

GUIDE Classification and Regression Trees*

User Manual for Version 19.0

Wei-Yin Loh
Department of Statistics
University of Wisconsin–Madison

April 16, 2015

Contents

1	Warranty disclaimer	2
2	Introduction	3
2.1	Installation	6
2.2	L ^A T _E X	6
3	Program operation	7
3.1	Required files	7
3.2	Input file creation	9
4	Classification	10
4.1	Default: univariate splits	10
4.1.1	Input file creation with default options	10
4.1.2	Contents of <code>input.txt</code>	12
4.1.3	Executing the program and interpreting the output	13
4.2	Non-default options	18
4.2.1	Linear splits	18
4.2.2	Equal priors	20

*Based upon work partially supported by grants from the U.S. Army Research Office, the National Science Foundation and the National Institutes of Health.

4.2.3	Unequal misclassification costs: hepatitis data	21
4.2.4	Nearest-neighbor estimates: car data	23
4.2.5	Kernel density estimates: car data	32
5	Regression	45
5.1	Stepwise least-squares	45
5.2	Least-squares simple polynomial	51
5.3	ANCOVA models	55
5.4	Quantile regression	62
5.5	Least median of squares	68
5.6	Poisson regression with offset	76
5.7	Censored response data	85
5.8	Multi-response data	94
5.9	Longitudinal data	99
5.10	Subgroup identification	111
5.11	Differential item functioning	117
6	Tree ensembles	124
6.1	Bagged GUIDE	124
6.2	GUIDE forest	127
7	Importance scores	130
7.1	Baseball data example	130
8	Other features	134
8.1	Pruning with test samples	134
8.2	Prediction of test samples	134
8.3	GUIDE in R and in simulations	134
8.4	Generation of powers and products	135
8.5	Data formatting functions	136

1 Warranty disclaimer

Redistribution and use in binary forms, with or without modification, are permitted provided that the following condition is met:

Redistributions in binary form must reproduce the above copyright notice, this condition and the following disclaimer in the documentation and/or other materials provided with the distribution.

THIS SOFTWARE IS PROVIDED BY WEI-YIN LOH “AS IS” AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL WEI-YIN LOH BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

The views and conclusions contained in the software and documentation are those of the author and should not be interpreted as representing official policies, either expressed or implied, of the University of Wisconsin.

2 Introduction

GUIDE stands for *Generalized, Unbiased, Interaction Detection and Estimation*. It is the only classification and regression tree algorithm with all these features:

1. Unbiased variable selection.
2. Kernel and nearest-neighbor node models for classification trees.
3. Weighted least squares, least median of squares, quantile, Poisson, and relative risk (proportional hazards) regression models.
4. Univariate, multivariate, and longitudinal response variables.
5. Pairwise interaction detection at each node.
6. Linear splits on two variables at a time for classification trees.
7. Categorical variables for splitting only, or for both splitting and fitting (via 0-1 dummy variables), in regression tree models.
8. Ranking and scoring of predictor variables.
9. Tree ensembles (bagging and forests).

Table 1: Comparison of GUIDE, QUEST, CRUISE, CART, and C4.5 classification tree algorithms. Node models: S = simple, K = kernel, L = linear discriminant, N = nearest-neighbor.

	GUIDE	QUEST	CRUISE	CART	C4.5
Unbiased splits	Yes	Yes	Yes	No	No
Splits per node	2	2	≥ 2	2	2
Interaction detection	Yes	No	Yes	No	No
Importance ranking	Yes	No	No	Yes	No
Class priors	Yes	Yes	Yes	Yes	No
Misclassification costs	Yes	Yes	Yes	Yes	No
Linear splits	Yes	Yes	Yes	Yes	No
Categorical splits	Subsets	Subsets	Subsets	Subsets	Atoms
Node models	S, K, N	S	S, L	S	S
Missing values	Special	Imputation	Surrogate	Surrogate	Weights
Tree diagrams	Text and L ^A T _E X			Proprietary	Text
Bagging	Yes	No	No	No	No
Forests	Yes	No	No	No	No

Tables 1 and 2 compare the features of GUIDE with CRUISE (Kim and Loh, 2001, 2003), QUEST (Loh and Shih, 1997), C4.5 (Quinlan, 1993), RPART¹, and M5' (Quinlan, 1992; Witten and Frank, 2000).

The GUIDE algorithm is documented in Loh (2002) for regression trees and Loh (2009) for classification trees. Loh (2008a), Loh (2011) and Loh (2014) review the subject. Advanced features of the algorithm are reported in Chaudhuri and Loh (2002), Loh (2006b), Kim et al. (2007), Loh et al. (2007), and Loh (2008b). For a list of third-party applications of GUIDE, CRUISE, QUEST, and the logistic regression tree algorithm LOTUS (Chan and Loh, 2004; Loh, 2006a), see <http://www.stat.wisc.edu/~loh/apps.html>

This manual illustrates the use of the program and interpretation of the output.

¹RPART is an implementation of CART (Breiman et al., 1984) in R. CART is a registered trademark of California Statistical Software, Inc.

Table 2: Comparison of GUIDE, CART and M5' regression tree algorithms

	GUIDE	CART	M5'
Unbiased splits	Yes	No	No
Pairwise interaction detection	Yes	No	No
Importance scores	Yes	Yes	No
Loss functions	Weighted least squares, least median of squares, quantile, Poisson, proportional hazards	Least squares, least absolute deviations	Least squares only
Survival, longitudinal and multi-response data	Yes, yes, yes	No, no, no	No, no, no
Node models	Constant, multiple, stepwise linear, polynomial, ANCOVA	Constant only	Constant and stepwise
Linear models	Multiple or stepwise (forward and forward-backward)	N/A	Stepwise
Variable roles	Split only, fit only, both, neither, weight, censored, offset	Split only	Split and fit
Categorical variable splits	Subsets of categorical values	Subsets	0-1 variables
Tree selection	Pruning or stopping rules	Pruning only	Pruning only
Tree diagrams	Text and \LaTeX	Proprietary	PostScript
Operation modes	Interactive and batch	Interactive and batch	Interactive
Case weights	Yes	Yes	No
Transformations	Powers and products	No	No
Missing values in split variables	Missing values treated as a special category	Surrogate splits	Imputation
Missing values in linear predictors	Choice of separate constant models or mean imputation	N/A	Imputation
Bagging & forests	Yes & yes	No & no	No & no
Data conversions	ARFF, C4.5, Minitab, R, SAS, Statistica, Systat, CSV	No	No

2.1 Installation

GUIDE is available free from www.stat.wisc.edu/~loh/guide.html in the form of compiled 32- and 64-bit executables for Linux, Mac OS X, and Windows on Intel and compatible processors.

Mac OS X and Linux: Make the unzipped file executable by issuing this command in a **Terminal** application in the folder where the file is located: `chmod a+x guide`

Mac OS X only: The Mac OS X version requires **Xcode** and **gfortran** to be installed. To ensure that the gfortran libraries are placed in the right place, follow the steps:

1. Install **Xcode** from <https://developer.apple.com/xcode/downloads/>.
2. Go to <http://hpc.sourceforge.net> and download file `gcc-4.9-bin.tar.gz` to your Downloads folder. The direct link to the file is <http://prdownloads.sourceforge.net/hpc/gcc-4.9-bin.tar.gz?download>
3. Open a **Terminal** window and type (or copy and paste):
 - (a) `cd ~/Downloads`
 - (b) `gunzip gcc-4.9-bin.tar.gz`
 - (c) `sudo tar -xvf gcc-4.9-bin.tar -C /`

2.2 L^AT_EX

GUIDE uses the public-domain software L^AT_EX (<http://www.ctan.org>) to produce tree diagrams. The specific locations are:

Linux: TeX Live <http://www.tug.org/texlive/>

Mac: MacTeX <http://tug.org/mactex/>

Windows: proTeXt <http://www.tug.org/protext/>

After L^AT_EX is installed, a pdf file of a L^AT_EX file, called `diagram.tex` say, produced by GUIDE can be obtained by typing these three commands in a terminal window:

1. `latex diagram`
2. `dvips diagram`

3. ps2pdf diagram.ps

The first command produces a file called `diagram.dvi` which the second command uses to create a postscript file called `diagram.ps`. The latter can be viewed and printed if a postscript viewer (such as *Preview* for the Mac) is installed. If no postscript viewer is available, the last command can be used to convert the postscript file into a pdf file, which can be viewed and printed with *Adobe Reader*.

The file `diagram.tex` can be edited to change colors, node sizes, etc. See, e.g., <http://tug.org/PSTricks/main.cgi/>.

3 Program operation

3.1 Required files

The GUIDE program requires two text files for input.

Data file: This file contains the training sample. Each file record consists of observations on the response (i.e., dependent) variable, the predictor (i.e., X or independent) variables, and optional weight and time variables. Entries in each record are comma, space, or tab delimited (multiple spaces are treated as one space, but not for commas). A record can occupy more than one line in the file, but each record must begin on a new line.

Values of categorical variables can contain any ascii character except single and double quotation marks, which are used to enclose values that contain spaces and commas. Values can be up to 60 characters long. Class labels are truncated to 10 characters in tabular displays.

Description file: This provides information about the name and location of the data file, names and column positions of the variables, and their roles in the analysis. This file permits different models to be fitted by changing the roles of the variables. We use the files `irisdesc.txt` and `irisdata.txt` (both obtainable from <http://www.stat.wisc.edu/~loh/guide.html>) to illustrate. The data give the sepal lengths and widths and the petal lengths and widths of 150 iris flowers. The response variable is the type of iris flower. The contents of `irisdesc.txt` are:

```
irisdata.txt
"?"
```

```
column, varname, vartype
1 sepallen n
2 sepalwid n
3 petallen n
4 petalwid n
5 class d
```

The first line of the file `irisdsc.txt` gives the name of the training sample file. If the data file `irisdata.txt` is not in the folder where GUIDE is installed, its full path (such as "c:\data\irisdata.txt") is needed. The second line gives the the missing value code, which can be up to 80 characters long. If it contains non alphanumeric characters, it must be surrounded by quotation marks. A missing value code must appear in the second line of the file even if there are no missing values in the data (in that case any character string not present among the data values can be used). The third line contains three character strings to indicate the column headers of the subsequent lines. The position, name and role of each variable comes next (in that order), with one line for each variable.

Variable names must begin with an alphabet and be not more than 60 characters long. If a name contains non-alphanumeric characters, it must be enclosed in matching single or double quotes. Spaces and the four characters #, %, {, and } are replaced by dots (periods) if they appear in a name. Variable names are truncated to 10 characters in tabular output. Leading and trailing spaces are dropped.

The following roles for the variables are permitted. Lower and upper case letters are accepted.

- b** Categorical variable that is used both for splitting and for node modeling in regression. It is transformed to 0-1 dummy variables for node modeling. It is converted to **c** type for classification.
- c** Categorical variable used for splitting only.
- d** Dependent variable. Except for multi-response data (see Sec. 5.8), there can only be one such variable. In the case of relative risk models, this is the **death** indicator. The variable can take character string values for classification.
- f** Numerical variable used only for fitting the linear models in the nodes of the tree. It is not used for splitting the nodes and is disallowed in classification.

- n** Numerical variable used both for splitting the nodes and for fitting the node models. It is converted to type **s** in classification.
- r** Categorical treatment (**Rx**) variable used only for fitting the linear models in the nodes of the tree. It is not used for splitting the nodes. If this variable is present, all **n** variables are automatically changed to **s**.
- s** Numerical-valued variable only used for splitting the nodes. It is not used as a regressor in the linear models. This role is suitable for ordinal categorical variables if they are given numerical values that reflect the orderings.
- t** Survival **time** (for proportional hazards models) or observation **time** (for longitudinal models) variable.
- w** Weight variable for weighted least squares regression or for excluding observations in the training sample from tree construction. See section 8.2 for the latter. Except for longitudinal models, a record with a missing value in a **d**, **t**, or **z**-variable is automatically assigned zero weight.
- x** Excluded variable. This allows models to be fitted to different subsets of the variables without reformatting the data file.
- z** Offset variable used only in Poisson regression.

GUIDE runs within a **terminal window** of the computer operating system.

Do not double-click its icon!

Linux. Any terminal program will do.

Mac OS X. The program is called **Terminal**; it is in the **Applications Folder**.

Windows. The terminal program is started from the **Start button** by choosing **All Programs → Accessories → Command Prompt**

3.2 Input file creation

GUIDE is started by typing its (lowercase) name in a terminal. The preferred way is to create an input file (option 1 below) for subsequent execution. The input file may be edited if you wish to change some input parameters later. In the following, the sign (**>**) is the terminal prompt (not to be typed!).

```
C:\Users\Weiyin> guide
GUIDE Classification and Regression Trees and Forests
Version 19.0 (build date: March 6, 2015)
Copyright (c) 1997-2015 Wei-Yin Loh. All rights reserved.
```

This software is based upon work supported by the U.S. Army Research Office, the National Science Foundation and the National Institutes of Health.

Choose one of the following options:

- 0. Read the warranty disclaimer
 - 1. Create an input file for batch run
 - 2. Fit a model without creating input file
 - 3. Convert data to other formats
 - 4. Variable importance scoring and differential item functioning
- Input your choice:

The meanings of these options are:

- 0. Print the warranty disclaimer.
- 1. Create an input file for subsequent execution.
- 2. Run the program right away without creating an input file.
- 3. Convert the data file into a format suitable for importation into database, spreadsheet, or statistics software. See Table 2 for the statistical packages supported. Section 8.5 has an example.
- 4. Importance scoring of variables and identification of items with differential item functioning.

4 Classification

4.1 Default: univariate splits

We first show how to generate an input file to produce a classification tree from the data in the file `irisdata.txt`, using the default options. Whenever you are prompted for a selection, there is usually range of permissible values given within square brackets and a default choice (indicated by the symbol `<cr>=`). The default may be selected by pressing the ENTER or RETURN key. Annotations are printed in *blue italics* in this manual.

4.1.1 Input file creation with default options

- 0. Read the warranty disclaimer
- 1. Create an input file for batch run

```

2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: irisin.txt
  This file will store your answers to the prompts.
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for
data conversion ([1:3], <cr>=1):
  Press the ENTER or RETURN key to accept the default selection.
Name of batch output file: irisout.txt
  This file will contain the results when you apply the input file to GUIDE later.
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
  Option 2 is for bagging and random forest-type methods.
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):
  The default option will produce a traditional classification tree.
  Choose option 2 for more advanced features.
Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): irisdesc.txt
Reading data description file ...
Training sample file: irisdata.txt
  The name of the data set is read from the description file.
  Some information about the data are printed in the next few lines.
Missing value code: ?
Warning: N variables changed to S
  This warning is triggered because we are fitting a classification model.
Dependent variable is class
Length of longest data entry = 11
Total number of cases =      150
Number of classes =      3
Checking data ...
Class name      Num. cases  Proportion
Setosa          50          0.33333333
Versicolour     50          0.33333333
Virginica       50          0.33333333
  Total  #cases w/  #missing
  #cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
    150      0      0      0      0      0      4      0      0
No. cases used for training =      150
Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
Input 1, 2, or 3 ([1:3], <cr>=1):
  See below for examples of equal priors and specified priors.
Choose 1 for unit misclassification costs, 2 to input costs from a file
Input 1 or 2 ([1:2], <cr>=1):
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
  Choose option 2 if you do not want LaTeX code.

```

```

Input file name to store LaTeX code (use .tex as suffix): iristree.tex
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node IDs and fitted values: irisfit.txt
This file will contain the node number and predicted class for each observation.
Input file is created!
Run GUIDE with the command: guide < irisin.txt
Press ENTER or RETURN to quit

```

4.1.2 Contents of input.txt

Here are the contents of the input file:

```

123321      (do not edit this file unless you know what you are doing)
  19.0      (version of GUIDE that generated this file)
    1      (1=model fitting, 2=importance or DIF scoring, 3=data conversion)
"irisout.txt" (name of output file)
    1      (1=one tree, 2=ensemble)
    1      (1=classification, 2=regression)
    1      (1=simple model, 2=nearest-neighbor, 3=kernel)
    1      (0=linear first, 1=univariate first, 2=skip linear, 3=skip linear and interaction tests)
    1      (1=prune by CV, 2=by test sample, 3=no pruning)
"irisdsc.txt" (name of data description file)
    10     (number of cross-validations)
    1      (1=mean-based CV tree, 2=median-based CV tree)
    0.500   (SE number for pruning)
    1      (1=estimated priors, 2=equal priors, 3=other priors)
    1      (1=unit misclassification costs, 2=other)
    2      (1=split point from quantiles, 2=use exhaustive search)
    1      (1=default max number of split levels, 2=specify no. in next line)
    1      (1=default min node size, 2=specify node size in next line)
    1      (1=write latex, 2=skip latex)
"iristree.tex" (latex file name)
    1      (1=vertical tree, 2=sideways tree)
    1      (1=include node numbers, 2=exclude)
    1      (1=number all nodes, 2=only terminal nodes)
    1      (1=color terminal nodes, 2=no colors)
    1      (0=#errors, 1=class sizes in nodes, 2=nothing)
    1      (1=no storage, 2=store fit and split variables, 3=store split variables and values)
    2      (1=do not save individual fitted values and node IDs, 2=save in a file)
"irisfit.txt" (file name for fitted values and node IDs)
    1      (1=do not save terminal node IDs for importance scoring in a file, 2=save them)
    1      (1=do not write R function, 2=write R function)

```

GUIDE reads only the first entry in each line; the remainder of the line is for human consumption. Because each question depends on the answers you have given to pre-

vious questions, only some of the entries in the input file may be changed. Examples are the SE number, maximum number of split levels, minimum node size, vertical vs. sideways tree diagram, coloring terminal nodes, and printing of class sizes in the L^AT_EX tree diagram.

4.1.3 Executing the program and interpreting the output

Once the input file is generated, GUIDE can be executed with the command:

```
guide < irisin.txt
```

Following is an annotated copy of the contents of the output file.

```
Classification tree
Pruning by cross-validation
Data description file: irisdisc.txt
Training sample file: irisdata.txt
Missing value code: ?
Warning: N variables changed to S
Dependent variable is class
Length of longest data entry = 11
Number of classes = 3
Class name      Num. cases  Proportion
Setosa           50      0.33333333
Versicolour      50      0.33333333
Virginica        50      0.33333333
This gives the number of observations in each class.
Summary information (without x variables)
d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,
n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
For categorical variables, #categories include one for missing values
Column  Variable      Variable  Minimum      Maximum      Number of      Number
number  name            type      value         value         categories      missing
    1  sepallen          s  4.3000E+00  7.9000E+00
    2  sepalwid          s  2.0000E+00  4.4000E+00
    3  petallen          s  1.0000E+00  6.9000E+00
    4  petalwid          s  1.0000E-01  2.5000E+00
    5  class              d
                                     3
This shows the type and minimum and maximum values of each ordered variable.
Total #cases w/  #missing
#cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
    150      0      0      0      0      0      4      0      0
No. cases used for training = 150
This shows the number of each type of variable.
Univariate split highest priority
```

Interaction and linear splits 2nd and 3rd priorities
 Pruning by v-fold cross-validation, with v = 10
 Selected tree is based on mean of CV estimates
 Simple node models
 Estimated priors
 Unit misclassification costs
 Split values for N and S variables based on exhaustive search
 Max number of split levels = 10
 Minimum node size = 3
 Number of SE's for pruned tree = 5.0000E-01

Size and CV mean cost and SE of subtrees:

Tree	#Tnodes	Mean Cost	SE(Mean)	BSE(Mean)	Median Cost	BSE(Median)
0	6	6.000E-02	1.939E-02	6.091E-03	6.667E-02	0.000E+00
1	5	6.000E-02	1.939E-02	6.091E-03	6.667E-02	0.000E+00
2**	3	4.667E-02	1.722E-02	1.579E-02	3.333E-02	3.111E-02
3	2	3.333E-01	3.849E-02	0.000E+00	3.333E-01	0.000E+00
4	1	6.667E-01	3.849E-02	0.000E+00	6.667E-01	0.000E+00

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

* tree, ** tree, + tree, and ++ tree all the same

The tree with the smallest mean CV cost is marked with an asterisk.

The selected tree is marked with two asterisks; it is the smallest one having mean CV cost within the specified standard error (SE) bounds.

The mean CV costs and SEs are given in the 3rd and 4th columns.

The other columns are bootstrap estimates used for experimental purposes.

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Node cost is node misclassification cost divided by number of training cases

Node label	Total cases	Train cases	Predicted class	Node cost	Split variables	Interacting variable
1	150	150	Setosa	6.667E-01	petalwid	
2T	50	50	Setosa	0.000E+00	-	
3	100	100	Versicolour	5.000E-01	petalwid	
6T	54	54	Versicolour	9.259E-02	sepalwid :petalwid	
7T	46	46	Virginica	2.174E-02	-	

This shows the tree structure in tabular form. A node with label k has its left and right child nodes are labeled 2k and 2k+1, respectively. Terminal nodes are indicated with the symbol T. The notation ':petawid' in node 6 indicates that

the variable petalwid has an interaction with the split variable sepalwid.

Number of terminal nodes of final tree: 3

Total number of nodes of final tree: 5

Classification tree:

The tree structure is shown next in indented text form.

Node 1: petalwid <= 0.80000

Node 2: Setosa

Node 1: petalwid > 0.80000 or ?

Node 3: petalwid <= 1.75000 or ?

Node 6: Versicolour

Node 3: petalwid > 1.75000 and not ?

Node 7: Virginica

Node 1: Intermediate node

A case goes into Node 2 if petalwid <= 8.0000000E-01

petalwid mean = 1.1987E+00

ClassName	Number	ClassPrior
Setosa	50	0.3333
Versicolou	50	0.3333
Virginica	50	0.3333

Number of training cases misclassified = 100

Predicted class is Setosa

Node 2: Terminal node

ClassName	Number	ClassPrior
Setosa	50	1.0000
Versicolou	0	0.0000
Virginica	0	0.0000

Number of training cases misclassified = 0

Predicted class is Setosa

Node 3: Intermediate node

A case goes into Node 6 if petalwid <= 1.7500000E+00 or ?

petalwid mean = 1.6760E+00

ClassName	Number	ClassPrior
Setosa	0	0.0000
Versicolou	50	0.5000
Virginica	50	0.5000

Number of training cases misclassified = 50

Predicted class is Versicolour

Node 6: Terminal node

```

ClassName      Number ClassPrior
Setosa         0      0.0000
Versicolou     49      0.9074
Virginica      5      0.0926
Number of training cases misclassified = 5
Predicted class is Versicolour
-----

```

```

Node 7: Terminal node
ClassName      Number ClassPrior
Setosa         0      0.0000
Versicolou     1      0.0217
Virginica      45     0.9783
Number of training cases misclassified = 1
Predicted class is Virginica
-----

```

Classification matrix for training sample:

Predicted	True class		
class	Setosa	Versicolo	Virginica
Setosa	50	0	0
Versicolou	0	49	5
Virginica	0	1	45
Total	50	50	50

Number of cases used for tree construction = 150

Number misclassified = 6

Resubstitution est. of mean misclassification cost = 4.0000000000000001E-002

This is the mean misclassification cost estimated from the training sample.

Observed and fitted values are stored in irisfit.txt

LaTeX code for tree is in iristree.tex

Elapsed time in seconds: 2.25169994E-02

The left side of Figure 1 shows the classification tree drawn by LaTeX using the file `iristree.tex` and the top lines of the file `irisfit.txt` are shown below. The order of the lines correspond to the order of the observations in the training sample file. The first column (labeled `train`) indicates whether the observation is used (“y”) or not used (“n”) to fit the model. Since we used the entire data set to fit the model here, all the entries in the first column are y.

train	node	observed	predicted
y	2	"Setosa"	"Setosa"
y	2	"Setosa"	"Setosa"
y	2	"Setosa"	"Setosa"
y	2	"Setosa"	"Setosa"
y	2	"Setosa"	"Setosa"
y	2	"Setosa"	"Setosa"

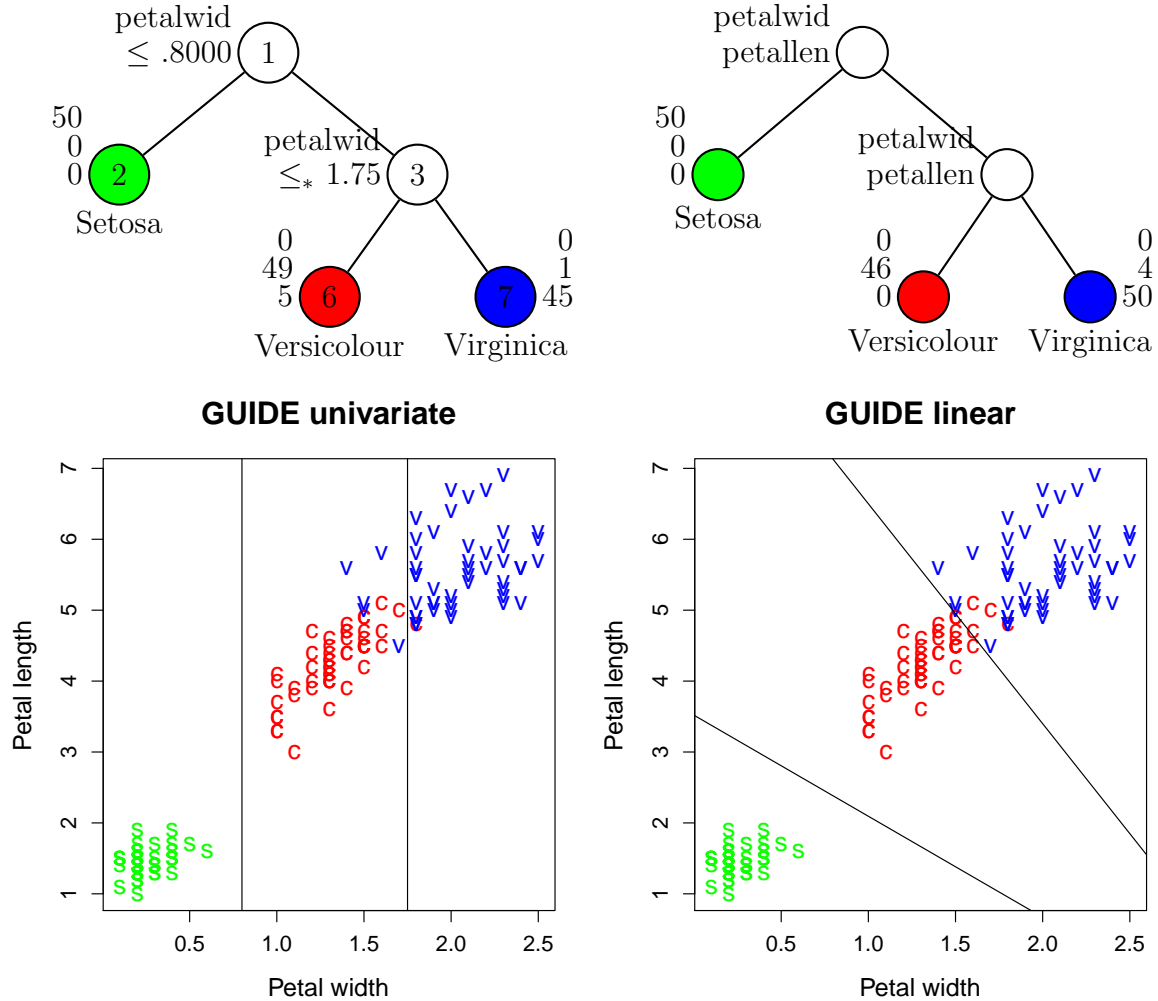


Figure 1: GUIDE classification trees and plots of the data partitions for the iris data using estimated priors and unit misclassification costs. The tree on the left uses univariate splits. At each intermediate node, an observation goes to the left branch if and only if the condition is satisfied. The symbol \leq_* denotes \leq or missing. The tree on the right uses linear splits on two variables at a time. Predicted classes (based on estimated misclassification cost) below terminal nodes; sample sizes for Setosa, Versicolour, and Virginica, respectively, beside nodes.

```
y          2    "Setosa"  "Setosa"
```

4.2 Non-default options

4.2.1 Linear splits

The above example uses the default options for classification trees. Other features are available with non-default options. We show how to obtain the linear splits shown in Figure 1 here.

```
0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: input2.txt
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: output2.txt
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1): 2
  Choosing 2 opens up the other options.
Input 1 for simple, 2 for nearest-neighbor, 3 for kernel method
([1:3], <cr>=1):
  Options 2 and 3 yield nearest-neighbor and kernel discriminant node models.
Input 0 for linear, interaction and univariate splits (in this order),
    1 for univariate, linear and interaction splits (in this order),
    2 to skip linear splits,
    3 to skip linear and interaction splits:
Input your choice ([0:3], <cr>=1): 0
  Option 1 is the default.
Input 1 to prune by CV, 2 by test sample, 3 for no pruning ([1:3], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): irisdisc.txt
Reading data description file ...
Training sample file: irisdata.txt
Missing value code: ?
Warning: N variables changed to S
Dependent variable is class
Length of longest data entry = 11
Total number of cases =      150
Number of classes =        3
Checking data ...
```

Class name	Num. cases	Proportion
Setosa	50	0.33333333
Versicolour	50	0.33333333
Virginica	50	0.33333333

Total #cases	#cases w/ miss. D	#missing ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
150	0	0	0	0	0	4	0	0

No. cases used for training = 150
 Default number of cross-validations = 10
 Input 1 to accept the default, 2 to change it ([1:2], <cr>=1):
 Best tree may be chosen based on mean or median CV estimate
 Input 1 for mean-based, 2 for median-based ([1:2], <cr>=1):
 Input number of SEs for pruning ([0.00:1000.00], <cr>=0.50):
 Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
 Input 1, 2, or 3 ([1:3], <cr>=1):
 Choose 1 for unit misclassification costs, 2 to input costs from a file
 Input 1 or 2 ([1:2], <cr>=1):
 Choose a split point selection method for numerical variables:
 Choose 1 to use faster method based on sample quantiles
 Choose 2 to use exhaustive search
 Input 1 or 2 ([1:2], <cr>=2):
 Default max number of split levels = 10
 Input 1 to accept this value, 2 to change it ([1:2], <cr>=1):
 Default minimum node sample size is 10
 Input 1 to use the default value, 2 to change it ([1:2], <cr>=1):
 Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
 Input file name to store LaTeX code (use .tex as suffix): tree2.tex
 Input 1 for a vertical tree, 2 for a sideways tree ([1:2], <cr>=1):
 Input 1 to include node numbers, 2 to omit them ([1:2], <cr>=1): 2
Choosing 2 will give a tree with no node labels.
 Input 1 to color terminal nodes, 2 otherwise ([1:2], <cr>=1):
 Choose amount of detail in nodes of LaTeX tree diagram
 Input 0 for #errors, 1 for class sizes, 2 for nothing ([0:2], <cr>=1):
Choose 2 for very large trees.
 You can store the variables and/or values used to split and fit in a file
 Choose 1 to skip this step, 2 to store split variables and their values
 Input your choice ([1:2], <cr>=1):
 Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1):
 Input 2 to save terminal node IDs for importance scoring; 1 otherwise ([1:2], <cr>=1):
 Input 2 to write R function for predicting new cases, 1 otherwise ([1:2], <cr>=1): 2
 Input file name: pred.r
 Input file is created!

