*Chapter 5*
# Build Decision Trees

## About the Tasks That You Will Perform

Now that you have verified the input data, it is time to build predictive models. You perform the following tasks to model the input data using nonparametric decision trees:

1.  You enable SAS Enterprise Miner to automatically train a full decision tree and to automatically prune the tree to an optimal size. When training the tree, you select split rules at each step to maximize the split decision logworth. Split decision logworth is a statistic that measures the effectiveness of a particular split decision at differentiating values of the target variable. For more information about logworth, see the SAS Enterprise Miner Help.

2.  You interactively train a decision tree. At each step, you select from a list of candidate rules to define the split rule that you deem to be the best.

3.  You use a Gradient Boosting node to generate a set of decision trees that form a single predictive model. Gradient boosting is a boosting approach that resamples the analysis data set several times to generate results that form a weighted average of the re-sampled data set.

## Automatically Train and Prune a Decision Tree

Decision tree models are advantageous because they are conceptually easy to understand, yet they readily accommodate nonlinear associations between input variables and one or more target variables. They also handle missing values without the need for imputation. Therefore, you decide to first model the data using decision trees. You will compare decision tree models to other models later in the example.

However, before you add and run the Decision Tree node, you will add a Control Point node. The Control Point node is used to simplify a process flow diagram by reducing the

number of connections between multiple interconnected nodes. By the end of this example, you will have created five different models of the input data set, and two Control Point nodes to connect these nodes. The first Control Point node, added here, will distribute the input data to each of these models. The second Control Point node will collect the models and send them to evaluation nodes.

To use the Control Point node:

1.  Select the **Utility** tab on the Toolbar.

2.  Select the **Control Point** node icon. Drag the node into the Diagram Workspace.

3.  Connect the **Replacement** node to the **Control Point** node.

SAS Enterprise Miner enables you to build a decision tree in two ways: automatically and interactively. You will begin by letting SAS Enterprise Miner automatically train and prune a tree.

To use the **Decision Tree** node to automatically train and prune a decision tree:

1.  Select the **Model** tab on the Toolbar.

2.  Select the **Decision Tree** node icon. Drag the node into the Diagram Workspace.

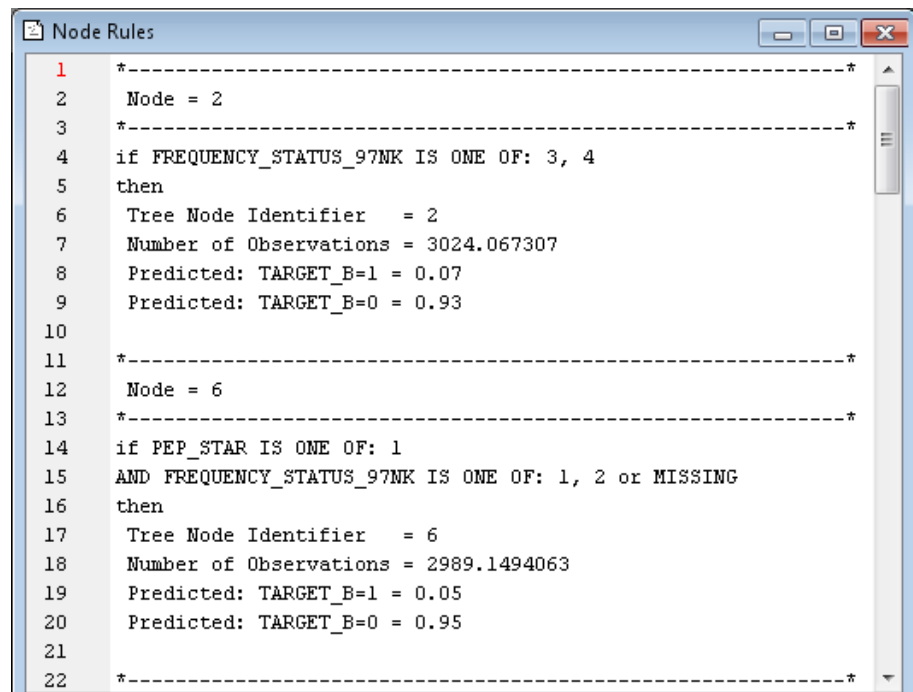3.  Connect the **Control Point** node to the **Decision Tree** node.



