


One solution to this problem is to stop the tree from growing once it reaches a certain number of decisions or if the decision nodes contain only a small number of examples. This is called early stopping or **pre-pruning** the decision tree. As the tree avoids doing needless work, this is an appealing strategy. However, one downside is that there is no way to know whether the tree will miss subtle, but important patterns that it would have learned had it grown to a larger size.

An alternative, called **post-pruning** involves growing a tree that is too large, then using pruning criteria based on the error rates at the nodes to reduce the size of the tree to a more appropriate level. This is often a more effective approach than pre-pruning because it is quite difficult to determine the optimal depth of a decision tree without growing it first. Pruning the tree later on allows the algorithm to be certain that all important data structures were discovered.

[ The implementation details of pruning operations are very technical and beyond the scope of this book. For a comparison of some of the available methods, see: *A comparative analysis of methods for pruning decision trees*. IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 19, pp. 476-491, by F. Esposito, D. Malerba, and G. Semeraro (1997).]

One of the benefits of the C5.0 algorithm is that it is opinionated about pruning – it takes care of many of the decisions, automatically using fairly reasonable defaults. Its overall strategy is to postprune the tree. It first grows a large tree that overfits the training data. Later, nodes and branches that have little effect on the classification errors are removed. In some cases, entire branches are moved further up the tree or replaced by simpler decisions. These processes of grafting branches are known as **subtree raising** and **subtree replacement**, respectively.

Balancing overfitting and underfitting a decision tree is a bit of an art, but if model accuracy is vital it may be worth investing some time with various pruning options to see if it improves performance on the test data. As you will soon see, one of the strengths of the C5.0 algorithm is that it is very easy to adjust the training options.

Example – identifying risky bank loans using C5.0 decision trees

The global financial crisis of 2007-2008 has highlighted the importance of transparency and rigor in banking practices. As the availability of credit has been limited, banks are increasingly tightening their lending systems and turning to machine learning to more accurately identify risky loans.

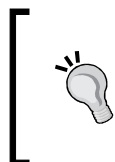
Decision trees are widely used in the banking industry due to their high accuracy and ability to formulate a statistical model in plain language. Since government organizations in many countries carefully monitor lending practices, executives must be able to explain why one applicant was rejected for a loan while others were approved. This information is also useful for customers hoping to determine why their credit rating is unsatisfactory.

It is likely that automated credit scoring models are employed for instantly approving credit applications on the telephone and the web. In this section, we will develop a simple credit approval model using C5.0 decision trees. We will also see how the results of the model can be tuned to minimize errors that result in a financial loss for the institution.

Step 1 – collecting data

The idea behind our credit model is to identify factors that make an applicant at higher risk of default. Therefore, we need to obtain data on a large number of past bank loans and whether the loan went into default, as well as information about the applicant.

Data with these characteristics are available in a dataset donated to the UCI Machine Learning Data Repository (<http://archive.ics.uci.edu/ml>) by *Hans Hofmann* of the University of Hamburg. They represent loans obtained from a credit agency in Germany.



The data presented in this chapter has been modified slightly from the original one for eliminating some preprocessing steps. To follow along with the examples, download the `credit.csv` file from Packt Publishing's website and save it to your R working directory.



The credit dataset includes 1,000 examples of loans, plus a combination of numeric and nominal features indicating characteristics of the loan and the loan applicant. A class variable indicates whether the loan went into default. Let's see if we can determine any patterns that predict this outcome.

Step 2 – exploring and preparing the data

As we have done previously, we will import the data using the `read.csv()` function. We will ignore the `stringsAsFactors` option (and therefore use the default value, `TRUE`) as the majority of features in the data are nominal. We'll also look at the structure of the credit data frame we created:

```
> credit <- read.csv("credit.csv")
> str(credit)
```

The first several lines of output from the `str()` function are as follows:

```
'data.frame':1000 obs. of 17 variables:
 $ checking_balance : Factor w/ 4 levels "< 0 DM", "> 200 DM", ..
 $ months_loan_duration: int 6 48 12 ...
 $ credit_history      : Factor w/ 5 levels "critical", "good", ...
 $ purpose            : Factor w/ 6 levels "business", "car", ...
 $ amount             : int 1169 5951 2096 ...
```

We see the expected 1,000 observations and 17 features, which are a combination of factor and integer data types.