Running the program with this input file yields the following results and the \LaTeX tree diagram and partitions on the right side of Figure 1.

Node 1: $7.0205078E-01 * \text{petallen} + \text{petalwid} \leq 2.4700244E+00$ or ?

```

Node 2: Setosa
Node 1: 7.0205078E-01 * petallen + petalwid > 2.4700244E+00
Node 3: 3.2242660E-01 * petallen + petalwid <= 3.0960117E+00 or ?
Node 6: Versicolour
Node 3: 3.2242660E-01 * petallen + petalwid > 3.0960117E+00
Node 7: Virginica

```

The R file `pred.r` contains this function

```

predicted <- function(){
  if(is.na(petalwid) | is.na(petallen) | 0.70205077945722660 * petallen + petalwid
    <= 2.4700244096702053){
    nodeid <- 2
    predict <- "Setosa"
  } else {
    if(is.na(petalwid) | is.na(petallen) | 0.32242659679310148 * petallen + petalwid
      <= 3.0960116541258524){
      nodeid <- 6
      predict <- "Versicolour"
    } else {
      nodeid <- 7
      predict <- "Virginica"
    }
  }
  return(c(nodeid,predict))
}

```

4.2.2 Equal priors

If a data set has one dominant class, a classification tree will often be null after pruning, because it is hard to beat the classifier that predicts every observation to belong to the dominant class. One way to obtain a non-null tree is to specify equal priors. We illustrate this with the hepatitis data set from

<http://archive.ics.uci.edu/ml/datasets/Hepatitis>. The files `hepdsc.txt` and `hepdatt.txt` are obtainable from <http://www.stat.wisc.edu/~loh/guide.html>. The data consist of observations from 155 individuals, of whom 32 are labeled “die” and 123 labeled “live”. The contents of the description file `hepdsc.txt` are:

```

hepdatt.txt
"?"
column, var, type
1 CLASS d
2 AGE n
3 SEX c
4 STEROID c
5 ANTIVIRALS c
6 FATIGUE c

```

```

7 MALAISE c
8 ANOREXIA c
9 BIGLIVER c
10 FIRMLIVER c
11 SPLEEN c
12 SPIDERS c
13 ASCITES c
14 VARICES c
15 BILIRUBIN n
16 ALKPHOSPHATE n
17 SGOT n
18 ALBUMIN n
19 PROTIME n
20 HISTOLOGY c

```

Using the default estimated priors yields a null tree with no splits. To obtain a nonnull tree, choose “2” for equal priors in this dialog step:

```

Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
Input 1, 2, or 3 ([1:3], <cr>=1): 2

```

The resulting tree in text form is:

```

Node 1: ASCITES = "yes"
  Node 2: die
Node 1: ASCITES /= "yes"
  Node 3: SPIDERS = "?", "yes"
    Node 6: die
  Node 3: SPIDERS /= "?", "yes"
    Node 7: live

```

The tree drawn by L^AT_EX is shown on the left of Figure 2. Nodes that predict the same class have the same color. Since the ratio of “die” to “live” classes is 32:123, the effect of equal priors is to treat one “die” observation as equivalent to $r = 123/32 = 3.84375$ “live” observations. Therefore a terminal node is classified as “die” if its ratio of “live” to “die” observations is less than r .

4.2.3 Unequal misclassification costs: hepatitis data

So far, we have assumed that the cost of misclassifying a “die” observation as “live” is the same as the opposite. Another way to obtain a nonnull tree for the hepatitis data is to use unequal misclassification costs. For example, if we think that the cost

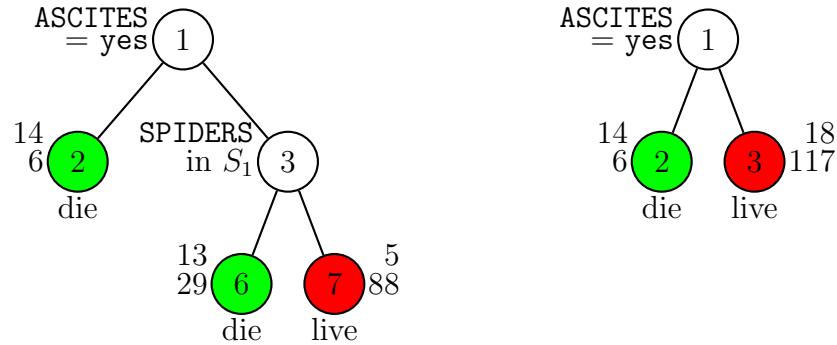


Figure 2: GUIDE 0.50-SE classification trees for predicting **CLASS** with equal priors (left) and unequal misclassification costs (right). At each split, an observation goes to the left branch if and only if the condition is satisfied. Set $S_1 = \{?, \text{yes}\}$. Predicted classes (based on estimated misclassification cost) printed below terminal nodes; sample sizes for die and live, respectively, beside nodes.

of misclassifying a “die” observation as “live” is four times that of the opposite, we will use the misclassification cost matrix

$$C = \begin{pmatrix} 0 & 1 \\ 4 & 0 \end{pmatrix}$$

where $C(i, j)$ denotes the cost of classifying an observation as class i given that it belongs to class j . Note that GUIDE sorts the class values in alphabetical order, so that “die” is treated as class 1 and “live” as class 2 here. This matrix is saved in the text file `cost.txt` which has these two lines:

```
0 1
4 0
```

The following lines in the input file generation step shows where this file is used:

```
Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
Input 1, 2, or 3 ([1:3], <cr>=1):
Choose 1 for unit misclassification costs, 2 to input costs from a file
Input 1 or 2 ([1:2], <cr>=1): 2
Input the name of a file containing the cost matrix C(i|j),
where C(i|j) is the cost of classifying class j as class i
The rows of the matrix must be in alphabetical order of the class names
Input name of file: cost.txt
```

The resulting tree is shown on the right of Figure 2.

4.2.4 Nearest-neighbor estimates: car data

The data file `drive.txt` gives the specifications and prices of 428 new cars in the 2004 model year. The data come from the *J. Statistics Education* website http://www.amstat.org/publications/jse/jse_data_archive.htm. Using the description file `drivedsc.txt` whose contents follow, the tree model for predicting Drive type misclassifies 90 of the 428 observations.

```
drive.txt
"*"
c1 c2 c3
1  Region x
2  Import x
3  Make c
4  Model x
5  Type c
6  Drive d
7  SC x
8  SUV x
9  Wagon x
10 Minivan x
11 Pickup x
12 Allwheel x
13 Rearwheel x
14 Rprice n
15 Dcost n
16 Enginsz n
17 Cylin n
18 Hp n
19 City n
20 Hwy n
21 Weight n
22 Whlbase n
23 Length n
24 Width n
```

In the examples so far, the observations in each terminal node of a classification tree are all predicted to belong to the class that minimizes the node misclassification cost. This method can be inefficient if the data are difficult to classify with a small number of splits. Alternatively, we can fit a classification model to the data in each node and use it to classify individual observations in the node. GUIDE has two means to achieve this: nearest-neighbor and kernel discrimination. For nearest-neighbor, an observation in a node is classified to the plurality class among observations within its neighborhood. The neighborhood is defined to be the whole node if the split variable is categorical. We illustrate this for the car data with the following input file generation log.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: drive.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: drive.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1): 2
Input 1 for simple, 2 for nearest-neighbor, 3 for kernel method ([1:3], <cr>=1): 2
  Choose nearest-neighbor method here.
Input 1 for univariate, 2 for bivariate preference ([1:2], <cr>=1):
  Default is univariate kernels.
Input 1 for interaction tests, 2 to skip them ([1:2], <cr>=1):
Input 1 to prune by CV, 2 by test sample, 3 for no pruning ([1:3], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): drivesc.txt
Reading data description file ...
Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Total number of cases =      428
Number of classes =      3
Col. no. Categorical variable   #levels   #missing values
      3 Make                    38         0
      5 Type                    6         0
Checking data ...
Class name      Num. cases   Proportion
4wd              94    0.21962617
fwd             224    0.52336449
rwd             110    0.25700935
      Total  #cases w/  #missing
      #cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
      428      0      0      10      0      0      11      0      2
No. cases used for training =      428
Default number of cross-validations =      10
Input 1 to accept the default, 2 to change it ([1:2], <cr>=1):
Best tree may be chosen based on mean or median CV estimate
Input 1 for mean-based, 2 for median-based ([1:2], <cr>=1):
Input number of SEs for pruning ([0.00:1000.00], <cr>=0.50):

```



```

Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
Input 1, 2, or 3 ([1:3], <cr>=1):
Choose 1 for unit misclassification costs, 2 to input costs from a file
Input 1 or 2 ([1:2], <cr>=1):
Choose a split point selection method for numerical variables:
Choose 1 to use faster method based on sample quantiles
Choose 2 to use exhaustive search
Input 1 or 2 ([1:2], <cr>=2):
Default max number of split levels =          10
Input 1 to accept this value, 2 to change it ([1:2], <cr>=1):
Default minimum node sample size is 10
Input 1 to use the default value, 2 to change it ([1:2], <cr>=1):
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): drive.tex
Input 1 for a vertical tree, 2 for a sideways tree ([1:2], <cr>=1):
Input 1 to include node numbers, 2 to omit them ([1:2], <cr>=1):
Input 1 to number all nodes, 2 to number leaves only ([1:2], <cr>=1):
Input 1 to color terminal nodes, 2 otherwise ([1:2], <cr>=1):
Choose amount of detail in nodes of LaTeX tree diagram
Input 0 for #errors, 1 for class sizes, 2 for nothing ([0:2], <cr>=1):
You can store the variables and/or values used to split and fit in a file
Choose 1 to skip this step, 2 to store split and fit variables,
3 to store split variables and their values
Input your choice ([1:3], <cr>=1):
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1):
Input 2 to save terminal node IDs for importance scoring; 1 otherwise ([1:2], <cr>=1):
Input file is created!

```

Results

```

Classification tree
Pruning by cross-validation
Data description file: drivedsc.txt
Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Number of classes = 3
Class      #Cases   Proportion
4wd         94     0.21962617
fwd        224     0.52336449
rwd        110     0.25700935

Summary information (without x variables)

```

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical, n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight

Column	Name		Minimum	Maximum	#Categories	#Missing
3	Make	c			38	
5	Type	c			6	
6	Drive	d			3	
14	Rprice	s	1.0280E+04	1.9246E+05		
15	Dcost	s	9.8750E+03	1.7356E+05		
16	Enginsz	s	1.3000E+00	8.3000E+00		
17	Cylin	s	-1.0000E+00	1.2000E+01		
18	Hp	s	7.3000E+01	5.0000E+02		
19	City	s	1.0000E+01	6.0000E+01		
20	Hwy	s	1.2000E+01	6.6000E+01		
21	Weight	s	1.8500E+03	7.1900E+03		
22	Whlbase	s	8.9000E+01	1.4400E+02		
23	Length	s	1.4300E+02	2.2800E+02		
24	Width	s	6.4000E+01	8.1000E+01		

Total	#cases w/	#missing							
#cases	miss. D	ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var	
428	0	0	10	0	0	11	0	2	

No. cases used for training: 428

Univariate split highest priority

Interaction splits 2nd priority; no linear splits

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Nearest-neighbor node models

Univariate preference

Estimated priors

Unit misclassification costs

Split values for N and S variables based on exhaustive search

Max number of split levels = 10

Minimum node size = 10

Number of SE's for pruned tree = 5.0000E-01

Size and CV mean cost and SE of subtrees:

Tree	#Tnodes	Mean Cost	SE(Mean)	BSE(Mean)	Median Cost	BSE(Median)
1	19	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
2	18	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
3	17	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
4	16	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
5	15	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
6	14	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
7	13	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
8	12	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02

9	11	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
10	10	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
11	9	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
12**	7	2.196E-01	2.001E-02	1.175E-02	2.118E-01	1.284E-02
13	6	2.453E-01	2.080E-02	9.596E-03	2.442E-01	1.558E-02
14	5	2.453E-01	2.080E-02	9.596E-03	2.442E-01	1.558E-02
15	3	2.547E-01	2.106E-02	1.652E-02	2.558E-01	2.811E-02
16	2	2.640E-01	2.131E-02	1.735E-02	2.674E-01	2.446E-02
17	1	3.808E-01	2.347E-02	1.384E-02	3.721E-01	1.732E-02

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

* tree, ** tree, + tree, and ++ tree all the same

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Node cost is node misclassification cost divided by number of training cases

Node label	Total cases	Train cases	Predicted class	Node cost	Split variable followed by (+)fit variable(s)
1	428	428	fwd	3.832E-01	Type +Type
2	73	73	rwd	3.425E-01	Type +Type
4T	24	24	4wd	2.083E-01	City +City
5T	49	49	rwd	1.020E-01	Make +Make
3	355	355	fwd	2.225E-01	Make +Make
6	105	105	rwd	2.762E-01	Make +Make
12T	33	33	4wd	2.121E-01	Width +Width
13T	72	72	rwd	1.944E-01	Type +Type
7	250	250	fwd	1.760E-01	Type +Type
14T	41	41	4wd	2.439E-01	City +City
15	209	209	fwd	1.053E-01	Hwy +Hwy
30T	32	32	fwd	2.500E-01	Length +Length
31T	177	177	fwd	2.825E-02	Cylin +Cylin

The variables preceded with a + sign are those used in the nearest neighbor models.

Number of terminal nodes of final tree: 7

Total number of nodes of final tree: 13

Classification tree:

Node 1: Type = "pickup", "sports"

```

Node 2: Type = "pickup"
Node 4: Mean cost = 2.08333E-01
Node 2: Type /= "pickup"
Node 5: Mean cost = 1.02041E-01
Node 1: Type /= "pickup", "sports"
Node 3: Make = "Audi", "BMW", "Hummer", "Infiniti", "Isuzu", "Jaguar",
          "Jeep", "Land-Rover", "Lexus", "Lincoln", "Mercedes",
          "Porsche", "Subaru"
Node 6: Make = "Audi", "Hummer", "Isuzu", "Jeep", "Land-Rover",
          "Porsche", "Subaru"
Node 12: Mean cost = 2.12121E-01
Node 6: Make /= "Audi", "Hummer", "Isuzu", "Jeep", "Land-Rover",
          "Porsche", "Subaru"
Node 13: Mean cost = 1.94444E-01
Node 3: Make /= "Audi", "BMW", "Hummer", "Infiniti", "Isuzu", "Jaguar",
          "Jeep", "Land-Rover", "Lexus", "Lincoln", "Mercedes",
          "Porsche", "Subaru"
Node 7: Type = "suv"
Node 14: Mean cost = 2.43902E-01
Node 7: Type /= "suv"
Node 15: Hwy <= 25.50000
Node 30: Mean cost = 2.50000E-01
Node 15: Hwy > 25.50000 or *
Node 31: Mean cost = 2.82486E-02

```

```

Node 1: Intermediate node
A case goes into Node 2 if Type =
  "pickup", "sports"
Nearest-neighbor K = 7
Type mode = car

```

Class	Number	ClassPrior	Fit variable Type
4wd	94	0.21963	
fwd	224	0.52336	
rwd	110	0.25701	

```

Number of training cases misclassified = 164
If node model is inapplicable due to missing values, predicted class =
fwd

```

Although the number of nearest neighbors is 7 in this node, the neighborhood is the entire node because the fit variable, Type, is categorical.

Node 2: Intermediate node

```

A case goes into Node 4 if Type = "pickup"
Nearest-neighbor K = 5
Type mode = sports
                                Fit variable
Class      Number  ClassPrior  Type
4wd         17      0.23288
fwd          8      0.10959
rwd         48      0.65753
Number of training cases misclassified = 25
If node model is inapplicable due to missing values, predicted class =
rwd
-----
Node 4: Terminal node
Nearest-neighbor K = 4
City mean = 1.6458E+01
                                Fit variable
Class      Number  ClassPrior  City
4wd         12      0.50000
fwd          0      0.00000
rwd         12      0.50000
Number of training cases misclassified = 5
If node model is inapplicable due to missing values, predicted class =
4wd
-----
Node 5: Terminal node
Nearest-neighbor K = 4
Make mode = Porsche
                                Fit variable
Class      Number  ClassPrior  Make
4wd         5      0.10204
fwd          8      0.16327
rwd         36      0.73469
Number of training cases misclassified = 5
If node model is inapplicable due to missing values, predicted class =
rwd
-----
Node 3: Intermediate node
A case goes into Node 6 if Make =
"Audi", "BMW", "Hummer", "Infiniti", "Isuzu", "Jaguar", "Jeep",
"Land-Rover", "Lexus", "Lincoln", "Mercedes", "Porsche", "Subaru"
Nearest-neighbor K = 6
Make mode = Toyota
                                Fit variable
Class      Number  ClassPrior  Make
4wd         77      0.21690
fwd         216      0.60845

```

```

rwd          62      0.17465
Number of training cases misclassified = 79
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 6: Intermediate node
A case goes into Node 12 if Make =
  "Audi", "Hummer", "Isuzu", "Jeep", "Land-Rover", "Porsche", "Subaru"
Nearest-neighbor K = 5
Make mode = Mercedes

                                Fit variable
Class      Number  ClassPrior  Make
4wd         45      0.42857
fwd         10      0.09524
rwd         50      0.47619
Number of training cases misclassified = 29
If node model is inapplicable due to missing values, predicted class =
rwd
-----
Node 12: Terminal node
Nearest-neighbor K = 4
Width mean = 7.1091E+01

                                Fit variable
Class      Number  ClassPrior  Width
4wd         26      0.78788
fwd          7      0.21212
rwd          0      0.00000
Number of training cases misclassified = 7
If node model is inapplicable due to missing values, predicted class =
4wd
-----
Node 13: Terminal node
Nearest-neighbor K = 5
Type mode = car

                                Fit variable
Class      Number  ClassPrior  Type
4wd         19      0.26389
fwd          3      0.04167
rwd         50      0.69444
Number of training cases misclassified = 14
If node model is inapplicable due to missing values, predicted class =
rwd
-----
Node 7: Intermediate node
A case goes into Node 14 if Type = "suv"
Nearest-neighbor K = 6

```

```

Type mode = car

                                Fit variable
Class      Number  ClassPrior  Type
4wd         32      0.12800
fwd         206     0.82400
rwd          12     0.04800
Number of training cases misclassified = 44
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 14: Terminal node
Nearest-neighbor K = 4
City mean = 1.6707E+01

                                Fit variable
Class      Number  ClassPrior  City
4wd         22      0.53659
fwd         19      0.46341
rwd          0      0.00000
Number of training cases misclassified = 10
If node model is inapplicable due to missing values, predicted class =
4wd
-----
Node 15: Intermediate node
A case goes into Node 30 if Hwy <= 2.5500000E+01
Nearest-neighbor K = 6
Hwy mean = 3.0148E+01

                                Fit variable
Class      Number  ClassPrior  Hwy
4wd         10      0.04785
fwd         187     0.89474
rwd          12     0.05742
Number of training cases misclassified = 22
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 30: Terminal node
Nearest-neighbor K = 4
Length mean = 1.9938E+02

                                Fit variable
Class      Number  ClassPrior  Length
4wd         6      0.18750
fwd         15     0.46875
rwd         11     0.34375
Number of training cases misclassified = 8
If node model is inapplicable due to missing values, predicted class =
fwd

```

```

-----
Node 31: Terminal node
Nearest-neighbor K = 6
Cylin mean = 4.8192E+00
                                Fit variable
Class      Number  ClassPrior  Cylin
4wd         4      0.02260
fwd        172     0.97175
rwd         1      0.00565
Number of training cases misclassified = 5
If node model is inapplicable due to missing values, predicted class =
fwd
-----

Classification matrix for training sample:
Predicted      True class
class          4wd      fwd      rwd
4wd             67      17       3
fwd             12     202       2
rwd             15       5     105
Total           94     224     110

Number of cases used for tree construction = 428
Number misclassified = 54
Resubstitution est. of mean misclassification cost = 0.12616822429906541

LaTeX code for tree is in drivenn.tex

```

Figure 3 shows the tree model, which misclassifies 54 observations. Figure 4 shows the observed and predicted values of `Drive` in node 30 of the tree.

4.2.5 Kernel density estimates: car data

An alternative to nearest-neighbor models is kernel discrimination models, where classification is based on maximum likelihood with class densities estimated by the kernel method. Unlike nearest-neighbor, however, this option also yields an estimated class probability vector for each observation. Therefore it can serve as a nonparametric alternative to multinomial logistic regression. Empirical evidence indicates that the nearest-neighbor and kernel methods possess similar prediction accuracy. See Loh (2009) for more details. Following is a log of the input file generation step for the kernel method.

```

0. Read the warranty disclaimer
1. Create an input file for batch run

```

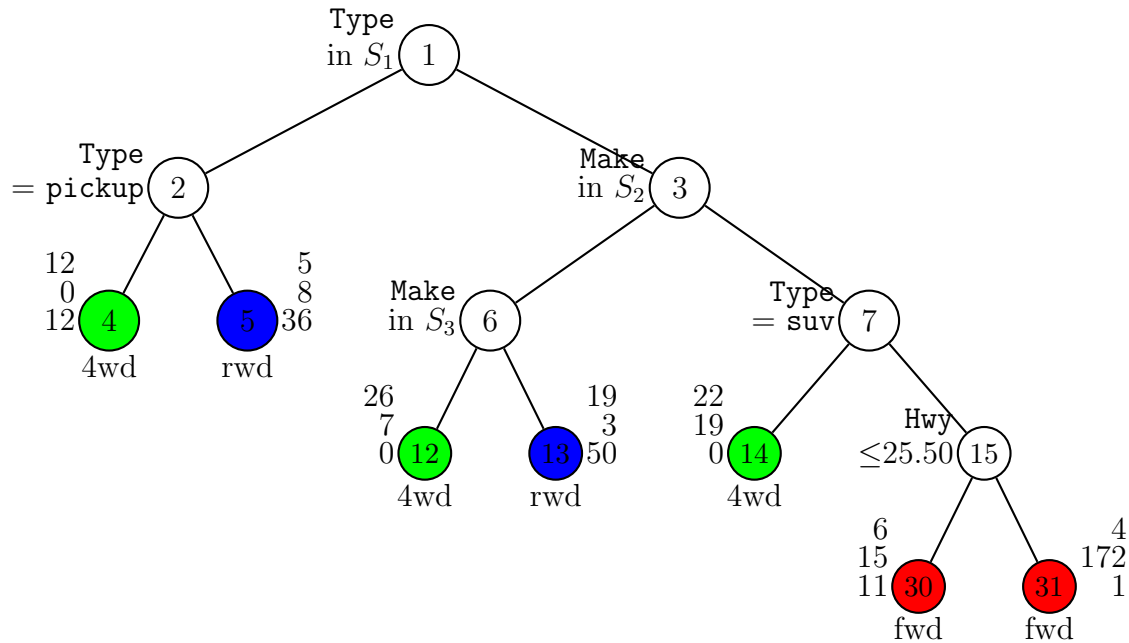



Figure 3: GUIDE 0.50-SE classification tree for predicting **Drive** with univariate nearest-neighbor node models, estimated priors and unit misclassification costs. At each split, an observation goes to the left branch if and only if the condition is satisfied. Set $S_1 = \{\text{pickup, sports}\}$. Set $S_2 = \{\text{Audi, BMW, Hummer, Infiniti, Isuzu, Jaguar, Jeep, Land-Rover, Lexus, Lincoln, Mercedes, Porsche, Subaru}\}$. Set $S_3 = \{\text{Audi, Hummer, Isuzu, Jeep, Land-Rover, Porsche, Subaru}\}$. Predicted classes (based on estimated misclassification cost) printed below terminal nodes; sample sizes for 4wd, fwd, and rwd, respectively, beside nodes.

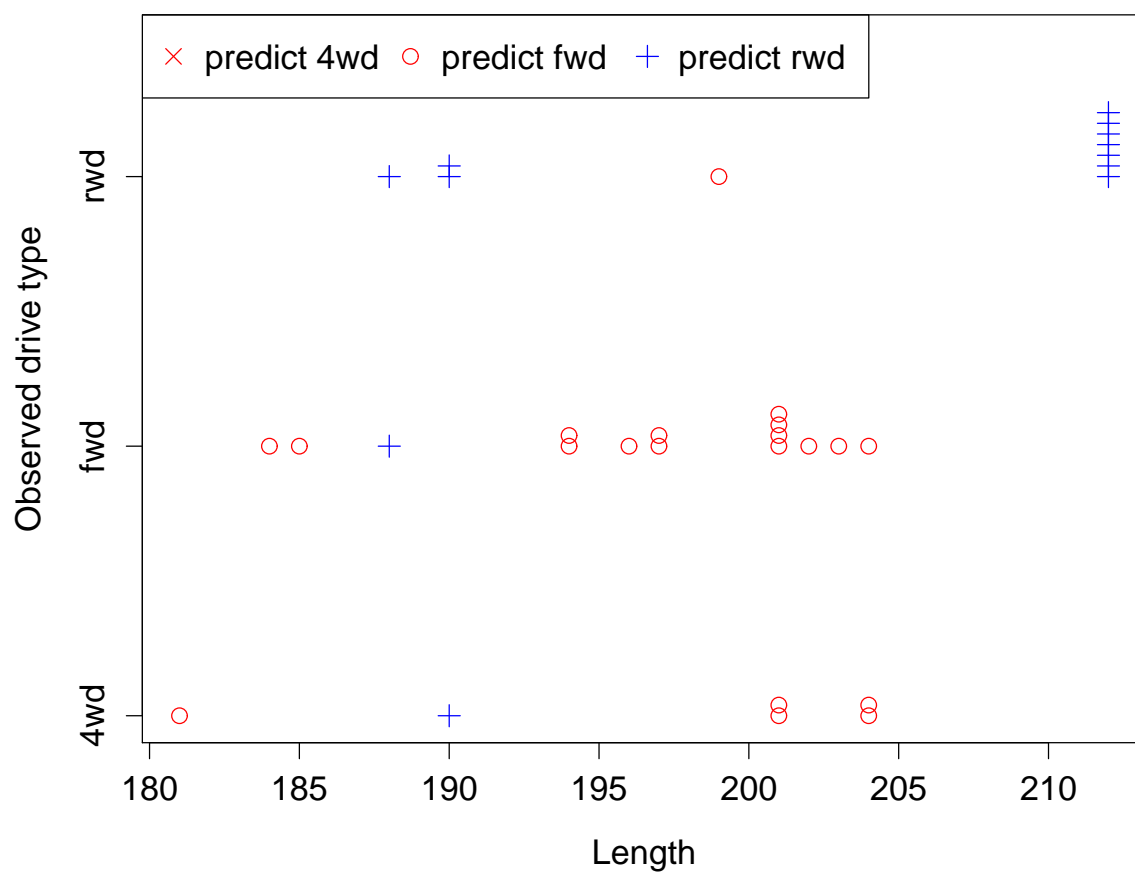


Figure 4: Observed and predicted values of drive type in node 30 of tree in Figure 3

```

2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: driveker.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: driveker.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1): 2
Input 1 for simple, 2 for nearest-neighbor, 3 for kernel method ([1:3], <cr>=1): 3
  This is where you choose kernel density estimation.
Input 1 for univariate, 2 for bivariate preference ([1:2], <cr>=1):
Input 1 for interaction tests, 2 to skip them ([1:2], <cr>=1):
Input 1 to prune by CV, 2 by test sample, 3 for no pruning ([1:3], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters):drivedsc.txt
Reading data description file ...
Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Total number of cases =      428
Number of classes =      3
Col. no. Categorical variable  #levels  #missing values
      3 Make                   38         0
      5 Type                   6         0
Checking data ...
Class name      Num. cases  Proportion
4wd              94    0.21962617
fwd             224    0.52336449
rwd             110    0.25700935
      Total  #cases w/  #missing
      #cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
      428      0      0      10      0      0      11      0      2
No. cases used for training =      428
Default number of cross-validations =      10
Input 1 to accept the default, 2 to change it ([1:2], <cr>=1):
Best tree may be chosen based on mean or median CV estimate
Input 1 for mean-based, 2 for median-based ([1:2], <cr>=1):
Input number of SEs for pruning ([0.00:1000.00], <cr>=0.50):
Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
Input 1, 2, or 3 ([1:3], <cr>=1):
Choose 1 for unit misclassification costs, 2 to input costs from a file

```

```

Input 1 or 2 ([1:2], <cr>=1):
Choose a split point selection method for numerical variables:
Choose 1 to use faster method based on sample quantiles
Choose 2 to use exhaustive search
Input 1 or 2 ([1:2], <cr>=2):
Default max number of split levels =          10
Input 1 to accept this value, 2 to change it ([1:2], <cr>=1):
Default minimum node sample size is 10
Input 1 to use the default value, 2 to change it ([1:2], <cr>=1):
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): driveker.tex
Input 1 for a vertical tree, 2 for a sideways tree ([1:2], <cr>=1):
Input 1 to include node numbers, 2 to omit them ([1:2], <cr>=1):
Input 1 to number all nodes, 2 to number leaves only ([1:2], <cr>=1):
Input 1 to color terminal nodes, 2 otherwise ([1:2], <cr>=1):
Choose amount of detail in nodes of LaTeX tree diagram
Input 0 for #errors, 1 for class sizes, 2 for nothing ([0:2], <cr>=1):
You can store the variables and/or values used to split and fit in a file
Choose 1 to skip this step, 2 to store split and fit variables,
3 to store split variables and their values
Input your choice ([1:3], <cr>=1):
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node IDs and fitted values: driveker.fit
This file contains the predicted class and terminal node label for each observation.
Input 2 to save terminal node IDs for importance scoring; 1 otherwise ([1:2], <cr>=1):
Input name of file to store predicted class and probability: driveker.pro
This file contains the estimated class probabilities for each observation.
Input file is created!
Run GUIDE with the command: guide < driveker.in

```

The results in the output file are given next.