4.  Select the **Decision Tree** node. In the Properties Panel, scroll down to view the **Train** properties:

    •   Click on the value of the **Maximum Depth** splitting rule property, and enter **10**. This specification enables SAS Enterprise Miner to train a tree that includes up to ten generations of the root node. The final tree in this example, however, will have fewer generations due to pruning.

    •   Click on the value of the **Leaf Size** node property, and enter **8**. This specification constrains the minimum number of training observations in any leaf to eight.

- Click on the value of the **Number of Surrogate Rules** node property, and enter
  **4**. This specification enables SAS Enterprise Miner to use up to four surrogate
  rules in each non-leaf node if the main splitting rule relies on an input whose
  value is missing.

*Note:* The **Assessment Measure** subtree property is automatically set to **Decision**
because you defined a profit matrix in "Create a Data Source" on page 11.
Accordingly, the Decision Tree node will build a tree that maximizes profit in the
validation data.

5. In the Diagram Workspace, right-click the Decision Tree node, and select **Run** from
   the resulting menu. Click **Yes** in the Confirmation window that opens.

6. In the window that appears when processing completes, click **Results**. The Results
   window appears.

   a. On the **View** menu, select **Model ⇨ Node Rules**. The English Rules window
      appears.

   b. Expand the Node Rules window. This window contains the IF-THEN logic that
      distributes observations into each leaf node of the decision tree.



   In the Output window, the **Tree Leaf Report** indicates that there are seven leaf
   nodes in this tree. For each leaf node, the following information is listed:

   - node number

   - number of training observations in the node

   - percentage of training observations in the node with TARGET_B=1 (did
     donate), adjusted for prior probabilities

   - percentage of training observations in the node with TARGET_B=0 (did not
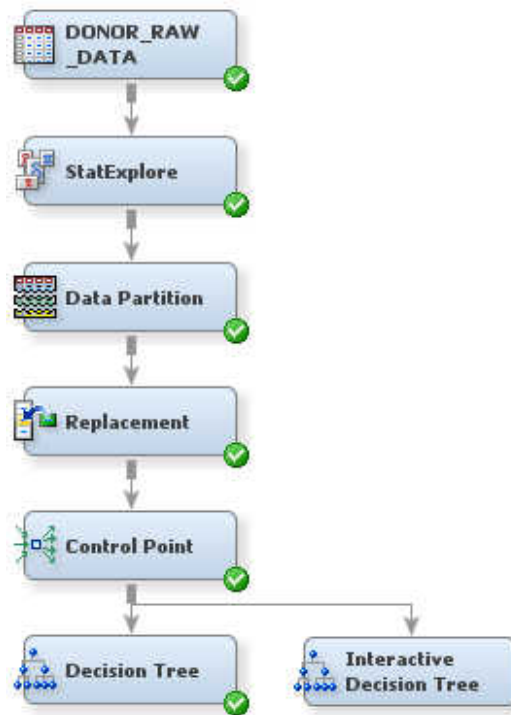     donate), adjusted for prior probabilities

