

# An Introduction to Recursive Partitioning Using the RPART Routines

Terry M. Therneau  
Elizabeth J. Atkinson  
Mayo Foundation

September 3, 1997

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Notation</b>	<b>4</b>
<b>3</b>	<b>Building the tree</b>	<b>5</b>
3.1	Splitting criteria . . . . .	5
3.2	Incorporating losses . . . . .	7
3.2.1	Generalized Gini index . . . . .	7
3.2.2	Altered priors . . . . .	8
3.3	Example: Stage C prostate cancer ( <code>class</code> method) . . . . .	9
<b>4</b>	<b>Pruning the tree</b>	<b>12</b>
4.1	Definitions . . . . .	12
4.2	Cross-validation . . . . .	13
4.3	Example: The Stochastic Digit Recognition Problem . . . . .	14
<b>5</b>	<b>Missing data</b>	<b>17</b>
5.1	Choosing the split . . . . .	17
5.2	Surrogate variables . . . . .	17
5.3	Example: Stage C prostate cancer (cont.) . . . . .	18
<b>6</b>	<b>Further options</b>	<b>20</b>
6.1	Program options . . . . .	20
6.2	Example: Consumer Report Auto Data . . . . .	22
6.3	Example: Kyphosis data . . . . .	26

<b>7</b>	<b>Regression</b>	<b>28</b>
7.1	Definition . . . . .	28
7.2	Example: Consumer Report Auto data (cont.) . . . . .	29
7.3	Example: Stage C prostate cancer ( <code>anova</code> method) . . . . .	33
<b>8</b>	<b>Poisson regression</b>	<b>35</b>
8.1	Definition . . . . .	35
8.2	Improving the method . . . . .	36
8.3	Example: solder data . . . . .	37
8.4	Example: Stage C prostate cancer ( <code>survival</code> method) . . . . .	41
<b>9</b>	<b>Plotting options</b>	<b>46</b>
<b>10</b>	<b>Other functions</b>	<b>49</b>
<b>11</b>	<b>Relation to other programs</b>	<b>50</b>
11.1	CART . . . . .	50
11.2	Tree . . . . .	51
<b>12</b>	<b>Source</b>	<b>52</b>

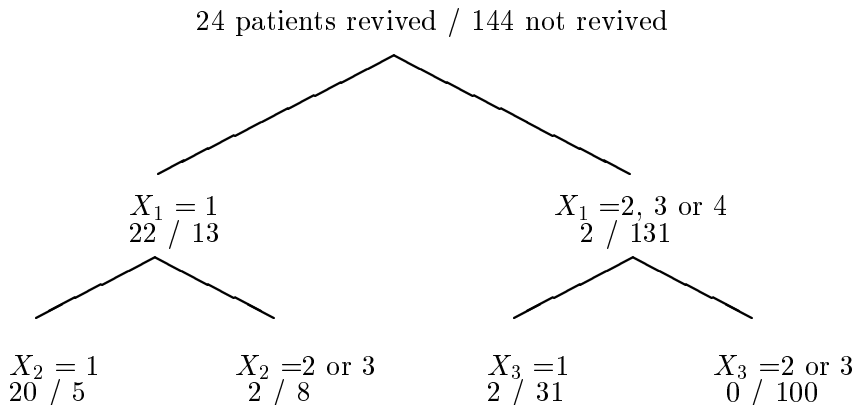


Figure 1: Revival data

## 1 Introduction

This document is an update of a technical report written several years ago at Stanford [5], and is intended to give a short overview of the methods found in the **rpart** routines, which implement many of the ideas found in the CART (Classification and Regression Trees) book and programs of Breiman, Friedman, Olshen and Stone [1]. Because CART is the trademarked name of a particular software implementation of these ideas, and *tree* has been used for the S-plus routines of Clark and Pregibon~[2] a different acronym — Recursive PARTitioning or **rpart** — was chosen.

The **rpart** programs build classification or regression models of a very general structure using a two stage procedure; the resulting models can be represented as binary trees. An example is some preliminary data gathered at Stanford on revival of cardiac arrest patients by paramedics. The goal is to predict which patients are revivable in the field on the basis of fourteen variables known at or near the time of paramedic arrival, e.g., sex, age, time from arrest to first care, etc. Since some patients who are not revived on site are later successfully resuscitated at the hospital, early identification of these “recalcitrant” cases is of considerable clinical interest.

The resultant model separated the patients into four groups as shown in figure 1, where

$X_1$  = initial heart rhythm  
 1= VF/VT, 2=EMD, 3=Asystole, 4=Other

$X_2$  = initial response to defibrillation

1=Improved, 2=No change, 3=Worse

$X_3$  = initial response to drugs

1=Improved, 2=No change, 3=Worse

The other 11 variables did not appear in the final model. This procedure seems to work especially well for variables such as  $X_1$ , where there is a definite ordering, but spacings are not necessarily equal.

The tree is built by the following process: first the single variable is found which best splits the data into two groups ('best' will be defined later). The data is separated, and then this process is applied *separately* to each sub-group, and so on recursively until the subgroups either reach a minimum size (5 for this data) or until no improvement can be made.

The resultant model is, with certainty, too complex, and the question arises as it does with all stepwise procedures of when to stop. The second stage of the procedure consists of using cross-validation to trim back the full tree. In the medical example above the full tree had ten terminal regions. A cross validated estimate of risk was computed for a nested set of subtrees; this final model, presented in figure 1, is the subtree with the lowest estimate of risk.

## 2 Notation

The partitioning method can be applied to many different kinds of data. We will start by looking at the classification problem, which is one of the more instructive cases (but also has the most complex equations). The sample population consists of  $n$  observations from  $C$  classes. A given model will break these observations into  $k$  terminal groups; to each of these groups is assigned a predicted class (this will be the response variable). In an actual application, most parameters will be estimated from the data, such estimates are given by  $\approx$  formulae.

$\pi_i$                      $i = 1, 2, \dots, C$             Prior probabilities of each class.

$L(i, j)$              $i = 1, 2, \dots, C$     Loss matrix for incorrectly classifying  
an  $i$  as a  $j$ .  $L(i, i) \equiv 0$ .

$A$                     Some node of the tree.  
Note that  $A$  represents both a set of individuals in  
the sample data, and, via the tree that produced it,  
a classification rule for future data.

$\tau(x)$	True class of an observation $x$ , where $x$ is the vector of predictor variables.
$\tau(A)$	The class assigned to $A$ , if $A$ were to be taken as a final node.
$n_i, n_A$	Number of observations in the sample that are class $i$ , number of obs in node $A$ .
$n_{iA}$	Number of observations in the sample that are class $i$ and node $A$ .
$P(A)$	Probability of $A$ (for future observations). $= \sum_{i=1}^C \pi_i P\{x \in A \mid \tau(x) = i\}$ $\approx \sum_{i=1}^C \pi_i n_{iA}/n_i$
$p(i A)$	$P\{\tau(x) = i \mid x \in A\}$ (for future observations) $= \pi_i P\{x \in A \mid \tau(x) = i\} / P\{x \in A\}$ $\approx \pi_i (n_{iA}/n_i) / \sum \pi_i (n_{iA}/n_i)$
$R(A)$	Risk of $A$ $= \sum_{i=1}^C p(i A) L(i, \tau(A))$ where $\tau(A)$ is chosen to minimize this risk.
$R(T)$	Risk of a model (or tree) $T$ $= \sum_{j=1}^k P(A_j) R(A_j)$ where $A_j$ are the terminal nodes of the tree.

If  $L(i, j) = 1$  for all  $i \neq j$ , and we set the prior probabilities  $\pi$  equal to the observed class frequencies in the sample then  $p(i|A) = n_{iA}/n_A$  and  $R(T)$  is the proportion misclassified.

## 3 Building the tree

### 3.1 Splitting criteria

If we split a node  $A$  into two sons  $A_L$  and  $A_R$  (left and right sons), we will have

$$P(A_L)R(A_L) + P(A_R)R(A_R) \leq P(A)R(A)$$

(this is proven in [1]). Using this, one obvious way to build a tree is to choose that split which maximizes  $\Delta R$ , the decrease in risk. There are defects with this, however, as the following example shows.

Suppose losses are equal and that the data is 80% class 1's, and that some trial split results in  $A_L$  being 54% class 1's and  $A_R$  being 100% class 1's. Class 1's versus class 0's are the outcome variable in this example. Since the minimum risk prediction for both the left and right son is  $\tau(A_L) = \tau(A_R) = 1$ , this split will have  $\Delta R = 0$ , yet scientifically this is a very informative division of the sample. In real data with such a majority, the first few splits very often can do no better than this.

A more serious defect with maximizing  $\Delta R$  is that the risk reduction is essentially linear. If there were two competing splits, one separating the data into groups of 85% and 50% purity respectively, and the other into 70%-70%, we would usually prefer the former, if for no other reason than because it better sets things up for the next splits.

One way around both of these problems is to use lookahead rules; but these are computationally very expensive. Instead **rpart** uses one of several measures of impurity, or diversity, of a node. Let  $f$  be some impurity function and define the impurity of a node  $A$  as

$$I(A) = \sum_{i=1}^C f(p_{iA})$$

where  $p_{iA}$  is the proportion of those in  $A$  that belong to class  $i$  for future samples. Since we would like  $I(A) = 0$  when  $A$  is pure,  $f$  must be concave with  $f(0) = f(1) = 0$ .

Two candidates for  $f$  are the information index  $f(p) = -p \log(p)$  and the Gini index  $f(p) = p(1 - p)$ . We then use that split with maximal impurity reduction

$$\Delta I = p(A)I(A) - p(A_L)I(A_L) - p(A_R)I(A_R)$$

The two impurity functions are plotted in figure (2), with the second plot scaled so that the maximum for both measures is at 1. For the two class problem the measures differ only slightly, and will nearly always choose the same split point.

Another convex criteria not quite of the above class is twoing for which

$$I(A) = \min_{C_1 C_2} [f(p_{C_1}) + f(p_{C_2})]$$

where  $C_1, C_2$  is some partition of the  $C$  classes into two disjoint sets. If  $C = 2$  twoing is equivalent to the usual impurity index for  $f$ . Surprisingly, twoing can be calculated almost as efficiently as the usual impurity index. One potential advantage of twoing

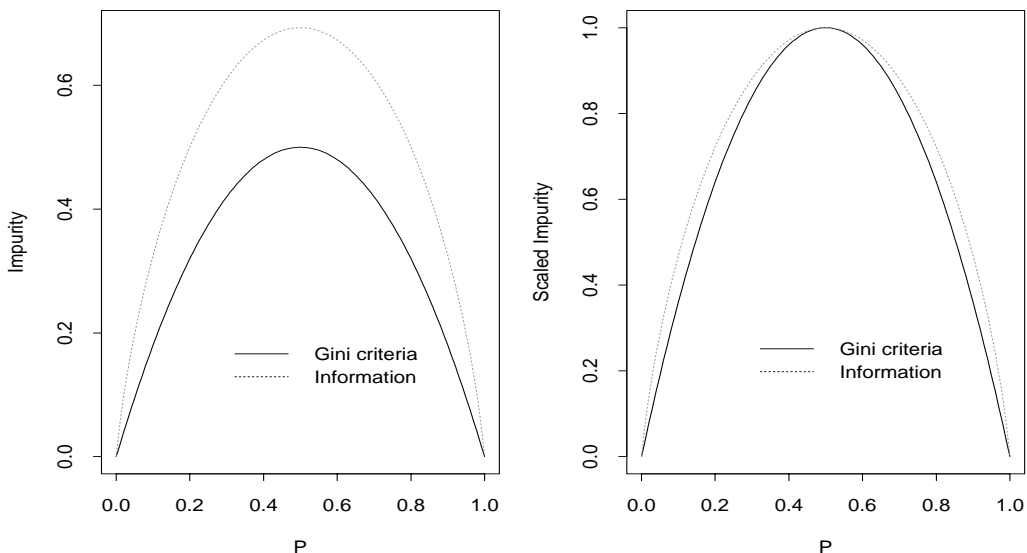


Figure 2: Comparison of Gini and Information indices

is that the output may give the user additional insight concerning the structure of the data. It can be viewed as the partition of  $C$  into two superclasses which are in some sense the most dissimilar for those observations in  $A$ . For certain problems there may be a natural ordering of the response categories (e.g. level of education), in which case ordered twoing can be naturally defined, by restricting  $C_1$  to be an interval  $[1, 2, \dots, k]$  of classes. Twoing is not part of **rpart**.

## 3.2 Incorporating losses

One salutatory aspect of the risk reduction criteria not found in the impurity measures is inclusion of the loss function. Two different ways of extending the impurity criteria to also include losses are implemented in CART, the generalized Gini index and altered priors. The **rpart** software implements only the altered priors method.

### 3.2.1 Generalized Gini index

The Gini index has the following interesting interpretation. Suppose an object is selected at random from one of  $C$  classes according to the probabilities  $(p_1, p_2, \dots, p_C)$  and is randomly assigned to a class using the same distribution. The probability of

misclassification is

$$\sum_i \sum_{j \neq i} p_i p_j = \sum_i \sum_j p_i p_j - \sum_i p_i^2 = \sum_i 1 - p_i^2 = \text{Gini index for } p$$

Let  $L(i, j)$  be the loss of assigning class  $j$  to an object which actually belongs to class  $i$ . The expected cost of misclassification is  $\sum_i \sum_j L(i, j) p_i p_j$ . This suggests defining a *generalized Gini* index of impurity by

$$G(p) = \sum_i \sum_j L(i, j) p_i p_j$$

The corresponding splitting criterion appears to be promising for applications involving variable misclassification costs. But there are several reasonable objections to it. First,  $G(p)$  is not necessarily a concave function of  $p$ , which was the motivating factor behind impurity measures. More seriously,  $G$  symmetrizes the loss matrix before using it. To see this note that

$$G(p) = (1/2) \sum_i \sum_j [L(i, j) + L(j, i)] p_i p_j$$

In particular, for two-class problems,  $G$  in effect ignores the loss matrix.