```

Classification tree
Pruning by cross-validation
Data description file: drivedsc.txt
Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Number of classes = 3
Class      #Cases   Proportion
4wd         94     0.21962617
fwd        224     0.52336449
rwd        110     0.25700935

```

Summary information (without x variables)

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical, n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight

Column	Name		Minimum	Maximum	#Categories	#Missing
3	Make	c			38	
5	Type	c			6	
6	Drive	d			3	
14	Rprice	s	1.0280E+04	1.9246E+05		
15	Dcost	s	9.8750E+03	1.7356E+05		
16	Enginsz	s	1.3000E+00	8.3000E+00		
17	Cylin	s	-1.0000E+00	1.2000E+01		
18	Hp	s	7.3000E+01	5.0000E+02		
19	City	s	1.0000E+01	6.0000E+01		
20	Hwy	s	1.2000E+01	6.6000E+01		
21	Weight	s	1.8500E+03	7.1900E+03		
22	Whlbase	s	8.9000E+01	1.4400E+02		
23	Length	s	1.4300E+02	2.2800E+02		
24	Width	s	6.4000E+01	8.1000E+01		

Total	#cases w/	#missing							
#cases	miss. D	ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var	
428	0	0	10	0	0	11	0	2	

No. cases used for training: 428

Univariate split highest priority

Interaction splits 2nd priority; no linear splits

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Kernel density node models

Univariate preference

Estimated priors

Unit misclassification costs

Split values for N and S variables based on exhaustive search

Max number of split levels = 10

Minimum node size = 10

Number of SE's for pruned tree = 5.0000E-01

Size and CV mean cost and SE of subtrees:

Tree	#Tnodes	Mean Cost	SE(Mean)	BSE(Mean)	Median Cost	BSE(Median)
1	19	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
2	18	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
3	17	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
4	16	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
5	15	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
6	14	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
7	13	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02

8	12	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
9--	11	2.243E-01	2.016E-02	1.076E-02	2.209E-01	1.614E-02
10	10	2.336E-01	2.045E-02	9.339E-03	2.326E-01	9.545E-03
11	9	2.360E-01	2.052E-02	1.098E-02	2.326E-01	9.655E-03
12**	8	2.336E-01	2.045E-02	1.218E-02	2.234E-01	1.628E-02
13	7	2.407E-01	2.066E-02	1.431E-02	2.326E-01	1.458E-02
14	5	2.407E-01	2.066E-02	1.293E-02	2.326E-01	1.709E-02
15	3	2.477E-01	2.086E-02	1.556E-02	2.558E-01	2.694E-02
16	2	2.617E-01	2.125E-02	1.851E-02	2.674E-01	2.446E-02
17	1	3.808E-01	2.347E-02	1.384E-02	3.721E-01	1.732E-02

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

* tree same as + tree

** tree same as ++ tree

* tree same as -- tree

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Node cost is node misclassification cost divided by number of training cases

Node label	Total cases	Train cases	Predicted class	Node cost	Split variable followed by (+)fit variable(s)
1	428	428	fwd	3.832E-01	Type +Type
2	73	73	rwd	3.425E-01	Type +Type
4T	24	24	4wd	3.750E-01	City +City
5T	49	49	rwd	1.020E-01	Make +Make
3	355	355	fwd	2.225E-01	Make +Make
6	105	105	rwd	2.762E-01	Make +Make
12T	33	33	4wd	2.121E-01	Width +Width
13T	72	72	rwd	1.944E-01	Type +Type
7	250	250	fwd	1.760E-01	Type +Type
14	41	41	4wd	3.902E-01	City +City
28T	16	16	4wd	3.125E-01	-
29T	25	25	fwd	2.800E-01	Whlbase +Whlbase
15	209	209	fwd	1.053E-01	Hwy +Hwy
30T	32	32	fwd	3.125E-01	Length +Length
31T	177	177	fwd	2.825E-02	Cylin +Cylin

In the above, ‘split variable’ refers to the variable selected to split the node and ‘fit variable(s)’ refers to the one(s) used to estimate the class kernel densities. Fit variables are indicated with a preceding + sign. In this example,

*the split and fit variables are the same in every node.
 If a categorical variable (e.g., Type) is selected for fitting, discrete kernel density estimates are used. A dash (-) indicates that a node is not split, usually due to sample size being too small, in which case all the observations in the node are predicted as belonging to the class that minimizes the misclassification cost.*

Number of terminal nodes of final tree: 8
 Total number of nodes of final tree: 15

Classification tree:

```

Node 1: Type = "pickup", "sports"
  Node 2: Type = "pickup"
    Node 4: Mean cost = 3.75000E-01
    Node 2: Type /= "pickup"
      Node 5: Mean cost = 1.02041E-01
  Node 1: Type /= "pickup", "sports"
    Node 3: Make = "Audi", "BMW", "Hummer", "Infiniti", "Isuzu", "Jaguar",
      "Jeep", "Land-Rover", "Lexus", "Lincoln", "Mercedes",
      "Porsche", "Subaru"
      Node 6: Make = "Audi", "Hummer", "Isuzu", "Jeep", "Land-Rover",
        "Porsche", "Subaru"
        Node 12: Mean cost = 2.12121E-01
        Node 6: Make /= "Audi", "Hummer", "Isuzu", "Jeep", "Land-Rover",
          "Porsche", "Subaru"
          Node 13: Mean cost = 1.94444E-01
    Node 3: Make /= "Audi", "BMW", "Hummer", "Infiniti", "Isuzu", "Jaguar",
      "Jeep", "Land-Rover", "Lexus", "Lincoln", "Mercedes",
      "Porsche", "Subaru"
    Node 7: Type = "suv"
      Node 14: City <= 15.50000
        Node 28: Mean cost = 3.12500E-01
        Node 14: City > 15.50000 or *
          Node 29: Mean cost = 2.80000E-01
      Node 7: Type /= "suv"
        Node 15: Hwy <= 25.50000
          Node 30: Mean cost = 3.12500E-01
          Node 15: Hwy > 25.50000 or *
            Node 31: Mean cost = 2.82486E-02
  
```

```

Node 1: Intermediate node
A case goes into Node 2 if Type =
  "pickup", "sports"
  
```

```
Type mode = car
```

Class	Number	ClassPrior	Bandwidth Type
4wd	94	0.21963	
fwd	224	0.52336	
rwd	110	0.25701	

Number of training cases misclassified = 164
 If node model is inapplicable due to missing values, predicted class = fwd

Categorical variables, such as Type, do not have bandwidths. Their kernel density estimates are the sample cell frequencies.

```
-----
Node 2: Intermediate node
A case goes into Node 4 if Type = "pickup"
Type mode = sports
```

Class	Number	ClassPrior	Bandwidth Type
4wd	17	0.23288	
fwd	8	0.10959	
rwd	48	0.65753	

Number of training cases misclassified = 25
 If node model is inapplicable due to missing values, predicted class = rwd

```
-----
Node 4: Terminal node
City mean = 1.6458E+01
```

Class	Number	ClassPrior	Bandwidth City
4wd	12	0.50000	3.3823E+00
fwd	0	0.00000	0.0000E+00
rwd	12	0.50000	5.1881E+00

Number of training cases misclassified = 9
 If node model is inapplicable due to missing values, predicted class = 4wd

The numbers in the last column give the kernel density bandwidth for each class.

```
-----
Node 5: Terminal node
Make mode = Porsche
```

Class	Number	ClassPrior	Fit variable Make
4wd	5	0.10204	
fwd	8	0.16327	
rwd	36	0.73469	


```

Number of training cases misclassified = 5
If node model is inapplicable due to missing values, predicted class =
rwd
-----
Node 3: Intermediate node
A case goes into Node 6 if Make =
"Audi", "BMW", "Hummer", "Infiniti", "Isuzu", "Jaguar", "Jeep",
"Land-Rover", "Lexus", "Lincoln", "Mercedes", "Porsche", "Subaru"
Make mode = Toyota

```

Class	Number	ClassPrior	Bandwidth Make
4wd	77	0.21690	
fwd	216	0.60845	
rwd	62	0.17465	

```

Number of training cases misclassified = 79
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 6: Intermediate node
A case goes into Node 12 if Make =
"Audi", "Hummer", "Isuzu", "Jeep", "Land-Rover", "Porsche", "Subaru"
Make mode = Mercedes

```

Class	Number	ClassPrior	Bandwidth Make
4wd	45	0.42857	
fwd	10	0.09524	
rwd	50	0.47619	

```

Number of training cases misclassified = 29
If node model is inapplicable due to missing values, predicted class =
rwd
-----
Node 12: Terminal node
Width mean = 7.1091E+01

```

Class	Number	ClassPrior	Bandwidth Width
4wd	26	0.78788	3.3807E+00
fwd	7	0.21212	1.2558E+00
rwd	0	0.00000	0.0000E+00

```

Number of training cases misclassified = 7
If node model is inapplicable due to missing values, predicted class =
4wd
-----
Node 13: Terminal node
Type mode = car

```

Class	Number	ClassPrior	Fit variable Type
-------	--------	------------	----------------------

```

4wd          19      0.26389
fwd           3      0.04167
rwd          50      0.69444
Number of training cases misclassified = 14
If node model is inapplicable due to missing values, predicted class =
rwd
-----
Node 7: Intermediate node
A case goes into Node 14 if Type = "suv"
Type mode = car

```

Class	Number	ClassPrior	Bandwidth Type
4wd	32	0.12800	
fwd	206	0.82400	
rwd	12	0.04800	

```

Number of training cases misclassified = 44
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 14: Intermediate node
A case goes into Node 28 if City <= 1.5500000E+01
City mean = 1.6707E+01

```

Class	Number	ClassPrior	Bandwidth City
4wd	22	0.53659	3.9949E+00
fwd	19	0.46341	3.0269E+00
rwd	0	0.00000	0.0000E+00

```

Number of training cases misclassified = 16
If node model is inapplicable due to missing values, predicted class =
4wd
-----
Node 28: Terminal node
Class      Number  ClassPrior
4wd         11     0.68750
fwd          5     0.31250
rwd          0     0.00000
Number of training cases misclassified = 5
Predicted class is 4wd
-----
Node 29: Terminal node
Whlbase mean = 1.0644E+02

```

Class	Number	ClassPrior	Bandwidth Whlbase
4wd	11	0.44000	5.6690E+00
fwd	14	0.56000	7.9672E+00
rwd	0	0.00000	0.0000E+00

```

Number of training cases misclassified = 7
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 15: Intermediate node
A case goes into Node 30 if Hwy <= 2.5500000E+01
Hwy mean = 3.0148E+01

```

Class	Number	ClassPrior	Bandwidth Hwy
4wd	10	0.04785	5.8465E+00
fwd	187	0.89474	3.9058E+00
rwd	12	0.05742	5.6846E+00

```

Number of training cases misclassified = 22
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 30: Terminal node
Length mean = 1.9938E+02

```

Class	Number	ClassPrior	Bandwidth Length
4wd	6	0.18750	1.1008E+01
fwd	15	0.46875	7.0085E+00
rwd	11	0.34375	1.6353E+01

```

Number of training cases misclassified = 10
If node model is inapplicable due to missing values, predicted class =
fwd
-----
Node 31: Terminal node
Cylin mean = 4.8192E+00

```

Class	Number	ClassPrior	Bandwidth Cylin
4wd	4	0.02260	9.4732E-01
fwd	172	0.97175	9.4317E-01
rwd	1	0.00565	1.9453E-01

```

Number of training cases misclassified = 5
If node model is inapplicable due to missing values, predicted class =
fwd
-----

```

Classification matrix for training sample:

Predicted	True class		
class	4wd	fwd	rwd
4wd	68	19	8
fwd	13	201	5
rwd	13	4	97

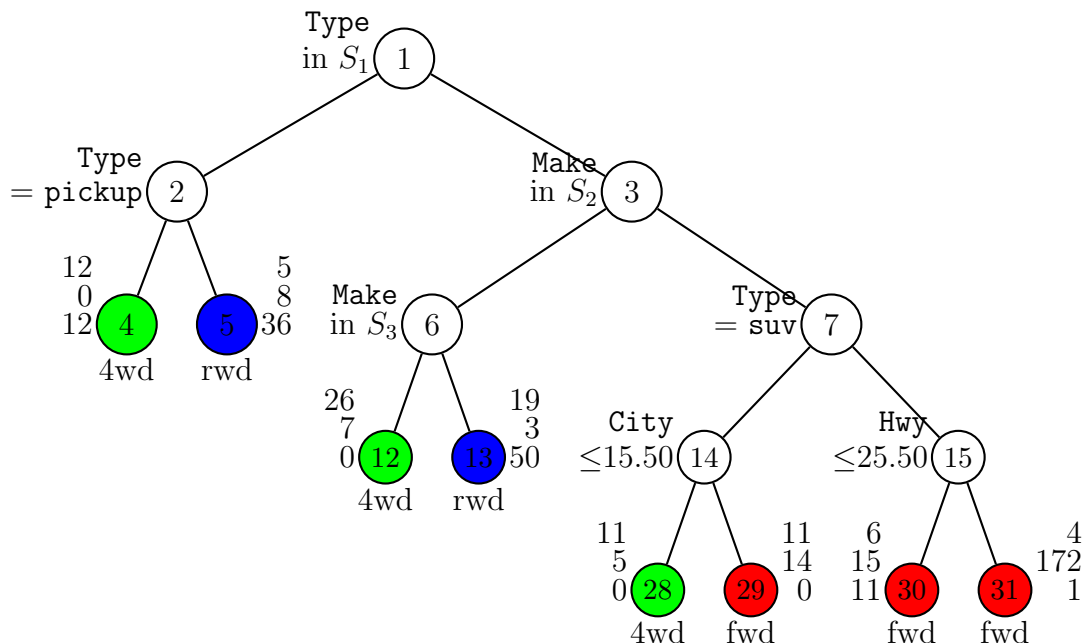


Figure 5: GUIDE 0.50-SE classification tree for predicting **Drive** with univariate kernel discrimination node models, estimated priors and unit misclassification costs. At each split, an observation goes to the left branch if and only if the condition is satisfied. Set $S_1 = \{\text{pickup, sports}\}$. Set $S_2 = \{\text{Audi, BMW, Hummer, Infiniti, Isuzu, Jaguar, Jeep, Land-Rover, Lexus, Lincoln, Mercedes, Porsche, Subaru}\}$. Set $S_3 = \{\text{Audi, Hummer, Isuzu, Jeep, Land-Rover, Porsche, Subaru}\}$. Predicted classes (based on estimated misclassification cost) printed below terminal nodes; sample sizes for 4wd, fwd, and rwd, respectively, beside nodes.

```

Total          94      224      110

Number of cases used for tree construction =  428
Number misclassified =  62
Resubstitution est. of mean misclassification cost =  0.14485981308411214

Predicted class probability estimates are stored in driveker.pro
Observed and fitted values are stored in driveker.fit
LaTeX code for tree is in driveker.tex
Elapsed time in seconds:  0.266952991

```

Figure 5 shows the \LaTeX tree diagram. The top several lines of the file `driveker.fit`, which contains the terminal node label and predicted class for each observation in the training sample, are:

train	node	observed	predicted
y	31	"fwd"	"fwd"
y	31	"fwd"	"fwd"
y	31	"fwd"	"fwd"
y	31	"fwd"	"fwd"
y	31	"fwd"	"fwd"
y	31	"fwd"	"fwd"
y	31	"fwd"	"fwd"
y	31	"fwd"	"fwd"

The corresponding lines of the file `driveker.pro`, giving the estimated class probabilities for each observation, are:

"4wd"	"fwd"	"rwd"	predicted	observed
0.02302	0.93575	0.04123	"fwd"	"fwd"
0.02302	0.93575	0.04123	"fwd"	"fwd"
0.02302	0.93575	0.04123	"fwd"	"fwd"
0.02302	0.93575	0.04123	"fwd"	"fwd"
0.02302	0.93575	0.04123	"fwd"	"fwd"
0.02302	0.93575	0.04123	"fwd"	"fwd"
0.02302	0.93575	0.04123	"fwd"	"fwd"

5 Regression

5.1 Stepwise least-squares

We use the baseball dataset `bbdat.txt` to show the results for regression trees when there are no missing values. The data give the log-salary and performance measures of 263 professional baseball players ([Hoaglin and Velleman, 1995](#)). The response variable is the logarithm of salary (`Logsalary`). The data description file `bbdsc.txt` consists of the following lines:

```
bbdat.txt
NA
column, varname, vartype
1 Id x
2 Name x
3 Bat86 n
4 Hit86 n
5 Hr86 n
6 Run86 n
7 Rb86 n
8 Wlk86 n
```

```

9 Yrs n
10 Batcr n
11 Hitcr n
12 Hrcr n
13 Runcr n
14 Rbcr n
15 Wlkr n
16 League86 b
17 Div86 b
18 Team86 c
19 Pos86 b
20 Puto86 n
21 Asst86 n
22 Err86 n
23 Salary x
24 League87 b
25 Team87 c
26 Logsalary d

```

Notice that there are four variables having the “b” variable type. This means that 0-1 dummy variables will be created for them in fitting the node linear models. The following shows how the input file is created.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: stepin.txt
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: stepout.txt
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1):
Input 1 for least squares, 2 least median of squares ([1:2], <cr>=1):
Choose complexity of model to use at each node:
0: stepwise linear, 1: multiple linear, 2: best polynomial, 3: constant,
4: stepwise simple ANCOVA ([0:4], <cr>=0):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): bbdsc.txt

```

```

Reading data description file ...
Training sample file: bbdatt.txt
Missing value code: NA
Dependent variable is Logsalary
Length of longest data entry = 17
Total number of cases =      263
Col. no. Categorical variable    #levels    #missing values
      16 Leag86                  2          0
      17 Div86                  2          0
      18 Team86                 24          0
      19 Pos86                  23          0
      24 Leag87                 2          0
      25 Team87                 24          0
Checking data ...
The program will try to create the variables in the description file.
If it is unsuccessful, please create the columns yourself...
Number of dummy variables created:      25
      Total  #cases w/  #missing
      #cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
      263      0      0      3      16      0      0      4      2
No weight variable in data file
No. cases used for training =      263
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): step.tex
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node IDs and fitted values: step.fit
Input file is created!

```

The contents from the file `stepout.txt` follow. They show a tree with two terminal nodes and give the regression coefficients, sample means of the dependent and predictor variables, MSE and R^2 values, and names of the split variables in each node.

```

Least squares regression tree
Predictions truncated at global min and max of D sample values
  The predicted values are truncated at the minimum and
  maximum values of the training sample by default.
Pruning by cross-validation
Data description file: bbdsc.txt
Training sample file: bbdatt.txt
Missing value code: NA
Dependent variable is Logsalary
Piecewise forward and backward stepwise regression
F-to-enter and F-to-delete =    4.0000000000000000    3.9900000000000002
  These default F values are the same as those used in SAS.

```

Using as many variables as needed
 Length of longest data entry = 17
 Number of dummy variables created = 25

Summary information (without x variables)

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical, n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
 For categorical variables, #categories include one for missing values

Column number	Variable name	Variable type	Minimum value	Maximum value	Number of categories	Number missing
3	Bat86	n	1.2700E+02	6.8700E+02		
4	Hit86	n	3.2000E+01	2.3800E+02		
5	Hr86	n	0.0000E+00	4.0000E+01		
6	Run86	n	1.3000E+01	1.3000E+02		
7	Rb86	n	8.0000E+00	1.2100E+02		
8	Wlk86	n	3.0000E+00	1.0500E+02		
9	Yrs	n	1.0000E+00	2.4000E+01		
10	Batcr	n	1.8100E+02	1.4053E+04		
11	Hitcr	n	4.2000E+01	4.2560E+03		
12	Hrcr	n	0.0000E+00	5.4800E+02		
13	Runcr	n	1.8000E+01	2.1650E+03		
14	Rbcr	n	9.0000E+00	1.6590E+03		
15	Wlkcr	n	8.0000E+00	1.5660E+03		
16	Leag86	b			2	
17	Div86	b			2	
18	Team86	c			24	
19	Pos86	b			23	
20	Puto86	n	0.0000E+00	1.3770E+03		
21	Asst86	n	0.0000E+00	4.9200E+02		
22	Err86	n	0.0000E+00	3.2000E+01		
24	Leag87	b			2	
25	Team87	c			24	
26	Logsalary	d	4.2121E+00	7.8079E+00		

===== Constructed variables =====

The F variables below this line are dummy variables constructed from the B variables.

27	Leag8.N	f	0.0000E+00	1.0000E+00
28	Div86.W	f	0.0000E+00	1.0000E+00
29	Pos86.10	f	0.0000E+00	1.0000E+00
30	Pos86.23	f	0.0000E+00	1.0000E+00
31	Pos86.2B	f	0.0000E+00	1.0000E+00
32	Pos86.2S	f	0.0000E+00	1.0000E+00
33	Pos86.32	f	0.0000E+00	1.0000E+00
34	Pos86.3B	f	0.0000E+00	1.0000E+00
35	Pos86.30	f	0.0000E+00	1.0000E+00
36	Pos86.3S	f	0.0000E+00	1.0000E+00
37	Pos86.C	f	0.0000E+00	1.0000E+00

38	Pos86.CD	f	0.0000E+00	1.0000E+00
39	Pos86.CF	f	0.0000E+00	1.0000E+00
40	Pos86.DH	f	0.0000E+00	1.0000E+00
41	Pos86.DO	f	0.0000E+00	1.0000E+00
42	Pos86.LF	f	0.0000E+00	1.0000E+00
43	Pos86.O1	f	0.0000E+00	1.0000E+00
44	Pos86.OD	f	0.0000E+00	1.0000E+00
45	Pos86.OF	f	0.0000E+00	1.0000E+00
46	Pos86.OS	f	0.0000E+00	1.0000E+00
47	Pos86.RF	f	0.0000E+00	1.0000E+00
48	Pos86.S3	f	0.0000E+00	1.0000E+00
49	Pos86.SS	f	0.0000E+00	1.0000E+00
50	Pos86.UT	f	0.0000E+00	1.0000E+00
51	Leag8.N	f	0.0000E+00	1.0000E+00

#cases	#cases w/ miss.	D	#missing ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
263	0		0	3	16	0	0	4	2

No weight variable in data file

No. cases used for training = 263

Missing N and F values imputed with node means

The default method of handling missing values is node-mean imputation.

Interaction tests on all variables

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Fraction of cases used for splitting each node = 1.0000

Max number of split levels = 10

Minimum node size = 13

Number of SE's for pruned tree = 5.0000E-01

Size and CV MSE and SE of subtrees:

Tree	#Tnodes	Mean MSE	SE(Mean)	BSE(Mean)	Median MSE	BSE(Median)
1	14	2.001E-01	2.068E-02	2.408E-02	1.866E-01	2.982E-02
2	13	2.001E-01	2.068E-02	2.408E-02	1.866E-01	2.982E-02
3	12	2.001E-01	2.068E-02	2.408E-02	1.866E-01	2.982E-02
4	11	1.768E-01	1.920E-02	2.206E-02	1.440E-01	3.984E-02
5	10	1.744E-01	1.935E-02	2.221E-02	1.440E-01	4.151E-02
6	8	1.790E-01	1.937E-02	2.107E-02	1.486E-01	3.760E-02
7	7	1.824E-01	1.979E-02	2.121E-02	1.492E-01	3.492E-02
8	4	1.574E-01	1.882E-02	1.869E-02	1.309E-01	3.241E-02
9	3	1.518E-01	1.765E-02	1.752E-02	1.296E-01	3.044E-02
10**	2	1.208E-01	1.439E-02	1.378E-02	1.166E-01	1.884E-02
11	1	3.469E-01	2.575E-02	2.224E-02	3.456E-01	3.877E-02

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +
 Selected-SE tree based on mean using naive SE is marked with **
 Selected-SE tree based on mean using bootstrap SE is marked with --
 Selected-SE tree based on median and bootstrap SE is marked with ++
 * tree, ** tree, + tree, and ++ tree all the same

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

D-mean is mean of Logsalary in the node

Cases fit give the number of cases used to fit node

MSE and R² are based on all cases in node

Node label	Total cases	Cases Matrix fit rank	Node D-mean	Node MSE	Node R ²	Split variable	Other variables
1	263	263 9	5.945E+00	2.907E-01	0.6391	Yrs	
2T	143	143 7	5.506E+00	8.336E-02	0.8907	Yrs	
3T	120	120 6	6.469E+00	1.258E-01	0.6456	Bat86	

Number of terminal nodes of final tree: 2

Total number of nodes of final tree: 3

Regression tree:

Node 1: Yrs <= 6.00000 or NA

Node 2: Logsalary-mean = 5.50632

Node 1: Yrs > 6.00000 and not NA

Node 3: Logsalary-mean = 6.46866

Node 1: Intermediate node

A case goes into Node 2 if Yrs <= 6.000000E+00 or NA

Yrs mean = 7.3802E+00

Coefficients of least squares regression function:

Regressor	Coefficient	t-stat	p-val	Min	Mean	Max
Constant	4.1321E+00	28.07	0.0000			
Bat86	-2.4528E-03	-2.73	0.0068	1.2700E+02	4.0829E+02	6.8700E+02
Hit86	1.4558E-02	5.16	0.0000	3.2000E+01	1.0916E+02	2.3800E+02
Wlk86	1.0020E-02	3.82	0.0002	3.0000E+00	4.1722E+01	1.0500E+02
Yrs	7.0500E-02	4.24	0.0000	1.0000E+00	7.3802E+00	2.4000E+01
Runcr	1.1939E-03	3.11	0.0021	1.8000E+01	3.6808E+02	2.1650E+03
Wlkcr	-9.6467E-04	-2.20	0.0291	8.0000E+00	2.6655E+02	1.5660E+03
Leag8.N	1.4081E-01	2.09	0.0373	0.0000E+00	4.7148E-01	1.0000E+00
Pos86.C	3.3729E-01	3.05	0.0025	0.0000E+00	1.1407E-01	1.0000E+00

```

Mean of Logsalary = 5.9454102795235446
Predicted values truncated at 4.2121275978784798 & 7.8079166289264101
-----
Node 2: Terminal node
Coefficients of least squares regression functions:
Regressor    Coefficient    t-stat  p-val      Min      Mean      Max
Constant     4.1385E+00     39.36  0.0000
Bat86        -1.8387E-03     -4.25  0.0000    1.5100E+02  4.1345E+02  6.8700E+02
Run86         1.5359E-02      6.46  0.0000    1.3000E+01  5.6699E+01  1.1900E+02
Yrs           1.1659E-01      4.86  0.0000    1.0000E+00  3.8042E+00  6.0000E+00
Batcr         5.0610E-04      5.44  0.0000    1.8100E+02  1.2046E+03  3.3740E+03
Rbcr          1.5954E-03      2.86  0.0049    9.0000E+00  1.4315E+02  4.7500E+02
Pos86.CF      -2.3212E-01     -2.74  0.0070    0.0000E+00  1.0490E-01  1.0000E+00
Mean of Logsalary = 5.5063156054013094
Predicted values truncated at 4.2121275978784798 & 7.8079166289264101
The truncation limits are the minimum and maximum values of the entire training sample.
-----
Node 3: Terminal node
Coefficients of least squares regression functions:
Regressor    Coefficient    t-stat  p-val      Min      Mean      Max
Constant     6.2274E+00     33.81  0.0000
Hit86         4.7082E-03      4.73  0.0000    3.2000E+01  1.0759E+02  2.0000E+02
Yrs           -1.0408E-01     -6.16  0.0000    7.0000E+00  1.1642E+01  2.4000E+01
Runcr         7.9765E-04      3.26  0.0015    6.7000E+01  6.1369E+02  2.1650E+03
Rbcr          5.9958E-04      2.56  0.0118    8.2000E+01  5.6998E+02  1.6590E+03
Puto86        4.1481E-04      3.43  0.0008    0.0000E+00  2.7723E+02  1.3140E+03
Mean of Logsalary = 6.4686647661858760
Predicted values truncated at 4.2121275978784798 & 7.8079166289264101
-----

Observed and predicted values are in file step.fit
Proportion of variance (R-squared) explained by tree model = .8745

Fitted values are stored in step.fit
LaTeX code for tree is in step.tex
Elapsed time in seconds: 1.98530900

```

The \LaTeX tree produced by the file `step.tex` is shown in Figure 6.