Let's take a look at some of the `table()` output for a couple of features of loans that seem likely to predict a default. The `checking_balance` and `savings_balance` features indicate the applicant's checking and savings account balance, and are recorded as categorical variables:

```
> table(credit$checking_balance)
  < 0 DM  > 200 DM 1 - 200 DM  unknown
    274      63      269      394

> table(credit$savings_balance)
  < 100 DM > 1000 DM 100 - 500 DM 500 - 1000 DM  unknown
    603      48      103      63      183
```

Since the loan data was obtained from Germany, the currency is recorded in Deutsche Marks (DM). It seems like a safe assumption that larger checking and savings account balances should be related to a reduced chance of loan default.

Some of the loan's features are numeric, such as its term (`months_loan_duration`), and the amount of credit requested (`amount`).

```
> summary(credit$months_loan_duration)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
   4.0    12.0    18.0    20.9    24.0    72.0

> summary(credit$amount)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  250    1366    2320    3271    3972   18420
```

The loan amounts ranged from 250 DM to 18,420 DM across terms of 4 to 72 months, with a median duration of 18 months and amount of 2,320 DM.

The `default` variable indicates whether the loan applicant was unable to meet the agreed payment terms and went into default. A total of 30 percent of the loans went into default:

```
> table(credit$default)
 no yes
700 300
```


A high rate of `default` is undesirable for a bank because it means that the bank is unlikely to fully recover its investment. If we are successful, our model will identify applicants that are likely to default, so that this number can be reduced.

Data preparation – creating random training and test datasets

As we have done in previous chapters, we will split our data into two portions: a training dataset to build the decision tree and a test dataset to evaluate the performance of the model on new data. We will use 90 percent of the data for training and 10 percent for testing, which will provide us with 100 records to simulate new applicants.

As prior chapters used data that had been sorted in a random order, we simply divided the dataset into two portions by taking the first 90 percent of records for training, and the remaining 10 percent for testing. In contrast, our data here is not randomly ordered. Suppose that the bank had sorted the data by the loan amount, with the largest loans at the end of the file. If we use the first 90 percent for training and the remaining 10 percent for testing, we would be building a model on only the small loans and testing the model on the big loans. Obviously, this could be problematic.


We'll solve this problem by randomly ordering our credit data frame prior to splitting. The `order()` function is used to rearrange a list of items in ascending or descending order. If we combine this with a function to generate a list of random numbers, we can generate a randomly-ordered list. For random number generation, we'll use the `runif()` function, which by default generates a sequence of random numbers between 0 and 1.

 If you're trying to figure out where the `runif()` function gets its name, the answer is due to the fact that it chooses numbers from a uniform distribution, which we learned about in *Chapter 2, Managing and Understanding Data*.

The following command creates a randomly-ordered `credit` data frame. The `set.seed()` function is used to generate random numbers in a predefined sequence, starting from a position known as a **seed** (set here to the arbitrary value 12345). It may seem that this defeats the purpose of generating random numbers, but there is a good reason for doing it this way. The `set.seed()` function ensures that if the analysis is repeated, an identical result is obtained.

```
> set.seed(12345)
> credit_rand <- credit[order(runif(1000)), ]
```

The `runif(1000)` command generates a list of 1,000 random numbers. We need exactly 1,000 random numbers because there are 1,000 records in the `credit` data frame. The `order()` function then returns a vector of numbers indicating the sorted position of the 1,000 random numbers. We then use these positions to select rows in the `credit` data frame and store in a new data frame named `credit_rand`.

 To better understand how this function works, note that `order(c(0.5, 0.25, 0.75, 0.1))` returns the sequence 4 1 2 3 because the smallest number (0.1) appears fourth, the second smallest (0.25) appears first, and so on.


To confirm that we have the same data frame sorted differently, we'll compare values on the `amount` feature across the two data frames. The following code shows the summary statistics:

```
> summary(credit$amount)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
   250   1366   2320   3271   3972   18420
> summary(credit_rand$amount)
   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
   250   1366   2320   3271   3972   18420
```

We can use the `head()` function to examine the first few values in each data frame:

```
> head(credit$amount)
[1] 1169 5951 2096 7882 4870 9055
> head(credit_rand$amount)
[1] 1199 2576 1103 4020 1501 1568
```

Since the summary statistics are identical while the first few values are different, this suggests that our random shuffle worked correctly.

 If your results do not match exactly with the previous ones, ensure that you run the command `set.seed(214805)` immediately prior to creating the `credit_rand` data frame.