### 3.2.2 Altered priors

Remember the definition of  $R(A)$

$$\begin{aligned} R(A) &\equiv \sum_{i=1}^C p_{iA} L(i, \tau(A)) \\ &= \sum_{i=1}^C \pi_i L(i, \tau(A)) (n_{iA}/n_i) (n/n_A) \end{aligned}$$

Assume there exists  $\tilde{\pi}$  and  $\tilde{L}$  be such that

$$\tilde{\pi}_i \tilde{L}(i, j) = \pi_i L(i, j) \quad \forall i, j \in C$$

Then  $R(A)$  is unchanged under the new losses and priors. If  $\tilde{L}$  is proportional to the zero-one loss matrix then the priors  $\tilde{\pi}$  should be used in the splitting criteria. This is possible only if  $L$  is of the form

$$L(i, j) = \begin{cases} L_i & i \neq j \\ 0 & i = j \end{cases}$$



in which case

$$\tilde{\pi}_i = \frac{\pi_i L_i}{\sum_j \pi_j L_j}$$

This is always possible when  $C = 2$ , and hence altered priors are exact for the two class problem. For arbitrary loss matrix of dimension  $C > 2$ , **rpart** uses the above formula with  $L_i = \sum_j L(i, j)$ .

A second justification for altered priors is this. An impurity index  $I(A) = \sum f(p_i)$  has its maximum at  $p_1 = p_2 = \dots = p_C = 1/C$ . If a problem had, for instance, a misclassification loss for class 1 which was twice the loss for a class 2 or 3 observation, one would wish  $I(A)$  to have its maximum at  $p_1 = 1/5$ ,  $p_2 = p_3 = 2/5$ , since this is the worst possible set of proportions on which to decide a node's class. The altered priors technique does exactly this, by shifting the  $p_i$ .

Two final notes

- When altered priors are used, they affect only the choice of split. The ordinary losses and priors are used to compute the risk of the node. The altered priors simply help the impurity rule choose splits that are likely to be “good” in terms of the risk.
- The argument for altered priors is valid for both the gini and information splitting rules.

### 3.3 Example: Stage C prostate cancer (`class` method)

This first example is based on a data set of 146 stage C prostate cancer patients [4]. The main clinical endpoint of interest is whether the disease recurs after initial surgical removal of the prostate, and the time interval to that progression (if any). The endpoint in this example is **status**, which takes on the value 1 if the disease has progressed and 0 if not. Later we'll analyze the data using the exponential (**exp**) method, which will take into account time to progression. A short description of each of the variables is listed below. The main predictor variable of interest in this study was DNA ploidy, as determined by flow cytometry. For diploid and tetraploid tumors, the flow cytometric method was also able to estimate the percent of tumor cells in a  $G_2$  (growth) stage of their cell cycle;  $G_2\%$  is systematically missing for most aneuploid tumors.

The variables in the data set are

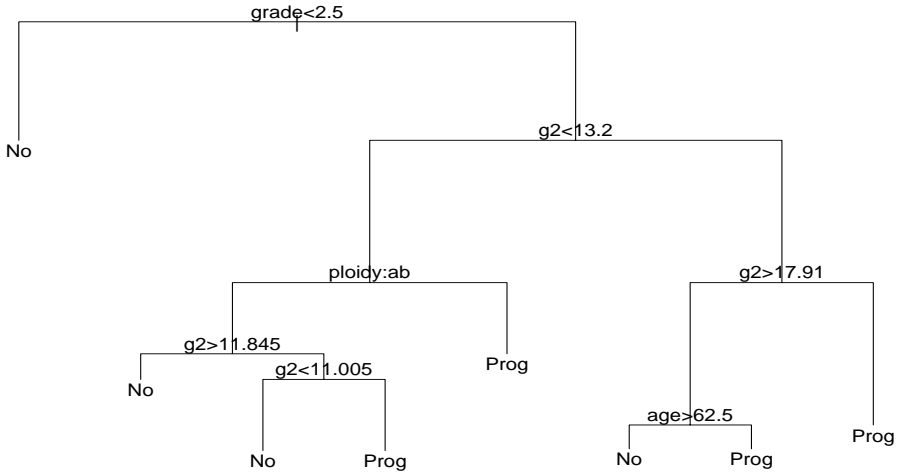


Figure 3: Classification tree for the Stage C data

pgtime    time to progression, or last follow-up free of progression  
 pgstat    status at last follow-up (1=progressed, 0=censored)  
 age       age at diagnosis  
 eet       early endocrine therapy (1=no, 0=yes)  
 ploidy    diploid/tetraploid/aneuploid DNA pattern  
 g2        % of cells in  $G_2$  phase  
 grade     tumor grade (1-4)  
 gleason   Gleason grade (3-10)

The model is fit by using the **rpart** function. The first argument of the function is a model formula, with the  $\sim$  symbol standing for “is modeled as”. The **print** function gives an abbreviated output, as for other S models. The **plot** and **text** command plot the tree and then label the plot, the result is shown in figure 3.

```

> progstat <- factor(stagec$pgstat, levels=0:1, labels=c("No", "Prog"))
> cfit <- rpart(progstat ~ age + eet + g2 + grade + gleason + ploidy,
  data=stagec, method='class')
> print(cfit)

node), split, n, loss, yval, (yprob)
  * denotes terminal node

1) root 146 54 No ( 0.6301 0.3699 )

```

```

2) grade<2.5 61 9 No ( 0.8525 0.1475 ) *
3) grade>2.5 85 40 Prog ( 0.4706 0.5294 )
6) g2<13.2 40 17 No ( 0.5750 0.4250 )
  12) ploidy:diploid,tetraploid 31 11 No ( 0.6452 0.3548 )
    24) g2>11.845 7 1 No ( 0.8571 0.1429 ) *
    25) g2<11.845 24 10 No ( 0.5833 0.4167 )
      50) g2<11.005 17 5 No ( 0.7059 0.2941 ) *
      51) g2>11.005 7 2 Prog ( 0.2857 0.7143 ) *
  13) ploidy:aneuploid 9 3 Prog ( 0.3333 0.6667 ) *
7) g2>13.2 45 17 Prog ( 0.3778 0.6222 )
  14) g2>17.91 22 8 No ( 0.6364 0.3636 )
    28) age>62.5 15 4 No ( 0.7333 0.2667 ) *
    29) age<62.5 7 3 Prog ( 0.4286 0.5714 ) *
  15) g2<17.91 23 3 Prog ( 0.1304 0.8696 ) *

```

```

> plot(cfit)
> text(cfit)

```

- The creation of a labeled factor variable as the response improves the labeling of the printout.
- We have explicitly directed the routine to treat `prognosis` as a categorical variable by asking for `method='class'`. (Since `prognosis` is a factor this would have been the default choice). Since no optional classification parameters are specified the routine will use the Gini rule for splitting, prior probabilities that are proportional to the observed data frequencies, and 0/1 losses.
- The child nodes of node  $x$  are always numbered  $2x$  (left) and  $2x + 1$  (right), to help in navigating the printout (compare the printout to figure 3).
- Other items in the list are the definition of the variable and split used to create a node, the number of subjects at the node, the loss or error at the node (for this example, with proportional priors and unit losses this will be the number misclassified), the classification of the node, and the predicted class for the node.
- \* indicates that the node is terminal.
- Grades 1 and 2 go to the left, grades 3 and 4 go to the right. The tree is arranged so that the branches with the largest “average class” go to the right.

## 4 Pruning the tree

### 4.1 Definitions

We have built a complete tree, possibly quite large and/or complex, and must now decide how much of that model to retain. In forward stepwise regression, for instance, this issue is addressed sequentially and no additional variables are added when the F-test for the remaining variables fails to achieve some level  $\alpha$ .

Let  $T_1, T_2, \dots, T_k$  be the terminal nodes of a tree  $T$ . Define

$$\begin{aligned} |T| &= \text{number of terminal nodes} \\ \text{risk of } T &= R(T) = \sum_{i=1}^k P(T_i)R(T_i) \end{aligned}$$

In comparison to regression,  $|T|$  is analogous to the model degrees of freedom and  $R(T)$  to the residual sum of squares.

Now let  $\alpha$  be some number between 0 and  $\infty$  which measures the 'cost' of adding another variable to the model;  $\alpha$  will be called a complexity parameter. Let  $R(T_0)$  be the risk for the zero split tree. Define

$$R_\alpha(T) = R(T) + \alpha|T|$$

to be the cost for the tree, and define  $T_\alpha$  to be that subtree of the full model which has minimal cost. Obviously  $T_0$  = the full model and  $T_\infty$  = the model with no splits at all. The following results are shown in [1].

1. If  $T_1$  and  $T_2$  are subtrees of  $T$  with  $R_\alpha(T_1) = R_\alpha(T_2)$ , then either  $T_1$  is a subtree of  $T_2$  or  $T_2$  is a subtree of  $T_1$ ; hence either  $|T_1| < |T_2|$  or  $|T_2| < |T_1|$ .
2. If  $\alpha > \beta$  then either  $T_\alpha = T_\beta$  or  $T_\alpha$  is a strict subtree of  $T_\beta$ .
3. Given some set of numbers  $\alpha_1, \alpha_2, \dots, \alpha_m$ ; both  $T_{\alpha_1}, \dots, T_{\alpha_m}$  and  $R(T_{\alpha_1}), \dots, R(T_{\alpha_m})$  can be computed efficiently.

Using the first result, we can uniquely define  $T_\alpha$  as the smallest tree  $T$  for which  $R_\alpha(T)$  is minimized.

Since any sequence of nested trees based on  $T$  has at most  $|T|$  members, result 2 implies that all possible values of  $\alpha$  can be grouped into  $m$  intervals,  $m \leq |T|$

$$\begin{aligned} I_1 &= [0, \alpha_1] \\ I_2 &= (\alpha_1, \alpha_2] \\ &\vdots \\ I_m &= (\alpha_{m-1}, \infty] \end{aligned}$$

where all  $\alpha \in I_i$  share the same minimizing subtree.

## 4.2 Cross-validation

Cross-validation is used to choose a best value for  $\alpha$  by the following steps:

1. Fit the full model on the data set  
 compute  $I_1, I_2, \dots, I_m$   
 set  $\beta_1 = 0$   
 $\beta_2 = \sqrt{\alpha_1 \alpha_2}$   
 $\beta_3 = \sqrt{\alpha_2 \alpha_3}$   
 $\vdots$   
 $\beta_{m-1} = \sqrt{\alpha_{m-2} \alpha_{m-1}}$   
 $\beta_m = \infty$   
 each  $\beta_i$  is a ‘typical value’ for its  $I_i$
2. Divide the data set into  $s$  groups  $G_1, G_2, \dots, G_s$  each of size  $s/n$ , and for each group separately:
  - fit a full model on the data set ‘everyone except  $G_i$ ’ and determine  $T_{\beta_1}, T_{\beta_2}, \dots, T_{\beta_m}$  for this reduced data set,
  - compute the predicted class for each observation in  $G_i$ , under each of the models  $T_{\beta_j}$  for  $1 \leq j \leq m$ ,
  - from this compute the risk for each subject.
3. Sum over the  $G_i$  to get an estimate of risk for each  $\beta_j$ . For that  $\beta$  (complexity parameter) with smallest risk compute  $T_\beta$  for the full data set, this is chosen as the best trimmed tree.

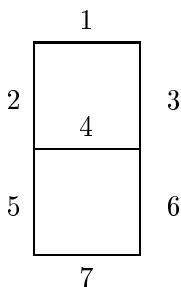
In actual practice, we may use instead the 1-SE rule. A plot of  $\beta$  versus risk often has an initial sharp drop followed by a relatively flat plateau and then a slow rise. The choice of  $\beta$  among those models on the plateau can be essentially random. To avoid this, both an estimate of the risk and its standard error are computed during the cross-validation. Any risk within one standard error of the achieved minimum is marked as being equivalent to the minimum (i.e. considered to be part of the flat plateau). Then the simplest model, among all those “tied” on the plateau, is chosen.

In the usual definition of cross-validation we would have taken  $s = n$  above, i.e., each of the  $G_i$  would contain exactly one observation, but for moderate  $n$  this is computationally prohibitive. A value of  $s = 10$  has been found to be sufficient, but users can vary this if they wish.

In Monte-Carlo trials, this method of pruning has proven very reliable for screening out ‘pure noise’ variables in the data set.

### 4.3 Example: The Stochastic Digit Recognition Problem

This example is found in section 2.6 of [1], and used as a running example throughout much of their book. Consider the segments of an unreliable digital readout



where each light is correct with probability 0.9, e.g., if the true digit is a 2, the lights 1, 3, 4, 5, and 7 are on with probability 0.9 and lights 2 and 6 are on with probability 0.1. Construct test data where  $Y \in \{0, 1, \dots, 9\}$ , each with proportion 1/10 and the  $X_i$ ,  $i = 1, \dots, 7$  are i.i.d. Bernoulli variables with parameter depending on  $Y$ .  $X_8 - X_{24}$  are generated as i.i.d. Bernoulli  $P\{X_i = 1\} = .5$ , and are independent of  $Y$ . They correspond to imbedding the readout in a larger rectangle of random lights.