5.2 Least-squares simple polynomial

Often it is useful to be able to visualize the fitted regression function and the data simultaneously. This can be accomplished by fitting a piecewise simple linear model, where the best single regressor is selected to fit a straight line in each node, as follows.

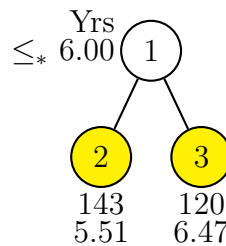


Figure 6: GUIDE piecewise linear least-squares regression tree with stepwise variable selection for predicting Logsalary. At each intermediate node, an observation goes to the left branch if and only if the condition is satisfied. The symbol ‘ \leq^* ’ stands for ‘ \leq or missing’. Sample sizes and means of Logsalary printed below nodes.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: linin.txt
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: linout.txt
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1):
Input 1 for least squares, 2 least median of squares ([1:2], <cr>=1):
Choose complexity of model to use at each node:
0: stepwise linear, 1: multiple linear, 2: best polynomial, 3: constant,
4: stepwise simple ANCOVA ([0:4], <cr>=0): 2
  The default degree of the polynomial is 1.
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): bbdsc.txt
Reading data description file ...
Training sample file: bbdatt.txt
Missing value code: NA
Warning: B variables changed to C
  This warning is triggered because the description file contains some variables
  with the B designation, which is not allowed in piecewise polynomial regression.

```

```

Dependent variable is Logsalary
Length of longest data entry = 17
Total number of cases =      263
Col. no. Categorical variable   #levels   #missing values
    16 Leag86                   2         0
    17 Div86                    2         0
    18 Team86                   24         0
    19 Pos86                    23         0
    24 Leag87                   2         0
    25 Team87                   24         0
Checking data ...
      Total #cases w/   #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
        263      0      0      3      16      0      0      0      6
No weight variable in data file
No. cases used for training =      263
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): lin.tex
A file by that name already exists
Input 1 to overwrite it, 2 to choose another name ([1:2], <cr>=1):
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node IDs and fitted values: lin.fit
Input file is created!

```

Partial contents of linout.txt

```

Size and CV MSE and SE of subtrees:
Tree   #Tnodes   Mean MSE   SE(Mean)   BSE(Mean)   Median MSE   BSE(Median)
  1      13    1.816E-01  2.760E-02  2.322E-02   1.594E-01   2.882E-02
  2      12    1.816E-01  2.760E-02  2.322E-02   1.594E-01   2.882E-02
  3      11    1.809E-01  2.761E-02  2.336E-02   1.594E-01   2.990E-02
  4      10    1.806E-01  2.763E-02  2.347E-02   1.599E-01   3.034E-02
  5       9    1.787E-01  2.757E-02  2.241E-02   1.599E-01   3.108E-02
  6       8    1.790E-01  2.757E-02  2.256E-02   1.599E-01   3.108E-02
  7       7    1.803E-01  2.754E-02  2.351E-02   1.811E-01   3.253E-02
  8       6    1.709E-01  2.654E-02  2.325E-02   1.610E-01   2.358E-02
  9       5    1.782E-01  2.683E-02  2.331E-02   1.651E-01   2.792E-02
10+      4    1.729E-01  2.115E-02  2.233E-02   1.575E-01   2.575E-02
11**     3    1.677E-01  2.150E-02  2.285E-02   1.643E-01   3.291E-02
12       2    1.887E-01  2.313E-02  2.491E-02   1.733E-01   4.050E-02
13       1    4.436E-01  3.247E-02  3.208E-02   4.574E-01   4.953E-02

```

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **
 Selected-SE tree based on mean using bootstrap SE is marked with --
 Selected-SE tree based on median and bootstrap SE is marked with ++
 ** tree same as -- tree
 * tree same as ** tree
 * tree same as -- tree

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

D-mean is mean of Logsalary in the node

Cases fit give the number of cases used to fit node

MSE and R² are based on all cases in node

Node label	Total cases	Cases fit	Matrix rank	Node D-mean	Node MSE	Node R ²	Split variable	Other variables
1	263	263	2	5.945E+00	4.502E-01	0.4257	Yrs +Hitcr	
2	143	143	2	5.506E+00	1.284E-01	0.8254	Hitcr +Batcr	
4T	110	110	2	5.146E+00	9.210E-02	0.7371	Wlkr +Hitcr	
5T	33	33	2	6.706E+00	7.500E-02	0.4395	Wlkr86 +Rbcr	
3T	120	120	2	6.469E+00	1.943E-01	0.4331	Wlkr +Hit86	

The last column, labeled 'Fit variables', give the regressor variable names and the signs of their regression coefficients.

Number of terminal nodes of final tree: 3

Total number of nodes of final tree: 5

Regression tree:

Node 1: Yrs <= 6.00000 or NA
 Node 2: Hitcr <= 4.59500E+02 or NA
 Node 4: Logsalary-mean = 5.14642
 Node 2: Hitcr > 4.59500E+02 and not NA
 Node 5: Logsalary-mean = 6.70595
 Node 1: Yrs > 6.00000 and not NA
 Node 3: Logsalary-mean = 6.46866

Node 1: Intermediate node

A case goes into Node 2 if Yrs <= 6.0000000E+00 or NA

Yrs mean = 7.3802E+00

Coefficients of least squares regression function:

Regressor	Coefficient	t-stat	p-val	Min	Mean	Max
Constant	5.2933E+00	84.65	0.0000			
Hitcr	8.8852E-04	13.91	0.0000	4.2000E+01	7.3392E+02	4.2560E+03

```

Mean of Logsalary = 5.9454102795235446
Predicted values truncated at 4.2121275978784798 & 7.8079166289264101
-----
Node 2: Intermediate node
A case goes into Node 4 if Hitcr <= 4.5950000E+02 or NA
Hitcr mean = 3.2022E+02
-----
Node 4: Terminal node
Coefficients of least squares regression functions:
Regressor    Coefficient    t-stat  p-val      Min      Mean      Max
Constant     4.2487E+00      71.82  0.0000
Hitcr        4.1935E-03      17.40  0.0000    4.2000E+01  2.1408E+02  4.5700E+02
Mean of Logsalary = 5.1464245329622393
Predicted values truncated at 4.2121275978784798 & 7.8079166289264101
-----
Node 5: Terminal node
Coefficients of least squares regression functions:
Regressor    Coefficient    t-stat  p-val      Min      Mean      Max
Constant     5.9484E+00      36.98  0.0000
Rbcr         2.5219E-03       4.93  0.0000    1.0300E+02  3.0039E+02  4.7500E+02
Mean of Logsalary = 6.7059525135315550
Predicted values truncated at 4.2121275978784798 & 7.8079166289264101
-----
Node 3: Terminal node
Coefficients of least squares regression functions:
Regressor    Coefficient    t-stat  p-val      Min      Mean      Max
Constant     5.4482E+00      47.48  0.0000
Hit86        9.4846E-03       9.50  0.0000    3.2000E+01  1.0759E+02  2.0000E+02
Mean of Logsalary = 6.4686647661858760
Predicted values truncated at 4.2121275978784798 & 7.8079166289264101
-----

Proportion of variance (R-squared) explained by tree model = .8279

Observed and fitted values are stored in lin.fit
LaTeX code for tree is in lin.tex
Elapsed time in seconds: 0.238200992

```

The \LaTeX tree is shown in Figure 7.

5.3 ANCOVA models

Besides, multiple linear, stepwise linear, and best simple polynomial regression, GUIDE can also fit a a best ANCOVA model in each node. The ANCOVA model

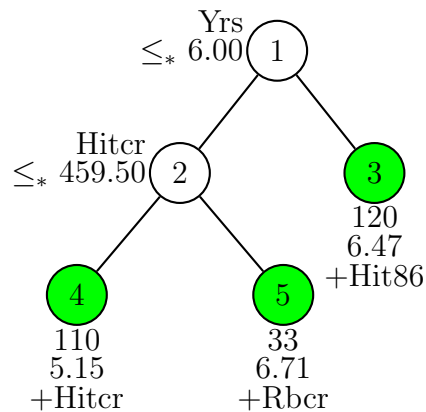


Figure 7: GUIDE piecewise simple linear least-squares regression tree for predicting Logsalary. At each intermediate node, an observation goes to the left branch if and only if the condition is satisfied. The symbol ' \leq_* ' stands for ' \leq or missing'. Sample sizes, means of Logsalary, and signs and names of regressor variable printed below nodes. Terminal nodes with negative, zero, and positive slopes are colored red, yellow, and green, respectively.

uses stepwise regression to find the best single linear regressor and the best subset of dummy (indicator) variables constructed from any B variables.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: ancova.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: ancova.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1):
Input 1 for least squares, 2 least median of squares ([1:2], <cr>=1):
Choose complexity of model to use at each node:
0: stepwise linear, 1: multiple linear, 2: best polynomial, 3: constant,
4: stepwise simple ANCOVA ([0:4], <cr>=0): 4
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):
  
```



```

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): bbdsc.txt
Reading data description file ...
Training sample file: bbdatt.txt
Missing value code: NA
Dependent variable is Logsalary
Length of longest data entry = 17
Total number of cases =      263
Col. no. Categorical variable    #levels    #missing values
    16 Leag86                    2          0
    17 Div86                     2          0
    18 Team86                    24          0
    19 Pos86                     23          0
    24 Leag87                    2          0
    25 Team87                    24          0
Checking data ...
The program will try to create the variables in the description file.
If it is unsuccessful, please create the columns yourself...
Number of dummy variables created:      25
      Total #cases w/ #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
        263      0      0      3      16      0      0      4      2
No weight variable in data file
No. cases used for training =      263
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): ancova.tex
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node IDs and fitted values: ancova.fit
Input file is created!

```

Results

```

Least squares regression tree
Predictions truncated at global min and max of D sample values
Pruning by cross-validation
Data description file: bbdsc.txt
Training sample file: bbdatt.txt
Missing value code: NA
Dependent variable is Logsalary
Piecewise simple linear ANCOVA model
F-to-enter and F-to-delete =  4.000 3.990
Length of longest data entry = 17
Number of dummy variables created =  25

```

Summary information (without x variables)

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical, n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
For categorical variables, #categories include one for missing values

Column number	Variable name	Variable type	Minimum value	Maximum value	Number of categories	Number missing
3	Bat86	n	1.2700E+02	6.8700E+02		
4	Hit86	n	3.2000E+01	2.3800E+02		
5	Hr86	n	0.0000E+00	4.0000E+01		
6	Run86	n	1.3000E+01	1.3000E+02		
7	Rb86	n	8.0000E+00	1.2100E+02		
8	Wlk86	n	3.0000E+00	1.0500E+02		
9	Yrs	n	1.0000E+00	2.4000E+01		
10	Batcr	n	1.8100E+02	1.4053E+04		
11	Hitcr	n	4.2000E+01	4.2560E+03		
12	Hrcr	n	0.0000E+00	5.4800E+02		
13	Runcr	n	1.8000E+01	2.1650E+03		
14	Rbcr	n	9.0000E+00	1.6590E+03		
15	Wlkcr	n	8.0000E+00	1.5660E+03		
16	Leag86	b			2	
17	Div86	b			2	
18	Team86	c			24	
19	Pos86	b			23	
20	Puto86	n	0.0000E+00	1.3770E+03		
21	Asst86	n	0.0000E+00	4.9200E+02		
22	Err86	n	0.0000E+00	3.2000E+01		
24	Leag87	b			2	
25	Team87	c			24	
26	Logsalary	d	4.2121E+00	7.8079E+00		
===== Constructed variables =====						
27	Leag8.N	f	0.0000E+00	1.0000E+00		
28	Div86.W	f	0.0000E+00	1.0000E+00		
29	Pos86.10	f	0.0000E+00	1.0000E+00		
30	Pos86.23	f	0.0000E+00	1.0000E+00		
31	Pos86.2B	f	0.0000E+00	1.0000E+00		
32	Pos86.2S	f	0.0000E+00	1.0000E+00		
33	Pos86.32	f	0.0000E+00	1.0000E+00		
34	Pos86.3B	f	0.0000E+00	1.0000E+00		
35	Pos86.30	f	0.0000E+00	1.0000E+00		
36	Pos86.3S	f	0.0000E+00	1.0000E+00		
37	Pos86.C	f	0.0000E+00	1.0000E+00		
38	Pos86.CD	f	0.0000E+00	1.0000E+00		
39	Pos86.CF	f	0.0000E+00	1.0000E+00		
40	Pos86.DH	f	0.0000E+00	1.0000E+00		
41	Pos86.DO	f	0.0000E+00	1.0000E+00		
42	Pos86.LF	f	0.0000E+00	1.0000E+00		

43	Pos86.01	f	0.0000E+00	1.0000E+00
44	Pos86.0D	f	0.0000E+00	1.0000E+00
45	Pos86.0F	f	0.0000E+00	1.0000E+00
46	Pos86.0S	f	0.0000E+00	1.0000E+00
47	Pos86.RF	f	0.0000E+00	1.0000E+00
48	Pos86.S3	f	0.0000E+00	1.0000E+00
49	Pos86.SS	f	0.0000E+00	1.0000E+00
50	Pos86.UT	f	0.0000E+00	1.0000E+00
51	Leag8.N	f	0.0000E+00	1.0000E+00

Total #cases	#cases w/ miss.	D ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
263	0	0	3	16	0	0	4	2

No weight variable in data file

No. cases used for training = 263

Missing N and F values imputed with node means

Interaction tests on all variables

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Fraction of cases used for splitting each node = 1.0000

Max number of split levels = 10

Minimum node size = 13

Number of SE's for pruned tree = 5.0000E-01

Size and CV MSE and SE of subtrees:

Tree	#Tnodes	Mean MSE	SE(Mean)	BSE(Mean)	Median MSE	BSE(Median)
1	16	2.427E-01	2.723E-02	3.309E-02	2.140E-01	4.431E-02
2	15	2.427E-01	2.723E-02	3.309E-02	2.140E-01	4.431E-02
3	14	2.478E-01	2.765E-02	3.205E-02	2.218E-01	3.847E-02
4	13	2.409E-01	2.882E-02	3.097E-02	2.218E-01	3.172E-02
5	12	2.397E-01	2.885E-02	3.124E-02	2.212E-01	3.288E-02
6	11	2.391E-01	2.883E-02	3.119E-02	2.201E-01	3.471E-02
7	10	2.493E-01	3.000E-02	3.331E-02	2.201E-01	5.040E-02
8	9	2.280E-01	2.802E-02	3.414E-02	2.025E-01	4.236E-02
9	6	2.253E-01	2.788E-02	3.359E-02	2.025E-01	3.815E-02
10	4	2.143E-01	2.656E-02	2.850E-02	2.021E-01	3.741E-02
11	3	1.931E-01	2.503E-02	2.633E-02	1.959E-01	3.924E-02
12**	2	1.845E-01	2.378E-02	2.321E-02	1.658E-01	3.923E-02
13	1	4.259E-01	3.050E-02	3.442E-02	4.367E-01	5.145E-02

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

* tree, ** tree, + tree, and ++ tree all the same

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

D-mean is mean of Logsalary in the node

Cases fit give the number of cases used to fit node

MSE and R² are based on all cases in node

Node label	Total cases	Cases fit	Matrix rank	Node D-mean	Node MSE	Node R ²	Split variable	Other variables
1	263	263	5	5.945E+00	4.160E-01	0.4755	Yrs +Hitcr	
2T	143	143	4	5.506E+00	1.184E-01	0.8413	Hitcr +Runcr	
3T	120	120	4	6.469E+00	1.838E-01	0.4729	Wlk86 +Hit86	

The linear predictor is the one with a + sign under the 'Other variable' column.

Number of terminal nodes of final tree: 2

Total number of nodes of final tree: 3

Regression tree:

Node 1: Yrs <= 6.00000 or NA

Node 2: Logsalary-mean = 5.50632

Node 1: Yrs > 6.00000 and not NA

Node 3: Logsalary-mean = 6.46866

Node 1: Intermediate node

A case goes into Node 2 if Yrs <= 6.0000000E+00 or NA

Yrs mean = 7.3802E+00

Coefficients of least squares regression function:

Regressor	Coefficient	t-stat	p-val	Min	Mean	Max
Constant	5.3659E+00	72.24	0.0000			
Hitcr	8.9484E-04	14.45	0.0000	4.2000E+01	7.3392E+02	4.2560E+03
Div86.W	-1.8665E-01	-2.34	0.0000	0.0000E+00	5.0951E-01	1.0000E+00
Pos86.RF	4.6442E-01	3.22	0.0000	0.0000E+00	8.3650E-02	1.0000E+00
Pos86.UT	-5.5260E-01	-2.63	0.0000	0.0000E+00	3.8023E-02	1.0000E+00

Mean of Logsalary = 5.9454102795235446

Predicted values truncated at 4.2121275978784798 & 7.8079166289264101

Node 2: Terminal node

Coefficients of least squares regression functions:

Regressor	Coefficient	t-stat	p-val	Min	Mean	Max
Constant	4.4528E+00	89.49	0.0000			
Runcr	6.8396E-03	26.94	0.0000	1.8000E+01	1.6197E+02	5.2900E+02

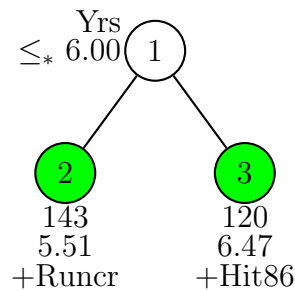


Figure 8: GUIDE piecewise ANCOVA regression tree for predicting Logsalary. At each intermediate node, an observation goes to the left branch if and only if the condition is satisfied. The symbol ‘ \leq_* ’ stands for ‘ \leq or missing’. Sample sizes, means of Logsalary, and signs and names of regressor variable printed below nodes. Terminal nodes with negative, zero, and positive slopes are colored red, yellow, and green, respectively.

```

Pos86.2B      -2.8035E-01      -2.72 0.0000      0.0000E+00      9.0909E-02      1.0000E+00
Pos86.CF      -2.7430E-01      -2.89 0.0000      0.0000E+00      1.0490E-01      1.0000E+00
Mean of Logsalary = 5.5063156054013094
Predicted values truncated at 4.2121275978784798      &      7.8079166289264101
-----
Node 3: Terminal node
Coefficients of least squares regression functions:
Regressor      Coefficient      t-stat p-val      Min      Mean      Max
Constant      5.5298E+00      46.67 0.0000
Hit86          9.2208E-03      9.08 0.0000      3.2000E+01      1.0759E+02      2.0000E+02
Pos86.2B      -2.7882E-01      -2.19 0.0000      0.0000E+00      1.0833E-01      1.0000E+00
Pos86.UT      -3.4458E-01      -2.11 0.0000      0.0000E+00      6.6667E-02      1.0000E+00
Mean of Logsalary = 6.4686647661858760
Predicted values truncated at 4.2121275978784798      &      7.8079166289264101
-----

Proportion of variance (R-squared) explained by tree model = .8172

Observed and fitted values are stored in ancova.fit
LaTeX code for tree is in ancova.tex
Elapsed time in seconds: 3.97944403
  
```

The \LaTeX tree, shown in Figure 8.

5.4 Quantile regression

Instead of estimating the conditional mean, we can estimate conditional quantiles (Chaudhuri and Loh, 2002; Koenker and Bassett, 1978). We demonstrate this with the data set `tuitiondat.txt`, which gives information on tuition and other variables for U.S. colleges. The data description file `tuitiondsc.txt` is:

```
tuitiondat.dat
NA
col_num var_name var_type
1 FICE x
2 CollName x
3 State x
4 PubPriv c
5 MathSAT x
6 VerbsAT x
7 CombsAT n
8 ACT x
9 Q1MSAT x
10 Q3MSAT x
11 Q1VSAT x
12 Q3VSAT x
13 Q1ACT x
14 Q3ACT x
15 AppsRec n
16 AppsAcc n
17 NewEnrol n
18 Top10 n
19 Top25 n
20 FUgrad n
21 PUgrad x
22 InTuition x
23 OutTuition d
24 RnBcost n
25 RmCost x
26 BrdCost x
27 AddFees x
28 BookCost x
29 PerSpend x
30 PFacPhD n
31 PFacTerm x
32 StudFac n
33 PAlDonate x
34 InstExp n
35 GradRate n
36 Type c
37 FullPSal n
```

```

38 AssocPSal x
39 AsstPSal x
40 AveSal x
41 FullPComp x
42 AssocPComp x
43 AsstPComp x
44 AveComp x
45 NFullProf n
46 NAssocProf x
47 NAsstProf x
48 NInstr x
49 NAllFac x

```

Following is a session log to create an input file for constructing a piecewise simple linear tree for the 90th percentile of out-of-state tuition (`OutTuition`).

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: quant.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: quant.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 2
Choose complexity of model to use at each node:
1: multiple linear, 2: best polynomial, 3: constant ([1:3], <cr>=1): 2
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):
Input 1 for 1 quantile, 2 for 2 quantiles ([1:2], <cr>=1):
  Option 2 allows simultaneous modeling of a pair of quantile values (e.g., 0.1 and 0.9)
Input quantile probability ([0.00:1.00], <cr>=0.50): 0.9

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): tuitiondsc.txt
Reading data description file ...
Training sample file: tuitiondat.txt
Missing value code: NA
Warning: B variables changed to C
Dependent variable is OutTuition

```

```

Length of longest data entry = 20
Total number of cases = 1134
Col. no. Categorical variable #levels #missing values
      4 PubPriv 2 0
     36 Type 3 0
Checking data ...
      Total #cases w/ #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
      1134 13 621 32 14 0 0 0 2
No. cases used for training = 1121
No. cases excluded due to 0 weight or missing D = 13
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): quant.tex
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node IDs and fitted values: quant.fit
Input file is created!

```

Results

Quantile regression tree with quantile probability .9000

No truncation of predicted values

Pruning by cross-validation

Data description file: tuitiondsc.txt

Training sample file: tuitiondat.txt

Missing value code: NA

Warning: B variables changed to C

Dependent variable is OutTuition

Piecewise simple linear or constant model

Length of longest data entry = 20

Summary information (without x variables)

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,

n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight

For categorical variables, #categories include one for missing values

Column number	Variable name	Variable type	Minimum value	Maximum value	Number of categories	Number missing
4	PubPriv	c			2	
7	CombSAT	n	6.0000E+02	1.4100E+03		471
15	AppsRec	n	5.7000E+01	4.8094E+04		9
16	AppsAcc	n	4.4000E+01	2.6330E+04		9
17	NewEnrol	n	2.1000E+01	7.4250E+03		5
18	Top10	n	1.0000E+00	9.8000E+01		183
19	Top25	n	1.1000E+01	1.0000E+02		155
20	FUgrad	n	1.1800E+02	3.1643E+04		3
23	OutTuition	d	1.0440E+03	2.5750E+04		13

24	RnBcost	n	1.3060E+03	8.7000E+03	57
30	PFacPhD	n	8.0000E+00	1.0500E+02	29
32	StudFac	n	2.5000E+00	4.2600E+01	2
34	InstExp	n	1.8340E+03	6.2469E+04	24
35	GradRate	n	8.0000E+00	1.1800E+02	69
36	Type	c			3
37	FullPSal	n	2.7000E+02	1.0090E+03	61
45	NFullProf	n	0.0000E+00	9.9700E+02	

Total #cases	w/ miss.	#missing D	ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
1134	13	621	32	14	0	0	0	0	2

No. cases used for training = 1121

No. cases excluded due to 0 weight or missing D = 13

Interaction tests on all variables

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Fraction of cases used for splitting each node = 1.0000

Max number of split levels = 11

Minimum node size = 55

Number of SE's for pruned tree = 5.0000E-01

Size and CV Loss and SE of subtrees:

Tree	#Tnodes	Mean Loss	SE(Mean)	BSE(Mean)	Median Loss	BSE(Median)
1	13	3.637E+02	1.442E+01	1.879E+01	3.605E+02	1.738E+01
2	12	3.641E+02	1.443E+01	1.883E+01	3.605E+02	1.783E+01
3+	11	3.644E+02	1.443E+01	1.881E+01	3.588E+02	1.708E+01
4	10	3.648E+02	1.447E+01	1.882E+01	3.608E+02	1.719E+01
5	9	3.612E+02	1.400E+01	1.990E+01	3.608E+02	2.034E+01
6++	8	3.568E+02	1.393E+01	1.794E+01	3.594E+02	2.100E+01
7	7	3.667E+02	1.407E+01	1.598E+01	3.680E+02	2.113E+01
8	6	3.667E+02	1.397E+01	1.489E+01	3.680E+02	1.910E+01
9**	5	3.576E+02	1.311E+01	1.272E+01	3.718E+02	1.420E+01
10	2	3.830E+02	1.368E+01	6.555E+00	3.862E+02	8.025E+00
11	1	4.625E+02	1.439E+01	1.121E+01	4.616E+02	2.334E+01

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

** tree same as -- tree

* tree same as ++ tree

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

D-quant is quantile of OutTuition in the node

Cases fit give the number of cases used to fit node

Node label	Total cases	Cases fit	Matrix rank	Node D-quant	Split variable	Other variables
1	1121	1098	2	1.591E+04	PubPriv	
2T	434	432	2	8.820E+03	RnBcost	
3	687	672	2	1.745E+04	RnBcost	
6T	209	203	2	1.160E+04	FUgrad	
7	478	469	2	1.848E+04	AppsAcc	
14T	220	214	2	1.468E+04	FullPSal	
15	258	255	2	1.900E+04	StudFac	
30T	108	105	2	1.970E+04	-	
31T	150	147	2	1.615E+04	AppsRec	

Number of terminal nodes of final tree: 5

Total number of nodes of final tree: 9

Regression tree:

Node 1: PubPriv = "Public"

Node 2: OutTuition sample quantile = 8.82000E+03

Node 1: PubPriv /= "Public"

Node 3: RnBcost <= 3.89050E+03

Node 6: OutTuition sample quantile = 1.16000E+04

Node 3: RnBcost > 3.89050E+03 or NA

Node 7: AppsAcc <= 9.26500E+02

Node 14: OutTuition sample quantile = 1.46750E+04

Node 7: AppsAcc > 9.26500E+02 or NA

Node 15: StudFac <= 12.10000

Node 30: OutTuition sample quantile = 1.97000E+04

Node 15: StudFac > 12.10000 or NA

Node 31: OutTuition sample quantile = 1.61545E+04

In the following the predictor node mean is mean of complete cases

Regression coefficients are computed from the complete cases

Node 1: Intermediate node

A case goes into Node 2 if PubPriv = "Public"

PubPriv mode = "Private"

Coefficients of quantile regression function:

Regressor	Coefficient	Min	Mean	Max
Constant	3.7443E+03			

```

InstExp      1.0271E+00    1.8340E+03    9.0272E+03    6.2469E+04
If regression function is inapplicable due to missing values, predicted quantile =
  10430.000000000000
-----

Node 2: Terminal node
Coefficients of quantile regression function:
Regressor    Coefficient      Min          Mean          Max
Constant     -8.9699E+02
FullPSal     1.6871E+01    3.5900E+02    5.4766E+02    8.9300E+02
If regression function is inapplicable due to missing values, predicted quantile =
  6297.000000000000
-----

Node 3: Intermediate node
A case goes into Node 6 if RnBcost <=  3.8905000E+03
RnBcost mean =  4.5535E+03
-----

Node 6: Terminal node
Coefficients of quantile regression function:
Regressor    Coefficient      Min          Mean          Max
Constant     5.8780E+03
InstExp      7.1975E-01    3.3650E+03    7.6059E+03    4.2926E+04
If regression function is inapplicable due to missing values, predicted quantile =
  10430.000000000000
-----

Node 7: Intermediate node
A case goes into Node 14 if AppsAcc <=  9.2650000E+02
AppsAcc mean =  1.4274E+03
-----

Node 14: Terminal node
Coefficients of quantile regression function:
Regressor    Coefficient      Min          Mean          Max
Constant     6.1379E+03
InstExp      7.7936E-01    3.0190E+03    9.6447E+03    6.2469E+04
If regression function is inapplicable due to missing values, predicted quantile =
  12700.000000000000
-----

Node 15: Intermediate node
A case goes into Node 30 if StudFac <=  1.2100000E+01
StudFac mean =  1.2708E+01
-----

Node 30: Terminal node
Coefficients of quantile regression function:
Regressor    Coefficient      Min          Mean          Max
Constant     1.8756E+04
RnBcost      1.3557E-01    3.8910E+03    5.6088E+03    8.1240E+03
If regression function is inapplicable due to missing values, predicted quantile =

```

```

25750.000000000000
-----
Node 31: Terminal node
Coefficients of quantile regression function:
Regressor    Coefficient      Min      Mean      Max
Constant      5.9828E+03
InstExp       9.1899E-01    2.5890E+03  9.4031E+03  2.2704E+04
If regression function is inapplicable due to missing values, predicted quantile =
10258.000000000000
-----

Observed and fitted values are stored in quant.fit
LaTeX code for tree is in quant.tex
Elapsed time in seconds: 55.0259514

```

The \LaTeX tree is shown in Figure 9 and plots of the data in the terminal nodes of the tree are given in Figure 10.

5.5 Least median of squares

Although median regression may be preferred to least-squares regression if there are large outliers in a data set, an alternative that is even more robust to outliers is *least median of squares* regression (Rousseeuw and Leroy, 1987). GUIDE can construct tree models using this criterion. We use the college tuition data for illustration. A session log of the input file generation is below, followed by the results and the \LaTeX tree diagram in Figure 11.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: lms.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: lms.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1):
Input 1 for least squares, 2 least median of squares ([1:2], <cr>=1): 2

```

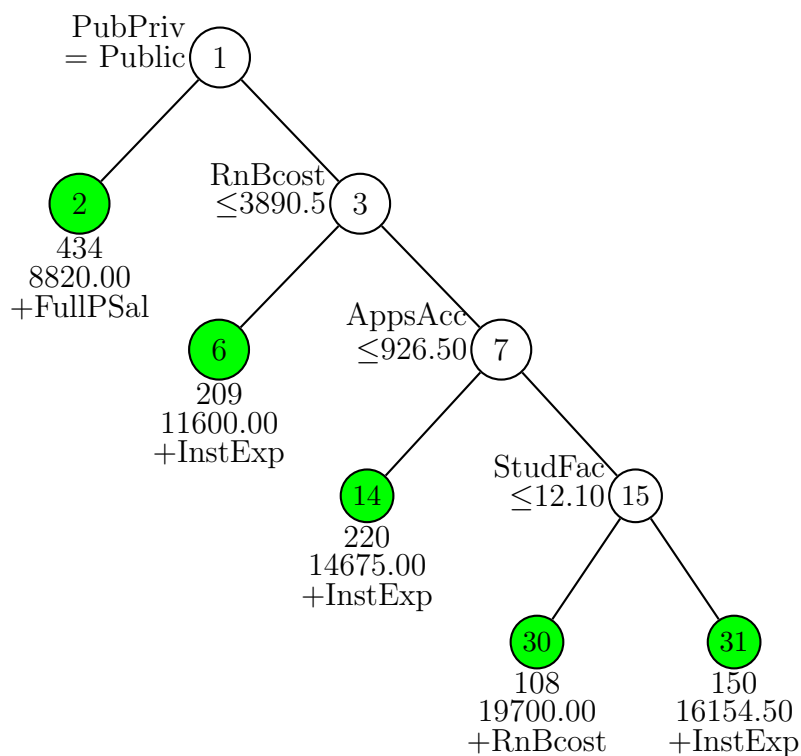


Figure 9: GUIDE piecewise simple linear 0.900-quantile regression tree for predicting OutTuition. At each intermediate node, an observation goes to the left branch if and only if the condition is satisfied. Sample sizes, 0.900-quantiles of OutTuition, and signs and names of best regressor printed below nodes. Terminal nodes with negative, zero, and positive slopes are colored red, yellow, and green, respectively.