   This tree has been automatically pruned to an optimal size. Therefore, the node
   numbers that appear in the final tree are not sequential. In fact, they reflect the
   positions of the nodes in the full tree, before pruning.
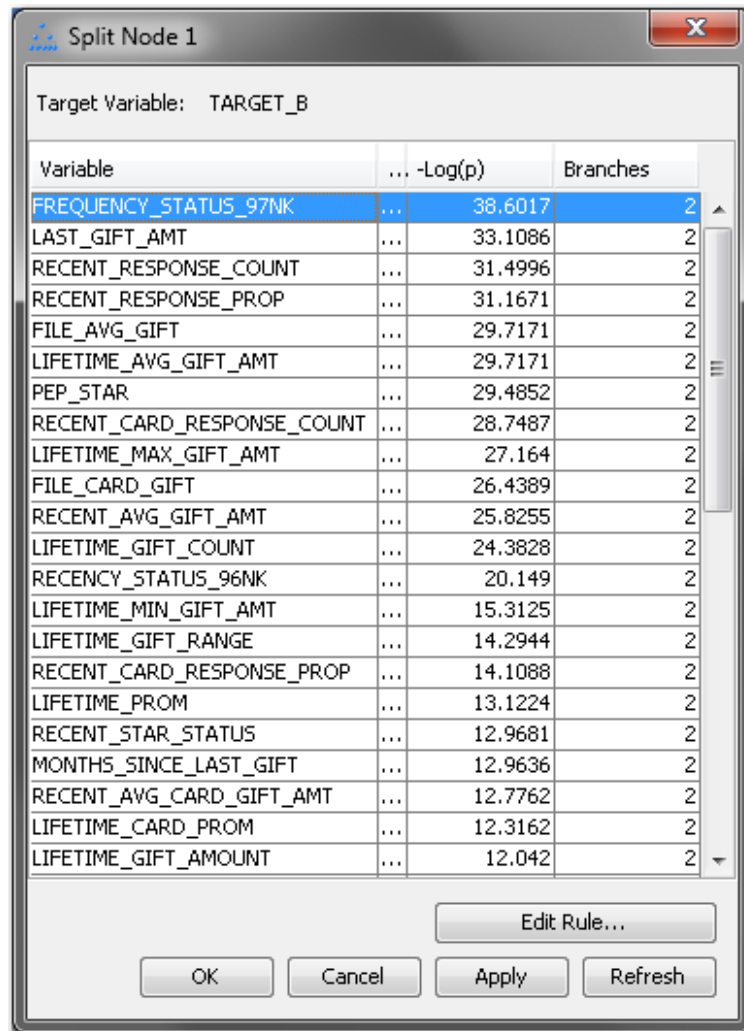
7. Close the Results window.

## Interactively Train a Decision Tree

To use the Decision Tree node to interactively train and prune a decision tree:

1. From the **Model** tab on the Toolbar, select the **Decision Tree** node icon. Drag the node into the Diagram Workspace.

2. In the Diagram Workspace, right-click the **Decision Tree** node, and select **Rename** from the resulting menu. Enter **Interactive Decision Tree** and then click **OK** in the window that opens.

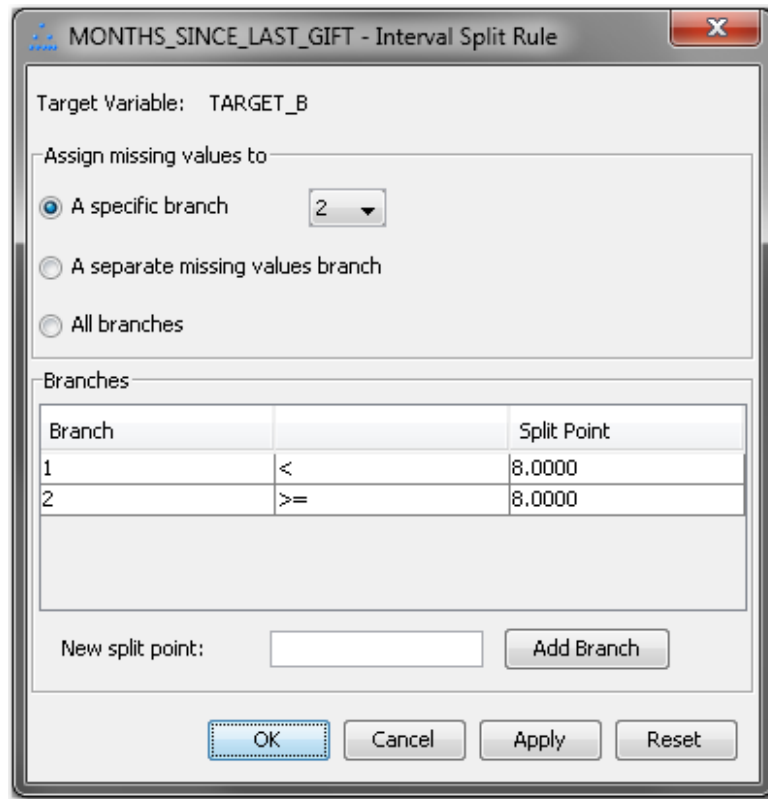3. Connect the Control Point node to the Interactive Decision Tree node.