Now, we can split into training (90 percent or 900 records), and test data (10 percent or 100 records) as we have done in previous analyses:

```
> credit_train <- credit_rand[1:900, ]
> credit_test <- credit_rand[901:1000, ]
```

If all went well, we should have about 30 percent of defaulted loans in each of the datasets.

```
> prop.table(table(credit_train$default))
      no      yes
0.7022222 0.2977778
> prop.table(table(credit_test$default))
      no      yes
0.68 0.32
```

This appears to be a fairly equal split, so we can now build our decision tree.

Step 3 – training a model on the data

We will use the C5.0 algorithm in the `C50` package for training our decision tree model. If you have not done so already, install the package with `install.packages("C50")` and load it to your R session using `library(C50)`.

The following syntax box lists some of the most commonly used commands for building decision trees. Compared to the machine learning approaches we have used previously, the C5.0 algorithm offers many more ways to tailor the model to a particular learning problem, but even more options are available. The `?C5.0Control` command displays the help page for more details on how to finely-tune the algorithm.

C5.0 decision tree syntax

using the `C5.0()` function in the `C50` package

Building the classifier:

```
m <- C5.0(train, class, trials = 1, costs = NULL)
```

- `train` is a data frame containing training data
- `class` is a factor vector with the class for each row in the training data
- `trials` is an optional number to control the number of boosting iterations (by default, 1)
- `costs` is an optional matrix specifying costs associated with types of errors

The function will return a C5.0 model object that can be used to make predictions.

Making predictions:

```
p <- predict(m, test, type = "class")
```

- `m` is a model trained by the `C5.0()` function
- `test` is a data frame containing test data with the same features as the training data used to build the classifier.
- `type` is either `"class"` or `"prob"` and specifies whether the predictions should be the most likely class value or the raw predicted probabilities

The function will return a vector of predicted class values or raw predicted probabilities depending upon the value of the `type` parameter.

Example:

```
credit_model <- C5.0(credit_train, loan_default)
credit_prediction <- predict(credit_model, credit_test)
```

For the first iteration of our credit approval model, we'll use the default C5.0 configuration, as shown in the following code. The 17th column in `credit_train` is the class variable, `default`, so we need to exclude it from the training data frame as an independent variable, but supply it as the target factor vector for classification:

```
> credit_model <- C5.0(credit_train[-17], credit_train$default)
```

The `credit_model` object now contains a C5.0 decision tree object. We can see some basic data about the tree by typing its name:

```
> credit_model
```

Call:

```
C5.0.default(x = credit_train[-17], y = credit_train$default)
```

Classification Tree

Number of samples: 900

Number of predictors: 16

Tree size: 67

The preceding text shows some simple facts about the tree, including the function call that generated it, the number of features (that is, `predictors`), and examples (that is, `samples`) used to grow the tree. Also listed is the tree size of 67, which indicates that the tree is 67 decisions deep—quite a bit larger than the trees we've looked at so far!

To see the decisions, we can call the `summary()` function on the model:

```
> summary(credit_model)
```

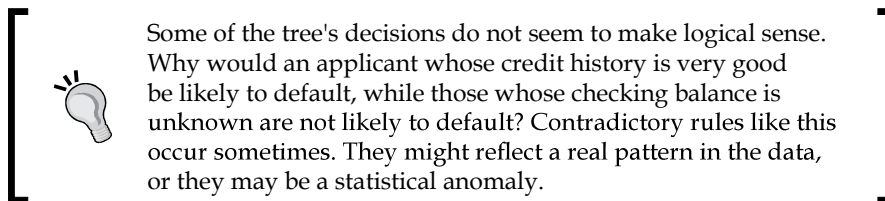
This results in the following output:

```
C5.0 [Release 2.07 GPL Edition]
-----
Class specified by attribute 'outcome'
Read 900 cases (17 attributes) from undefined.data
Decision tree:
checking_balance = unknown: no (358/44)
checking_balance in {< 0 DM,> 200 DM,1 - 200 DM}:
...credit_history in {perfect,very good}:
...dependents > 1: yes (10/1)
:
:   dependents <= 1:
:   ...savings_balance = < 100 DM: yes (39/11)
:   :
:   :   savings_balance in {> 1000 DM,500 - 1000 DM,unknown}: no (8/1)
:   :   savings_balance = 100 - 500 DM:
:   :   ...checking_balance = < 0 DM: no (1)
:   :   checking_balance in {> 200 DM,1 - 200 DM}: yes (5/1)
```


The preceding output shows some of the first branches in the decision tree. The first four lines could be represented in plain language as:

1. If the checking account balance is unknown, then classify as **not likely to default**.
2. Otherwise, if the checking account balance is less than zero DM, between one and 200 DM, or greater than 200 DM and...
3. The credit history is very good or perfect, and...
4. There is more than one dependent, then classify as **likely to default**.

