak, M.P. and E.L. LaDue. "Stabilizing and Extending, Qualitative and Quantitative Measure in Multiperiod Agricultural Credit Evaluation Model." *Agr. Fin. Rev.* ( forthcoming, 1997)
- Oltman, A.W. "Aggregate Loan Quality Assessment in the Search for Related Credit-Scoring Model." *Agr. Fin. Rev.* 54 (1994):94-107.
- Smith, S.F., W.A. Knoblauch, and L.D. Putnam. "Dairy Farm Management Business Summary, New York State, 1993" Department of Agricultural, Resource, and Managerial Economics, Cornell University, Ithaca, NY. September 1994. R.B.94-07.
- Splett, N.S., P.J. Barry, B.L. Dixon, and P.N. Ellinger. "A Joint Experience and Statistical Approach to Credit Scoring." *Agr. Fin. Rev.* 54 (1994):39-54.
- Srinivansan, V., and Y.H. Kim. "Credit Granting: A Comparative Analysis of Classification Procedures." *Journal of Finance* 42 (1987):665-81.

Steinberg, D. and P. Colla. CART Tree-structured Non-Parametric Data Analysis. San Diego, CA: Salford Systems, 1995.

Turvey, C.G. "Credit Scoring for Agricultural Loans: A Review with Application". *Agr. Fin. Rev.* 51 (1991):43-54.

Turvey, C.G. and R. Brown. "Credit Scoring for Federal Lending Institutions: The Case of Canada's Farm Credit Corporations." *Agr. Fin. Rev.* 50(1990):47-57.

Ziari, H.A., D.J. Leatham, and Calum G. Turvey. "Application of Mathematical Programming Techniques in Credit Scoring of Agricultural Loans." *Agr. Fin. Rev.* 55(1995):74-88.

Zmijewski, M.E. "Methodological Issues Related to the Estimation of Financial Distress Prediction Models." *Journal of Accounting Research Supplement* 22 (1994):59-86.

Table 1.

**Mean Value of the Sixteen FFCS Recommended Financial Ratios  
and Measures, 70 New York Dairy Farms , 1985-93.**

<u>Ratio/Measure</u>	<u>1985</u>	<u>1986</u>	<u>1987</u>	<u>1988</u>	<u>1989</u>	<u>1990</u>	<u>1991</u>	<u>1992</u>	<u>1993</u>
<u>Liquidity</u>									
Current Ratio	2.89	2.94	3.06	3.25	3.48	2.90	2.77	2.59	2.50
Working Capital (\$)	52,711	49,111	63,799	70,272	84,755	65,891	53,295	57,443	40,148
<u>Solvency</u>									
Debt/Asset Ratio	0.34	0.34	0.31	0.30	0.27	0.28	0.29	0.29	0.29
Equity/Asset Ratio	0.66	0.66	0.69	0.70	0.73	0.72	0.71	0.71	0.71
Debt/Equity Ratio	-0.58	2.17	0.73	0.63	0.51	0.51	0.56	0.56	0.53
<u>Profitability</u>									
Rate of Return on Assets	0.09	0.09	0.10	0.09	0.11	0.09	0.07	0.07	0.06
Rate of Return on Equity	0.10	0.10	0.11	0.10	0.13	0.10	0.06	0.08	0.07
Operating Profit Margin Ratio	0.21	0.20	0.22	0.21	0.23	0.20	0.16	0.17	0.15
Net Farm Income (\$)	105,352	100,588	122,144	129,899	148,560	133,232	105,790	133,809	130,104
<u>Debt Repayment Capacity</u>									
TDACLCR <sup>1</sup>	2.70	3.28	3.73	3.40	3.76	3.59	3.29	2.73	2.32
CRATDRM <sup>2</sup> (\$)	79,199	70,967	95,968	86,035	106,381	76,021	55,275	83,920	69,612
<u>Financial Efficiency</u>									
Asset Turnover Ratio	0.41	0.42	0.43	0.43	0.46	0.46	0.40	0.43	0.43
Operating Expense Ratio	0.60	0.62	0.60	0.62	0.60	0.64	0.68	0.67	0.69
Depreciation Expense Ratio	0.13	0.12	0.11	0.10	0.10	0.09	0.09	0.09	0.09
Interest Expense Ratio	0.08	0.07	0.06	0.06	0.05	0.05	0.06	0.05	0.05
Net Farm Income from Operation Ratio	0.27	0.26	0.29	0.28	0.30	0.27	0.23	0.24	0.23



Table 2.