A sample of size 200 was generated accordingly and the procedure applied using the gini index (see 3.2.1) to build the tree. The S-plus code to compute the simulated data and the fit are shown below.

```
> n <- 200
> temp <- c(1,1,1,0,1,1,1,
            0,0,1,0,0,1,0,
            1,0,1,1,1,0,1,
            1,0,1,1,0,1,1,
            0,1,1,1,0,1,0,
            1,1,0,1,0,1,1,
            0,1,0,1,1,1,1,
            1,0,1,0,0,1,0,
            1,1,1,1,1,1,1,
            1,1,1,1,0,1,0)

> lights <- matrix(temp, 10, 7, byrow=T) # The true light pattern 0-9
> temp1 <- matrix(rbinom(n*7, 1, .9), n, 7) # Noisy lights
> temp1 <- ifelse(lights[y+1, ]==1, temp1, 1-temp1)
> temp2 <- matrix(rbinom(n*17, 1, .5), n, 17) #Random lights
> x <- cbind(temp1, temp2) #x is the matrix of predictors
```

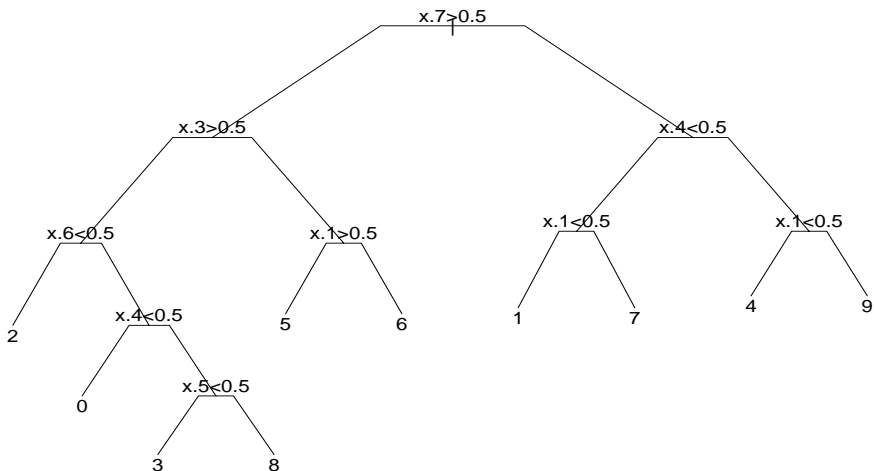


Figure 4: Optimally pruned tree for the stochastic digit recognition data

```
> y <- rep(0:9, length=200)
```

The particular data set of this example can be replicated by setting `.Random.seed` to `c(21, 14, 49, 32, 43, 1, 32, 22, 36, 23, 28, 3)` before the call to `rbinom`. Now we fit the model:

```
> temp3 <- rpart.control(xval=10, minbucket=2, minsplit=4, cp=0)
> dfit <- rpart(y ~ x, method='class', control=temp3)
> printcp(dfit)
```

Classification tree:

```
rpart(formula = y ~ x, method = "class", control = temp3)
```

Variables actually used in tree construction:

```
[1] x.1 x.10 x.12 x.13 x.15 x.19 x.2 x.20 x.22 x.3 x.4 x.5 x.6 x.7 x
```

Root node error: 180/200 = 0.9

	CP	nsplit	rel error	xerror	xstd
1	0.1055556	0	1.00000	1.09444	0.0095501
2	0.0888889	2	0.79444	1.01667	0.0219110
3	0.0777778	3	0.70556	0.90556	0.0305075
4	0.0666667	5	0.55556	0.75000	0.0367990

5	0.0555556	8	0.361111	0.561111	0.0392817
6	0.0166667	9	0.305556	0.361111	0.0367990
7	0.0111111	11	0.272222	0.37778	0.0372181
8	0.0083333	12	0.261111	0.361111	0.0367990
9	0.0055556	16	0.22778	0.355556	0.0366498
10	0.0027778	27	0.16667	0.34444	0.0363369
11	0.0013889	31	0.15556	0.36667	0.0369434
12	0.0000000	35	0.15000	0.36667	0.0369434

```
> fit9 <- prune(dfit, cp=.02)
> plot(fit9, branch=.3, compress=T)
> text(fit9)
```

The `cp` table differs from that in section 3.5 of [1] in several ways, the last two of which are somewhat important.

- The actual values are different, of course, because of different random number generators in the two runs.
- The table is printed from the smallest tree (no splits) to the largest one (35 splits). We find it easier to compare one tree to another when they start at the same place.
- The number of splits is listed, rather than the number of nodes. The number of nodes is always  $1 +$  the number of splits.
- For easier reading, the error columns have been scaled so that the first node has an error of 1. Since in this example the model with no splits must make 180/200 misclassifications, multiply columns 3-5 by 180 to get a result in terms of absolute error. (Computations are done on the absolute error scale, and printed on relative scale).
- The complexity parameter column (`cp`) has been similarly scaled.

Looking at the `cp` table, we see that the best tree has 10 terminal nodes (9 splits), based on cross-validation (using 1-SE rule of  $0.3444 + 0.0363369$ ). This subtree is extracted with call to `prune` and saved in `fit9`. The pruned tree is shown in figure 4. Two options have been used in the plot. The `compress` option tries to narrow the printout by vertically overlapping portions of the plot. (It has only a small effect on this particular dendrogram). The `branch` option controls the shape of the branches that connect a node to its children. The section on plotting (9) will discuss this and other options in more detail.

The largest tree, with 36 terminal nodes, correctly classifies  $170/200 = 85\%$  ( $1 - 0.15$ ) of the observations, but uses several of the random predictors in doing



so and seriously overfits the data. If the number of observations per terminal node (minbucket) had been set to 1 instead of 2, then every observation would be classified correctly in the final model, many in terminal nodes of size 1.

## 5 Missing data

### 5.1 Choosing the split

Missing values are one of the curses of statistical models and analysis. Most procedures deal with them by refusing to deal with them – incomplete observations are tossed out. **Rpart** is somewhat more ambitious. Any observation with values for the dependent variable and at least one independent variable will participate in the modeling.

The quantity to be maximized is still

$$\Delta I = p(A)I(A) - p(A_L)I(A_L) - p(A_R)I(A_R)$$

The leading term is the same for all variables and splits irrespective of missing data, but the right two terms are somewhat modified. Firstly, the impurity indices  $I(A_R)$  and  $I(A_L)$  are calculated only over the observations which are not missing a particular predictor. Secondly, the two probabilities  $p(A_L)$  and  $p(A_R)$  are also calculated only over the relevant observations, but they are then adjusted so that they sum to  $p(A)$ . This entails some extra bookkeeping as the tree is built, but ensures that the terminal node probabilities sum to 1.

In the extreme case of a variable for which only 2 observations are non-missing, the impurity of the two sons will both be zero when splitting on that variable. Hence  $\Delta I$  will be maximal, and this ‘almost all missing’ coordinate is guaranteed to be chosen as best; the method is certainly flawed in this extreme case. It is difficult to say whether this bias toward missing coordinates carries through to the non-extreme cases, however, since a more complete variable also affords for itself more possible values at which to split.

### 5.2 Surrogate variables

Once a splitting variable and a split point for it have been decided, what *is* to be done with observations missing that variable? One approach is to estimate the missing datum using the other independent variables; **rpart** uses a variation of this to define *surrogate* variables.

As an example, assume that the split ( $\text{age} < 40$ ,  $\text{age} \geq 40$ ) has been chosen. The surrogate variables are found by re-applying the partitioning algorithm (without

recursion) to predict the two categories ‘age <40’ vs. ‘age ≥40’ using the other independent variables.

For each predictor an optimal split point and a misclassification error are computed. (Losses and priors do not enter in — none are defined for the age groups — so the risk is simply  $\# \text{misclassified} / n$ .) Also evaluated is the blind rule ‘go with the majority’ which has misclassification error  $\min(p, 1 - p)$  where

$$p = (\# \text{ in } A \text{ with age} < 40) / n_A.$$

The surrogates are ranked, and any variables which do no better than the blind rule are discarded from the list.

Assume that the majority of subjects have age  $\leq 40$  and that there is another variable  $x$  which is uncorrelated to age; however, the subject with the largest value of  $x$  is also over 40 years of age. Then the surrogate variable  $x < \max$  versus  $x \geq \max$  will have one less error than the blind rule, sending 1 subject to the right and  $n - 1$  to the left. A continuous variable that is completely unrelated to age has probability  $1 - p^2$  of generating such a trim-one-end surrogate by chance alone. For this reason the **rpart** routines impose one more constraint during the construction of the surrogates: a candidate split must send at least 2 observations to the left and at least 2 to the right.

Any observation which is missing the split variable is then classified using the first surrogate variable, or if missing that, the second surrogate is used, and etc. If an observation is missing all the surrogates the blind rule is used. Other strategies for these ‘missing everything’ observations can be convincingly argued, but there should be few or no observations of this type (we hope).

### 5.3 Example: Stage C prostate cancer (cont.)

Let us return to the stage C prostate cancer data of the earlier example. For a more detailed listing of the **rpart** object, we use the **summary** function. It includes the information from the CP table (not repeated below), plus information about each node. It is easy to print a subtree based on a different **cp** value using the **cp** option. Any value between 0.0555 and 0.1049 would produce the same result as is listed below, that is, the tree with 3 splits. Because the **printout** is long, the **file** option of **summary.rpart** is often useful.

```
> printcp(fit)
```

```
Classification tree:
```

```
rpart(formula = prostat ~ age + eet + g2 + grade + gleason + ploidy,
      data = stagec)
```

Variables actually used in tree construction:

[1] age g2 grade ploidy

Root node error: 54/146 = 0.36986

	CP	nsplit	rel error	xerror	xstd
1	0.104938	0	1.00000	1.0000	0.10802
2	0.055556	3	0.68519	1.1852	0.11103
3	0.027778	4	0.62963	1.0556	0.10916
4	0.018519	6	0.57407	1.0556	0.10916
5	0.010000	7	0.55556	1.0556	0.10916

> summary(cfit,cp=.06)

Node number 1: 146 observations, complexity param=0.1049

predicted class= No expected loss= 0.3699

class counts: 92 54

probabilities: 0.6301 0.3699

left son=2 (61 obs) right son=3 (85 obs)

Primary splits:

grade < 2.5 to the left, improve=10.360, (0 missing)

gleason < 5.5 to the left, improve= 8.400, (3 missing)

ploidy splits as LRR, improve= 7.657, (0 missing)

g2 < 13.2 to the left, improve= 7.187, (7 missing)

age < 58.5 to the right, improve= 1.388, (0 missing)

Surrogate splits:

gleason < 5.5 to the left, agree=0.8630, (0 split)

ploidy splits as LRR, agree=0.6438, (0 split)

g2 < 9.945 to the left, agree=0.6301, (0 split)

age < 66.5 to the right, agree=0.5890, (0 split)

Node number 2: 61 observations

predicted class= No expected loss= 0.1475

class counts: 52 9

probabilities: 0.8525 0.1475

Node number 3: 85 observations, complexity param=0.1049

predicted class= Prog expected loss= 0.4706

class counts: 40 45

probabilities: 0.4706 0.5294

left son=6 (40 obs) right son=7 (45 obs)

Primary splits:

g2 < 13.2 to the left, improve=2.1780, (6 missing)

ploidy splits as LRR, improve=1.9830, (0 missing)

age < 56.5 to the right, improve=1.6600, (0 missing)

```

gleason < 8.5    to the left,  improve=1.6390, (0 missing)
eet    < 1.5    to the right, improve=0.1086, (1 missing)
Surrogate splits:
ploidy splits as LRL, agree=0.9620, (6 split)
age    < 68.5   to the right, agree=0.6076, (0 split)
gleason < 6.5   to the left,  agree=0.5823, (0 split)
.
.
.
```

- There are 54 progressions (class 1) and 92 non-progressions, so the first node has an expected loss of  $54/146 \approx 0.37$ . (The computation is this simple only for the default priors and losses).
- Grades 1 and 2 go to the left, grades 3 and 4 to the right. The tree is arranged so that the “more severe” nodes go to the right.
- The improvement is  $n$  times the change in impurity index. In this instance, the largest improvement is for the variable `grade`, with an improvement of 10.36. The next best choice is Gleason score, with an improvement of 8.4. The actual values of the improvement are not so important, but their relative size gives an indication of the comparative utility of the variables.
- Ploidy is a categorical variable, with values of diploid, tetraploid, and aneuploid, in that order. (To check the order, type `table(stagec$ploidy)`). All three possible splits were attempted: aneuploid+diploid vs. tetraploid, aneuploid+tetraploid vs. diploid, and aneuploid vs. diploid + tetraploid. The best split sends diploid to the right and the others to the left (node 6, see figure (3)).
- For node 3, the primary split variable is missing on 6 subjects. All 6 are split based on the first surrogate, ploidy. Diploid and aneuploid tumors are sent to the left, tetraploid to the right.

	g2 < 13.2	g2 > 13.2	NA
Diploid/aneuploid	33	2	5
Tetraploid	1	43	1

## 6 Further options

### 6.1 Program options

The central fitting function is `rpart`, whose main arguments are

- **formula**: the model formula, as in `lm` and other S model fitting functions. The right hand side may contain both continuous and categorical (factor) terms. If the outcome  $y$  has more than two levels, then categorical predictors must be fit by exhaustive enumeration, which can take a very long time.
- **data**, **weights**, **subset**: as for other S models. Weights are not yet supported, and will be ignored if present.
- **method**: the type of splitting rule to use. Options at this point are classification, anova, Poisson, and exponential.
- **parms**: a list of method specific optional parameters. For classification, the list can contain any of: the vector of prior probabilities (component **prior**), the loss matrix (component **loss**) or the splitting index (component **split**). The priors must be positive and sum to 1. The loss matrix must have zeros on the diagonal and positive off-diagonal elements. The splitting index can be ‘gini’ or ‘information’.
- **na.action**: the action for missing values. The default action for `rpart` is `na.rpart`, this default is not overridden by the `options(na.action)` global option. The default action removes only those rows for which either the response  $y$  or *all* of the predictors are missing. This ability to retain partially missing observations is perhaps the single most useful feature of `rpart` models.
- **control**: a list of control parameters, usually the result of the `rpart.control` function. The list must contain
  - **minsplit**: The minimum number of observations in a node for which the routine will even try to compute a split. The default is 20. This parameter can save computation time, since smaller nodes are almost always pruned away by cross-validation.
  - **minbucket**: The minimum number of observations in a terminal node. This defaults to `minsplit/3`.
  - **maxcompete**: It is often useful in the printout to see not only the variable that gave the best split at a node, but also the second, third, etc best. This parameter controls the number that will be printed. It has no effect on computational time, and a small effect on the amount of memory used. The default is 5.
  - **xval**: The number of cross-validations to be done. Usually set to zero during exploratory phases of the analysis. A value of 10, for instance, increases the compute time to 11-fold over a value of 0.