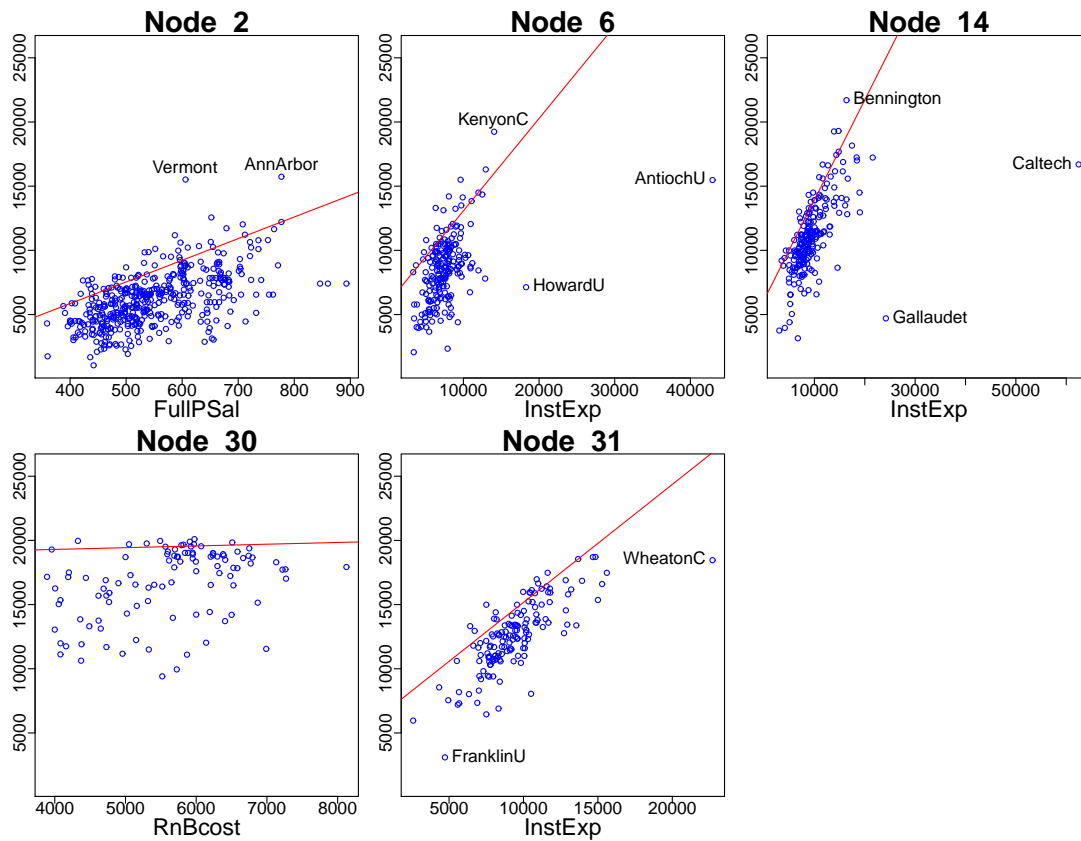


Figure 10: Plots of data and estimated linear 0.9-quantile fits in terminal nodes of tree in Figure 9.

This is where the option for least median of squares is selected.

Choose complexity of model to use at each node:

1: multiple linear, 2: best simple linear, 3: constant ([1:3], <cr>=2):

Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes if it contains spaces or non alphanumeric characters): tuitiondsc.txt

Reading data description file ...

Training sample file: tuitiondat.txt

Missing value code: NA

Warning: B variables changed to C

Dependent variable is OutTuition

Length of longest data entry = 20

Total number of cases = 1134

Col. no.	Categorical variable	#levels	#missing values
4	PubPriv	2	0
36	Type	3	0

Checking data ...

#cases	Total #cases w/ miss. D	ord. vals	#missing	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
1134	13	621		32	14	0	0	0	2

No weight variable in data file

No. cases used for training = 1121

No. cases excluded due to 0 weight or missing D = 13

Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):

Input file name to store LaTeX code (use .tex as suffix): lms.tex

Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2

Input name of file to store node IDs and fitted values: lms.fit

Input file is created!

Results

Least median of squares regression tree

Predictions truncated at global min and max of D sample values

Pruning by cross-validation

Data description file: tuitiondsc.txt

Training sample file: tuitiondat.txt

Missing value code: NA

Warning: B variables changed to C

Dependent variable is OutTuition

Piecewise simple linear or constant model

Length of longest data entry = 20

Summary information (without x variables)

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,

```

n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
Column Name Minimum Maximum #Categories #Missing
  4 PubPriv c 2
  7 CombSAT n 6.0000E+02 1.4100E+03 471
 15 AppsRec n 5.7000E+01 4.8094E+04 9
 16 AppsAcc n 4.4000E+01 2.6330E+04 9
 17 NewEnrol n 2.1000E+01 7.4250E+03 5
 18 Top10 n 1.0000E+00 9.8000E+01 183
 19 Top25 n 1.1000E+01 1.0000E+02 155
 20 FUgrad n 1.1800E+02 3.1643E+04 3
 23 OutTuition d 1.0440E+03 2.5750E+04 13
 24 RnBcost n 1.3060E+03 8.7000E+03 57
 30 PFacPhD n 8.0000E+00 1.0500E+02 29
 32 StudFac n 2.5000E+00 4.2600E+01 2
 34 InstExp n 1.8340E+03 6.2469E+04 24
 35 GradRate n 8.0000E+00 1.1800E+02 69
 36 Type c 3
 37 FullPSal n 2.7000E+02 1.0090E+03 61
 45 NFullProf n 0.0000E+00 9.9700E+02

```

```

Total #cases w/ #missing
#cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
 1134      13      621      32      14      0      0      0      2

```

No weight variable in data file

No. cases used for training: 1121

No. cases excluded due to 0 weight or missing D: 13

Missing N and F values imputed with node means

Interaction tests on all variables

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Fraction of cases used for splitting each node = 1.0000

Max number of split levels = 11

Minimum node size = 55

Number of SE's for pruned tree = 5.0000E-01

Size and CV median absolute residual (MAR) and SE of subtrees:

Tree	#Tnodes	Mean MAR	BSE(Mean)	Median MAR	BSE(Median)
1	14	1.344E+05	5.452E+01	1.176E+03	8.878E+01
2+	13	1.325E+05	4.928E+01	1.148E+03	7.454E+01
3	12	1.333E+05	4.793E+01	1.161E+03	6.814E+01
4	10	1.353E+05	5.546E+01	1.161E+03	8.151E+01
5	9	1.367E+05	4.932E+01	1.222E+03	6.692E+01
6	8	1.364E+05	4.972E+01	1.230E+03	7.063E+01
7--	7	1.313E+05	2.490E+01	1.181E+03	4.526E+01
8	6	1.340E+05	2.932E+01	1.223E+03	3.248E+01

9	5	1.388E+05	2.292E+01	1.248E+03	2.001E+01
10	4	1.406E+05	3.343E+01	1.264E+03	3.146E+01
11	3	1.476E+05	5.183E+01	1.266E+03	7.042E+01
12	2	1.752E+05	8.364E+01	1.553E+03	1.203E+02
13	1	1.901E+05	8.738E+01	1.716E+03	9.536E+01

The selected tree is marked by two dashes.

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

++ tree same as -- tree

* tree same as ++ tree

* tree same as -- tree

Following tree is based on mean CV with bootstrap SE estimate (--).

Structure of final tree. Each terminal node is marked with a T.

D-mean is mean of OutTuition in the node

Cases fit give the number of cases used to fit node

MAR is median of absolute residuals

Node label	Total cases	Cases fit	Matrix rank	Node D-median	Node MAR	Split variable	Other variables
1	1121	1121	2	8.820E+03	1.687E+03	PubPriv	+InstExp
2	434	434	2	6.114E+03	9.854E+02	NFullProf	+RnBcost
4T	198	198	2	5.130E+03	7.377E+02	StudFac	+RnBcost
5T	236	236	2	6.832E+03	9.981E+02	InstExp	+FullPSal
3	687	687	2	1.096E+04	1.485E+03	InstExp	+InstExp
6	577	577	2	1.042E+04	1.314E+03	RnBcost	+InstExp
12	242	242	2	8.444E+03	1.268E+03	InstExp	+InstExp
24T	179	179	2	7.896E+03	1.178E+03	RnBcost	+InstExp
25T	63	63	2	9.950E+03	7.434E+02	-	+FUgrad
13	335	335	2	1.155E+04	1.072E+03	NewEnrol	+InstExp
26T	157	157	2	1.086E+04	8.482E+02	RnBcost	+InstExp
27T	178	178	2	1.259E+04	1.006E+03	Type	+InstExp
7T	110	110	2	1.788E+04	8.988E+02	-	+GradRate

Number of terminal nodes of final tree: 7

Total number of nodes of final tree: 13

Regression tree:

Node 1: PubPriv = "Public"

Node 2: NFullProf <= 93.00000

Node 4: OutTuition-mean = 5.13000E+03

```

Node 2: NFullProf > 93.00000 or NA
Node 5: OutTuition-mean = 6.83150E+03
Node 1: PubPriv /= "Public"
Node 3: InstExp <= 1.38525E+04 or NA
Node 6: RnBcost <= 3.99950E+03 or NA
Node 12: InstExp <= 8.33000E+03 or NA
Node 24: OutTuition-mean = 7.89600E+03
Node 12: InstExp > 8.33000E+03 and not NA
Node 25: OutTuition-mean = 9.95000E+03
Node 6: RnBcost > 3.99950E+03 and not NA
Node 13: NewEnrol <= 3.19500E+02
Node 26: OutTuition-mean = 1.08600E+04
Node 13: NewEnrol > 3.19500E+02 or NA
Node 27: OutTuition-mean = 1.25860E+04
Node 3: InstExp > 1.38525E+04 and not NA
Node 7: OutTuition-mean = 1.78825E+04

```

In the following the predictor node mean is mean of complete cases
Regression coefficients are computed from the complete cases

Node 1: Intermediate node

A case goes into Node 2 if PubPriv = "Public"

PubPriv mode = "Private"

Coefficients of least median of squares regression function:

Regressor	Coefficient	Minimum	Mean	Maximum
Constant	973.9978			

InstExp	1.0571E+00	1834.00*****	6.2469E+04	
---------	------------	--------------	------------	--

Mean of OutTuition = 8820.000000000000

Predicted values truncated at 1044.000000000000 & 25750.000000000000

Node 2: Intermediate node

A case goes into Node 4 if NFullProf <= 9.3000000E+01

NFullProf mean = 1.7104E+02

Node 4: Terminal node

Coefficients of least median of squares regression function:

Regressor	Coefficient	Minimum	Mean	Maximum
Constant	1.3867E+03			

RnBcost	1.2073E+00	1306.00*****	5.5440E+03	
---------	------------	--------------	------------	--

Mean of OutTuition = 5130.000000000000

Predicted values truncated at 1044.000000000000 & 25750.000000000000

Node 5: Terminal node

Coefficients of least median of squares regression function:

```

Regressor Coefficient      Minimum      Mean      Maximum
Constant   -330.8814
FullPSal    1.2814E+01    440.00***** 8.9300E+02
Mean of OutTuition = 6831.500000000000000
Predicted values truncated at 1044.000000000000000 & 25750.0000000000000
-----

Node 3: Intermediate node
A case goes into Node 6 if InstExp <= 1.3852500E+04 or NA
InstExp mean = 1.0340E+04
-----

Node 6: Intermediate node
A case goes into Node 12 if RnBcost <= 3.9995000E+03 or NA
RnBcost mean = 4.3413E+03
-----

Node 12: Intermediate node
A case goes into Node 24 if InstExp <= 8.3300000E+03 or NA
InstExp mean = 7.3339E+03
-----

Node 24: Terminal node
Coefficients of least median of squares regression function:
Regressor Coefficient      Minimum      Mean      Maximum
Constant   -1.3502E+03
InstExp     1.3199E+00    2589.00***** 8.3240E+03
Mean of OutTuition = 7896.000000000000000
Predicted values truncated at 1044.000000000000000 & 25750.0000000000000
-----

Node 25: Terminal node
Coefficients of least median of squares regression function:
Regressor Coefficient      Minimum      Mean      Maximum
Constant     8.1732E+03
FUgrad       2.3909E+00    139.00***** 5.0640E+03
Mean of OutTuition = 9950.000000000000000
Predicted values truncated at 1044.000000000000000 & 25750.0000000000000
-----

Node 13: Intermediate node
A case goes into Node 26 if NewEnrol <= 3.1950000E+02
NewEnrol mean = 4.1792E+02
-----

Node 26: Terminal node
Coefficients of least median of squares regression function:
Regressor Coefficient      Minimum      Mean      Maximum
Constant     6.7469E+03
InstExp       3.9716E-01    4054.00***** 1.3844E+04
Mean of OutTuition = 10860.000000000000000
Predicted values truncated at 1044.000000000000000 & 25750.0000000000000
-----

```

```

Node 27: Terminal node
Coefficients of least median of squares regression function:
Regressor Coefficient      Minimum      Mean      Maximum
Constant      792.8736
InstExp      1.3078E+00  3480.00*****  1.3706E+04
Mean of OutTuition = 12586.000000000000
Predicted values truncated at 1044.0000000000000 & 25750.000000000000
-----

Node 7: Terminal node
Coefficients of least median of squares regression function:
Regressor Coefficient      Minimum      Mean      Maximum
Constant      1.3307E+04
GradRate      6.1667E+01  46.0082.1468  1.0000E+02
Mean of OutTuition = 17882.500000000000
Predicted values truncated at 1044.0000000000000 & 25750.000000000000
-----

Proportion of deviance explained by tree model = 0.64638321073675831

Observed and fitted values are stored in lms.fit
LaTeX code for tree is in lms.tex
Elapsed time in seconds: 2.67295408

```

5.6 Poisson regression with offset

We use a data set from www.statsci.org/data/general/motorins.html on motor insurance claims in Sweden for the year 1977 (Andrews and Herzberg, 1985; Hallin and Ingenbleek, 1983). The description and data files are `swedendsc.txt` and `swendendat.txt`. The dependent variable is the number of claims. The other variables are `mileagegp` with ordered values 1–5, `zone` with 7 unordered values, `bonus` which is the number of years plus one since last claim, `make` of car with 9 unordered values, `insured` which is the number of insured in policy-years, and `payment` which is the total value of payments in Skr. We ignore `insured` and `payments` by giving them the `x` designation. To fit a Poisson regression model for the claim rate, we created an offset variable `lninsured` which is the log of `insured` and designate as `z`. Because `mileagegp` is an ordered categorical variable, we designate it as `s` to prevent it from being used as a linear predictor for fitting the Poisson node models. We designate `bonus` as `f` so that it is only used as a linear predictor and not for splitting the nodes.

```

swendendat.txt
?
column  variable  type
1 mileagegp s

```

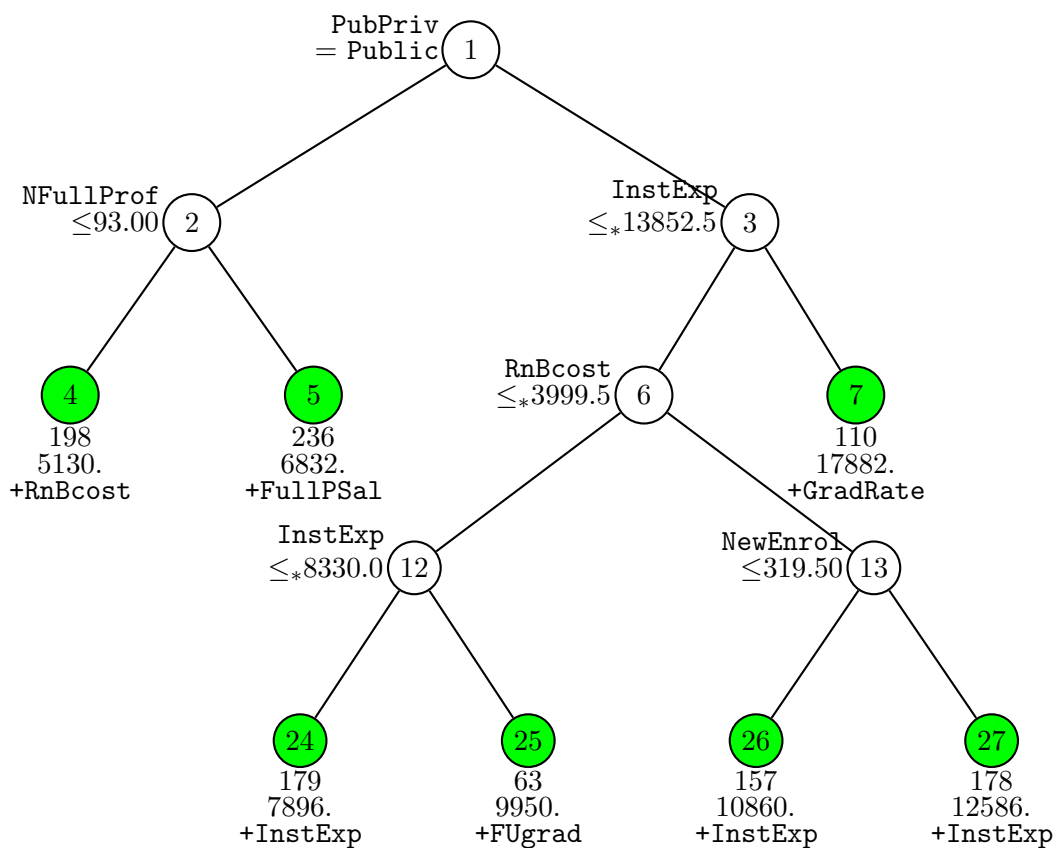


Figure 11: GUIDE 0.50-SE piecewise simple linear least-median-of-squares regression tree for predicting `OutTuition`. At each split, an observation goes to the left branch if and only if the condition is satisfied. The symbol ' \leq_* ' stands for ' \leq or missing'. Sample sizes, means of `OutTuition`, and signs and names of regressor variable are printed below nodes. Terminal nodes with negative, zero, and positive slopes are colored red, yellow, and green, respectively.

```

2 zone c
3 bonus f
4 make c
5 insured x
6 lninsured z
7 claims d
8 payments x

```

Since there is only one linear predictor, the multiple linear Poisson model is the same as the best simple linear Poisson model.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: input.txt
File input.txt exists
Input 1 to overwrite it, 2 to choose another name ([1:2], <cr>=1):
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: output.txt
File output.txt exists
Input 1 to overwrite it, 2 to choose another name ([1:2], <cr>=1):
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 3
Choose complexity of model to use at each node:
1: multiple linear, 2: best polynomial, 3: constant ([1:3], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): swedendsc.txt
Reading data description file ...
Training sample file: swedendat.txt
Missing value code: ?
Dependent variable is claims
Length of longest data entry = 19
Total number of cases =      2182
Col. no. Categorical variable    #levels    #missing values
      2 zone                     7             0
      4 make                     9             0
Checking data ...
Number of cases with positive D values =      1797

```

```

      Total #cases w/ #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
      2182      0      0      2      0      1      1      0      2
Offset variable in column:      6
No. cases used for training =      2182
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): sweden.tex
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1):
Input file is created!

```

Results

```

Poisson regression tree
No truncation of predicted values
Pruning by cross-validation
Data description file: swedendsc.txt
Training sample file: swendendat.txt
Missing value code: ?
Dependent variable is claims
Piecewise linear model
Length of longest data entry = 19
Number of cases with positive D values: 1797

```

```

Summary information (without x variables)
d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,
n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
z=offset variable

```

Column	Name		Minimum	Maximum	#Categories	#Missing
1	mileagegp	s	1.0000E+00	5.0000E+00		
2	zone	c			7	
3	bonus	f	1.0000E+00	7.0000E+00		
4	make	c			9	
6	lninsured	z	-4.6052E+00	1.1757E+01		
7	claims	d	0.0000E+00	3.3380E+03		

```

      Total #cases w/ #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
      2182      0      0      2      0      1      1      0      2
Offset variable in column 6
No. cases used for training: 2182

```

```

Interaction tests on all variables
Pruning by v-fold cross-validation, with v = 10
Selected tree is based on mean of CV estimates
Fraction of cases used for splitting each node = 1.0000

```

Max number of split levels = 12
 Minimum node size = 109
 100 bootstrap calibration replicates
 Number of SE's for pruned tree = 5.0000E-01

Size and CV Loss and SE of subtrees:

Tree	#Tnodes	Mean Loss	SE(Mean)	BSE(Mean)	Median Loss	BSE(Median)
1*	10	3.217E+00	2.507E-01	2.542E-01	2.870E+00	2.166E-01
2	9	3.229E+00	2.510E-01	2.538E-01	2.964E+00	2.098E-01
3**	7	3.301E+00	2.591E-01	2.974E-01	2.964E+00	2.109E-01
4	6	3.458E+00	2.812E-01	2.940E-01	3.196E+00	2.070E-01
5	5	3.584E+00	2.871E-01	2.760E-01	3.355E+00	2.192E-01
6	4	3.566E+00	2.974E-01	2.813E-01	3.523E+00	2.159E-01
7	3	4.243E+00	3.823E-01	3.722E-01	4.207E+00	5.342E-01
8	1	6.519E+00	6.365E-01	5.622E-01	6.609E+00	7.858E-01

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

* tree same as + tree

** tree same as ++ tree

** tree same as -- tree

++ tree same as -- tree

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Rate is mean of Y/exp(offset)

Cases fit give the number of cases used to fit node

Deviance is mean residual deviance for all cases in node

Node label	Total cases	Cases fit	Matrix rank	Node rate	Node deviance	Split variable	Other variables
1	2182	2182	2	4.749E-02	6.408E+00	make	
2T	482	482	2	3.453E-02	2.253E+00	zone	
3	1700	1700	2	4.864E-02	6.663E+00	zone	
6	490	490	2	6.375E-02	6.030E+00	zone	
12T	245	245	2	5.627E-02	4.450E+00	-	
13T	245	245	2	7.258E-02	5.188E+00	-	
7	1210	1210	2	4.209E-02	3.927E+00	mileagegp	
14T	243	243	2	3.613E-02	4.163E+00	-	
15	967	967	2	4.497E-02	2.558E+00	mileagegp	
30	490	490	2	4.368E-02	2.986E+00	zone	

60T	196	196	2	5.011E-02	2.371E+00	-
61T	294	294	2	4.056E-02	2.316E+00	make
31T	477	477	2	5.024E-02	1.229E+00	mileagegp

Number of terminal nodes of final tree: 7

Total number of nodes of final tree: 13

Regression tree:

Node 1: make = "4", "6"

Node 2: claims sample rate = 3.45272E-02

Node 1: make /= "4", "6"

Node 3: zone = "1", "2"

Node 6: zone = "2"

Node 12: claims sample rate = 5.62714E-02

Node 6: zone /= "2"

Node 13: claims sample rate = 7.25828E-02

Node 3: zone /= "1", "2"

Node 7: mileagegp <= 1.00000

Node 14: claims sample rate = 3.61274E-02

Node 7: mileagegp > 1.00000 or ?

Node 15: mileagegp <= 3.00000 or ?

Node 30: zone = "3", "5"

Node 60: claims sample rate = 5.01083E-02

Node 30: zone /= "3", "5"

Node 61: claims sample rate = 4.05632E-02

Node 15: mileagegp > 3.00000 and not ?

Node 31: claims sample rate = 5.02372E-02

Node 1: Intermediate node

A case goes into Node 2 if make =

"4", "6"

make mode = "1"

Coefficients of loglinear regression function:

Regressor	Coefficient	t-stat	p-val	Minimum	Mean	Maximum
Constant	-2.0572E+00	-305.02	0.0000			
bonus	-1.8651E-01	-148.88	0.0000	1.0000E+00	4.0151E+00	7.0000E+00

Node mean for offset variable = 4.4928E+00

If regression function is inapplicable due to missing values, predicted rate =

4.7487588464521328E-002

Node 2: Terminal node

Coefficients of loglinear regression function:

```

Regressor Coefficient      t-stat  p-val      Minimum      Mean      Maximum
Constant -2.5849E+00      -104.45 0.0000
bonus     -1.6243E-01      -32.72 0.0000  1.0000E+00  4.0249E+00  7.0000E+00
Node mean for offset variable =      4.3366E+00
If regression function is inapplicable due to missing values, predicted rate =
  3.4527165629107126E-002
-----
Node 3: Intermediate node
A case goes into Node 6 if zone =
  "1", "2"
zone mode = "1"
-----
Node 6: Intermediate node
A case goes into Node 12 if zone = "2"
zone mode = "1"
-----
Node 12: Terminal node
Coefficients of loglinear regression function:
Regressor Coefficient      t-stat  p-val      Minimum      Mean      Maximum
Constant -1.8885E+00      -117.98 0.0000
bonus     -1.8707E-01      -62.66 0.0000  1.0000E+00  4.0000E+00  7.0000E+00
Node mean for offset variable =      5.1038E+00
If regression function is inapplicable due to missing values, predicted rate =
  5.6271427885877690E-002
-----
Node 13: Terminal node
Coefficients of loglinear regression function:
Regressor Coefficient      t-stat  p-val      Minimum      Mean      Maximum
Constant -1.5438E+00      -107.39 0.0000
bonus     -2.1116E-01      -75.60 0.0000  1.0000E+00  4.0000E+00  7.0000E+00
Node mean for offset variable =      4.9811E+00
If regression function is inapplicable due to missing values, predicted rate =
  7.2582822726790452E-002
-----
Node 7: Intermediate node
A case goes into Node 14 if mileagegp <=  1.0000000E+00
mileagegp mean =  2.9860E+00
-----
Node 14: Terminal node
Coefficients of loglinear regression function:
Regressor Coefficient      t-stat  p-val      Minimum      Mean      Maximum
Constant -2.3934E+00      -152.64 0.0000
bonus     -1.8117E-01      -60.43 0.0000  1.0000E+00  4.0123E+00  7.0000E+00
Node mean for offset variable =      4.6679E+00
If regression function is inapplicable due to missing values, predicted rate =
  3.6127394092405472E-002

```