Number of Annual, Two-Year Average, and Three-Year  
Average Creditworthy Farms, 70 New York Dairy Farms, 1985-93

<u>One Year</u>									
Cut-off Value	1985	1986	1987	1988	1989	1990	1991	1992	1993
1.00	64	65	65	60	66	61	52	54	50
<u>Two-Year Average</u>									
Cut-off Value	1985-86		1987-88		1989-90		1991-92		
1.00	66		65		63		57		
<u>Three-Year Average</u>									
Cut-off Value	1985-86-87			1988-89-90			1991-92-93		
1.00	68			65			57		

Table 3.

## Logistic Parameter Estimates of Creditworthiness Models

Variables	Annual	Two-Year Average	Three-Year Average
Intercept	2.02 (0.01) <sup>a</sup>	0.70 (0.59)	0.39 (0.09)
Debt/Asset Ratio	-1.90 (0.03)	-1.72 (0.26)	-0.92 (0.73)
Current Ratio	0.03 (0.78)	0.15 (0.51)	0.13 (0.72)
Lagged Dependent Variable	0.96 (0.05)	2.26 (0.01)	2.36 (0.21)
Model X <sup>2</sup>	14.26	18.71	6.16
Prior Probabilities	0.85	0.90	0.90

a) P-values are reported in parenthesis.

Table 4

## Number of Borrowers Misclassified using RPA Models

Within-Sample									
Cost <sup>1</sup>	<u>Annual</u>			<u>Two-Year Average</u>			<u>Three-Year Average</u>		
	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>
1:1 <sup>2</sup>	33	0	(33)	12	0	(12)	0	1	(1)
2:1	33	0	(33)	1	15	(16)	0	1	(1)
3:1	11	61	(72)	1	15	(16)	0	1	(1)
4:1	11	61	(72)	1	15	(16)	0	1	(1)
5:1	11	61	(72)	1	15	(16)	0	1	(1)
Out-of-Sample									
Cost	<u>1991</u>			<u>1991-92</u>			<u>1991-93</u>		
	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>
1:1	18	0	(18)	13	0	(13)	7	2	(9)
2:1	18	0	(18)	8	7	(15)	7	2	(9)
3:1	6	10	(16)	8	7	(15)	7	2	(9)
4:1	6	10	(16)	8	7	(15)	7	2	(9)
5:1	6	10	(16)	8	7	(15)	7	2	(9)
Cost	<u>1992</u>			<u>1993</u>					
	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>
1:1	16	0	(16)						
2:1	16	0	(16)						
3:1	4	16	(20)						
4:1	4	16	(20)						
5:1	4	16	(20)						
1:1	20	0	(20)						
2:1	20	0	(20)						
3:1	6	13	(19)						
4:1	6	13	(19)						
5:1	6	13	(19)						

1) Relative cost of type I misclassification (a less creditworthy borrower classified as creditworthy) to a type II misclassification (a creditworthy borrower classified as less creditworthy).

2) Summarized for each relative cost model is the number of type I misclassifications, number of type II misclassifications and, in parentheses, the total number of misclassifications.

Table 5

Number of Borrowers Misclassified Using the  
Logistic Regression Models

Within-Sample									
Cost <sup>1</sup>	<u>Annual</u>			<u>Two-Year Average</u>			<u>Three-Year Average</u>		
	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>
1:1 <sup>2</sup>	23	35	(58)	6	12	(18)	3	4	(7)
2:1	23	35	(58)	6	12	(18)	3	4	(7)
3:1	23	35	(58)	6	12	(18)	3	4	(7)
4:1	23	35	(58)	6	12	(18)	3	4	(7)
5:1	23	35	(58)	6	12	(18)	3	4	(7)
Out-of-Sample									
Cost	<u>1991</u>			<u>1991-92</u>			<u>1991-93</u>		
	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>	<u>I</u>	<u>II</u>	<u>T</u>
1:1	15	5	(20)	7	4	(11)	8	2	(10)
2:1	15	5	(20)	7	4	(11)	8	2	(10)
3:1	15	5	(20)	7	4	(11)	8	2	(10)
4:1	15	5	(20)	7	4	(11)	8	2	(10)
5:1	15	5	(20)	7	4	(11)	8	2	(10)
	<u>1992</u>								
	<u>I</u>	<u>II</u>	<u>T</u>						
1:1	5	9	(14)						
2:1	5	9	(14)						
3:1	5	9	(14)						
4:1	5	9	(14)						
5:1	5	9	(14)						
	<u>1993</u>								
	<u>I</u>	<u>II</u>	<u>T</u>						
1:1	11	4	(15)						
2:1	11	4	(15)						
3:1	11	4	(15)						
4:1	11	4	(15)						
5:1	11	4	(15)						