- **maxsurrogate**: The maximum number of surrogate variables to retain at each node. (No surrogate that does worse than “go with the majority” is printed or used). Setting this to zero will cut the computation time in half, and set **usesurrogate** to zero. The default is 5. Surrogates give different information than competitor splits. The competitor list asks “which other splits would have as many correct classifications”, surrogates ask “which other splits would classify the same subjects in the same way”, which is a harsher criteria.
- **usesurrogate**: A value of **usesurrogate**=2, the default, splits subjects in the way described previously. This is similar to CART. If the value is 0, then a subject who is missing the primary split variable does not progress further down the tree. A value of 1 is intermediate: all surrogate variables except “go with the majority” are used to send a case further down the tree.
- **cp**: The threshold complexity parameter.

The complexity parameter **cp** is, like **minsplit**, an advisory parameter, but is considerably more useful. It is specified according to the formula

$$R_{cp}(T) \equiv R(T) + cp * |T| * R(T_0)$$

where  $T_0$  is the tree with no splits. This scaled version is much more user friendly than the original CART formula (4.1) since it is unitless. A value of **cp**=1 will always result in a tree with no splits. For regression models (see next section) the scaled **cp** has a very direct interpretation: if any split does not increase the overall  $R^2$  of the model by at least **cp** (where  $R^2$  is the usual linear-models definition) then that split is decreed to be, a priori, not worth pursuing. The program does not split said branch any further, and saves considerable computational effort. The default value of .01 has been reasonably successful at ‘pre-pruning’ trees so that the cross-validation step need only remove 1 or 2 layers, but it sometimes overprunes, particularly for large data sets.

## 6.2 Example: Consumer Report Auto Data

A second example using the **class** method demonstrates the outcome for a response with multiple ( $> 2$ ) categories. We also explore the difference between Gini and information splitting rules. The dataset **cu.summary** contains a collection of variables from the April, 1990 Consumer Reports summary on 117 cars. For our purposes,

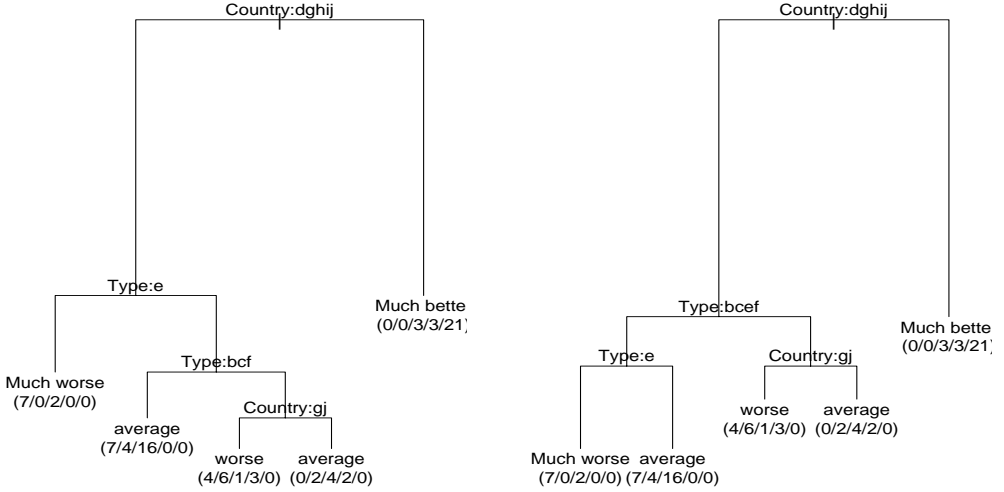


Figure 5: Displays the *rpart*-based model relating automobile Reliability to car type, price, and country of origin. The figure on the left uses the *gini* splitting index and the figure on the right uses the *information* splitting index.

car reliability will be treated as the response. The variables are:

Reliability	an ordered factor (contains NAs): Much worse < worse < average < better < Much Better
Price	numeric: list price in dollars, with standard equipment
Country	factor: country where car manufactured (Brazil, England, France, Germany, Japan, Japan/USA, Korea, Mexico, Sweden, USA)
Mileage	numeric: gas mileage in miles/gallon, contains NAs
Type	factor: Small, Sporty, Compact, Medium, Large, Van

In the analysis we are treating reliability as an unordered outcome. Nodes potentially can be classified as much worse, worse, average, better, or much better, though there are none that are labelled as just “better”. The 32 cars with missing response (listed as NA) were not used in the analysis. Two fits are made, one using the Gini index and the other the information index.

```
> fit1 <- rpart(Reliability ~ Price + Country + Mileage + Type,
  data=cu.summary, parms=list(split='gini'))
> fit2 <- rpart(Reliability ~ Price + Country + Mileage + Type,
```

```

data=cu.summary, parms=list(split='information'))

> par(mfrow=c(1,2))
> plot(fit1); text(fit1,use.n=T,cex=.9)
> plot(fit2); text(fit2,use.n=T,cex=.9)

```

The first two nodes from the Gini tree are

```

Node number 1: 85 observations,    complexity param=0.3051
predicted class= average expected loss= 0.6941
  class counts:  18 12 26  8 21
  probabilities: 0.2118 0.1412 0.3059 0.0941 0.2471
left son=2 (58 obs) right son=3 (27 obs)
Primary splits:
  Country splits as ---LRRLLLL, improve=15.220, (0 missing)
  Type splits as RLLRLL, improve= 4.288, (0 missing)
  Price < 11970 to the right, improve= 3.200, (0 missing)
  Mileage < 24.5 to the left, improve= 2.476, (36 missing)

```

```

Node number 2: 58 observations,    complexity param=0.08475
predicted class= average expected loss= 0.6034
  class counts:  18 12 23  5  0
  probabilities: 0.3103 0.2069 0.3966 0.0862 0.0000
left son=4 (9 obs) right son=5 (49 obs)
Primary splits:
  Type splits as RRRRLR, improve=3.187, (0 missing)
  Price < 11230 to the left, improve=2.564, (0 missing)
  Mileage < 24.5 to the left, improve=1.802, (30 missing)
  Country splits as ---L--RLRL, improve=1.329, (0 missing)

```

The fit for the information splitting rule is

```

Node number 1: 85 observations,    complexity param=0.3051
predicted class= average expected loss= 0.6941
  class counts:  18 12 26  8 21
  probabilities: 0.2118 0.1412 0.3059 0.0941 0.2471
left son=2 (58 obs) right son=3 (27 obs)
Primary splits:
  Country splits as ---LRRLLLL, improve=38.540, (0 missing)
  Type splits as RLLRLL, improve=11.330, (0 missing)
  Price < 11970 to the right, improve= 6.241, (0 missing)
  Mileage < 24.5 to the left, improve= 5.548, (36 missing)

```

```

Node number 2: 58 observations,    complexity param=0.0678
predicted class= average expected loss= 0.6034
  class counts:  18 12 23  5  0

```



```

probabilities: 0.3103 0.2069 0.3966 0.0862 0.0000
left son=4 (36 obs) right son=5 (22 obs)
Primary splits:
  Type      splits as RLLRLL, improve=9.281, (0 missing)
  Price    < 11230 to the left, improve=5.609, (0 missing)
  Mileage < 24.5 to the left, improve=5.594, (30 missing)
  Country splits as ---L--RRRL, improve=2.891, (0 missing)
Surrogate splits:
  Price    < 10970 to the right, agree=0.8793, (0 split)
  Country splits as ---R--RRRL, agree=0.7931, (0 split)

```

The first 3 countries (Brazil, England, France) had only one or two cars in the listing, all of which were missing the reliability variable. There are no entries for these countries in the first node, leading to the – symbol for the rule. The information measure has larger “improvements”, consistent with the difference in scaling between the information and Gini criteria shown in figure 2, but the relative merits of different splits are fairly stable.

The two rules do not choose the same primary split at node 2. The data at this point are

	Compact	Large	Medium	Small	Sporty	Van
Much worse	2	2	4	2	7	1
worse	5	0	4	3	0	0
average	3	5	8	2	2	3
better	2	0	0	3	0	0
Much better	0	0	0	0	0	0

Since there are 6 different categories, all  $2^5 = 32$  different combinations were explored, and as it turns out there are several with a nearly identical improvement. The Gini and information criteria make different “random” choices from this set of near ties. For the Gini index, *Sporty vs others* *Compact/Small vs others* have improvements of 37.19 and 37.20, respectively. For the information index, the improvements are 67.3 versus 64.2, respectively. Interestingly, the two splitting criteria arrive at exactly the same final nodes, for the full tree, although by different paths. (Compare the class counts of the terminal nodes).

We have said that for a categorical predictor with  $m$  levels, all  $2^{m-1}$  different possible splits are tested.. When there are a large number of categories for the predictor, the computational burden of evaluating all of these subsets can become large. For instance, the call `rpart(Reliabilty ~ ., data=car.all)` does not return for a *long*, long time: one of the predictors in that data set is a factor with 79 levels! Luckily, for any ordered outcome there is a computational shortcut that allows the program to find the best split using only  $m - 1$  comparisons. This includes the

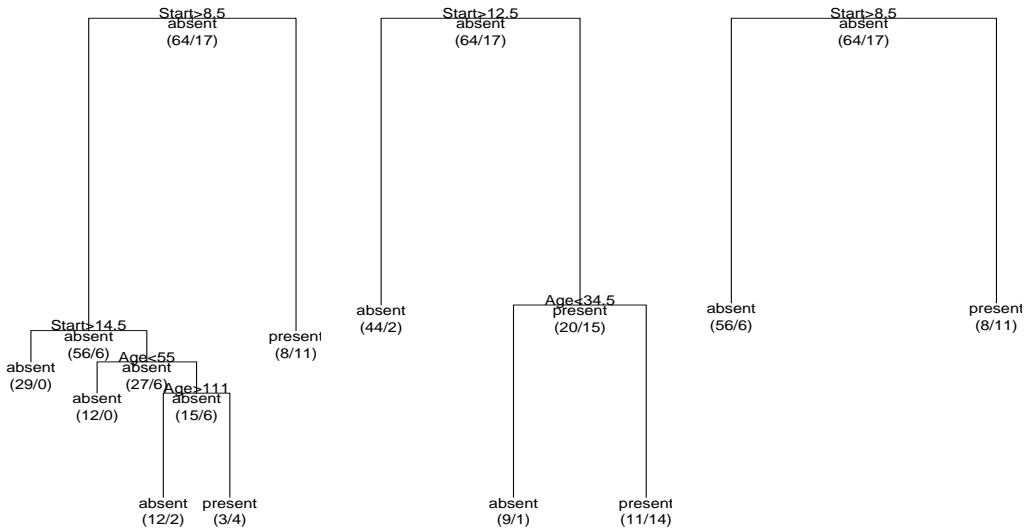


Figure 6: Displays the *rpart*-based models for the (presence/absence) of kyphosis. The figure on the left uses the default prior  $(0.79, 0.21)$  and loss; the middle figure uses the user-defined prior  $(0.65, 0.35)$  and default loss; and the third figure uses the default prior and the user-defined loss  $L(1, 2) = 3$ ,  $L(2, 1) = 4$ .

classification method when there are only two categories, along with the anova and Poisson methods to be introduced later.

### 6.3 Example: Kyphosis data

A third `class` method example explores the parameters `prior` and `loss`. The dataset `kyphosis` has 81 rows representing data on 81 children who have had corrective spinal surgery. The variables are:

Kyphosis	factor: lists if postoperative deformity is present/absent
Age	numeric: age of child in months
Number	numeric: number of vertebrae involved in operation
Start	numeric: beginning of the range of vertebrae involved

```
> lmat <- matrix(c(0,4,3,0), nrow=2, ncol=2, byrow=F)
> fit1 <- rpart(Kyphosis ~ Age + Number + Start, data=kyphosis)
> fit2 <- rpart(Kyphosis ~ Age + Number + Start, data=kyphosis,
```

```

      parms=list(prior=c(.65,.35)))
> fit3 <- rpart(Kyphosis ~ Age + Number + Start,data=kyphosis,
      parms=list(loss=lmat))

> par(mfrow=c(1,3))
> plot(fit1); text(fit1,use.n=T,all=T)
> plot(fit2); text(fit2,use.n=T,all=T)
> plot(fit3); text(fit3,use.n=T,all=T)

```

This example shows how even the initial split changes depending on the prior and loss that are specified. The first and third fits have the same initial split (Start < 8.5), but the improvement differs. The second fit splits Start at 12.5 which moves 46 people to the left instead of 62.

Looking at the leftmost tree, we see that the sequence of splits on the left hand branch yields only a single node classified as *present*. For any loss greater than 4 to 3, the routine will instead classify this node as *absent*, and the entire left side of the tree collapses, as seen in the right hand figure. This is not unusual — the most common effect of alternate loss matrices is to change the amount of pruning in the tree, more in some branches and less in others, rather than to change the choice of splits.