4. Select the **Interactive Decision Tree** node. In the Properties Panel, in the **Train** properties group, click on the ellipses that represent the value of **Interactive**. The Interactive Decision Tree window appears.

   a. Select the root node (at this point, the only node in the tree), and then from the **Action** menu, select **Split Node**. The Split Node window appears that lists the candidate splitting rules ranked by logworth (-Log(p)). The FREQUENCY_STATUS_97NK rule has the highest logworth. Ensure that this row is selected, and click **OK**.

b.  The tree now has two additional nodes. Select the lower left node (where
    FREQUENCY_STATUS_97NK is 3 or 4), and then from the **Action** menu,
    select **Split Node**. In the Split Node window that opens, select
    MONTHS_SINCE_LAST_GIFT, which ranks second in logworth, and click
    **Edit Rule** to manually specify the split point for this rule. The Interval Split Rule
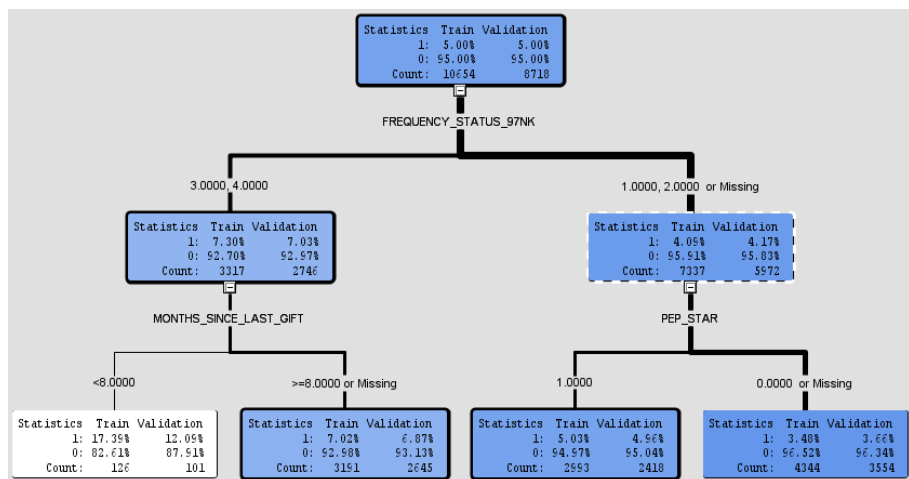    window appears.

    Enter **8** as the **New split point**, and click **Add Branch**. Then, select Branch 3
    (>= 8.5) and click **Remove Branch**. Click **OK**.

Ensure that MONTHS_SINCE_LAST_GIFT is selected in the Split Node window, and click **OK**.

c. Select the first generation node that you have not yet split (where FREQUENCY_STATUS_97NK is 1, 2, or Missing). From the **Action** menu, select **Split Node**. In the Split Node window that opens, ensure that PEP_STAR is selected, and click **OK**.

The tree now has seven nodes, four of which are leaf nodes. The nodes are colored from light to dark, corresponding to low to high percentages of correctly classified observations.



d. Select the lower right node (where FREQUENCY_STATUS_97NK is 1, 2, or Missing and PEP_STAR is 0 or Missing). From the **Action** menu, select **Train Node**. This selection causes SAS Enterprise Miner to continue adding

generations of this node until a stopping criterion is met. For more information about stopping criteria for decision trees, see the SAS Enterprise Miner Help.

*Note:* In the Interactive Decision Tree window, you can prune decision trees. However, in this example, you will leave the tree in its current state.