1) Relative cost of type I misclassification (a less creditworthy borrower classified as creditworthy) to a type II misclassification to a type II misclassification (a creditworthy borrower classified as less creditworthy).

2) Summarized for each relative cost model is the number of type I misclassifications, number of type II misclassifications and, in parentheses, the total number of misclassifications. Also, note since the development of the logistic regression is not explicitly affected by misclassification costs, the misclassification results do not change for different relative misclassification costs. This is done for comparison purposes.

Table 6

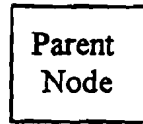
Expected Cost of Misclassification<sup>1</sup> for the RPA and Logistic Regression Models

Within-Sample							
RPA				Logistic Regression			
Cost	<u>One-Year</u>	<u>Two-Year</u>	<u>Three-Year</u>	Cost <sup>2</sup>	<u>One-Year</u>	<u>Two-Year</u>	<u>Three-Year</u>
1:1	0.150	0.100	0.014	1:1	0.198	0.134	0.110
2:1	0.300	0.122	0.014	2:1	0.303	0.184	0.164
3:1	0.314	0.131	0.014	3:1	0.408	0.234	0.218
4:1	0.364	0.139	0.014	4:1	0.512	0.284	0.272
5:1	0.414	0.147	0.014	5:1	0.617	0.334	0.326
Out-of-Sample							
RPA				Logistic Regression			
Cost	<u>1991</u>	<u>1991-92</u>	<u>1991-93</u>	Cost <sup>2</sup>	<u>1991</u>	<u>1991-92</u>	<u>1991-93</u>
1:1	0.150	0.100	0.080	1:1	0.207	0.117	0.087
2:1	0.300	0.234	0.129	2:1	0.332	0.171	0.143
3:1	0.314	0.295	0.177	3:1	0.457	0.225	0.198
4:1	0.364	0.357	0.226	4:1	0.582	0.279	0.254
5:1	0.414	0.418	0.274	5:1	0.707	0.332	0.309
	<u>1992</u>				<u>1992</u>		
1:1	0.150			1:1	0.189		
2:1	0.300			2:1	0.235		
3:1	0.338			3:1	0.282		
4:1	0.366			4:1	0.329		
5:1	0.395			5:1	0.376		
	<u>1993</u>				<u>1993</u>		
1:1	0.150			1:1	0.151		
2:1	0.300			2:1	0.233		
3:1	0.356			3:1	0.316		
4:1	0.401			4:1	0.398		
5:1	0.446			5:1	0.481		

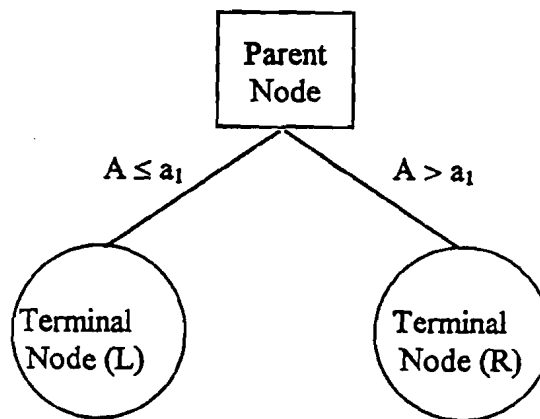
1) Expected cost of misclassification =  $\pi_i c_{ij} n_{ij}(T)/N_j + \pi_j c_{ji} n_{ji}(T)/N_i$ , where  $n_{ij}(T)$  = total number from group  $j$  observations misclassified as group  $i$ , similar for  $n_{ji}$ , and  $N_i$  = sample size of group  $i$ , similarly for  $N_j$ .

2) The logistic regression does not explicitly account for cost of misclassification during the development of the model, however for comparative purposes the relative costs are varied. In other words, the same number of borrowers are misclassified for each model and relative costs, except the misclassifications are weighted differently.