The first node from the default tree is

```

Node number 1: 81 observations,      complexity param=0.1765
predicted class= absent expected loss= 0.2099
  class counts:  64 17
  probabilities: 0.7901 0.2099
left son=2 (62 obs) right son=3 (19 obs)
Primary splits:
  Start < 8.5  to the right, improve=6.762, (0 missing)
  Number < 5.5 to the left,  improve=2.867, (0 missing)
  Age    < 39.5 to the left,  improve=2.250, (0 missing)
Surrogate splits:
  Number < 6.5 to the left,  agree=0.8025, (0 split)

```

The fit using the prior (0.65,0.35) is

```

Node number 1: 81 observations,      complexity param=0.302
predicted class= absent expected loss= 0.35
  class counts:  64 17
  probabilities: 0.65 0.35
left son=2 (46 obs) right son=3 (35 obs)
Primary splits:
  Start < 12.5 to the right, improve=10.900, (0 missing)

```

```

    Number < 4.5  to the left,  improve= 5.087, (0 missing)
    Age    < 39.5 to the left,  improve= 4.635, (0 missing)
Surrogate splits:
    Number < 3.5  to the left,  agree=0.6667, (0 split)

```

And first split under 4/3 losses is

```

Node number 1: 81 observations,      complexity param=0.01961
  predicted class= absent  expected loss= 0.6296
    class counts:  64 17
    probabilities:  0.7901 0.2099
  left son=2 (62 obs) right son=3 (19 obs)
Primary splits:
    Start < 8.5  to the right, improve=5.077, (0 missing)
    Number < 5.5 to the left,  improve=2.165, (0 missing)
    Age    < 39.5 to the left,  improve=1.535, (0 missing)
Surrogate splits:
    Number < 6.5 to the left,  agree=0.8025, (0 split)

```

## 7 Regression

### 7.1 Definition

Up to this point the classification problem has been used to define and motivate our formulae. However, the partitioning procedure is quite general and can be extended by specifying 5 “ingredients”:

- A splitting criterion, which is used to decide which variable gives the best split. For classification this was either the Gini or log-likelihood function. In the anova method the splitting criteria is  $SS_T - (SS_L + SS_R)$ , where  $SS_T = \sum (y_i - \bar{y})^2$  is the sum of squares for the node, and  $SS_R$ ,  $SS_L$  are the sums of squares for the right and left son, respectively. This is equivalent to choosing the split to maximize the between-groups sum-of-squares in a simple analysis of variance. This rule is identical to the regression option for `tree`.
- A summary statistic or vector, which is used to describe a node. The first element of the vector is considered to be the fitted value. For the anova method this is the mean of the node; for classification the response is the predicted class followed by the vector of class probabilities.
- The error of a node. This will be the variance of  $y$  for anova, and the predicted loss for classification.

- The prediction error for a new observation, assigned to the node. For anova this is  $(y_{new} - \bar{y})$ .
- Any necessary initialization.

The anova method leads to regression trees; it is the default method if  $y$  a simple numeric vector, i.e., not a factor, matrix, or survival object.

## 7.2 Example: Consumer Report Auto data (cont.)

The dataset `car.all` contains a collection of variables from the April, 1990 Consumer Reports; it has 36 variables on 111 cars. Documentation may be found in the S-Plus manual. We will work with a subset of 23 of the variables, created by the first two lines of the example below. We will use `Price` as the response. This data set is a good example of the usefulness of the missing value logic in `rpart`: most of the variables are missing on only 3-5 observations, but only 42/111 have a complete subset.

```
> cars <- car.all[, c(1:12, 15:17, 21, 28, 32:36)]
> cars$Eng.Rev <- as.numeric(as.character(car.all$Eng.Rev2))
> fit3 <- rpart(Price ~ ., data=cars)
> fit3
node), split, n, deviance, yval
    * denotes terminal node

1) root 105 7118.00 15.810
  2) Disp.<156 70 1492.00 11.860
    4) Country:Brazil,Japan,Japan/USA,Korea,Mexico,USA 58 421.20 10.320
      8) Type:Small 21 50.31 7.629 *
      9) Type:Compact,Medium,Sporty,Van 37 132.80 11.840 *
    5) Country:France,Germany,Sweden 12 270.70 19.290 *
  3) Disp.>156 35 2351.00 23.700
    6) HP.revs<5550 24 980.30 20.390
      12) Disp.<267.5 16 396.00 17.820 *
      13) Disp.>267.5 8 267.60 25.530 *
    7) HP.revs>5550 11 531.60 30.940 *

> printcp(fit3)

Regression tree:
rpart(formula = Price ~ ., data = cars)

Variables actually used in tree construction:
[1] Country Disp. HP.revs Type
```

Root node error: 7.1183e9/105 = 6.7793e7

	CP	nsplit	rel error	xerror	xstd
1	0.460146	0	1.00000	1.02413	0.16411
2	0.117905	1	0.53985	0.79225	0.11481
3	0.044491	3	0.30961	0.60042	0.10809
4	0.033449	4	0.26511	0.58892	0.10621
5	0.010000	5	0.23166	0.57062	0.11782

Only 4 of 22 predictors were actually used in the fit: engine displacement in cubic inches, country of origin, type of vehicle, and the revolutions for maximum horsepower (the “red line” on a tachometer).

- The relative error is  $1 - R^2$ , similar to linear regression. The xerror is related to the PRESS statistic. The first split appears to improve the fit the most. The last split adds little improvement to the apparent error.
- The 1-SE rule would choose a tree with 3 splits.
- This is a case where the default cp value of .01 may have overpruned the tree, since the cross-validated error is not yet at a minimum. A rerun with the cp threshold at .002 gave a maximum tree size of 8 splits, with a minimum cross-validated error for the 5 split model.
- For any CP value between 0.46015 and 0.11791 the best model has one split; for any CP value between 0.11791 and 0.04449 the best model is with 3 splits; and so on.

The print command also recognizes the cp option, which allows the user to see which splits are the most important.

```
> print(fit3,cp=.10)
node), split, n, deviance, yval
    * denotes terminal node
```

```
1) root 105 7.118e+09 15810
  2) Disp.<156 70 1.492e+09 11860
    4) Country:Brazil,Japan,Japan/USA,Korea,Mexico,USA 58 4.212e+08 10320 *
    5) Country:France,Germany,Sweden 12 2.707e+08 19290 *
  3) Disp.>156 35 2.351e+09 23700
    6) HP.revs<5550 24 9.803e+08 20390 *
    7) HP.revs>5550 11 5.316e+08 30940 *
```

The first split on displacement partitions the 105 observations into groups of 70 and 35 (nodes 2 and 3) with mean prices of 11,860 and 23,700. The deviance (corrected sum-of-squares) at these 2 nodes are  $1.49 \times 10^9$  and  $2.35 \times 10^9$ , respectively. More detailed summarization of the splits is again obtained by using the function `summary.rpart`.

```
> summary(fit3, cp=.10)
Node number 1: 105 observations,      complexity param=0.4601
mean=15810 , SS/n=67790000
left son=2 (70 obs) right son=3 (35 obs)
Primary splits:
  Disp.      < 156   to the left,   improve=0.4601, (0 missing)
  HP         < 154   to the left,   improve=0.4549, (0 missing)
  Tank       < 17.8  to the left,   improve=0.4431, (0 missing)
  Weight     < 2890  to the left,   improve=0.3912, (0 missing)
  Wheel.base < 104.5 to the left,   improve=0.3067, (0 missing)
Surrogate splits:
  Weight     < 3095  to the left,   agree=0.9143, (0 split)
  HP         < 139   to the left,   agree=0.8952, (0 split)
  Tank       < 17.95 to the left,   agree=0.8952, (0 split)
  Wheel.base < 105.5 to the left,   agree=0.8571, (0 split)
  Length     < 185.5 to the left,   agree=0.8381, (0 split)

Node number 2: 70 observations,      complexity param=0.1123
mean=11860 , SS/n=21310000
left son=4 (58 obs) right son=5 (12 obs)
Primary splits:
  Country splits as L-RRLLLLRL, improve=0.5361, (0 missing)
  Tank    < 15.65 to the left,   improve=0.3805, (0 missing)
  Weight  < 2568  to the left,   improve=0.3691, (0 missing)
  Type    splits as R-RLRR, improve=0.3650, (0 missing)
  HP      < 105.5 to the left,   improve=0.3578, (0 missing)
Surrogate splits:
  Tank          < 17.8  to the left,   agree=0.8571, (0 split)
  Rear.Seating < 28.75 to the left,   agree=0.8429, (0 split)
  .
  .
  .
```

- The improvement listed is the percent change in sums of squares (SS) for this split, i.e.,  $1 - (SS_{right} + SS_{left})/SS_{parent}$ .
- The weight and displacement are very closely related, as shown by the surrogate split agreement of 91%.

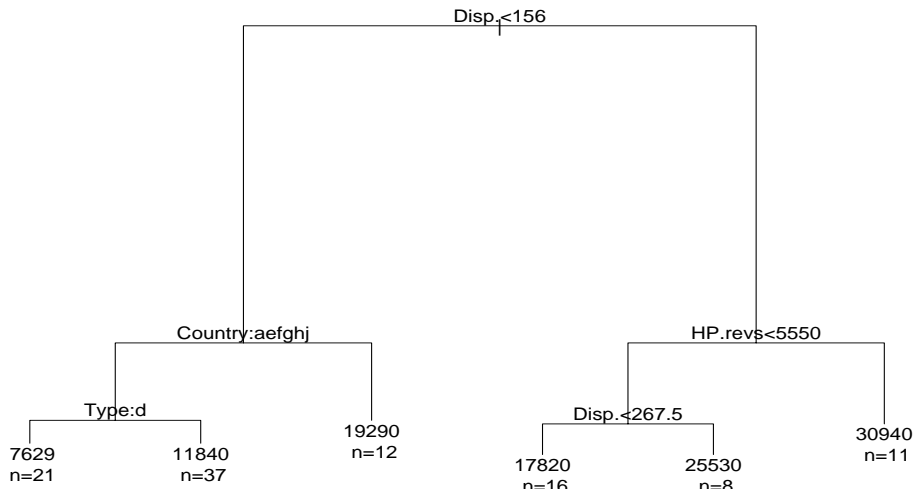


Figure 7: A *anova* tree for the *car.test.frame* dataset. The label of each node indicates the mean *Price* for the cars in that node.

- Not all types are represented in node 2, e.g., there are no representatives from England (the second category). This is indicated by a - in the list of split directions.

```

> plot(fit3)
> text(fit3,use.n=T)

```

As always, a plot of the fit is useful for understanding the *rpart* object. In this plot, we use the option `use.n=T` to add the number of cars in each node. (The default is for only the mean of the response variable to appear). Each individual split is ordered to send the less expensive cars to the left.

Other plots can be used to help determine the best *cp* value for this model. The function `rsq.rpart` plots the jackknifed error versus the number of splits. Of interest is the smallest error, but any number of splits within the “error bars” (1-SE rule) are considered a reasonable number of splits (in this case, 1 or 3 splits seem to be sufficient). As is often true with modelling, simpler is usually better. Another useful plot is the  $R^2$  versus number of splits. The (1 - apparent error) and (1 - relative error) show how much is gained with additional splits. This plot highlights the differences between the  $R^2$  values (figure 8).

Finally, it is possible to look at the residuals from this model, just as with a regular linear regression fit, as shown in the following figure.



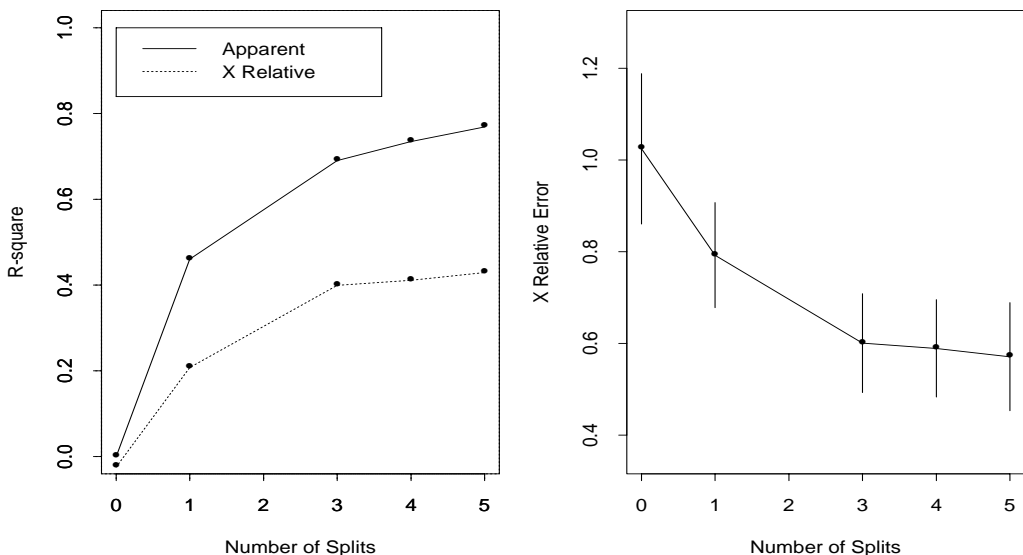


Figure 8: Both plots were obtained using the function `rsq.rpart(fit3)`. The figure on the left shows that the first split offers the most information. The figure on the right suggests that the tree should be pruned to include only 1 or 2 splits.

```
> plot(predict(fit3), resid(fit3))
> axis(3, at=fit3$frame$yval[fit3$frame$var=='<leaf>'],
      labels=row.names(fit3$frame)[fit3$frame$var=='<leaf>'])
> mtext('leaf number', side=3, line=3)
> abline(h=0)
```

### 7.3 Example: Stage C prostate cancer (anova method)

The stage C prostate cancer data of the earlier section can also be fit using the anova method, by treating the status variable as though it were continuous.

```
> cfit2 <- rpart(pgstat ~ age + eet + g2 + grade + gleason + ploidy,
                 data=stagec)

> printcp(cfit2)
Regression tree:
rpart(formula = pgstat ~ age + eet + g2 + grade + gleason + ploidy, data =
      stagec)

Variables actually used in tree construction:
[1] age    g2    grade ploidy
```