```

-----
Node 15: Intermediate node
A case goes into Node 30 if mileagegp <= 3.0000000E+00 or ?
mileagegp mean = 3.4850E+00
-----
Node 30: Intermediate node
A case goes into Node 60 if zone =
"3", "5"
zone mode = "3"
-----
Node 60: Terminal node
Coefficients of loglinear regression function:
Regressor Coefficient    t-stat  p-val    Minimum    Mean    Maximum
Constant -1.9094E+00        -93.79 0.0000
bonus     -2.0022E-01       -53.60 0.0000  1.0000E+00 4.0000E+00 7.0000E+00
Node mean for offset variable = 5.1513E+00
If regression function is inapplicable due to missing values, predicted rate =
5.0108317504291920E-002
-----
Node 61: Terminal node
Coefficients of loglinear regression function:
Regressor Coefficient    t-stat  p-val    Minimum    Mean    Maximum
Constant -2.0908E+00       -131.72 0.0000
bonus     -2.0319E-01       -70.53 0.0000  1.0000E+00 4.0000E+00 7.0000E+00
Node mean for offset variable = 4.7057E+00
If regression function is inapplicable due to missing values, predicted rate =
4.0563186633902397E-002
-----
Node 31: Terminal node
Coefficients of loglinear regression function:
Regressor Coefficient    t-stat  p-val    Minimum    Mean    Maximum
Constant -2.0319E+00       -64.33 0.0000
bonus     -1.6109E-01       -30.64 0.0000  1.0000E+00 4.0377E+00 7.0000E+00
Node mean for offset variable = 3.5950E+00
If regression function is inapplicable due to missing values, predicted rate =
5.0237213688715164E-002
-----

LaTeX code for tree is in sweden.tex
Elapsed time in seconds: 3.71040082

```

The tree is shown in Figure 12.

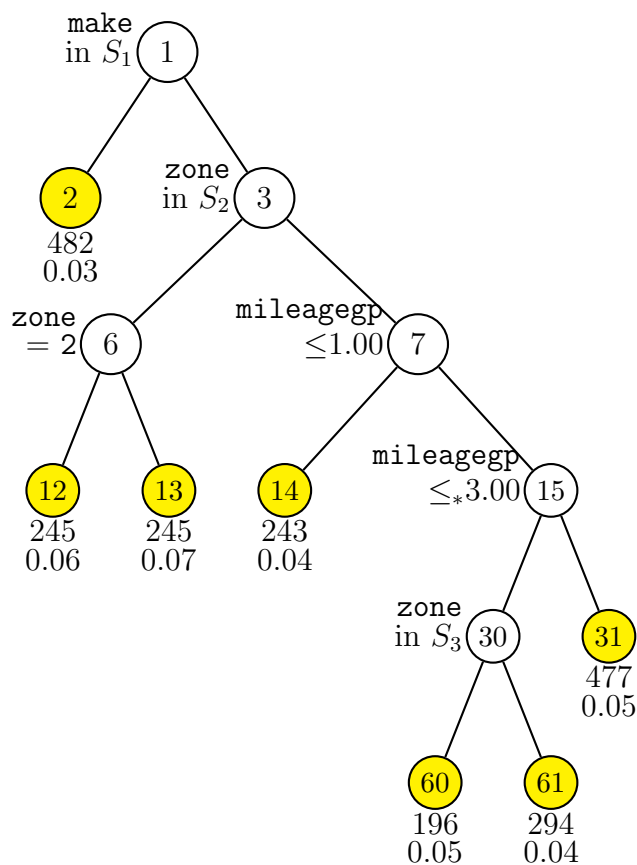


Figure 12: GUIDE 0.50-SE Poisson regression tree for predicting `claims`. At each split, an observation goes to the left branch if and only if the condition is satisfied. The symbol ' \leq_* ' stands for ' \leq or missing'. Set $S_1 = \{4, 6\}$. Set $S_2 = \{1, 2\}$. Set $S_3 = \{3, 5\}$. Sample sizes, sample rates, and names of regressor are printed below nodes.

5.7 Censored response data

GUIDE can fit a piecewise-constant, piecewise-simple linear, or piecewise multiple linear proportional hazards regression model to censored response data. Using usual notation, let $\lambda(\mathbf{x}, t)$ denote the hazard rate at time t for a subject with covariate vector \mathbf{x} . In a proportional hazards model, the hazard rate can be factored as $\lambda(\mathbf{x}, t) = \lambda_0(t)f(\mathbf{x}, \boldsymbol{\beta})$, where $\lambda_0(t)$ is a “baseline” hazard rate that is independent of the covariates and $f(\mathbf{x}, \boldsymbol{\beta})$ is a function of \mathbf{x} and some coefficients $\boldsymbol{\beta}$, independent of t . The Cox proportional hazards model uses $\lambda(\mathbf{x}, t) = \lambda_0(t) \exp(\boldsymbol{\beta}'\mathbf{x})$. GUIDE fits the more general model

$$\lambda(\mathbf{x}, t) = \lambda_0(t) \sum_i I(\mathbf{x} \in S_i) \exp(\boldsymbol{\beta}_i'\mathbf{x}),$$

where S_i is a set corresponding node i and $\boldsymbol{\beta}_i$ is its associated coefficient vector. See [Loh et al. \(2015\)](#) for more details.

We illustrate the piecewise-constant model $\lambda(\mathbf{x}, t) = \lambda_0(t) \sum_i I(\mathbf{x} \in S_i) \exp(\beta_{i0})$ with a data set from the Worcester Heart Attack Study analyzed in [Hosmer et al. \(2008\)](#). The data are in the file `whas500.csv` and the description file in `whas500.dsc` whose contents are repeated below.

```
whas500.csv
NA
c1 c2 c3
1 id x
2 age n
3 gender c
4 hr n
5 sysbp n
6 diasbp n
7 bmi n
8 cvd c
9 afb c
10 sho c
11 chf c
12 av3 c
13 miord c
14 mitype c
15 year c
16 admitdate x
17 disdate x
18 fdate x
19 los n
20 dstat x
21 lenfol t
```

22 fstat d

The goal of the study is to observe survival rates following hospital admission for acute myocardial infarction. The response variable is `lenfol`, which stands for total length of follow-up in days. Variable `fstat` is status at last follow-up (0=alive, 1=dead) and variable `chf` is congestive heart complications (0=no, 1=yes).

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: whas500.in
Input 1 for model fitting, 2 for importance or DIF scoring,
    3 for data conversion ([1:3], <cr>=1):
Name of batch output file: whas500.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
    1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
    5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 4
Choose complexity of model to use at each node:
1: multiple linear, 2: best simple linear, 3: constant ([1:3], <cr>=3): 3
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (max 100 characters);
enclose with matching quotes if it has spaces: whas500.dsc
Reading data description file ...
Training sample file: whas500.csv
Missing value code: NA
Warning: N variables changed to S
Dependent variable is fstat
Length of longest data entry = 10
Total number of cases: 500
Col. no. Categorical variable    #levels    #missing values
      3 gender                  2              0
      8 cvd                    2              0
      9 afb                    2              0
     10 sho                    2              0
     11 chf                    2              0
     12 av3                    2              0
     13 miord                  2              0
     14 mitype                 2              0
     15 year                   3              0

```

```

Checking data ...
Smallest uncensored T: 1.0000
No. complete cases excluding censored T < smallest uncensored T: 500
No. cases used to compute baseline hazard: 500
No. cases with D=1 and T >= smallest uncensored: 215
Total #cases w/ #missing
#cases   miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
    500      0      0      5      0      0      6      0      9
Survival time variable in column: 21
Censoring indicator variable in column: 22
Proportion of uncensored among nonmissing T and D variables = 0.430
No. cases used for training: 500
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): whas500.tex
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=2):
Input name of file to store node IDs and fitted values: whas500.fit
Input file is created!

```

Results

```

Proportional hazards regression with relative risk estimates
Pruning by cross-validation
Data description file: whas500.dsc
Training sample file: whas500.csv
Missing value code: NA
Warning: N variables changed to S
Dependent variable is fstat
Piecewise constant model
Length of longest data entry = 10
Smallest uncensored T: 1.0000
No. complete cases excluding censored T < smallest uncensored T: 500
No. cases used to compute baseline hazard: 500
No. cases with D=1 and T >= smallest uncensored: 215

Summary information (without x variables)
d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,
n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
t=survival time variable

```

Column	Name		Minimum	Maximum	#Categories	#Missing
2	age	s	3.0000E+01	1.0400E+02		
3	gender	c			2	
4	hr	s	3.5000E+01	1.8600E+02		
5	sysbp	s	5.7000E+01	2.4400E+02		
6	diasbp	s	6.0000E+00	1.9800E+02		
7	bmi	s	1.3045E+01	4.4839E+01		
8	cvd	c			2	

```

    9  afb      c      2
   10  sho      c      2
   11  chf      c      2
   12  av3      c      2
   13  miord    c      2
   14  mitype   c      2
   15  year     c      3
   19  los      s  0.0000E+00  4.7000E+01
   21  lenfol   t  1.0000E+00  2.3580E+03
   22  fstat    d  0.0000E+00  1.0000E+00
===== Constructed variables =====
   23  lnbasehaz z -4.1352E+00  9.7549E-01

Total  #cases w/  #missing
#cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
    500      0      0      5      0      0      6      0      9
Survival time variable in column: 21
Censoring indicator variable in column: 22
Proportion of uncensored among nonmissing T and D variables: 0.430
No. cases used for training: 500

Interaction tests on all variables
Pruning by v-fold cross-validation, with v = 10
Selected tree is based on mean of CV estimates
Fraction of cases used for splitting each node = 1.0000
Max number of split levels = 10
Minimum node size = 3
Number of iterations = 5
Number of SE's for pruned tree = 5.0000E-01

Size and CV Loss and SE of subtrees:
Tree  #Tnodes  Mean Loss  SE(Mean)  BSE(Mean)  Median Loss  BSE(Median)
    1      60  1.579E+00  1.039E-01  7.117E-02  1.630E+00  8.633E-02
    2      59  1.579E+00  1.039E-01  7.132E-02  1.632E+00  8.726E-02
    3      58  1.574E+00  1.037E-01  7.088E-02  1.629E+00  8.560E-02
    4      56  1.575E+00  1.037E-01  7.134E-02  1.629E+00  8.698E-02
    5      55  1.576E+00  1.036E-01  7.085E-02  1.623E+00  8.370E-02
    6      54  1.565E+00  1.021E-01  7.149E-02  1.609E+00  8.418E-02
    7      53  1.570E+00  1.023E-01  6.965E-02  1.609E+00  8.278E-02
    8      52  1.573E+00  1.024E-01  6.999E-02  1.609E+00  8.370E-02
    9      49  1.578E+00  1.023E-01  7.098E-02  1.609E+00  8.634E-02
   10      48  1.574E+00  1.023E-01  7.246E-02  1.609E+00  8.744E-02
   11      46  1.574E+00  1.023E-01  7.246E-02  1.609E+00  8.744E-02
   12      45  1.569E+00  1.021E-01  7.195E-02  1.609E+00  8.400E-02
   13      44  1.569E+00  1.018E-01  6.850E-02  1.609E+00  8.305E-02
   14      43  1.561E+00  1.014E-01  6.161E-02  1.604E+00  8.005E-02

```


15	42	1.560E+00	1.015E-01	6.124E-02	1.604E+00	7.960E-02
16	38	1.551E+00	1.013E-01	6.017E-02	1.593E+00	6.745E-02
17	37	1.556E+00	1.015E-01	6.525E-02	1.558E+00	8.787E-02
18	33	1.554E+00	1.014E-01	6.493E-02	1.558E+00	8.577E-02
19	30	1.553E+00	1.014E-01	6.458E-02	1.558E+00	8.543E-02
20	29	1.520E+00	9.861E-02	7.153E-02	1.558E+00	8.354E-02
21	28	1.508E+00	9.803E-02	7.192E-02	1.519E+00	6.852E-02
22	26	1.504E+00	9.723E-02	7.143E-02	1.516E+00	6.417E-02
23	25	1.501E+00	9.733E-02	7.324E-02	1.516E+00	6.428E-02
24	23	1.483E+00	9.532E-02	7.663E-02	1.465E+00	6.276E-02
25	21	1.478E+00	9.519E-02	7.655E-02	1.454E+00	6.210E-02
26	20	1.472E+00	9.488E-02	7.686E-02	1.452E+00	6.764E-02
27	19	1.437E+00	9.287E-02	8.809E-02	1.452E+00	8.111E-02
28	18	1.389E+00	8.985E-02	8.560E-02	1.376E+00	7.902E-02
29	12	1.382E+00	8.916E-02	8.459E-02	1.376E+00	7.450E-02
30	11	1.378E+00	9.014E-02	8.603E-02	1.397E+00	8.716E-02
31	9	1.329E+00	8.411E-02	6.501E-02	1.422E+00	7.183E-02
32	7	1.265E+00	7.544E-02	4.533E-02	1.309E+00	6.645E-02
33*	6	1.238E+00	7.306E-02	3.864E-02	1.257E+00	6.255E-02
34**	5	1.243E+00	7.171E-02	3.991E-02	1.263E+00	5.840E-02
35++	4	1.281E+00	7.154E-02	4.064E-02	1.270E+00	4.950E-02
36	3	1.300E+00	7.015E-02	3.646E-02	1.319E+00	3.915E-02
37	2	1.325E+00	6.506E-02	2.839E-02	1.313E+00	3.556E-02
38	1	1.503E+00	5.698E-02	2.544E-02	1.484E+00	3.699E-02

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

* tree same as + tree

** tree same as -- tree

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Rel. risk is mean risk relative to sample average ignoring covariates

Cases fit give the number of cases used to fit node

Deviance is mean residual deviance for all cases in node

Node label	Total cases	Cases Matrix fit rank	Node rel.risk	Node deviance	Split variable	Other variables
1	500	500 1	1.000E+00	1.505E+00	age	
2	244	244 1	3.726E-01	9.913E-01	chf	

4T	49	49	1	1.110E+00	1.413E+00	miord
5T	195	195	1	2.124E-01	7.383E-01	year
3	256	256	1	1.890E+00	1.526E+00	chf
6T	106	106	1	3.028E+00	1.372E+00	sho
7	150	150	1	1.365E+00	1.469E+00	age
14T	120	120	1	1.063E+00	1.360E+00	los
15T	30	30	1	3.322E+00	1.278E+00	year

Number of terminal nodes of final tree: 5

Total number of nodes of final tree: 9

Regression tree:

Node 1: age <= 71.00000

Node 2: chf = "1"

Node 4: Risk relative to sample average ignoring covariates = 1.10956

Node 2: chf /= "1"

Node 5: Risk relative to sample average ignoring covariates = 0.21235

Node 1: age > 71.00000 or NA

Node 3: chf = "1"

Node 6: Risk relative to sample average ignoring covariates = 3.02760

Node 3: chf /= "1"

Node 7: age <= 85.00000 or NA

Node 14: Risk relative to sample average ignoring covariates = 1.06334

Node 7: age > 85.00000 and not NA

Node 15: Risk relative to sample average ignoring covariates = 3.32215

Node 1: Intermediate node

A case goes into Node 2 if age <= 7.1000000E+01

age mean = 6.9846E+01

Coefficients of log-relative risk function:

Regressor	Coefficient	t-stat	p-val
Constant	-3.5381E-02	-0.52	0.6041

Predicted relative risk = 1.0000000000000000

Node 2: Intermediate node

A case goes into Node 4 if chf = "1"

chf mode = "0"

Node 4: Terminal node

Coefficients of log-relative risk function:

Regressor	Coefficient	t-stat	p-val
Constant	6.8580E-02	0.34	0.7332

```

Predicted relative risk = 1.1095574995429367
-----
Node 5: Terminal node
Coefficients of log-relative risk function:
Regressor Coefficient      t-stat  p-val
Constant -1.5849E+00      -7.43  0.0000
Predicted relative risk = 0.21235168812700622
-----
Node 3: Intermediate node
A case goes into Node 6 if chf = "1"
chf mode = "0"
-----
Node 6: Terminal node
Coefficients of log-relative risk function:
Regressor Coefficient      t-stat  p-val
Constant 1.0724E+00       9.89  0.0000
Predicted relative risk = 3.0276015801608627
-----
Node 7: Intermediate node
A case goes into Node 14 if age <= 8.5000000E+01 or NA
age mean = 8.0667E+01
-----
Node 14: Terminal node
Coefficients of log-relative risk function:
Regressor Coefficient      t-stat  p-val
Constant 2.6029E-02       0.19  0.8459
Predicted relative risk = 1.0633351387096228
-----
Node 15: Terminal node
Coefficients of log-relative risk function:
Regressor Coefficient      t-stat  p-val
Constant 1.1652E+00       6.05  0.0000
Predicted relative risk = 3.3221527399879980
-----

Observed and fitted values are stored in whas500.fit
LaTeX code for tree is in whas500.tex
Elapsed time in seconds: 7.85750771

```

The tree model, given in Figure 13, shows that risk of death is lowest (0.21 relative to the sample average for the whole data set) for those younger than 72 with no congestive heart complications. The groups with the highest risks (3.0–3.32 relative to average) are those older than 85 and those between 72 and 85 with congestive heart complications.

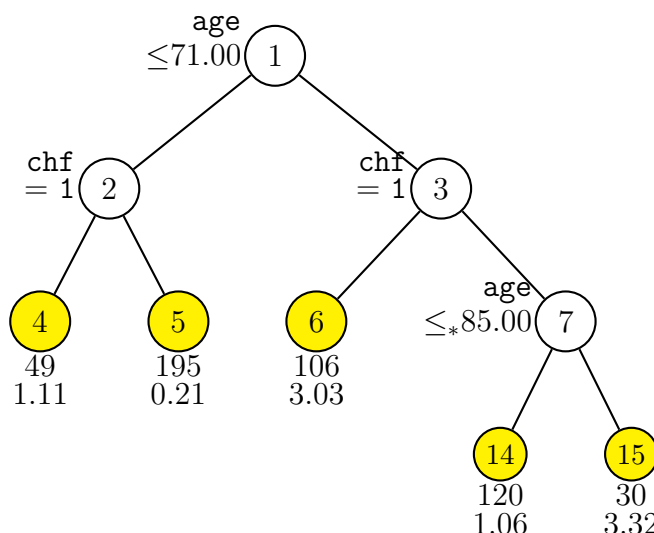


Figure 13: GUIDE 0.50-SE piecewise constant relative risk regression tree for predicting `fstat`. At each split, an observation goes to the left branch if and only if the condition is satisfied. The symbol ‘ \leq_* ’ stands for ‘ \leq or missing’. Sample sizes and mean relative risks (relative to sample average ignoring covariates) are printed below nodes.

The top 8 lines of the file `whas500.fit` and its column definitions are:

train	node	survivaltime	logbasecumhaz	relativerisk	survivalprob	mediansurvtime
y	14	2.178000E+03	-7.667985E-02	1.063335E+00	3.865048E-01	1.553841E+03
y	5	2.172000E+03	-7.667985E-02	2.123517E-01	8.270912E-01	2.354932E+03
y	5	2.190000E+03	-7.667985E-02	2.123517E-01	8.270912E-01	2.354932E+03
y	4	2.970000E+02	-1.320296E+00	1.109557E+00	7.512523E-01	1.534972E+03
y	5	2.131000E+03	-2.213734E-01	2.123517E-01	8.485159E-01	2.354932E+03
y	5	1.000000E+00	-4.352824E+00	2.123517E-01	9.973654E-01	2.354932E+03
y	5	2.122000E+03	-2.213734E-01	2.123517E-01	8.485159E-01	2.354932E+03
y	5	1.496000E+03	-4.919833E-01	2.123517E-01	8.822135E-01	2.354932E+03

train: “y” if the observation is used for model fitting, “n” if not.

node: terminal node label of observation.

survivaltime: observed survival time t .

logbasecumhaz: log of the estimated baseline cumulative hazard function $\log \Lambda_0(t) = \log \int_0^t \lambda_0(u) du$ at observed time t .

relativerisk: $\exp(\beta' \mathbf{x} - \beta_*)$, risk of death relative to the average for the sample, where \mathbf{x} is the covariate vector of the observation, β is the estimated regression coefficient vector in the node, and β_* is the coefficient in the constant model $\lambda_0(t) \exp(\beta_*)$ fitted to all the training cases in the root node. Because a constant is fitted to each node here, $\beta_* = -0.035381$ is the value of β at the root node. For example, the first subject, which is in node 14, has $\beta = 0.026029$ and so $\text{relativerisk} = \exp(\beta - \beta_*) = \exp(0.026029 + 0.035381) = 1.063335$.

survivalprob: probability that the subject survives up to observed time t . For the first subject, this is

$$\begin{aligned} \exp\{-\Lambda_0(t) \exp(\beta' \mathbf{x})\} &= \exp\{-\exp(\beta_* + \text{logbasecumhaz}) \times \text{relativerisk}\} \\ &= \exp(-\exp(-0.035381 - 0.07667985) \times 1.063335) \\ &= 0.3865048. \end{aligned}$$

mediansurvtime: estimated median survival time t such that $\exp\{-\Lambda_0(t) \exp(\beta' \mathbf{x})\} = 0.5$, or, equivalently, $\Lambda_0(t) \exp(\beta' \mathbf{x}) = -\log(0.5)$, or $\text{logbasecumhaz}(t) = \log \log(2) - \beta' \mathbf{x}$, using linear interpolation of $\Lambda_0(t)$. Median survival times greater than the largest observed time have a trailing plus (+) sign. Figure 14 shows plots of $\log \Lambda_0(t)$ and $\Lambda_0(t)$ for this data set.

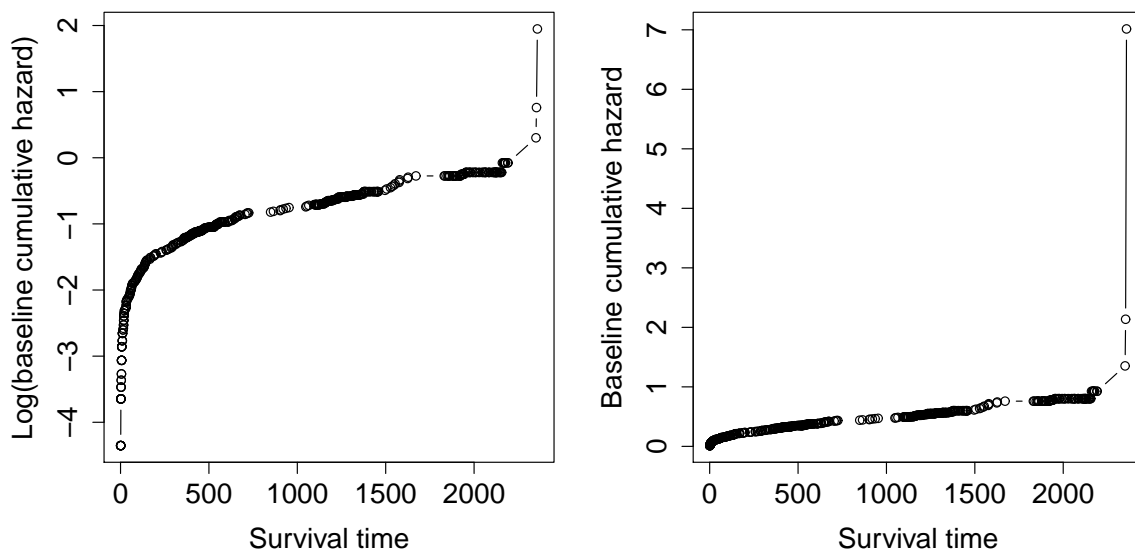


Figure 14: Plots of estimated baseline cumulative hazard function

5.8 Multi-response data

GUIDE can fit a piecewise-constant regression model for two or more dependent variables simultaneously. Following is an example from [Loh and Zheng \(2013\)](#) on estimating the strength and viscosity of concrete. The comma-delimited data file `concrete.csv` is from [Yeh \(2007\)](#). The data description file is below. Notice that there are three D variables. Our goal is to construct a single regression tree that predicts all three D variables simultaneously.

```
concrete.csv
NA
c1 c2 c3
1 No x
2 Cement n
3 Slag n
4 FlyAsh n
5 Water n
6 SP n
7 CoarseAggr n
8 FineAggr n
9 Slump d
10 Flow d
11 Strength d
```

Following is an annotated log of the input file creation.

```
0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: concrete.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: concrete.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 5
  Option 5 is for multiresponse data.
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):
Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): concretedsc.txt
Reading data description file ...
Training sample file: concrete.csv
```

```

Missing value code: NA
Warning: N variables changed to S
Number of D variables = 3
D variables are:
Slump
Flow
Strength
Choose multivariate or univariate split variable selection:
Choose multivariate if fewer than 5 D variables, choose univariate otherwise
Input 1 for multivariate, 2 for univariate ([1:2], <cr>=1):
The D vector can be grouped into segments to look for patterns
Input 0 for no grouping, 1 for roughly equal groups, 2 for other choices
Input your selection ([0:2], <cr>=0):
Grouping is recommended if there are numerous D variables and they are clustered.
Here, Slump and Flow may be considered as belonging to one cluster and Strength
to another, but in this illustration, we treat them as separate.
Input 1 to normalize D variables, 2 for no normalization ([1:2], <cr>=1):
Normalization means scaling each variable to have variance 1.
Input 1 for equal, 2 for unequal weighting of D variables ([1:2], <cr>=1):
Length of longest data entry = 6
Total number of cases = 103
Checking data ...
#cases w/ miss. D = number of cases with all D values missing
      Total #cases w/ #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
      103      0      0      1      0      0      7      0      0
No. cases used for training = 103
Warning: interaction tests skipped
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): concrete.tex
Input 2 to save node IDs of cases, 1 otherwise ([1:2], <cr>=1): 1
Input 2 to save node fitted values; 1 otherwise ([1:2], <cr>=1): 1
Input file is created!

```

Results

```

Multi-response or longitudinal data without T variables
Pruning by cross-validation
Data description file: concretedsc.txt
Training sample file: concrete.csv
Missing value code: NA
Warning: N variables changed to S
Number of D variables = 3
Multivariate split variable selection method
Missing D values treated as below average for variable selection

```

No grouping of D variables

Segment boundaries are:

1.500 2.500

GUIDE labels the D variables as Y_1, Y_2, \dots (in the order of their appearance in the data file). The segment boundaries refer to the grouping of the indices. In this case, since there is no grouping, each variable is its own group. As a result, the indices are grouped into intervals (0.5, 1.5), (1.5, 2.5), (2.5, 3.5).

Mean-squared errors (MSE) are calculated from normalized D variables

This is a reminder that the D variables are normalized.