Tree  $T_0$



Tree  $T_1$



Tree  $T_2$

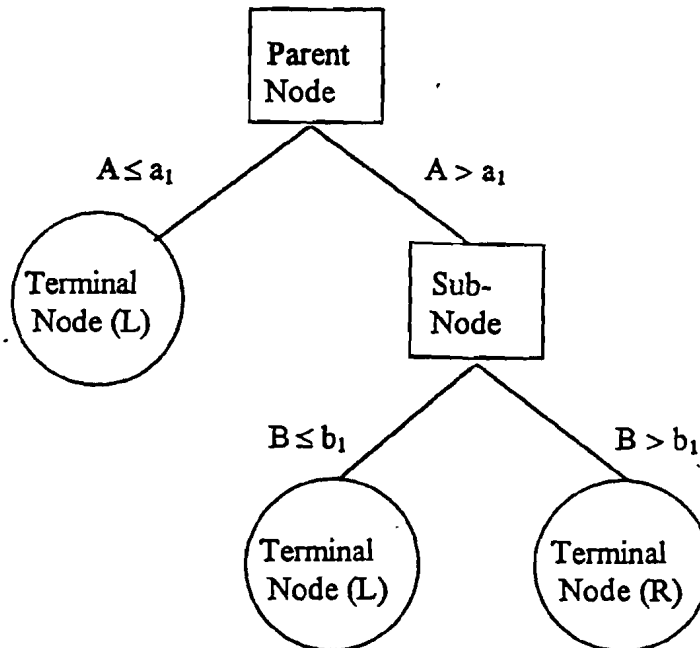


Figure 1. Hypothetical Recursive Partitioning Algorithm Tree

Tree T<sub>3</sub>

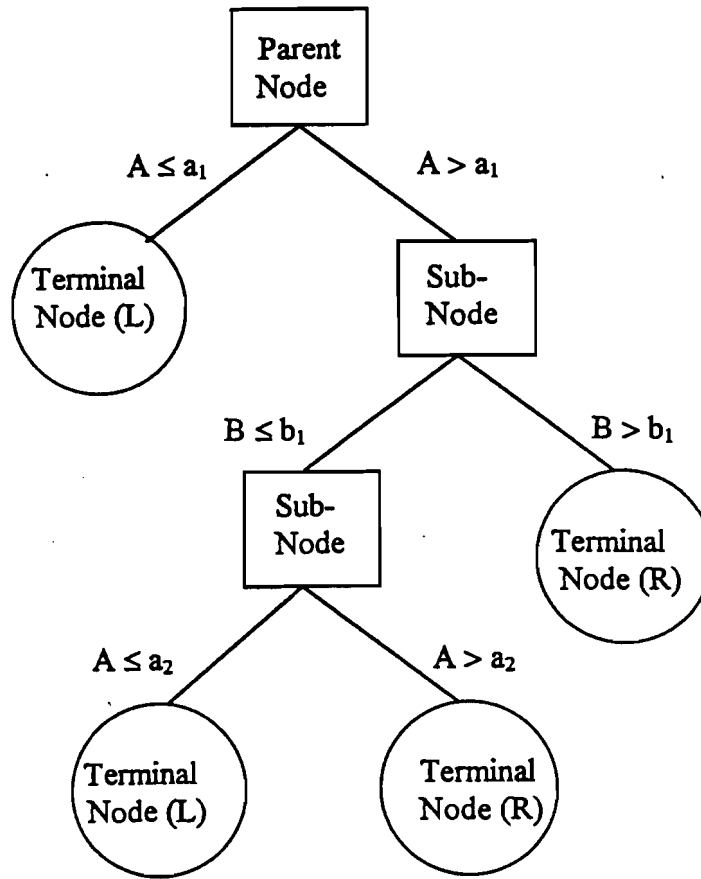


Figure 1 (Continued). Hypothetical Recursive Partitioning Algorithm Tree

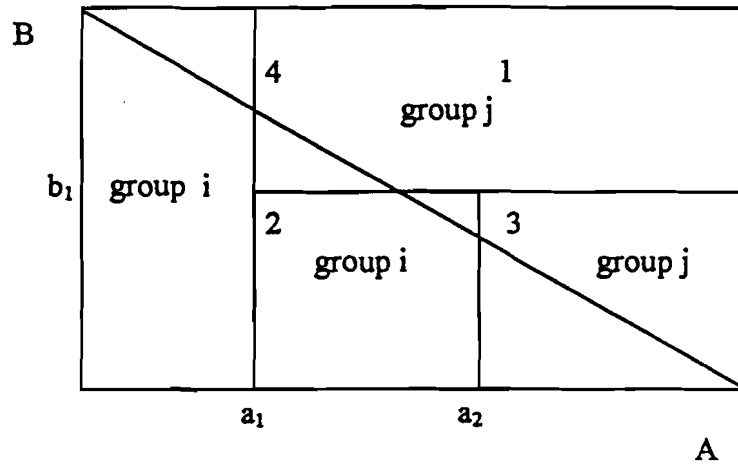
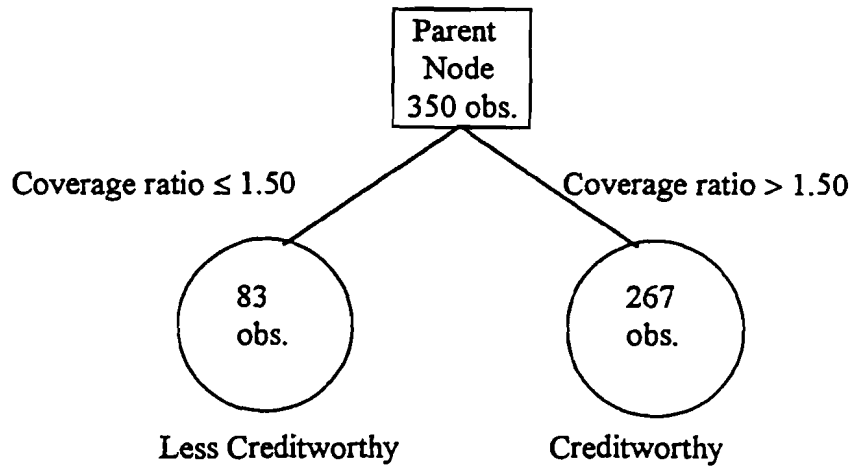


Figure 2. Observation Space





Surogate Variables

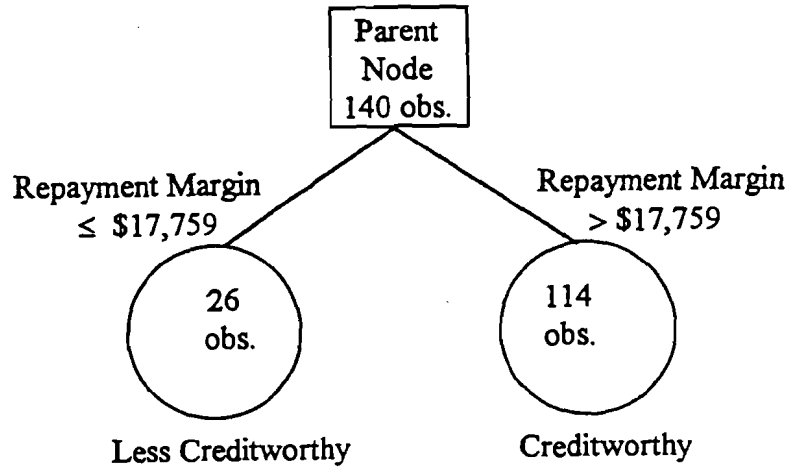
Split Values

1. Capital replacement and term debt repayment margin	\$18,552
2. Net farm income from operations ratio	0.181
3. Binary labbed dependent variable	0.500
4. Predicted probability of creditworthiness	0.837
5. Operating expense ratio	0.747

Competitor Variables

1. Capital replacement and term debt repayment margin	\$18,419
2. Debt/equity ratio	0.408
3. Debt/asset ratio	0.290
4. Operating expense ratio	0.640
5. Operating profit margin ratio	0.152