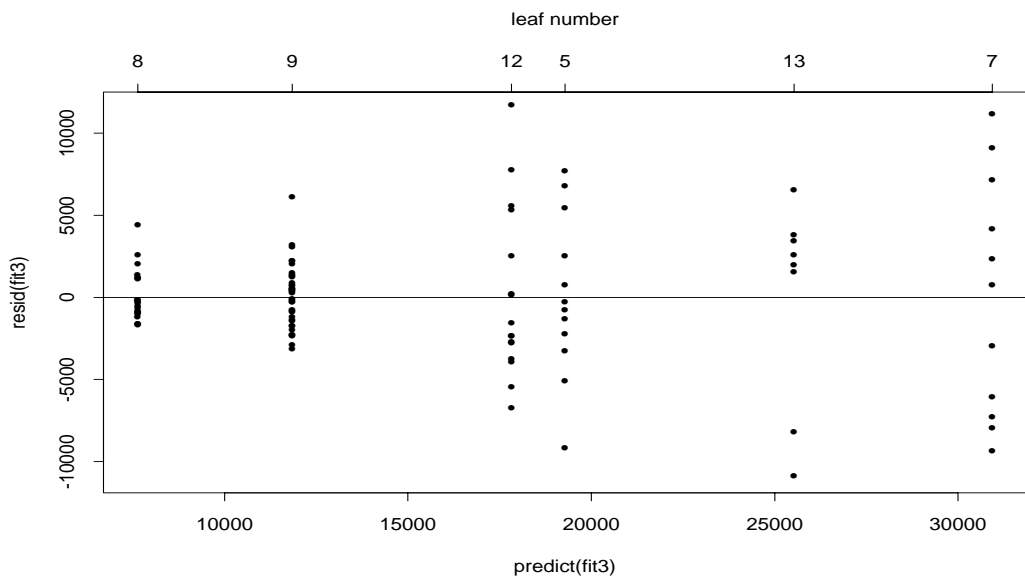


Figure 9: *This plot shows the (observed-expected) cost of cars versus the predicted cost of cars based on the nodes/leaves in which the cars landed. There appears to be more variability in node 7 than in some of the other leaves.*

Root node error:  $34.027/146 = 0.23306$

	CP	nsplit	rel error	xerror	xstd
1	0.152195	0	1.00000	1.01527	0.045470
2	0.054395	1	0.84781	0.86670	0.063447
3	0.032487	3	0.73901	0.86524	0.075460
4	0.019932	4	0.70653	0.95702	0.085390
5	0.013027	8	0.63144	1.05606	0.092566
6	0.010000	9	0.61841	1.07727	0.094466

```
> print(cfit2, cp=.03)
node), split, n, deviance, yval
    * denotes terminal node
```

- 1) root 146 34.030 0.3699
- 2) grade<2.5 61 7.672 0.1475
  - 4) g2<13.19 40 1.900 0.0500 \*
  - 5) g2>13.19 21 4.667 0.3333 \*
- 3) grade>2.5 85 21.180 0.5294
  - 6) g2<13.2 40 9.775 0.4250 \*

```

7) g2>13.2 45 10.580 0.6222
14) g2>17.91 22 5.091 0.3636 *
15) g2<17.91 23 2.609 0.8696 *

```

If this tree is compared to the earlier results, we see that it has chosen exactly the same variables and split points as before. The only addition is further splitting of node 2, the upper left “No” of figure 3. This is no accident, for the two class case the Gini splitting rule reduces to  $2p(1 - p)$ , which is the variance of a node.

The two methods differ in their evaluation and pruning, however. Note that nodes 4 and 5, the two children of node 2, contain 2/40 and 7/21 progressions, respectively. For classification purposes both nodes have the same predicted value (No) and the split will be discarded since the error (# of misclassifications) with and without the split is identical. In the regression context the two predicted values of .05 and .33 *are* different — the split has identified a nearly pure subgroup of significant size.

This setup is known as *odds regression*, and may be a more sensible way to evaluate a split when the emphasis of the model is on understanding/explanation rather than on prediction error per se. Extension of this rule to the multiple class problem is appealing, but has not yet been implimented in rpart.

## 8 Poisson regression

### 8.1 Definition

The Poisson splitting method attempts to extend rpart models to event rate data. The model in this case is

$$\lambda = f(x)$$

where  $\lambda$  is an event rate and  $x$  is some set of predictors. As an example consider hip fracture rates. For each county in the United States we can obtain

- number of fractures in patients age 65 or greater (from Medicare files)
- population of the county (US census data)
- potential predictors such as
  - socio-economic indicators
  - number of days below freezing
  - ethnic mix
  - physicians/1000 population

– etc.

Such data would usually be approached by using Poisson regression; can we find a tree based analogue? In adding criteria for rates regression to this ensemble, the guiding principle was the following: the between groups sum-of-squares is not a very robust measure, yet tree based regression works very well. So do the simplest thing possible.

Let  $c_i$  be the observed event count for observation  $i$ ,  $t_i$  be the observation time, and  $x_{ij}, j = 1, \dots, p$  be the predictors. The  $y$  variable for the program will be a 2 column matrix.

Splitting criterion: The likelihood ratio test for two Poisson groups

$$D_{\text{parent}} - (D_{\text{left son}} + D_{\text{right son}})$$

Summary statistics: The observed event rate and the number of events.

$$\hat{\lambda} = \frac{\# \text{ events}}{\text{total time}} = \frac{\sum c_i}{\sum t_i}$$

Error of a node: The within node deviance.

$$D = \sum \left[ c_i \log \left( \frac{c_i}{\hat{\lambda} t_i} \right) - (c_i - \hat{\lambda} t_i) \right]$$

Prediction error: The deviance contribution for a new observation, using  $\hat{\lambda}$  of the node as the predicted rate.

## 8.2 Improving the method

There is a problem with the criterion just proposed, however: cross-validation of a model often produces an infinite value for the deviance. The simplest case where this occurs is easy to understand. Assume that some terminal node of the tree has 20 subjects, but only 1 of the 20 has experienced any events. The cross-validated error (deviance) estimate for that node will have one subset — the one where the subject with an event is left out — which has  $\hat{\lambda} = 0$ . When we use the prediction for the 10% of subjects who were set aside, the deviance contribution of the subject with an event is

$$\dots + c_i \log(c_i/0) + \dots$$

which is infinite since  $c_i > 0$ . The problem is that when  $\hat{\lambda} = 0$  the occurrence of an event is infinitely improbable, and, using the deviance measure, the corresponding model is then infinitely bad.

One might expect this phenomenon to be fairly rare, but unfortunately it is not so. One given of tree-based modeling is that a right-sized model is arrived at by purposely over-fitting the data and then pruning back the branches. A program that aborts due to a numeric exception during the first stage is uninformative to say the least. Of more concern is that this edge effect does not seem to be limited to the pathologic case detailed above. Any near approach to the boundary value  $\lambda = 0$  leads to large values of the deviance, and the procedure tends to discourage any final node with a small number of events.

An ad hoc solution is to use the revised estimate

$$\hat{\lambda} = \max\left(\hat{\lambda}, \frac{k}{\sum t_i}\right)$$

where  $k$  is  $1/2$  or  $1/6$ . That is, pure nodes are given a partial event. (This is similar to the starting estimates used in the GLM program for a Poisson regression.) This is unsatisfying, however, and we propose instead using a shrinkage estimate.

Assume that the true rates  $\lambda_j$  for the leaves of the tree are random values from a  $\text{Gamma}(\mu, \sigma)$  distribution. Set  $\mu$  to the observed overall event rate  $\sum c_i / \sum t_i$ , and let the user choose as a prior the coefficient of variation  $k = \sigma / \mu$ . A value of  $k = 0$  represents extreme pessimism (“the leaf nodes will all give the same result”), whereas  $k = \infty$  represents extreme optimism. The Bayes estimate of the event rate for a node works out to be

$$\hat{\lambda}_k = \frac{\alpha + \sum c_i}{\beta + \sum t_i},$$

where  $\alpha = 1/k^2$  and  $\beta = \alpha / \hat{\lambda}$ .

This estimate is scale invariant, has a simple interpretation, and shrinks least those nodes with a large amount of information. In practice, a value of  $k = 10$  does essentially no shrinkage. For `method='poisson'`, the optional parameters list is the single number  $k$ , with a default value of 1. This corresponds to prior coefficient of variation of 1 for the estimated  $\lambda_j$ . We have not nearly enough experience to decide if this is a good value. (It does stop the `log(0)` message though).

Cross-validation does not work very well. The procedure gives very conservative results, and quite often declares the no-split tree to be the best. This may be another artifact of the edge effect.

### 8.3 Example: solder data

The solder data frame, as explained in the Splus help file, is a design object with 900 observations, which are the results of an experiment varying 5 factors relevant to the wave-soldering procedure for mounting components on printed circuit boards. The

response variable, skips, is a count of how many solder skips appeared to a visual inspection. The other variables are listed below:

Opening	factor: amount of clearance around the mounting pad (S < M < L)
Solder	factor: amount of solder used (Thin < Thick)
Mask	factor: Type of solder mask used (5 possible)
PadType	factor: Mounting pad used (10 possible)
Panel	factor: panel (1, 2 or 3) on board being counted

In this call, the `rpart.control` options are modified: `maxcompete = 2` means that only 2 other competing splits are listed (default is 4); `cp = .05` means that a smaller tree will be built initially (default is .01). The `y` variable for Poisson partitioning may be a two column matrix containing the observation time in column 1 and the number of events in column 2, or it may be a vector of event counts alone.

```
fit <- rpart(skips ~ Opening + Solder + Mask + PadType
            + Panel, data=solder, method='poisson',
            control=rpart.control(cp=.05, maxcompete=2))
```

The `print` command summarizes the tree, as in the previous examples.

```
node), split, n, deviance, yval
  * denotes terminal node

1) root 900 8788.0  5.530
 2) Opening:M,L 600 2957.0  2.553
   4) Mask:A1.5,A3,B3 420  874.4  1.033 *
   5) Mask:A6,B6 180  953.1  6.099 *
 3) Opening:S 300 3162.0 11.480
   6) Mask:A1.5,A3 150  652.6  4.535 *
   7) Mask:A6,B3,B6 150 1155.0 18.420 *
```

- The response value is the expected event rate (with a time variable), or in this case the expected number of skips. The values are shrunk towards the global estimate of 5.530 skips/observation.
- The deviance is the same as the null deviance (sometimes called the residual deviance) that you'd get when calculating a Poisson glm model for the given subset of data.

```
> summary(fit,cp=.10)
```

```

Call:
rpart(formula = skips ~ Opening + Solder + Mask + PadType + Panel, data =
      solder, method = "poisson", control = rpart.control(cp = 0.05,
      maxcompete = 2))

```

	CP	nsplit	rel error	xerror	xstd
1	0.3038	0	1.0000	1.0051	0.05248
2	0.1541	1	0.6962	0.7016	0.03299
3	0.1285	2	0.5421	0.5469	0.02544
4	0.0500	3	0.4137	0.4187	0.01962

```

Node number 1: 900 observations,      complexity param=0.3038
events=4977,  estimated rate=5.53 , deviance/n=9.765
left son=2 (600 obs) right son=3 (300 obs)
Primary splits:
  Opening splits as  RLL, improve=2670, (0 missing)
  Mask    splits as  LLRLR, improve=2162, (0 missing)
  Solder  splits as  RL, improve=1168, (0 missing)

```

```

Node number 2: 600 observations,      complexity param=0.1285
events=1531,  estimated rate=2.553 , deviance/n=4.928
left son=4 (420 obs) right son=5 (180 obs)
Primary splits:
  Mask    splits as  LLRLR, improve=1129.0, (0 missing)
  Opening splits as  -RL, improve= 250.8, (0 missing)
  Solder  splits as  RL, improve= 219.8, (0 missing)

```

```

Node number 3: 300 observations,      complexity param=0.1541
events=3446,  estimated rate=11.48 , deviance/n=10.54
left son=6 (150 obs) right son=7 (150 obs)
Primary splits:
  Mask    splits as  LLRRR, improve=1354.0, (0 missing)
  Solder  splits as  RL, improve= 976.9, (0 missing)
  PadType splits as  RRRRLRLRLL, improve= 313.2, (0 missing)
Surrogate splits:
  Solder splits as  RL, agree=0.6, (0 split)

```

```

Node number 4: 420 observations
events=433,  estimated rate=1.033 , deviance/n=2.082

```

```

Node number 5: 180 observations
events=1098,  estimated rate=6.099 , deviance/n=5.295

```

```

Node number 6: 150 observations
events=680,  estimated rate=4.535 , deviance/n=4.351

```

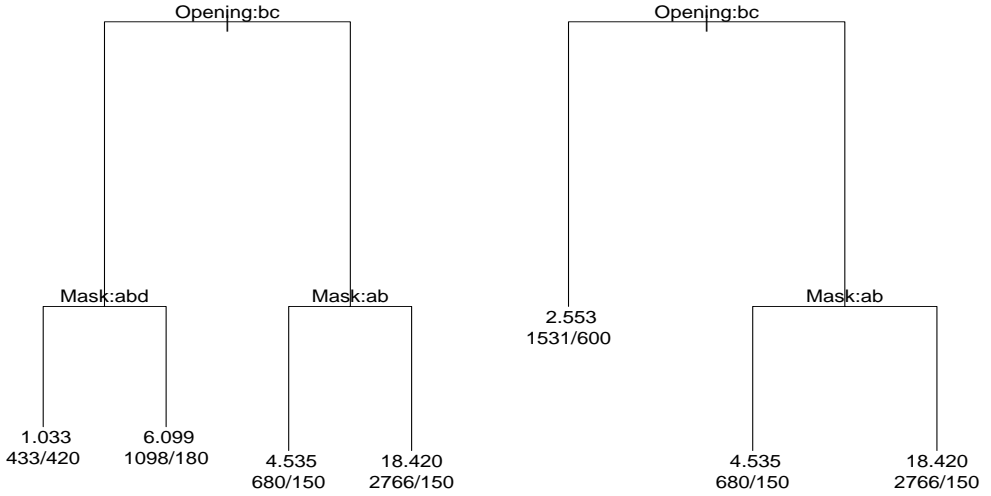


Figure 10: The first figure shows the solder data, fit with the *poisson* method, using a *cp* value of 0.05. The second figure shows the same fit, but with a *cp* value of 0.15. The function *prune.rpart* was used to produce the smaller tree.