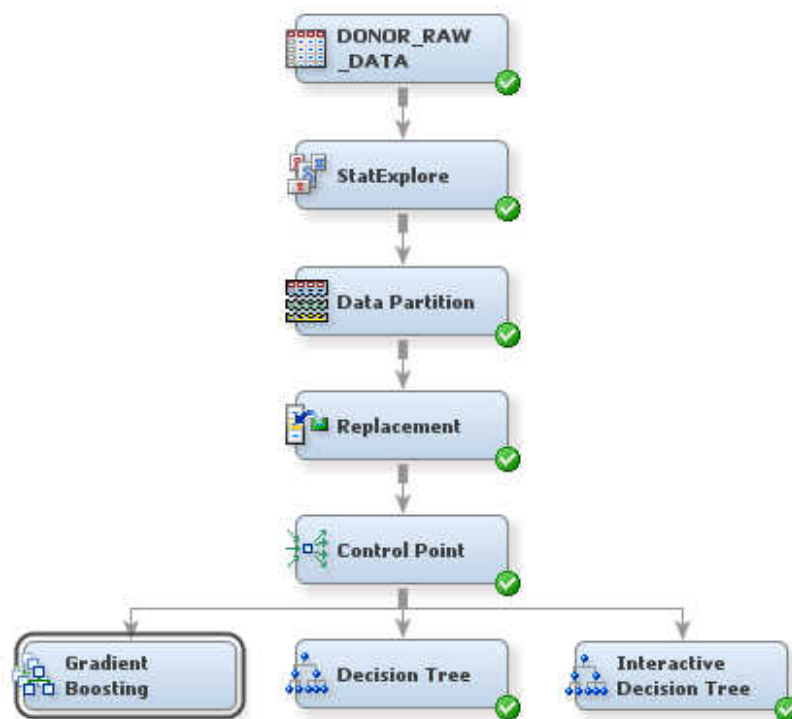e. Close the Interactive Decision Tree window.

# Create a Gradient Boosting Model of the Data

The Gradient Boosting node uses a partitioning algorithm to search for an optimal partition of the data for a single target variable. Gradient boosting is an approach that resamples the analysis data several times to generate results that form a weighted average of the resampled data set. Tree boosting creates a series of decision trees that form a single predictive model.

Like decision trees, boosting makes no assumptions about the distribution of the data. Boosting is less prone to overfit the data than a single decision tree. If a decision tree fits the data fairly well, then boosting often improves the fit. For more information about the Gradient Boosting node, see the SAS Enterprise Miner help documentation.

To create a gradient boosting model of the data:

1. Select the **Model** tab on the Toolbar.

2. Select the **Gradient Boosting** node icon. Drag the node into the Diagram Workspace.

3. Connect the **Control Point** node to the **Gradient Boosting** node.



4. Select the **Gradient Boosting** node. In the Properties Panel, set the following properties:

- Click on the value for the **Maximum Depth** property, in the **Splitting Rule** subgroup, and enter `10`. This property determines the number of generations in each decision tree created by the Gradient Boosting node.

- Click on the value for the **Number of Surrogate Rules** property, in the **Node** subgroup, and enter `2`. Surrogate rules are backup rules that are used in the event of missing data. For example, if your primary splitting rule sorts donors based on their ZIP codes, then a reasonable surrogate rule would sort based on the donor's city of residence.

5. In the Diagram Workspace, right-click the Gradient Boosting node, and select **Run** from the resulting menu. Click **Yes** in the Confirmation window that opens.

6. In the Run Status window, select **OK**.

*T I P*   The book "Decision Trees for Analytics Using SAS Enterprise Miner" offers additional information about alternative measures of the effectiveness of a split, options for training and pruning, suggestions for guiding tree growth, and examples of multiple tree and gradient boosting models.

*T I P*   "Predictive Modeling with SAS Enterprise Miner: Practical Solutions for Business Applications" offers examples that automatically and interactively train and prune decision tree models and examples that create gradient boosting models.