Figure 3. RPA Tree Using Annual Data



Surrogate Variables

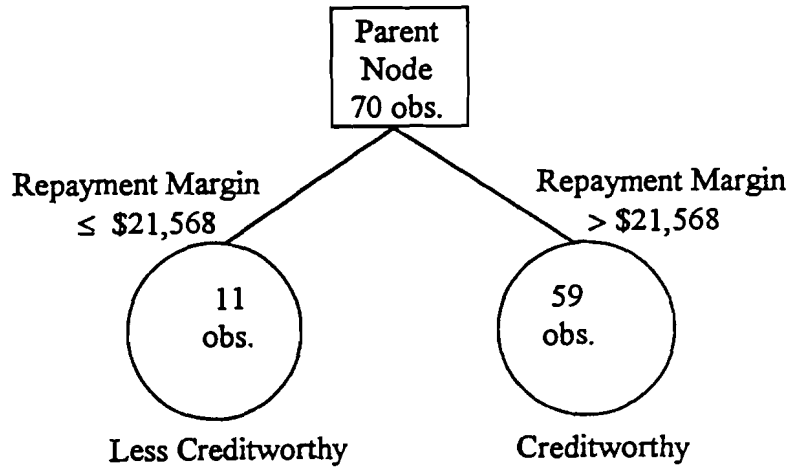
Split Values

1. Term debt and capital lease coverage ratio	1.405
2. Predicted probability of creditworthiness	0.818
3. Binary lagged dependent variable	0.500
4. Net farm income	..\$ 22,922
5. Interest expense ratio	0.158

Competitor Variables

1. Term debt and capital lease coverage ratio	1.698
2. Operating expense ratio	0.749
3. Predicted probability of creditworthiness	0.853
4. Rate of return on equity	0.013
5. Net farm income	\$69,172

Figure 4. RPA Tree Using Two-Year Average Data



Surrogate Variables

1. Term debt and capital lease coverage ratio
2. Operating expense ratio
3. Net farm income
4. Rate of return of assets
5. Current ratio

Split Values

- 1.429  
0.748  
\$22,265  
0.046  
0.856

Competitor Variables

1. Term debt and capital lease coverage ratio
2. Operating expense ratio
3. Rate of return on assets
4. Interest expense ratio
5. Operating profit margin ratio

- 1.663  
0.748  
0.046  
0.277  
0.158

Figure 5. RPA Tree Using Three-Year Average Data

**OTHER A.R.M.E. WORKING PAPERS**

<b><u>WP No</u></b>	<b><u>Title</u></b>	<b><u>Author(s)</u></b>
97-17	Trust in Japanese Interfirm Relations: Institutional Sanctions Matter	Hagen, J.M. and S. Choe
97-16	Effective Incentives and Chickpea Competitiveness in India	Rao, K. and S. Kyle
97-15	Can Hypothetical Questions Reveal True Values? A Laboratory Comparison of Dichotomous Choice and Open-Ended Contingent Values with Auction Values	Balistreri, E., G. McClelland, G. Poe and W. Schulze
97-14	Global Hunger: The Methodologies Underlying the Official Statistics	Poleman, T.T.
97-13	Agriculture in the Republic of Karakalpakstan and Khorezm Oblast of Uzbekistan	Kyle, S. and P. Chabot
97-12	Crop Budgets for the Western Region of Uzbekistan	Chabot, P. and S. Kyle
97-11	Farmer Participation in Reforestation Incentive Programs in Costa Rica	Thacher, T., D.R. Lee and J.W. Schelhas
97-10	Ecotourism Demand and Differential Pricing of National Park Entrance Fees in Costa Rica	Chase, L.C., D.R. Lee, W.D. Schulze and D.J. Anderson
97-09	The Private Provision of Public Goods: Tests of a Provision Point Mechanism for Funding Green Power Programs	Rose, S.K., J. Clark, G.L. Poe, D. Rondeau and W.D. Schulze
97-08	Nonrenewability in Forest Rotations: Implications for Economic and Ecosystem Sustainability	Erickson, J.D., D. Chapman, T. Fahey and M.J. Christ
97-07	Is There an Environmental Kuznets Curve for Energy? An Econometric Analysis	Agras, J. and D. Chapman
97-06	A Comparative Analysis of the Economic Development of Angola and Mozambique	Kyle, S.
97-05	Success in Maximizing Profits and Reasons for Profit Deviation on Dairy Farms	Tauer, L. and Z. Stefanides
97-04	A Monthly Cycle in Food Expenditure and Intake by Participants in the U.S. Food Stamp Program	Wilde, P. and C. Ranney
97-03	Estimating Individual Farm Supply and Demand Elasticities Using Nonparametric Production Analysis	Stefanides, Z. and L. Tauer