Node number 7: 150 observations  
 events=2766, estimated rate=18.42 , deviance/n=7.701

- The improvement is  $\text{Deviance}_{\text{parent}} - (\text{Deviance}_{\text{left}} + \text{Deviance}_{\text{right}})$ , which is the likelihood ratio test for comparing two Poisson samples.
- The cross-validated error has been found to be overly pessimistic when describing how much the error is improved by each split. This is likely an effect of the boundary effect mentioned earlier, but more research is needed.
- The variation *xstd* is not as useful, given the bias of *xerror*.

```
> plot(fit)
> text(fit,use.n=T)

> fit.prune <- prune(fit,cp=.15)
> plot(fit.prune)
> text(fit.prune,use.n=T)
```



The `use.n=T` option specifies that number of events / total N should be listed along with the predicted rate (number of events/person-years). The function `prune` trims the tree fit to the `cp` value 0.15. The same tree could have been created by specifying `cp = .15` in the original call to `rpart`.

## 8.4 Example: Stage C prostate cancer (survival method)

One special case of the Poisson model is of particular interest for medical consulting (such as the authors do). Assume that we have survival data, i.e., each subject has either 0 or 1 event. Further, assume that the time values have been pre-scaled so as to fit an exponential model. That is, stretch the time axis so that a Kaplan-Meier plot of the data will be a straight line when plotted on the logarithmic scale. An approximate way to do this is

```
temp <- coxph(Surv(time, status) ~1)
newtime <- predict(temp, type='expected')
```

and then do the analysis using the `newtime` variable. (This replaces each time value by  $\Lambda(t)$ , where  $\Lambda$  is the cumulative hazard function).

A slightly more sophisticated version of this which we will call *exponential scaling* gives a straight line curve for  $\log(\text{survival})$  under a parametric exponential model. The only difference from the approximate scaling above is that a subject who is censored between observed death times will receive “credit” for the intervening interval, i.e., we assume the baseline hazard to be linear between observed deaths. If the data is pre-scaled in this way, then the Poisson model above is equivalent to the *local full likelihood* tree model of LeBlanc and Crowley [3]. They show that this model is more efficient than the earlier suggestion of Therneau et. al. [6] to use the martingale residuals from a Cox model as input to a regression tree (anova method). Exponential scaling or `method='exp'` is the default if `y` is a `Surv` object.

Let us again return to the stage C cancer example. Besides the variables explained previously we will use `pgtime`, which is time to tumor progression.

```
> fit <- rpart(Surv(pgtime, pgstat) ~ age + eet + g2 + grade +
               gleason + ploidy, data=stagec)
> print(fit)
```

```
node), split, n, deviance, yval
    * denotes terminal node
```

```
1) root 146 195.30 1.0000
  2) grade<2.5 61  44.98 0.3617
    4) g2<11.36 33   9.13 0.1220 *
    5) g2>11.36 28  27.70 0.7341 *
```

```

3) grade>2.5 85 125.10 1.6230
6) age>56.5 75 104.00 1.4320
12) gleason<7.5 50 66.49 1.1490
24) g2<13.475 25 29.10 0.8817 *
25) g2>13.475 25 36.05 1.4080
50) g2>17.915 14 18.72 0.8795 *
51) g2<17.915 11 13.70 2.1830 *
13) gleason>7.5 25 34.13 2.0280
26) g2>15.29 10 11.81 1.2140 *
27) g2<15.29 15 19.36 2.7020 *
7) age<56.5 10 15.52 3.1980 *

```

```

> plot(fit, uniform=T, branch=.4, compress=T)
> text(fit, use.n=T)

```

Note that the primary split on grade is the same as when status was used as a dichotomous endpoint, but that the splits thereafter differ.

```

> summary(fit,cp=.02)

```

Call:

```

rpart(formula = Surv(pgtime, pgstat) ~ age + eet + g2 + grade +
      gleason + ploidy, data = stagec)

```

	CP	nsplit	rel error	xerror	xstd
1	0.12913	0	1.0000	1.0060	0.07389
2	0.04169	1	0.8709	0.8839	0.07584
3	0.02880	2	0.8292	0.9271	0.08196
4	0.01720	3	0.8004	0.9348	0.08326
5	0.01518	4	0.7832	0.9647	0.08259
6	0.01271	5	0.7680	0.9648	0.08258
7	0.01000	8	0.7311	0.9632	0.08480

```

Node number 1: 146 observations,    complexity param=0.1291
events=54, estimated rate=1 , deviance/n=1.338
left son=2 (61 obs) right son=3 (85 obs)

```

Primary splits:

```

grade < 2.5   to the left, improve=25.270, (0 missing)
gleason < 5.5 to the left, improve=21.630, (3 missing)
ploidy splits as LRR, improve=14.020, (0 missing)
g2 < 13.2    to the left, improve=12.580, (7 missing)
age < 58.5   to the right, improve= 2.796, (0 missing)

```

Surrogate splits:

```

gleason < 5.5 to the left, agree=0.8630, (0 split)
ploidy splits as LRR, agree=0.6438, (0 split)

```

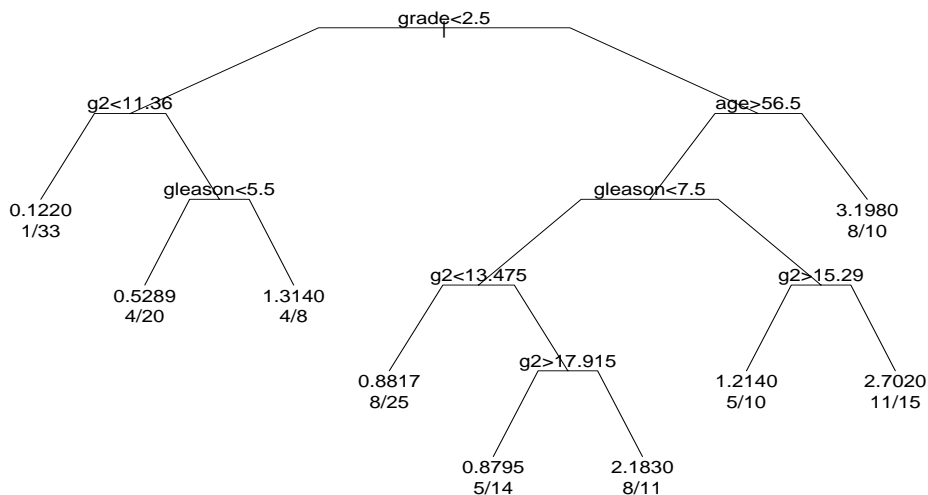


Figure 11: *The prostate cancer data as a survival tree*

```

g2      < 9.945 to the left,  agree=0.6301, (0 split)
age     < 66.5  to the right, agree=0.5890, (0 split)
.
.
.

```

Suppose that we wish to simplify this tree, so that only four terminal nodes remain. Looking at the table of complexity parameters, we see that `prune(fit, cp=.015)` would give the desired result. It is also possible to trim the figure interactively using `snip.rpart`. Point the mouse on a node and click with the left button to get some simple descriptives of the node. Double-clicking with the left button will ‘remove’ the sub-tree below, or one may click on another node. Multiple branches may be snipped off one by one; clicking with the right button will end interactive mode and return the pruned tree.

```

> plot(fit)
> text(fit,use.n=T)
> fit2 <- snip.rpart(fit)

```

```

node number: 6  n= 75
response= 1.432013 ( 37 )

```

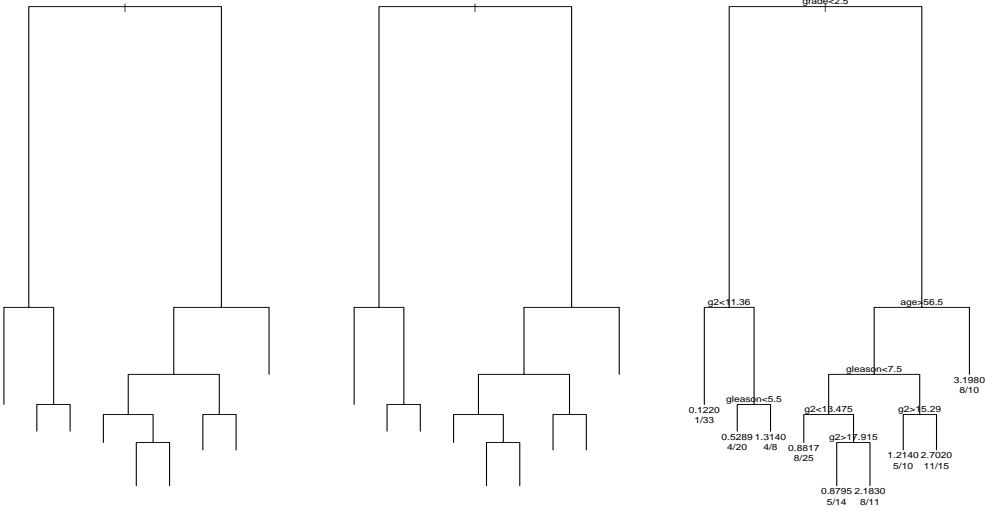


Figure 12: *An illustration of how `snip.rpart` works. The full tree is plotted in the first panel. After selecting node 6 with the mouse (double clicking on left button), the subtree disappears from the plot (shown in the second panel). Finally, the new subtree is redrawn to use all available space and it is labelled.*

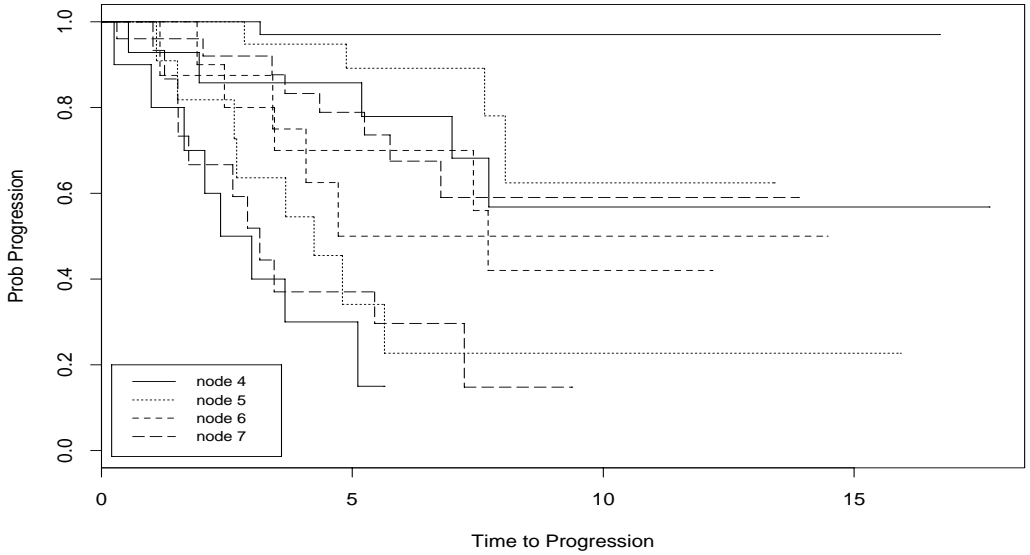


Figure 13: *Survival plot based on snipped rpart object. The probability of tumor progression is greatest in node 7, which has patients who are younger and have a higher initial tumor grade.*

```
Error (dev) = 103.9783
```

```
> plot(fit2)
> text(fit2,use.n=T)
```

For a final summary of the model, it can be helpful to plot the probability of survival based on the final bins in which the subjects landed. To create new variables based on the rpart groupings, use `where`. The nodes of `fit2` above are shown in the right hand panel of figure 12: node 4 has 1 event, 33 subjects, grade = 1-2 and  $g2 < 11.36$ ; node 5 has 8 events, 28 subjects, grade = 1-2 and  $g2 > 11.36$ ; node 6 has 37 events, 75 subjects, grade = 3-4, age  $> 56.5$ ; node 7 has 8 events, 10 subjects, grade = 3-4, age  $< 56.5$ . Patients who are younger and have a higher initial grade tend to have more rapid progression of disease.

```
> newgrp <- fit2$where
> plot(survfit(Surv(pgtime,pgstat) ~ newgrp, data=stagec),
      mark.time=F, lty=1:4)
> title(xlab='Time to Progression',ylab='Prob Progression')
> legend(.2,.2, legend=paste('node',c(4,5,6,7)), lty=1:4)
```

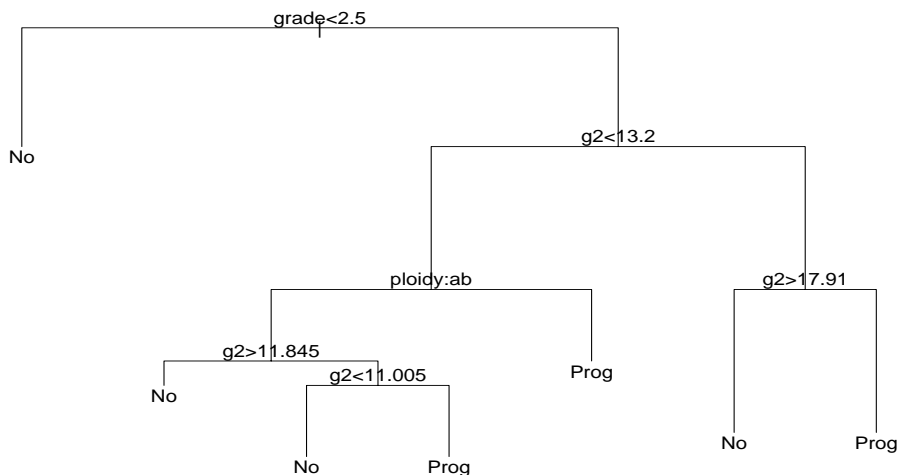


Figure 14: `plot(fit); text(fit)`

## 9 Plotting options

This section examines the various options that are available when plotting an `rpart` object. For simplicity, the same model (data from Example 1) will be used throughout.