The numbers in parentheses indicate the number of examples meeting the criteria for that decision, and the number incorrectly classified by the decision. For instance, on the first line, (358/44) indicates that of the 358 examples reaching the decision, 44 were incorrectly classified as `no`, that is, not likely to default. In other words, 44 applicants actually defaulted in spite of the model's prediction to the contrary.



After the tree output, the `summary(credit_model)` displays a confusion matrix, which is a cross-tabulation that indicates the model's incorrectly classified records in the training data:

Evaluation on training data (900 cases):

```
Decision Tree
-----
Size      Errors
 66  125 (13.9%)  <<

(a)   (b)   <-classified as
----  ----
609    23   (a): class no
102   166   (b): class yes
```

The `Errors` field notes that the model correctly classified all but 125 of the 900 training instances for an error rate of 13.9 percent. A total of 23 actual `no` values were incorrectly classified as `yes` (false positives), while 102 `yes` values were misclassified as `no` (false negatives).

Decision trees are known for having a tendency to overfit the model to the training data. For this reason, the error rate reported on training data may be overly optimistic, and it is especially important to evaluate decision trees on a test dataset.

Step 4 – evaluating model performance

To apply our decision tree to the test dataset, we use the `predict()` function as shown in the following line of code:

```
> credit_pred <- predict(credit_model, credit_test)
```

This creates a vector of predicted class values, which we can compare to the actual class values using the `CrossTable()` function in the `gmodels` package. Setting the `prop.c` and `prop.r` parameters to `FALSE` removes the column and row percentages from the table. The remaining percentage (`prop.t`) indicates the proportion of records in the cell out of the total number of records.

```
> library(gmodels)
> CrossTable(credit_test$default, credit_pred,
             prop.chisq = FALSE, prop.c = FALSE, prop.r = FALSE,
             dnn = c('actual default', 'predicted default'))
```

This results in the following table:

actual default	predicted default		Row Total
	no	yes	
no	57 0.570	11 0.110	68
yes	16 0.160	16 0.160	32
Column Total	73	27	100

Out of the 100 test loan application records, our model correctly predicted that 57 did not default and 16 did default, resulting in an accuracy of 73 percent and an error rate of 27 percent. This is somewhat worse than its performance on the training data, but not unexpected, given that a model's performance is often worse on unseen data. Also note that the model only correctly predicted 50 percent of the 32 loan defaults in the test data. Unfortunately, this type of error is a potentially very costly mistake. Let's see if we can improve the result with a bit more effort.

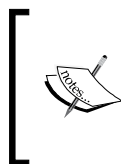
Step 5 – improving model performance

Our model's error rate is likely to be too high to deploy it in a real-time credit scoring application. In fact, if the model had predicted "no default" for every test case, it would have been correct 68 percent of the time – a result not much worse than our model, but requiring much less effort! Predicting loan defaults from 900 examples seems to be a challenging problem.

Making matters even worse, our model performed especially poorly at identifying applicants who default. Luckily, there are a couple of simple ways to adjust the C5.0 algorithm that may help to improve the performance of the model, both overall and for the more costly mistakes.

Boosting the accuracy of decision trees

One way the C5.0 algorithm improved upon the C4.5 algorithm was by adding **adaptive boosting**. This is a process in which many decision trees are built, and the trees vote on the best class for each example.



The idea of boosting is based largely upon research by Rob Schapire and Yoav Freund. For more information, try searching the web for their publications or their recent textbook: *Boosting: Foundations and Algorithms* (The MIT Press, 2012).

As boosting can be applied more generally to any machine learning algorithm, it is covered in more detail later in this book in *Chapter 11, Improving Model Performance*. For now, it suffices to say that boosting is rooted in the notion that by combining a number of weak performing learners, you can create a team that is much stronger than any one of the learners alone. Each of the models has a unique set of strengths and weaknesses, and may be better or worse at certain problems. Using a combination of several learners with complementary strengths and weaknesses can therefore dramatically improve the accuracy of a classifier.