D variables equally weighted

Piecewise constant model

Length of longest data entry = 6

Summary information (without x variables)

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,

n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight

For categorical variables, #categories include one for missing values

Column number	Variable name	Variable type	Minimum value	Maximum value	Number of categories	Number missing
2	Cement	s	1.3700E+02	3.7400E+02		
3	Slag	s	0.0000E+00	1.9300E+02		
4	FlyAsh	s	0.0000E+00	2.6000E+02		
5	Water	s	1.6000E+02	2.4000E+02		
6	SP	s	4.4000E+00	1.9000E+01		
7	CoarseAggr	s	7.0800E+02	1.0499E+03		
8	FineAggr	s	6.4060E+02	9.0200E+02		
9	Slump	d	0.0000E+00	2.9000E+01		
10	Flow	d	2.0000E+01	7.8000E+01		
11	Strength	d	1.7190E+01	5.8530E+01		

#cases w/ miss. D = number of cases with all D values missing

Total #cases w/ #missing

#cases	miss. D	ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
103	0	0	1	0	0	7	0	0

No. cases used for training = 103

Warning: interaction tests skipped

No interaction tests

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Split values for N and S variables based on exhaustive search

Max number of split levels = 10

Minimum node size = 10

Number of SE's for pruned tree = 5.0000E-01

Size and CV Loss and SE of subtrees:

Tree	#Tnodes	Mean Loss	SE(Mean)	BSE(Mean)	Median Loss	BSE(Median)
0*	7	7.112E-01	9.308E-02	8.141E-02	6.231E-01	1.582E-01
1**	6	7.213E-01	8.920E-02	6.935E-02	6.693E-01	1.214E-01
2	3	7.693E-01	9.331E-02	9.137E-02	7.043E-01	1.236E-01
3	2	8.717E-01	9.806E-02	8.897E-02	7.856E-01	1.209E-01
4	1	1.011E+00	9.866E-02	7.561E-02	9.799E-01	1.073E-01

0-SE tree based on mean is marked with *

0-SE tree based on median is marked with +

Selected-SE tree based on mean using naive SE is marked with **

Selected-SE tree based on mean using bootstrap SE is marked with --

Selected-SE tree based on median and bootstrap SE is marked with ++

* tree same as + tree

** tree same as ++ tree

** tree same as -- tree

++ tree same as -- tree

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Cases fit give the number of cases used to fit node

MSE is residual sum of squares divided by number of cases in node

Node	Total	Cases	Node	Split
label	cases	fit	MSE	variable
1	103	103	1.000E+00	Water
2	29	29	9.453E-01	Cement
4T	14	14	4.479E-01	-
5T	13	15	1.078E+00	-
3	74	74	6.728E-01	Slag
6	64	64	5.374E-01	Cement
12T	18	18	1.874E-01	-
13	46	46	5.620E-01	FlyAsh
26T	20	20	3.240E-01	-
27T	26	26	4.488E-01	FlyAsh
7T	10	10	1.068E+00	-

Number of terminal nodes of final tree: 6

Total number of nodes of final tree: 11

Regression tree for multi-response data:

Node 1: Water <= 1.82250E+02

Node 2: Cement <= 1.61000E+02

Node 4: Mean cost = 4.15924E-01

```

Node 2: Cement > 1.61000E+02 or NA
Node 5: Mean cost = 9.95311E-01
Node 1: Water > 1.82250E+02 or NA
Node 3: Slag <= 1.35000E+02 or NA
Node 6: Cement <= 1.57100E+02
Node 12: Mean cost = 1.76946E-01
Node 6: Cement > 1.57100E+02 or NA
Node 13: FlyAsh <= 1.17500E+02
Node 26: Mean cost = 3.07784E-01
Node 13: FlyAsh > 1.17500E+02 or NA
Node 27: Mean cost = 4.31499E-01
Node 3: Slag > 1.35000E+02 and not NA
Node 7: Mean cost = 9.61024E-01

*****

Node 1: Intermediate node
A case goes into Node 2 if Water <= 1.8225000E+02
Water mean = 1.9717E+02
Estimated D values are:
1.8049E+01 4.9611E+01 3.6039E+01
-----
Node 2: Intermediate node
A case goes into Node 4 if Cement <= 1.6100000E+02
Cement mean = 1.9275E+02
-----
Node 4: Terminal node
Estimated D values are:
1.5286E+01 4.0536E+01 3.5470E+01
The three numbers are the sample mean values of Slump, Flow and Strength, in
order of their appearance in the data file.
-----
Node 5: Terminal node
Estimated D values are:
5.6667E+00 2.5347E+01 4.4189E+01
-----
Node 3: Intermediate node
A case goes into Node 6 if Slag <= 1.3500000E+02 or NA
Slag mean = 7.8059E+01
-----
Node 6: Intermediate node
A case goes into Node 12 if Cement <= 1.5710000E+02
Cement mean = 2.5277E+02
-----
Node 12: Terminal node
Estimated D values are:

```

```

      2.2083E+01   5.8389E+01   2.9007E+01
-----
Node 13: Intermediate node
A case goes into Node 26 if FlyAsh <=  1.1750000E+02
FlyAsh mean =  1.1328E+02
-----
Node 26: Terminal node
Estimated D values are:
      2.0950E+01   5.4325E+01   3.1665E+01
-----
Node 27: Terminal node
Estimated D values are:
      2.2712E+01   6.1450E+01   4.2456E+01
-----
Node 7: Terminal node
Estimated D values are:
      1.5300E+01   4.2700E+01   2.9337E+01
-----

Case and node IDs are in file: node.txt
Node fitted values are in file: fit.txt
Observed and fitted values are stored in node.txt
LaTeX code for tree is in tree.tex
Elapsed time in seconds:   3.15099992E-02

```

The \LaTeX tree is shown in Figure 15.

5.9 Longitudinal data with irregular time points

The data come from a longitudinal study on the hourly wage of 888 male high-school dropouts (246 black, 204 Hispanic, 438 white), where the observation time points as well as their number (1–13) varied across individuals (Murnane et al., 1999; Singer and Willett, 2003). An earlier version of GUIDE was used to analyze these data in Loh and Zheng (2013).

The response variable is hourly wage (in 1990 dollars) and the predictor variables are `hgc` (highest grade completed; 6–12), `exper` (years in labor force; 0.001–12.7 yrs), and `race` (Black, Hispanic, and White). The data file `widewage.txt` is in *wide format*, where each record refers to one individual. The description file `wage.dsc` is given below. Note that observation time points are marked as `t`.

```

widewage.txt
NA
c1 c2 c3
1 id x

```

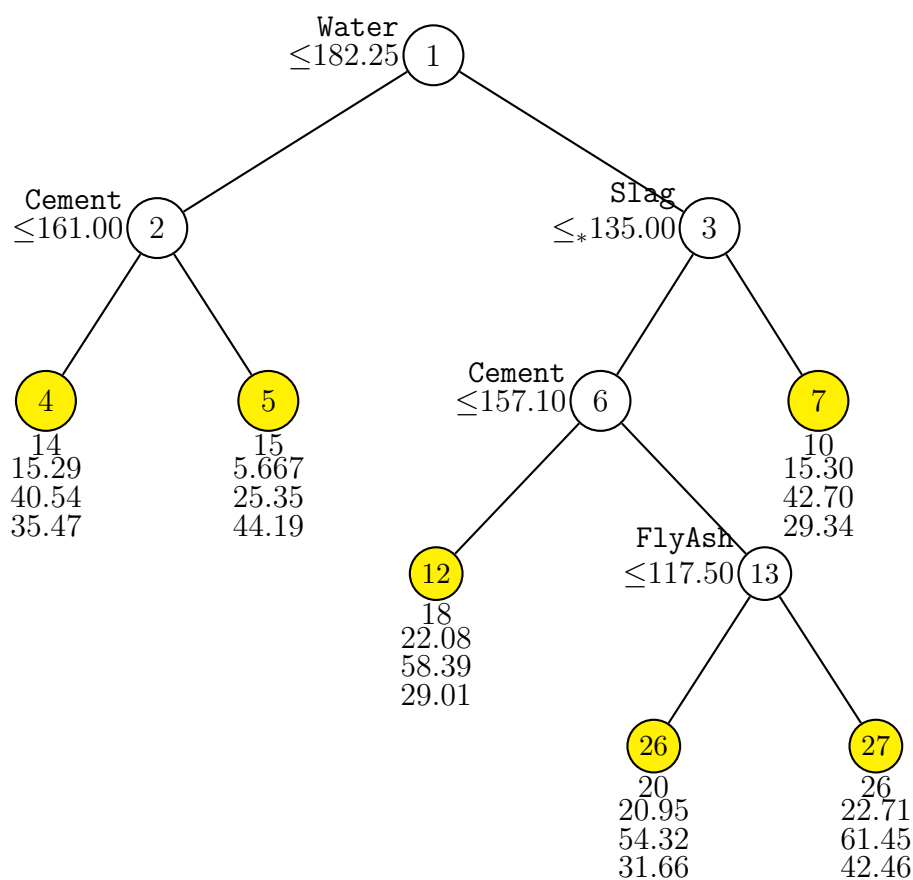


Figure 15: GUIDE 0.50-SE regression tree for predicting response variables **Slump**, **Flow**, and **Strength**. At each split, an observation goes to the left branch if and only if the condition is satisfied. The symbol ‘ \leq^* ’ stands for ‘ \leq or missing’. Sample sizes and predicted values of **Slump**, **Flow**, and **Strength** are printed below nodes.

```
2 hgc n
3 exper1 t
4 exper2 t
5 exper3 t
6 exper4 t
7 exper5 t
8 exper6 t
9 exper7 t
10 exper8 t
11 exper9 t
12 exper10 t
13 exper11 t
14 exper12 t
15 exper13 t
16 postexp1 x
17 postexp2 x
18 postexp3 x
19 postexp4 x
20 postexp5 x
21 postexp6 x
22 postexp7 x
23 postexp8 x
24 postexp9 x
25 postexp10 x
26 postexp11 x
27 postexp12 x
28 postexp13 x
29 wage1 d
30 wage2 d
31 wage3 d
32 wage4 d
33 wage5 d
34 wage6 d
35 wage7 d
36 wage8 d
37 wage9 d
38 wage10 d
39 wage11 d
40 wage12 d
41 wage13 d
42 ged1 x
43 ged2 x
44 ged3 x
45 ged4 x
46 ged5 x
47 ged6 x
```

```

48 ged7 x
49 ged8 x
50 ged9 x
51 ged10 x
52 ged11 x
53 ged12 x
54 ged13 x
55 uerate1 x
56 uerate2 x
57 uerate3 x
58 uerate4 x
59 uerate5 x
60 uerate6 x
61 uerate7 x
62 uerate8 x
63 uerate9 x
64 uerate10 x
65 uerate11 x
66 uerate12 x
67 uerate13 x
68 race c

```

Because the default 0.5-SE rule yields a trivial tree with no splits, we show how the options can be changed to produce a tree with the 0-SE rule. Following is a session log.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: wage.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: wage.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 6
Input 1 for lowess smoothing, 2 for spline smoothing ([1:2], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1): 2
  Choosing 1 will produce a 0.5-SE tree. We choose 2 to allow more options.
Input 1 for interaction tests, 2 to skip them ([1:2], <cr>=1):
Input 1 to prune by CV, 2 for no pruning ([1:2], <cr>=1):

```

```

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): wage.dsc
Reading data description file ...
Training sample file: widewage.txt
Missing value code: NA
Warning: N variables changed to S
Number of D variables = 13
Number of D variables =          13
D variables are:
wage1
wage2
wage3
wage4
wage5
wage6
wage7
wage8
wage9
wage10
wage11
wage12
wage13
T variables are:
exper1
exper2
exper3
exper4
exper5
exper6
exper7
exper8
exper9
exper10
exper11
exper12
exper13
The D variables can be grouped into segments to look for patterns
Input 1 for roughly equal-sized groups, 2 for customized groups ([1:2], <cr>=1):
Input number of roughly equal-sized groups ([2:9], <cr>=3):
Input number of interpolating points for prediction ([10:100], <cr>=31):
Length of longest data entry = 16
Total number of cases =          888
Col. no. Categorical variable    #levels    #missing values
      68 race                    3              0
Checking data ...

```

```

#cases w/ miss. D = number of cases with all D values missing
      Total #cases w/ #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
      888      0      0      40      0      0      1      0      1
No. cases used for training =      888
No. cases excluded due to 0 weight or missing D =      0
Warning: interaction tests skipped
Default number of cross-validations =      10
Input 1 to accept the default, 2 to change it ([1:2], <cr>=1):
Best tree may be chosen based on mean or median CV estimate
Input 1 for mean-based, 2 for median-based ([1:2], <cr>=1):
Input number of SEs for pruning ([0.00:1000.00], <cr>=0.50): 0
This is where we choose the 0-SE pruning rule.
Choose a split point selection method for numerical variables:
Choose 1 to use faster method based on sample quantiles
Choose 2 to use exhaustive search
Input 1 or 2 ([1:2], <cr>=2):
Default max number of split levels =      10
Input 1 to accept this value, 2 to change it ([1:2], <cr>=1):
Default minimum node sample size is 44
Input 1 to use the default value, 2 to change it ([1:2], <cr>=1):
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): wage.tex
Input 1 for a vertical tree, 2 for a sideways tree ([1:2], <cr>=1):
Input 1 to include node numbers, 2 to omit them ([1:2], <cr>=1):
Input 1 to number all nodes, 2 to number leaves only ([1:2], <cr>=1):
Choose a color for the terminal nodes:
(1) white
(2) lightgray
(3) gray
(4) darkgray
(5) black
(6) yellow
(7) red
(8) blue
(9) green
(10) magenta
(11) cyan
Input your choice ([1:11], <cr>=6):
You can store the variables and/or values used to split and fit in a file
Choose 1 to skip this step, 2 to store split and fit variables,
3 to store split variables and their values
Input your choice ([1:3], <cr>=1):
Input 2 to save node IDs of cases, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store terminal node IDs: wage.nid
Input 2 to save node fitted values; 1 otherwise ([1:2], <cr>=1): 2

```



```
Input name of file to store node fitted values: wage.fit
Input 2 to save terminal node IDs for importance scoring; 1 otherwise ([1:2], <cr>=1):
Input 2 to write R function for predicting new cases, 1 otherwise ([1:2], <cr>=1):
Input file is created!
```

Results

```
Lowess smoothing
Longitudinal data with T variables
Pruning by cross-validation
Data description file: wagedsc.txt
Training sample file: wagedat.txt
Missing value code: NA
Warning: N variables changed to S
Number of D variables = 13
Number of D variables = 13
D variables are:
wage1
wage2
wage3
wage4
wage5
wage6
wage7
wage8
wage9
wage10
wage11
wage12
wage13
T variables are:
exper1
exper2
exper3
exper4
exper5
exper6
exper7
exper8
exper9
exper10
exper11
exper12
exper13
Length of longest data entry = 16
```

Summary information (without x variables)

d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical, n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight

Column	Name		Minimum	Maximum	#Categories	#Missing
2	hgc	s	6.0000E+00	1.2000E+01		
3	exper1	t	1.0000E-03	5.6370E+00		
4	exper2	t	0.0000E+00	7.5840E+00		38
5	exper3	t	0.0000E+00	9.7770E+00		77
6	exper4	t	0.0000E+00	1.0815E+01		124
7	exper5	t	0.0000E+00	1.1777E+01		159
8	exper6	t	0.0000E+00	1.0587E+01		233
9	exper7	t	0.0000E+00	1.1279E+01		325
10	exper8	t	0.0000E+00	1.0582E+01		428
11	exper9	t	0.0000E+00	1.1621E+01		551
12	exper10	t	0.0000E+00	1.2260E+01		678
13	exper11	t	0.0000E+00	1.1980E+01		791
14	exper12	t	0.0000E+00	1.2558E+01		856
15	exper13	t	0.0000E+00	1.2700E+01		882
29	wage1	d	2.0299E+00	6.8649E+01		
30	wage2	d	2.0689E+00	5.0400E+01		38
31	wage3	d	2.0462E+00	3.4501E+01		77
32	wage4	d	2.1170E+00	3.3149E+01		124
33	wage5	d	2.1043E+00	4.9304E+01		159
34	wage6	d	2.2078E+00	7.3995E+01		233
35	wage7	d	2.1043E+00	4.7276E+01		325
36	wage8	d	2.3164E+00	3.7713E+01		428
37	wage9	d	2.5294E+00	4.6109E+01		551
38	wage10	d	2.9982E+00	5.6543E+01		678
39	wage11	d	4.0837E+00	2.2198E+01		791
40	wage12	d	3.4315E+00	4.6201E+01		856
41	wage13	d	4.5631E+00	7.7757E+00		882
68	race	c			3	

Total	#cases w/	#missing							
#cases	miss. D	ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var	
888	0	0	40	0	0	1	0	1	

No. cases used for training: 888

No. cases excluded due to 0 weight or missing D: 0

Warning: interaction tests skipped

No interaction tests

Pruning by v-fold cross-validation, with v = 10

Selected tree is based on mean of CV estimates

Split values for N and S variables based on exhaustive search

Max number of split levels = 10

Minimum node size = 44
 Number of SE's for pruned tree = 0.0000E+00

Size and CV Loss and SE of subtrees:

Tree	#Tnodes	Mean Loss	SE(Mean)	BSE(Mean)	Median Loss	BSE(Median)
0	10	1.257E+02	1.044E+01	8.502E+00	1.204E+02	1.525E+01
1	9	1.257E+02	1.044E+01	8.502E+00	1.204E+02	1.525E+01
2	8	1.257E+02	1.044E+01	8.502E+00	1.204E+02	1.525E+01
3	6	1.242E+02	1.049E+01	8.463E+00	1.181E+02	1.537E+01
4**	5	1.238E+02	1.058E+01	8.434E+00	1.175E+02	1.530E+01
5++	1	1.244E+02	1.064E+01	8.700E+00	1.157E+02	1.577E+01

0-SE tree based on mean is marked with *
 0-SE tree based on median is marked with +
 Selected-SE tree based on mean using naive SE is marked with **
 Selected-SE tree based on mean using bootstrap SE is marked with --
 Selected-SE tree based on median and bootstrap SE is marked with ++
 ** tree same as -- tree
 + tree same as ++ tree
 * tree same as ** tree
 * tree same as -- tree

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Cases fit give the number of cases used to fit node
 MSE is residual sum of squares divided by number of cases in node

Node label	Total cases	Cases fit	Node MSE	Split variable
1	888	888	1.222E+02	race
2T	246	246	1.111E+02	hgc
3	642	642	1.259E+02	race
6	204	204	1.278E+02	hgc
12T	127	127	1.085E+02	-
13T	77	77	1.514E+02	-
7	438	438	1.252E+02	hgc
14T	299	299	9.813E+01	hgc
15T	139	139	1.777E+02	hgc

Number of terminal nodes of final tree: 5
 Total number of nodes of final tree: 9

Regression tree for longitudinal data:

```

Node 1: race = "black"
Node 2: Mean cost = 1.10602E+02
Node 1: race /= "black"
Node 3: race = "hispanic"
Node 6: hgc <= 9.50000 or NA
Node 12: Mean cost = 1.07621E+02
Node 6: hgc > 9.50000 and not NA
Node 13: Mean cost = 1.49412E+02
Node 3: race /= "hispanic"
Node 7: hgc <= 9.50000 or NA
Node 14: Mean cost = 9.78002E+01
Node 7: hgc > 9.50000 and not NA
Node 15: Mean cost = 1.76394E+02

```

```

Node 1: Intermediate node
A case goes into Node 2 if race = "black"
race mode = "white"
-----
Node 2: Terminal node
-----
Node 3: Intermediate node
A case goes into Node 6 if race = "hispanic"
race mode = "white"
-----
Node 6: Intermediate node
A case goes into Node 12 if hgc <= 9.5000000E+00 or NA
hgc mean = 8.9118E+00
-----
Node 12: Terminal node
-----
Node 13: Terminal node
-----
Node 7: Intermediate node
A case goes into Node 14 if hgc <= 9.5000000E+00 or NA
hgc mean = 8.8973E+00
-----
Node 14: Terminal node
-----
Node 15: Terminal node
-----

```

```

Case and node IDs are in file: wage.nid
Node fitted values are in file: wage.fit

```

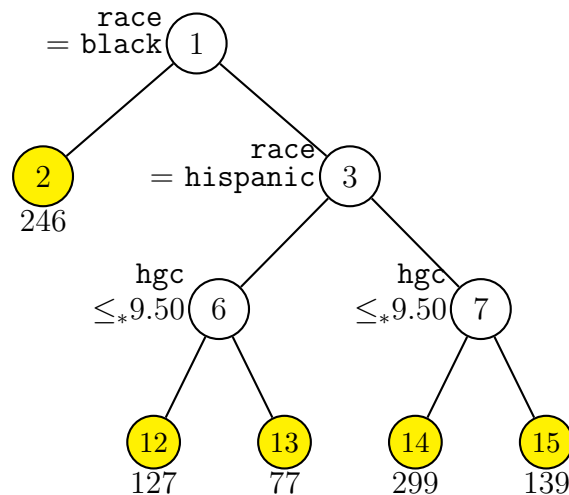


Figure 16: GUIDE 0.00-SE regression tree for predicting longitudinal variables `wage1`, `wage2`, etc. At each split, an observation goes to the left branch if and only if the condition is satisfied. The symbol ' \leq_* ' stands for ' \leq or missing'. Sample sizes are printed below nodes.

```

Observed and fitted values are stored in wage.nid
LaTeX code for tree is in wage.tex
Elapsed time in seconds:    2.98005509
  
```

Figure 16 shows the tree and Figure 17 plots lowess-smoothed curves of mean wage in the two terminal nodes. The plotting values are obtained from the result file `wage.fit` whose contents are given below. The first column gives the node number and the next two columns the start and end of the times at which fitted values are computed. The other columns give the fitted values equally spaced between the start and end times.

node	t.start	t.end	fitted1	fitted2	fitted3	fitted4	fitted5	fitted6	fitted7	fitted8	fitted9	fitted10
2	0.40000E-02	0.12558E+02	0.50794E+01	0.52623E+01	0.54112E+01	0.55477E+01	0.56649E+01	0.57721E+01	0.58793E+01	0.59865E+01	0.60937E+01	0.62009E+01
12	0.60000E-02	0.12535E+02	0.47994E+01	0.50688E+01	0.53388E+01	0.55076E+01	0.54340E+01	0.55076E+01	0.55764E+01	0.56452E+01	0.57140E+01	0.57828E+01
13	0.12200E+00	0.11990E+02	0.56361E+01	0.58877E+01	0.61037E+01	0.62417E+01	0.62780E+01	0.63143E+01	0.63506E+01	0.63869E+01	0.64232E+01	0.64595E+01
14	0.10000E-02	0.12700E+02	0.50837E+01	0.52324E+01	0.53638E+01	0.54764E+01	0.55631E+01	0.56498E+01	0.57365E+01	0.58232E+01	0.59099E+01	0.59966E+01
15	0.20000E-02	0.12045E+02	0.57487E+01	0.59247E+01	0.60944E+01	0.62441E+01	0.63143E+01	0.63845E+01	0.64547E+01	0.65249E+01	0.65951E+01	0.66653E+01

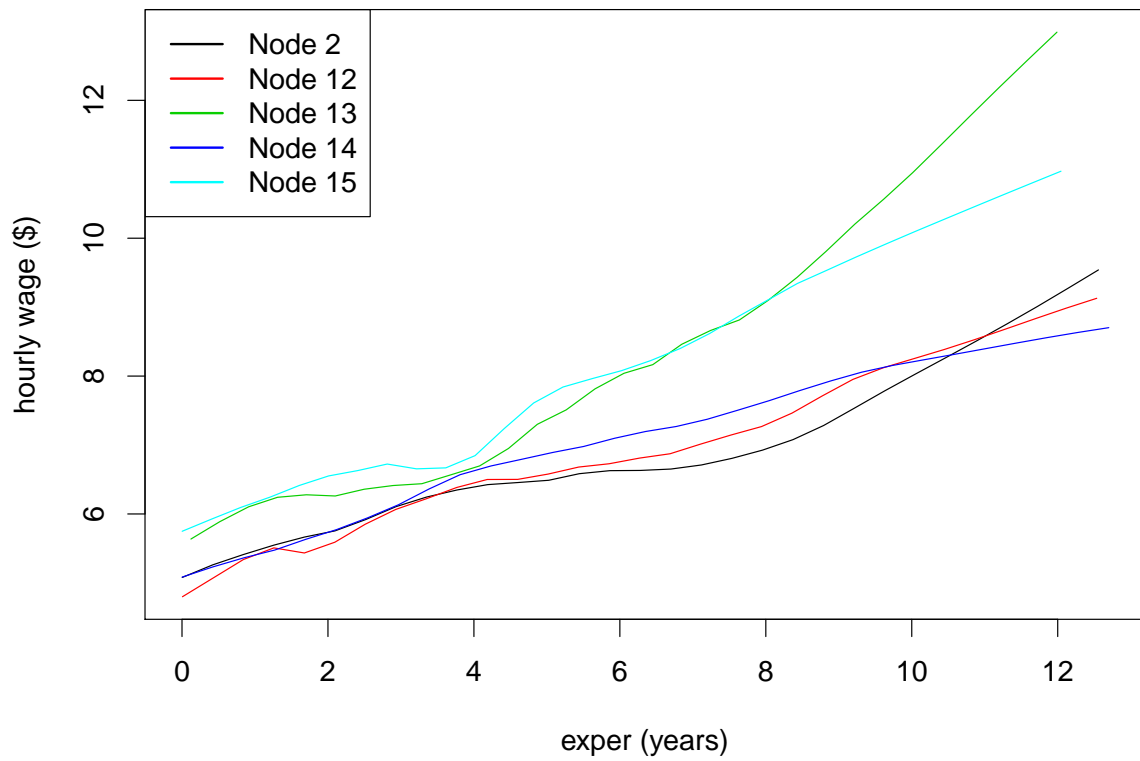


Figure 17: Lowess-smoothed mean wage curves in the terminal nodes of Figure 16.

5.10 Subgroup identification

If there is a treatment variable in the data, GUIDE can fit a tree model find subgroups with differential treatment effects. The dependent variable can be censored or not.

1. The treatment variable is designated as `R`.
2. If there is no censoring, the response variable is designated as `D` as usual.
3. If there is censoring, then the survival time is designated as `T` and the event (typically death) indicator is designated as `D`, taking value 1 for death and 0 for censored.

GUIDE has two methods for solving this problem, called `gi` and `gs`. The former is more sensitive to predictive variables (i.e., variables that interact with treatment) and the latter is equally sensitive to prognostic and predictive variables. The methods are documented in [Loh et al. \(2015\)](#).

We illustrate the `gi` method with data from a breast cancer trial ([Schmoor et al., 1996](#)). The data are in the file `cancer.txt` from the `ipred` R package ([Peters and Hothorn, 2012](#)). In the description file `cancerdsc.txt` below, the treatment variable is hormone therapy, `horTh`. The variable `time` is censored survival time and `death` is the event indicator (1=death, 0=censored).

```
cancer.txt
NA
c1 c2 c3
1 horTh r
2 age s
3 menostat c
4 tsize s
5 tgrade c
6 pnodes s
7 progrec s
8 estrec s
9 time t
10 death d
```

Input file generation

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning

```

Input your choice: 1
Name of batch input file: cancerin.txt
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: cancerout.txt
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 4
Choose complexity of model to use at each node:
1: multiple linear, 2: best simple linear, 3: constant ([1:3], <cr>=3):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):
  Whatever option chosen here will automatically be changed to multiple linear
  when the program detects the presence of a 'R' variable in the description file.

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): cancerdsc.txt
Reading data description file ...
Training sample file: cancer.txt
Missing value code: NA
  R variable present
  Warning: model fit changed to linear in treatment variable
Dependent variable is death
Length of longest data entry = 4
Total number of cases = 686
Col. no. Categorical variable    #levels    #missing values
      1 horTh                    2              0
      3 menostat                 2              0
      5 tgrade                   3              0
Checking data ...
Smallest uncensored T = 72.0000000000000000
No. cases dropped due to missing D or T or censored T < smallest uncensored T = 14
No. complete cases excluding censored T < smallest uncensored T = 672
No. cases used to compute baseline hazard = 672
No. cases with D=1 and T >= smallest uncensored = 299
The program will try to create the variables in the description file.
If it is unsuccessful, please create the columns yourself...
Number of dummy variables created: 1
Choose a subgroup identification method:
1 = Sum of chi-squares (Gs)
2 = Treatment interactions (Gi)
Input your choice: ([1:2], <cr>=2):
  Option 2 is generally more sensitive to detecting treatment interactions.
  Total  #cases w/  #missing
  #cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var  #R-var

```



```

        686          0          0          0          0          0          5          0          2          1
Survival time variable in column:          9
Censoring indicator variable in column:      10
Proportion of uncensored among nonmissing T and D variables =    0.445
No. cases used for training =          672
Warning: interaction tests skipped
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): cancer.tex
Input 2 to save fitted values and node IDs, 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node IDs and fitted values: cancer.fit
Input file is created!

```

Results The following results show that the tree splits once on variable progrec.