The simplest labelled plot (figure 14) is called by using `plot` and `text` without changing any of the defaults. This is useful for a first look, but sometimes you'll want more information about each of the nodes.

```

> fit <- rpart(progstat ~ age + eet + g2 + grade + gleason + ploidy,
               stagec, control=rpart.control(cp=.025))

> plot(fit)
> text(fit)

```

The next plot (figure 15) has uniform stem lengths (`uniform=T`), has extra information (`use.n=T`) specifying number of subjects at each node (here it lists how many are diseased and how many are not diseased), and has labels on all the nodes, not just the terminal nodes (`all=T`).

```

> plot(fit, uniform=T)

```

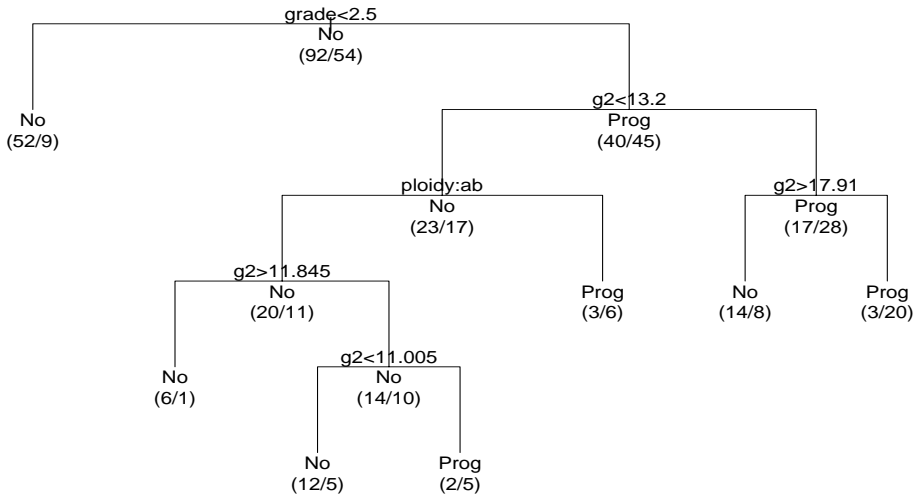


Figure 15: `plot(fit, uniform=T); text(fit,use.n=T,all=T)`

```
> text(fit, use.n=T, all=T)
```

Fancier plots can be created by modifying the `branch` option, which controls the shape of branches that connect a node to its children. The default for the plots is to have square shouldered trees (`branch = 1.0`). This can be taken to the other extreme with no shoulders at all (`branch=0`) as shown in figure 16.

```
> plot(fit, branch=0)
> text(fit, use.n=T)
```

These options can be combined with others to create the plot that fits your particular needs. The default plot may be inefficient in it's use of space: the terminal nodes will always lie at x-coordinates of 1,2,... . The `compress` option attempts to improve this by overlapping some nodes. It has little effect on figure 17, but in figure 4 it allows the lowest branch to “tuck under” those above. If you want to play around with the spacing with `compress`, try using `nospace` which regulates the space between a terminal node and a split.

```
> plot(fit,branch=.4,uniform=T,compress=T)
> text(fit,all=T,use.n=T)
```

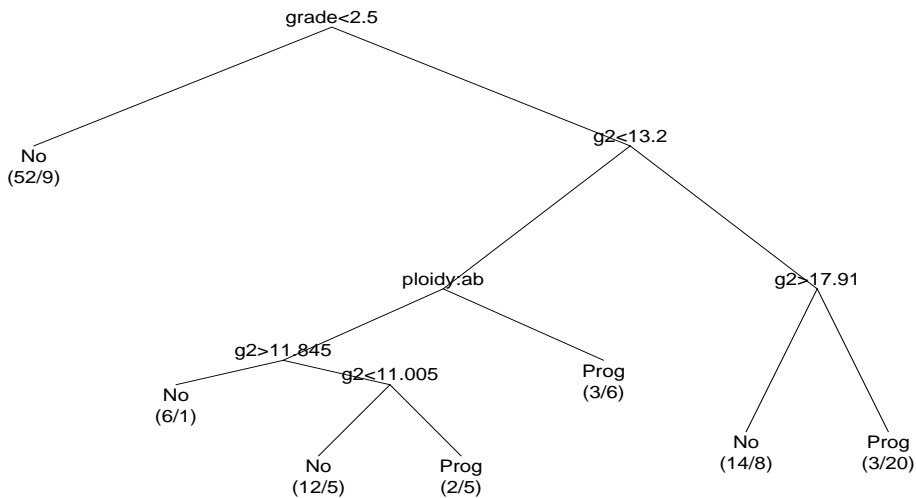


Figure 16: `plot(fit, branch=0); text(fit,use.n=T)`

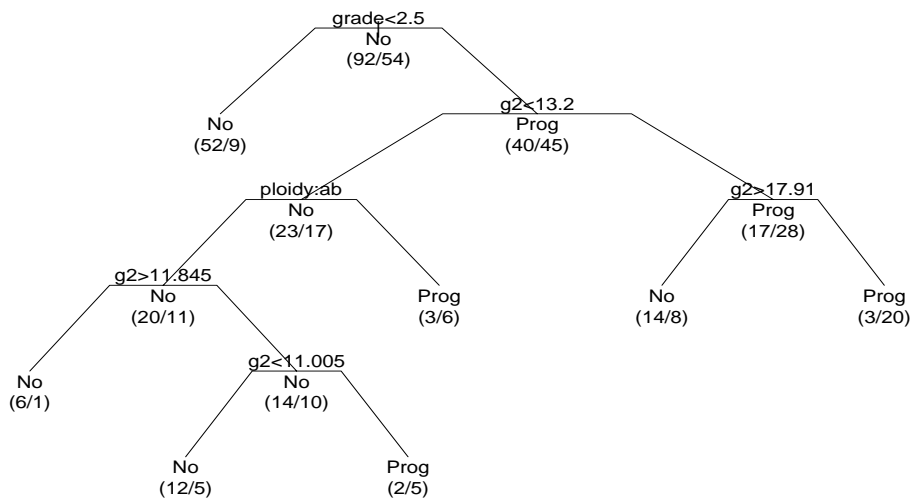


Figure 17: `plot(fit, branch=.4, uniform=T,compress=T)`



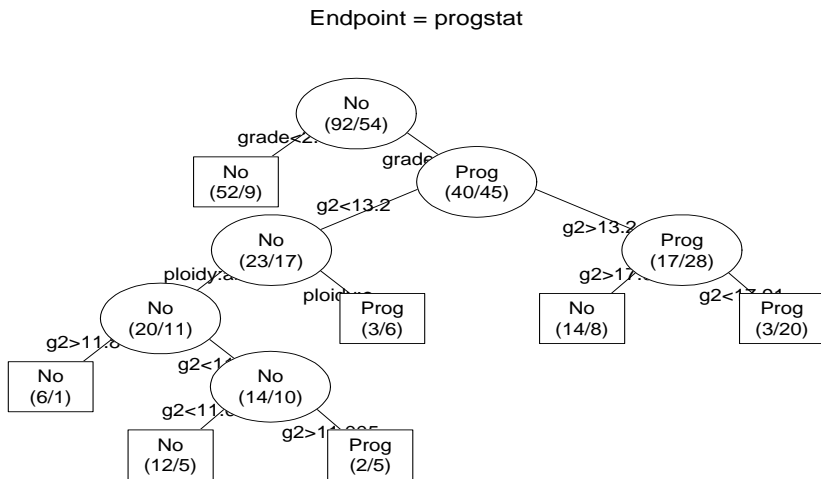


Figure 18: `post(fit)`

We have combined several of these options into a function called `post.rpart`. Results are shown in figure 18. The code is essentially

```
> plot(tree, uniform = T, branch = 0.2, compress = T, margin = 0.1)
> text(tree, all = T, use.n=T, fancy = T)
```

The `fancy` option of `text` creates the ellipses and rectangles, and moves the splitting rule to the midpoints of the branches. `Margin` shrinks the plotting region slightly so that the text boxes don't run over the edge of the plot. The `branch` option makes the lines exit the ellipse at a “good” angle. The call `post(fit)` will create a postscript file *fit.ps* in the current directory. The additional argument `file=""` will cause the plot to appear on the current active device. Note that `post.rpart` is just our choice of options to the functions `plot.rpart` and `text.rpart`. The reader is encouraged to try out other possibilities, as these specifications may not be the best choice in all situations.

## 10 Other functions

A more general approach to cross-validation can be gained using the `xpred.rpart` function. Given an `rpart` fit, a vector of  $k$  complexity parameters, and the desired number of cross-validation groups, this function returns an  $n$  by  $k$  matrix containing

the predicted value  $\hat{y}_{(-i)}$  for each subject, from the model that was fit without that subject. The *cp* vector defaults to the geometric mean of the *cp* sequence for the pruned tree on the full data set.

```
> fit <- rpart(pgtime ~ age + eet + g2 + grade + gleason + ploidy, stagec)
> fit$cpstable
      CP nsplit rel error   xerror   xstd
1 0.07572983      0 1.0000000 1.012323 0.1321533
2 0.02825076      2 0.8485403 1.094811 0.1547652
3 0.01789441      7 0.7219625 1.216495 0.1603581
4 0.01295145      8 0.7040681 1.223120 0.1610315
5 0.01000000      9 0.6911166 1.227213 0.1660616

> temp <- xpred.rpart(fit)
> err <- (stagec$pgtime - temp)^2
> sum.err <- apply(err,2,sum)
> sum.err / (fit$frame)$dev[1]
0.27519053 0.04625392 0.02248401 0.01522362 0.01138044
1.021901 1.038632 1.14714 1.179571 1.174196
```

The answer from `xpred.rpart` differs from the model's computation due to different random number seeds, which causes slightly different cross-validation groups.

## 11 Relation to other programs

### 11.1 CART

Almost all of the definitions in `rpart` are equivalent to those used in CART, and the output should usually be very similar. The printout given by `summary.rpart` was also strongly influenced by some early `cart` output. Some known differences are

- **Surrogate splits:** `cart` uses the percentage agreement between the surrogate and the primary split, and `rpart` uses the total number of agreements. When one of the surrogate variables has missing values this can lead to a different ordering. For instance, assume that the best surrogate based on  $x_1$  has 45/50 = 90% agreement (with 10 missing), and the best based on  $x_2$  has 46/60 (none missing). Then `rpart` will pick  $x_2$ . This is only a serious issue when there are a large number of missing values for one variable, and indeed the change was motivated by examples where a nearly-100% best surrogate due to perfect concordance with the primary split.
- **Computation:** Some versions of the `cart` code have been optimized for very large data problems, and include such features as subsampling from the larger nodes. Large data sets can be a problem in S-plus.

## 11.2 Tree

The user interface to **rpart** is almost identical to that of the **tree** functions. In fact, the **rpart** object was designed to be a simple superset of the **tree** class, and to inherit most methods from it (saving a lot of code writing). However, this close connection had to be abandoned. The **rpart** object is still very similar to **tree** objects, differing in 3 respects

- Addition of a **method** component. This was the single largest reason for divergence. In **tree**, splitting of a categorical variable results in a **yprob** element in the data structure, but regression does not. Most of the “downstream” functions then contain the code fragment “if the object contains a **yprob** component, then do A, else do B”. **rpart** has more than two methods, and this simple approach does not work. Rather, the method used is itself retained in the output.
- Additional components to describe the tree. This includes the **yval2** component, which contains further response information beyond the primary value. (For the gini method, for instance, the primary response value is the predicted class for a node, and the additional value is the complete vector of class counts. The predicted probability vector **yprob** is a function of these, the priors, and the tree topology.) Other additional components store the competitor and surrogate split information.
- The **xlevels** component in **rpart** is a list containing, for each factor variable, the list of levels for that factor. In **tree**, the list also contains **NULL** values for the non-factor predictors. In one problem with a very large (4096) number of predictors we found that the processing of this list consumed nearly as much time and memory as the problem itself. (The inefficiency in S-plus that caused this may have since been corrected).

Although the **rpart** structure does not inherit from class **tree**, some of the **tree** functions are used by the **rpart** methods, e.g. **tree.depth**. **Tree** methods that have not been implemented as of yet include **burl.tree**, **cv.tree**, **hist.tree**, **rug.tree**, and **tile.tree**.

Not all of the plotting functions available for **tree** objects have been implemented for **rpart** objects. However, many can be accessed by using the **as.tree** function, which reformats the components of an **rpart** object into a **tree** object. The resultant value may not work for all operations, however. The **tree** functions ignore the **method** component of the result, instead they *assume* that

- if *y* is a factor, then classification was performed based on the information index,

- otherwise, anova splitting (regression) was done.

Thus the result of `as.tree` from a Poisson fit will work reasonably for some plotting functions, e.g., `hist.tree`, but would not make sense for functions that do further computations such as `burl.tree`.

## 12 Source

The software exists in two forms: a stand-alone version, which can be found in statlib in the ‘general’ section, and an S version, also on statlib, but in the ‘S’ section of the library. The splitting rules and other computations exist in both versions, but the S version has been enhanced with several graphics options, most of which are modeled (copied actually) on the `tree` functions.

## References

- [1] L. Breiman, J.H. Friedman, R.A. Olshen, , and C.J Stone. *Classification and Regression Trees*. Wadsworth, Belmont, Ca, 1983.
- [2] L.A. Clark and D. Pregibon. Tree-based models. In J.M. Chambers and T.J. Hastie, editors, *Statistical Models in S*, chapter 9. Wadsworth and Brooks/Cole, Pacific Grove, Ca, 1992.
- [3] M. LeBlanc and J Crowley. Relative risk trees for censored survival data. *Biometrics*, 48:411–425, 1992.
- [4] O. Nativ, Y. Raz, H.Z. Winkler, Y. Hosaka, E.T. Boyle, T.M. Therneau, G.M. Farrow, R.P. Meyers, H. Zincke, and M.M Lieber. Prognostic value of flow cytometric nuclear DNA analysis in stage C prostate carcinoma. *Surgical Forum*, pages 685–687, 1988.
- [5] T.M. Therneau. A short introduction to recursive partitioning. Orion Technical Report 21, Stanford University, Department of Statistics, 1983.
- [6] T.M. Therneau, Grambsch P.M., and T.R. Fleming. Martingale based residuals for survival models. *Biometrika*, 77:147–160, 1990.