The `C5.0()` function makes it easy to add boosting to our C5.0 decision tree. We simply need to add an additional `trials` parameter indicating the number of separate decision trees to use in the boosted team. The `trials` parameter sets an upper limit; the algorithm will stop adding trees if it recognizes that additional trials do not seem to be improving the accuracy. We'll start with 10 trials—a number that has become the de facto standard, as research suggests that this reduces error rates on test data by about 25 percent.

```
> credit_boost10 <- C5.0(credit_train[-17], credit_train$default,
                           trials = 10)
```

While examining the resulting model, we can see that some additional lines have been added indicating the changes:

```
> credit_boost10
Number of boosting iterations: 10
Average tree size: 56
```

Across the 10 iterations, our tree size shrunk. If you would like, you can see all 10 trees by typing `summary(credit_boost10)` at the command prompt.

Let's take a look at the performance on our training data:

```
> summary(credit_boost10)

      (a)   (b)   <-classified as
-----
 626      6   (a): class no
 25    243   (b): class yes
```

The classifier made 31 mistakes on 900 training examples for an error rate of 3.4 percent. This is quite an improvement over the 13.9 percent training error rate we noted before adding boosting! However, it remains to be seen whether we see a similar improvement on the test data. Let's take a look:

```
> credit_boost_pred10 <- predict(credit_boost10, credit_test)
> CrossTable(credit_test$default, credit_boost_pred10,
              prop.chisq = FALSE, prop.c = FALSE, prop.r = FALSE,
              dnn = c('actual default', 'predicted default'))
```

The resulting table is as follows:

actual default	predicted default		Row Total
	no	yes	
no	60 0.600	8 0.080	68
yes	15 0.150	17 0.170	32
Column Total	75	25	100

Here, we reduced the total error rate from 27 percent prior to boosting down to 23 percent in the boosted model. It does not seem like a large gain, but it is reasonably close to the 25 percent reduction we hoped for. On the other hand, the model is still not doing well at predicting defaults, getting $15 / 32 = 47\%$ wrong. The lack of an even greater improvement may be a function of our relatively small training dataset, or it may just be a very difficult problem to solve.

That said, if boosting can be added this easily, why not apply it by default to every decision tree? The reason is twofold. First, if building a decision tree once takes a great deal of computation time, building many trees may be computationally impractical. Secondly, if the training data is very noisy, then boosting might not result in an improvement at all. Still, if greater accuracy is needed, it's worth giving it a try.

Making some mistakes more costly than others

Giving a loan out to an applicant who is likely to default can be an expensive mistake. One solution to reduce the number of false negatives may be to reject a larger number of borderline applicants. The few years' worth of interest that the bank would earn from a risky loan is far outweighed by the massive loss it would take if the money was never paid back at all.

The C5.0 algorithm allows us to assign a penalty to different types of errors in order to discourage a tree from making more costly mistakes. The penalties are designated in a **cost matrix**, which specifies how many times more costly each error is, relative to any other. Suppose we believe that a loan default costs the bank four times as much as a missed opportunity. Our cost matrix then could be defined as:

```
> error_cost <- matrix(c(0, 1, 4, 0), nrow = 2)
```

This creates a matrix with two rows and two columns, arranged somewhat differently than the confusion matrixes we have been working with. The value 1 indicates no and the value 2 indicates yes. Rows are for predicted values and columns are for actual values:

```
> error_cost
      [,1] [,2]
[1,]    0    4
[2,]    1    0
```

As defined by this matrix, there is no cost assigned when the algorithm classifies a no or yes correctly, but a false negative has a cost of 4 versus a false positive's cost of 1. To see how this impacts classification, let's apply it to our decision tree using the costs parameter of the `C5.0()` function. We'll otherwise use the same steps as before:

```
> credit_cost <- C5.0(credit_train[-17], credit_train$default,
                      costs = error_cost)
> credit_cost_pred <- predict(credit_cost, credit_test)
> CrossTable(credit_test$default, credit_cost_pred,
             prop.chisq = FALSE, prop.c = FALSE, prop.r = FALSE,
             dnn = c('actual default', 'predicted default'))
```

This produces the following confusion matrix:

actual default	predicted default		Row Total
	no	yes	
no	42 0.420	26 0.260	68
yes	6 0.060	26 0.260	32
column Total	48	52	100

Compared to our best boosted model, this version makes more mistakes overall: 32 percent here versus 23 percent in the boosted case. However, the types of mistakes vary dramatically. Where the previous models incorrectly classified nearly half of the defaults incorrectly, in this model, only 25 percent of the defaults were predicted to be non-defaults. This trade resulting in a reduction of false negatives at the expense of increasing false positives may be acceptable if our cost estimates were accurate.