```

Proportional hazards regression with relative risk estimates
Pruning by cross-validation
Data description file: cancerdsc.txt
Training sample file: cancer.txt
Missing value code: NA
R variable present
Warning: model fit changed to linear in treatment variable
Dependent variable is death
Piecewise linear model
Length of longest data entry =  4
Smallest uncensored T = 72.00
No. cases dropped due to missing D or T or censored T < smallest uncensored T =  14
No. complete cases excluding censored T < smallest uncensored T =  672
No. cases used to compute baseline hazard =  672
No. cases with D=1 and T >= smallest uncensored =  299
Number of dummy variables created =  1

Summary information (without x variables)
d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,
n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
t=survival time variable
For categorical variables, #categories include one for missing values

```

Column number	Variable name	Variable type	Minimum value	Maximum value	Number of categories	Number missing
1	horTh	r			2	
2	age	s	2.1000E+01	8.0000E+01		
3	menostat	c			2	
4	tsize	s	3.0000E+00	1.2000E+02		
5	tgrade	c			3	
6	pnodes	s	1.0000E+00	5.1000E+01		
7	progrec	s	0.0000E+00	2.3800E+03		

```

      8  estrec          s  0.0000E+00  1.1440E+03
      9  time           t  7.2000E+01  2.6590E+03
     10  death          d  0.0000E+00  1.0000E+00
===== Constructed variables =====
     11  lnbasehaz0     z  -6.5103E+00  5.8866E-02
     12  horTh.yes      f  0.0000E+00  1.0000E+00

Total  #cases w/  #missing
#cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
    686      0      0      0      0      0      0      5      0      2

Survival time variable in column 9
Censoring indicator variable in column 10
Proportion of uncensored among nonmissing T and D variables = 0.445
No. cases used for training = 672

Warning: interaction tests skipped
Treatment interactions (Gi)
No interaction tests
Pruning by v-fold cross-validation, with v = 10
Selected tree is based on mean of CV estimates
Fraction of cases used for splitting each node = 1.0000
Max number of split levels = 10
Minimum node size = 34
Number of iterations = 5
No calibration needed: no N variables
Number of SE's for pruned tree = 5.0000E-01

Size and CV Loss and SE of subtrees:
Tree  #Tnodes  Mean Loss  SE(Mean)  BSE(Mean)  Median Loss  BSE(Median)
  1      15  1.453E+00  5.891E-02  4.830E-02  1.408E+00  5.079E-02
  2      14  1.453E+00  5.885E-02  4.827E-02  1.408E+00  5.061E-02
  3      13  1.441E+00  5.676E-02  4.165E-02  1.395E+00  5.337E-02
  4      11  1.445E+00  5.640E-02  4.195E-02  1.405E+00  5.353E-02
  5       7  1.444E+00  5.632E-02  4.201E-02  1.405E+00  5.235E-02
  6       5  1.445E+00  5.627E-02  4.149E-02  1.405E+00  5.190E-02
 7**      2  1.408E+00  5.211E-02  3.686E-02  1.386E+00  3.980E-02
 8       1  1.442E+00  5.157E-02  1.216E-02  1.450E+00  1.474E-02

0-SE tree based on mean is marked with *
0-SE tree based on median is marked with +
Selected-SE tree based on mean using naive SE is marked with **
Selected-SE tree based on mean using bootstrap SE is marked with --
Selected-SE tree based on median and bootstrap SE is marked with ++
* tree, ** tree, + tree, and ++ tree all the same

```

Following tree is based on mean CV with naive SE estimate (**).

Structure of final tree. Each terminal node is marked with a T.

Rel. risk is mean risk relative to sample average ignoring covariates

Cases fit give the number of cases used to fit node

Deviance is mean residual deviance for all cases in node

Node	Total	Cases	Matrix	Node	Node	Split
label	cases	fit	rank	rel.risk	deviance	variable
1	672	672	1	1.000E+00	1.414E+00	progrec
2T	274	274	1	1.588E+00	1.584E+00	menostat
3T	398	398	1	7.095E-01	1.172E+00	menostat

Number of terminal nodes of final tree: 2

Total number of nodes of final tree: 3

Regression tree:

Node 1: progrec <= 21.00000

Node 2: Risk relative to sample average ignoring covariates = 1.58824

Node 1: progrec > 21.00000 or NA

Node 3: Risk relative to sample average ignoring covariates = 0.70947

Node 1: Intermediate node

A case goes into Node 2 if progrec <= 2.1000000E+01

progrec mean = 1.1092E+02

Coefficients of log-relative risk function:

Regressor	Coefficient	t-stat	p-val	Min	Mean	Max
Constant	1.2903E-01	1.85	0.0651			
horTh.yes	-3.6984E-01	-2.97	0.0031	0.0000E+00	3.6012E-01	1.0000E+00

Predicted relative risk = 1.0000000000000000

Node 2: Terminal node

Coefficients of log-relative risk function:

Regressor	Coefficient	t-stat	p-val	Min	Mean	Max
Constant	5.0439E-01	5.04	0.0000			
horTh.yes	-1.1775E-01	-0.71	0.4786	0.0000E+00	3.6131E-01	1.0000E+00

Predicted relative risk = 1.5882374130670138

Node 3: Terminal node

Coefficients of log-relative risk function:

Regressor	Coefficient	t-stat	p-val	Min	Mean	Max
Constant	-1.3184E-01	-1.35	0.1775			

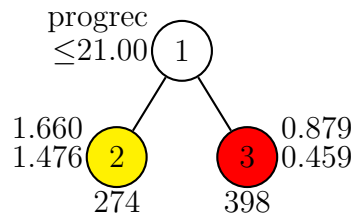


Figure 18: GUIDE Gi proportional hazards regression tree for differential treatment effects. At each intermediate node, an observation goes to the left branch if and only if the condition is satisfied. Numbers beside terminal nodes are estimated relative risks (relative to average for sample ignoring covariates) corresponding to treatment levels **no** and **yes**; numbers below are sample sizes.

```
horTh.yes      -6.5011E-01      -3.40 0.0007      0.0000E+00      3.5930E-01      1.0000E+00
Predicted relative risk = 0.70947399586773086
```

```
-----
Observed and fitted values are stored in cancer.fit
LaTeX code for tree is in cancer.tex
Elapsed time in seconds:      12.6887312
```

The \LaTeX tree diagram and the Kaplan-Meier survival functions estimated from the data in the terminal nodes of the tree are shown in Figures 18 and 19, respectively.

Estimated relative risks and survival probabilities The file `cancer.fit` gives the terminal node number, estimated survival time, log baseline cumulative hazard, relative risk (relative to the average for the data, ignoring covariates), survival probability, and median survival time of each observation in the training sample file `cancer.txt`. The results for the first few observations are shown below. See Section 5.7 for definitions of the terms.

train	node	survivaltime	logbasecumhaz	relativerisk	survivalprob	mediansurvtime
y	3	1.814000E+03	-3.317667E-01	8.787636E-01	5.331186E-01	2.014420E+03
y	3	2.018000E+03	-2.024282E-01	4.587030E-01	6.882035E-01	2.659000E+03+
y	3	7.120000E+02	-1.300331E+00	4.587030E-01	8.828100E-01	2.659000E+03+
y	3	1.807000E+03	-3.550694E-01	4.587030E-01	7.255880E-01	2.659000E+03+
y	3	7.720000E+02	-1.176558E+00	8.787636E-01	7.631865E-01	2.014420E+03
y	2	4.480000E+02	-2.105688E+00	1.660293E+00	8.173929E-01	1.038277E+03

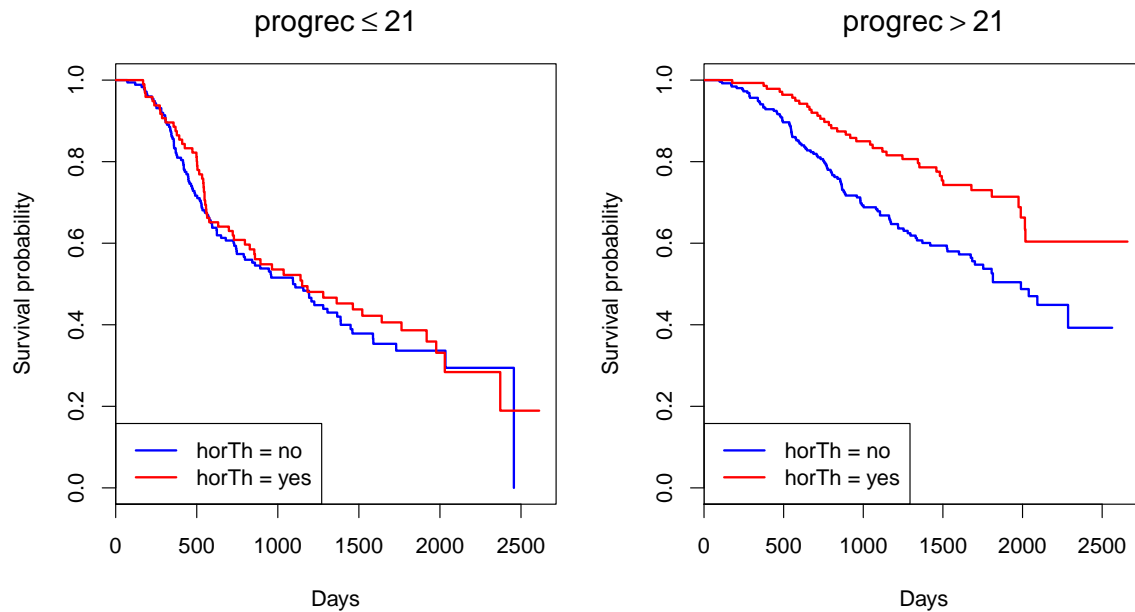


Figure 19: Estimated survival probability functions for breast cancer data

5.11 Differential item functioning

GUIDE has an experimental option to identify important predictor variables and items with differential item functioning (DIF) in a data set with two or more item (dependent variable) scores. We illustrate it with a data set from [Broekman et al. \(2011, 2008\)](#) and [Marc et al. \(2008\)](#). It consists of responses from 1978 subjects on 15 items. There are 3 predictor variables (age, education, and gender). The data and description files are `GDS.dat` and `GDS.dsc`. Although the item responses in this example are 0-1, GUIDE allows them to be in any ordinal (e.g., Likert) scale. The contents of `GDS.dsc` are:

```
GDS.dat
NA
c1 c2 c3
1 SATIS d
2 SPIRIT d
3 HAPPY d
4 ALIVE d
5 ENERGY d
6 DROP d
7 EMPTY d
8 BORED d
9 AFRAID d
10 HELP d
```

```

11 HOME d
12 MEMORY d
13 WORTH d
14 HOPE d
15 BETTER d
16 MALE c
17 EDUCATION n
18 AGE n

```

Here is the session log to create an input file for identifying DIF items and the important predictor variables:

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: GDSimp.in
Input 1 for model fitting, 2 for importance or DIF scoring,
    3 for data conversion ([1:3], <cr>=1): 2
Name of batch output file: GDSimp.out
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
    1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
    5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 5
    Choose option 5 for item response data.
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (max 100 characters);
enclose with matching quotes if it has spaces: GDS.dsc
Reading data description file ...
Training sample file: GDS.dat
Missing value code: NA
Warning: N variables changed to S
Number of D variables = 15
D variables are:
SATIS
SPIRIT
HAPPY
ALIVE
ENERGY
DROP
EMPTY
BORED

```

```

AFRAID
HELP
HOME
MEMORY
WORTH
HOPE
BETTER
Choose multivariate or univariate split variable selection:
Choose multivariate if there is an order among the D variables; otherwise choose univariate
Input 1 for multivariate, 2 for univariate ([1:2], <cr>=2):
Input 1 to normalize D variables, 2 for no normalization ([1:2], <cr>=2):
Input 1 for equal, 2 for unequal weighting of D variables ([1:2], <cr>=1):
Length of longest data entry = 2
Total number of cases: 1978
Col. no. Categorical variable    #levels    #missing values
      16 MALE                    2            0
Checking data ...
PCA can be used for variable selection
Do not use PCA if differential item functioning (DIF) scores are wanted
Input 1 to use PCA, 2 otherwise ([1:2], <cr>=2):
  Choose the default because DIF scoring is desired.

#cases w/ miss. D = number of cases with all D values missing
      Total  #cases w/  #missing
      #cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
      1978      0      0      0      0      0      0      2      0      1
No. cases used for training: 1978
Warning: interaction tests skipped
Input expected fraction of noise variables erroneously selected ([0.00:0.99], <cr>=0.01):
Input 1 to save p-value matrix for differential item functioning (DIF),
      2 otherwise ([1:2], <cr>=1):
Input file name to store DIF p-values: GDSimp.pv
You can create a description file with the selected variables included or excluded
Input 2 to create such a file, 1 otherwise ([1:2], <cr>=1): 2
Input 1 to keep only selected variables, 2 to exclude selected variables ([1:2], <cr>=1):
Input file name: GDSsub.dsc
You can also output the importance scores and variable names to a file
Input 1 to create such a file, 2 otherwise ([1:2], <cr>=1):
Input file name: GDSimp.scr
Input file is created!

```

The importance scores in the output file `GDSimp.scr` shows that all three predictor variables are distinguishable from noise (because their scores are above 1.0):

Rank	Score	Variable
1	9.702	AGE

2	3.501	MALE
3	1.598	EDUCATION

The last column of GDSimp.pv below shows that three items (#4, 10, 13) have DIF.

Item	Itemname	EDUCATION	AGE	GENDER	DIF
1	SATIS	0.310E-01	0.946E-01	0.476E-01	no
2	SPIRIT	0.988E+00	0.456E+00	0.437E-02	no
3	HAPPY	0.938E+00	0.930E-01	0.375E-01	no
4	ALIVE	0.129E+00	0.282E-01	0.382E+00	no
5	ENERGY	0.721E+00	0.845E+00	0.573E-06	yes
6	DROP	0.107E-01	0.117E+00	0.951E+00	no
7	EMPTY	0.369E-02	0.194E-02	0.315E-01	no
8	BORED	0.750E-07	0.166E+00	0.416E+00	yes
9	AFRAID	0.106E+00	0.323E-02	0.287E-02	no
10	HELP	0.928E-01	0.678E+00	0.148E-02	no
11	HOME	0.128E+00	0.826E+00	0.779E-03	no
12	MEMORY	0.434E+00	0.000E+00	0.440E-01	yes
13	WORTH	0.934E+00	0.573E+00	0.624E+00	no
14	HOPE	0.653E+00	0.799E+00	0.109E+00	no
15	BETTER	0.956E+00	0.525E+00	0.747E+00	no

The following output file GDSsub.dsc can be used to fit a model to the selected item and predictor variables:

```
"GDS.dat"
"NA"
colnumber  varname  vartype
1 SATIS x
2 SPIRIT x
3 HAPPY x
4 ALIVE x
5 ENERGY d
6 DROP x
7 EMPTY x
8 BORED d
9 AFRAID x
10 HELP x
11 HOME x
12 MEMORY d
13 WORTH x
14 HOPE x
15 BETTER x
16 MALE c
```



```
17 EDUCATION n
18 AGE n
```

Following is the input file creation log that uses GDSsub.dsc.

```
0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: GDSsub.in
Input 1 for model fitting, 2 for importance or DIF scoring,
    3 for data conversion ([1:3], <cr>=1):
Name of batch output file: GDSsub.out
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
    1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
    5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1): 5
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (max 100 characters);
enclose with matching quotes if it has spaces: GDSsub.dsc
Reading data description file ...
Training sample file: GDS.dat
Missing value code: NA
Warning: N variables changed to S
Number of D variables = 3
D variables are:
ENERGY
BORED
MEMORY
Choose multivariate or univariate split variable selection:
Choose multivariate if there is an order among the D variables; otherwise choose univariate
Input 1 for multivariate, 2 for univariate ([1:2], <cr>=1): 2
    Choose 2 because items are not ordered.
Input 1 to normalize D variables, 2 for no normalization ([1:2], <cr>=2):
Input 1 for equal, 2 for unequal weighting of D variables ([1:2], <cr>=1):
Length of longest data entry = 2
Total number of cases: 1978
Col. no. Categorical variable    #levels    #missing values
      16 MALE                    2              0
Checking data ...
PCA can be used for variable selection
```

```

Do not use PCA if differential item functioning (DIF) scores are wanted
Input 1 to use PCA, 2 otherwise ([1:2], <cr>=2):
#cases w/ miss. D = number of cases with all D values missing
      Total #cases w/ #missing
      #cases miss. D ord. vals #X-var #N-var #F-var #S-var #B-var #C-var
      1978      0      0      12      0      0      2      0      1
No. cases used for training: 1978
Warning: interaction tests skipped
Input 1 for LaTeX tree code, 2 to skip it ([1:2], <cr>=1):
Input file name to store LaTeX code (use .tex as suffix): GDSsub.tex
Input 2 to save node IDs of individual cases, 1 otherwise ([1:2], <cr>=2):
Input name of file to store terminal node ID of each case: GDSsub.nid
Input 2 to save fitted values at each terminal node; 1 otherwise ([1:2], <cr>=1): 2
Input name of file to store node fitted values: GDSsub.fit
Input file is created!

```

The result of running this input file produces a tree with one split on $\text{AGE} \leq 64$. The file `GDSsub.nid` gives the terminal node number that each case resides in the tree. The first 10 cases are shown below.

case	train	node
1	y	3
2	y	3
3	y	3
4	y	3
5	y	3
6	y	3
7	y	3
8	y	2
9	y	3
10	y	2

The file `GDSsub.fit` gives the sample mean values of the dependent (item) variables in each terminal node:

node	GDBORED	GDMEMORY	GDENERGY
2	0.13109E+00	0.64419E+00	0.75655E+00
3	0.70175E-01	0.37953E+00	0.77485E+00

Figure 20 shows the tree and the mean item responses in each terminal node by item.

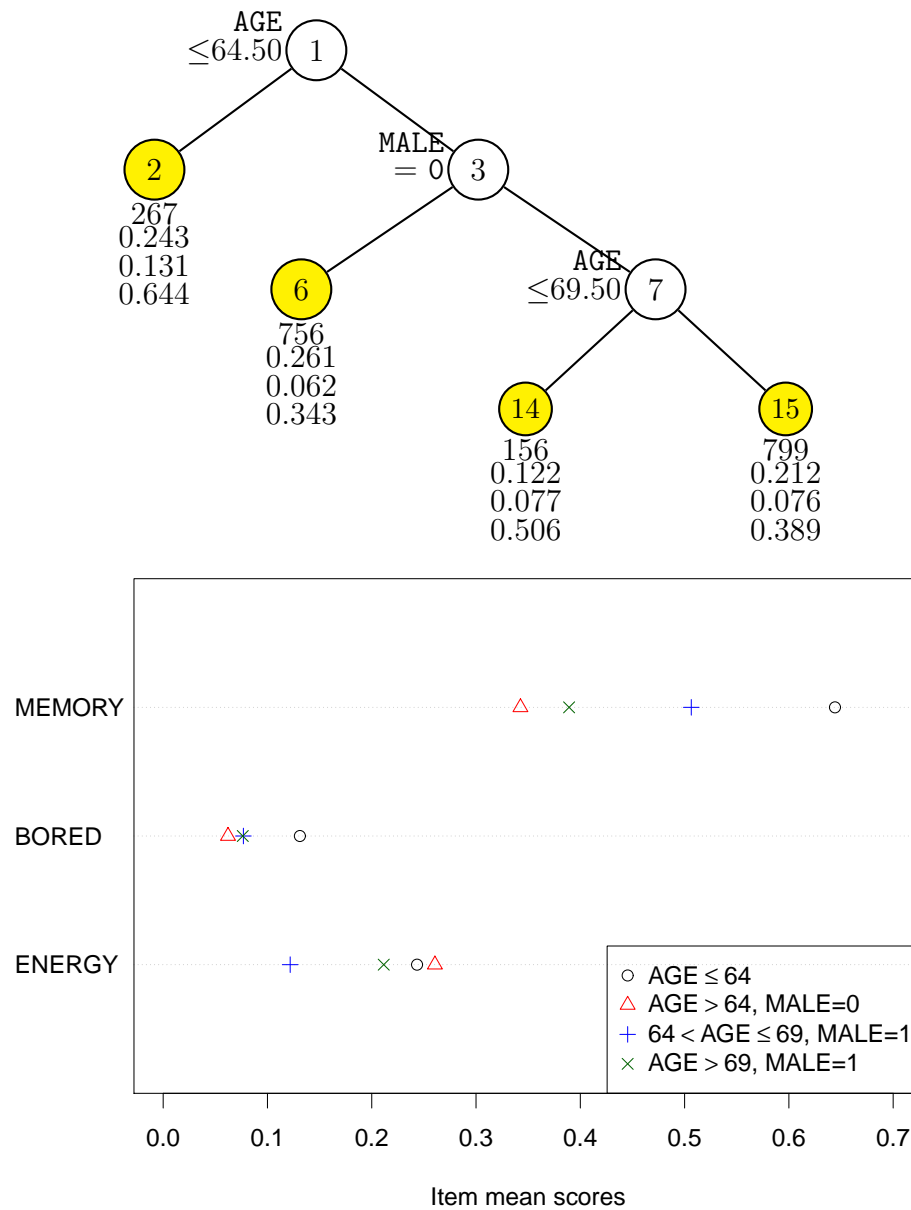


Figure 20: GUIDE 0.50-SE regression tree for predicting response variables **ENERGY**, **BORED**, and **MEMORY**. PCA not used. At each split, an observation goes to the left branch if and only if the condition is satisfied. Sample sizes and predicted values of **ENERGY**, **BORED**, and **MEMORY** are printed below nodes.

6 Tree ensembles

A tree ensemble is a collection of trees. GUIDE has two methods of constructing an ensemble. One is called “bagged GUIDE”, which fits *pruned* GUIDE trees to bootstrap samples of the training data (Breiman, 1996). The other is called “GUIDE forest”; it is similar to random forest (Breiman, 2001), which fits *unpruned* trees to bootstrap samples but randomly selects a small subset of variables for split selection at each node. There is some empirical evidence that, if there are many variables of which only a few are useful for prediction, bagged GUIDE tends to be more accurate than GUIDE forest (Loh, 2009, 2012). But GUIDE forest is computationally faster.

6.1 Bagged GUIDE

We first demonstrate bagged GUIDE on the car data.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: bagin.txt
Input 1 for model fitting, 2 for importance or DIF scoring,
    3 for data conversion ([1:3], <cr>=1):
Name of batch output file: bagout.txt
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1): 2
    This is where an ensemble method is selected.
Input 1 for bagging, 2 for rforest: ([1:2], <cr>=2): 1
    Option 1 is bagged GUIDE, option 2 is GUIDE forest.
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):
Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): drivesc.txt
Reading data description file ...
Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Total number of cases =      428
Number of classes =      3
Col. no. Categorical variable  #levels  #missing values
      3 Make                  38         0
      5 Type                  6         0

```

```

Checking data ...
Class name      Num. cases  Proportion
4wd              94      0.21962617
fwd             224      0.52336449
rwd             110      0.25700935

  Total  #cases w/  #missing
#cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
    428      0      0      10      0      0      11      0      2

No. cases used for training =      428
Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
Input 1, 2, or 3 ([1:3], <cr>=1):
Choose 1 for unit misclassification costs, 2 to input costs from a file
Input 1 or 2 ([1:2], <cr>=1):
Input name of file to store predicted class and probability: bagfit.txt
Input file is created!

```

Results

```

Ensemble of bagged classification trees
Pruning by cross-validation
Data description file: drivesc.txt
Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Number of classes = 3
Class      #Cases  Proportion
4wd         94      0.21962617
fwd        224      0.52336449
rwd        110      0.25700935

Summary information (without x variables)
d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,
n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
Column  Name      Minimum      Maximum  #Categories  #Missing
    3  Make      c                      38
    5  Type      c                      6
    6  Drive     d                      3
   14  Rprice    s   1.0280E+04  1.9246E+05
   15  Dcost     s   9.8750E+03  1.7356E+05
   16  Enginsz   s   1.3000E+00  8.3000E+00
   17  Cylin     s  -1.0000E+00  1.2000E+01
   18  Hp        s   7.3000E+01  5.0000E+02
   19  City      s   1.0000E+01  6.0000E+01

```

20	Hwy	s	1.2000E+01	6.6000E+01
21	Weight	s	1.8500E+03	7.1900E+03
22	Whlbase	s	8.9000E+01	1.4400E+02
23	Length	s	1.4300E+02	2.2800E+02
24	Width	s	6.4000E+01	8.1000E+01

Total #cases	#cases w/ miss.	D	#missing ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
428	0		0	10	0	0	11	0	2

No. cases used for training: 428

Univariate split highest priority
 Interaction splits 2nd priority; no linear splits
 Number of trees in ensemble = 100
 Pruning by v-fold cross-validation, with v = 5
 Selected tree is based on mean of CV estimates
 Simple node models
 Estimated priors
 Unit misclassification costs
 Fraction of cases used for splitting each node = 0.23364
 Max number of split levels = 7
 Minimum node size = 10
 Number of SE's for pruned tree = 5.0000E-01

Mean number of terminal nodes = 8.020

Classification matrix for training sample:

Predicted	True class		
class	4wd	fwd	rwd
4wd	61	9	2
fwd	19	212	10
rwd	14	3	98
Total	94	224	110

Number of cases used for tree construction = 428

Number misclassified = 57

Resubstitution est. of mean misclassification cost = 0.13317757009345793

Note: The above results will likely differ slightly from one run to another due to the randomness of bagging.

Predicted class probability estimates are stored in bagfit.txt

Elapsed time in seconds: 7.46616936

The following lines from the top of the file `bagfit.txt` give the estimated class probabilities and the predicted class of the observations.

"4wd"	"fwd"	"rwd"	predicted	observed		
0.50866E-01	0.90605E+00	0.43088E-01	"fwd"	"fwd"		
0.50866E-01	0.90605E+00	0.43088E-01	"fwd"	"fwd"		
0.50866E-01	0.90605E+00	0.43088E-01	"fwd"	"fwd"		
0.50866E-01	0.90605E+00	0.43088E-01	"fwd"	"fwd"		
0.50866E-01	0.90605E+00	0.43088E-01	"fwd"	"fwd"		
0.64426E-01	0.85457E+00	0.81001E-01	"fwd"	"fwd"		
0.64426E-01	0.85457E+00	0.81001E-01	"fwd"	"fwd"		

6.2 GUIDE forest

GUIDE forest differs from bagged GUIDE in two respects:

1. At each node a random subset of the variables is used for split selection.
2. The trees in GUIDE forest are not pruned.

These are the same principles in random forest. GUIDE forest differs from the latter in using GUIDE's unbiased variable selection method instead of greedy search. GUIDE forest typically requires many more trees than bagged GUIDE to achieve similar accuracy, but because the former does not prune the trees, the former is still faster to compute.

Input file creation

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: forestin.txt
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for data conversion ([1:3], <cr>=1):
Name of batch output file: forestout.txt
Input 1 for single tree, 2 for ensemble ([1:2], <cr>=1): 2
Input 1 for bagging, 2 for rforest: ([1:2], <cr>=2):
Input 1 for random splits of missing values, 2 for nonrandom: ([1:2], <cr>=2):
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1):
Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): drivedsc.txt
Reading data description file ...

```

```

Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Total number of cases =      428
Number of classes =      3
Col. no. Categorical variable  #levels  #missing values
      3 Make                   38        0
      5 Type                   6         0
Checking data ...
Class name      Num. cases  Proportion
4wd              94    0.21962617
fwd             224    0.52336449
rwd             110    0.25700935
Total  #cases w/  #missing
#cases  miss. D  ord. vals  #X-var  #N-var  #F-var  #S-var  #B-var  #C-var
    428      0      0      10      0      0      11      0      2
No. cases used for training =      428
Choose 1 for estimated priors, 2 for equal priors, 3 to input the priors from a file
Input 1, 2, or 3 ([1:3], <cr>=1):
Choose 1 for unit misclassification costs, 2 to input costs from a file
Input 1 or 2 ([1:2], <cr>=1):
Input name of file to store predicted class and probability: forest.fit
Input file is created!

```

Results

```

Random forest of classification trees
Data description file: drivesc.txt
Training sample file: drive.txt
Missing value code: *
Warning: N variables changed to S
Dependent variable is Drive
Length of longest data entry = 26
Number of classes = 3
Class      #Cases  Proportion
4wd         94    0.21962617
fwd        224    0.52336449
rwd        110    0.25700935

Summary information (without x variables)
d=dependent, b=split and fit cat variable using 0-1 dummies, c=split-only categorical,
n=split and fit numerical, f=fit-only numerical, s=split-only numerical, w=weight
Column Name      Minimum      Maximum  #Categories  #Missing

```


3	Make	c			38
5	Type	c			6
6	Drive	d			3
14	Rprice	s	1.0280E+04	1.9246E+05	
15	Dcost	s	9.8750E+03	1.7356E+05	
16	Enginsz	s	1.3000E+00	8.3000E+00	
17	Cylin	s	-1.0000E+00	1.2000E+01	
18	Hp	s	7.3000E+01	5.0000E+02	
19	City	s	1.0000E+01	6.0000E+01	
20	Hwy	s	1.2000E+01	6.6000E+01	
21	Weight	s	1.8500E+03	7.1900E+03	
22	Whlbase	s	8.9000E+01	1.4400E+02	
23	Length	s	1.4300E+02	2.2800E+02	
24	Width	s	6.4000E+01	8.1000E+01	

Total	#cases w/ #cases	miss. D	#missing ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
428	0	0	10	0	0	11	0	2	

No. cases used for training: 428

Univariate split highest priority
 No interaction and linear splits
 Number of trees in ensemble = 500
 Number of variables used for splitting = 5
 Simple node models
 Estimated priors
 Unit misclassification costs
 Fraction of cases used for splitting each node = 0.23364
 Max number of split levels = 10
 Minimum node size = 5
 Mean number of terminal nodes = 29.82

Classification matrix for training sample:

Predicted	True class		
class	4wd	fwd	rwd
4wd	80	6	0
fwd	4	216	3
rwd	10	2	107
Total	94	224	110

Number of cases used for tree construction = 428
 Number misclassified = 25
 Resubstitution est. of mean misclassification cost = 5.8411214953271021E-002

Predicted class probability estimates are stored in forest.fit
 Elapsed time in seconds: 3.62463903

The above results are not particularly useful because it is impossible to analyze the individual trees. The results mostly provide a record of the parameter values chosen to construct the forest. The most interesting results are the predicted values in the file `forest.fit`, the top few lines of which are shown below.

```
"4wd" "fwd" "rwd" predicted observed
0.62238E-02 0.98398E+00 0.97993E-02 "fwd" "fwd"
0.63508E-02 0.98035E+00 0.13298E-01 "fwd" "fwd"
0.55558E-02 0.98198E+00 0.12466E-01 "fwd" "fwd"
0.53931E-02 0.98608E+00 0.85241E-02 "fwd" "fwd"
0.56024E-02 0.98205E+00 0.12348E-01 "fwd" "fwd"
0.92692E-02 0.96686E+00 0.23874E-01 "fwd" "fwd"
0.92692E-02 0.96697E+00 0.23756E-01 "fwd" "fwd"
```

7 Importance scores

GUIDE can rank the variables in order of their importance for predicting the dependent variable. In addition, it provides a threshold score for distinguishing the important variables from the unimportant ones.

7.1 Baseball data example

We demonstrate this capability with the baseball data below.

```
0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 1
Name of batch input file: bbimp.in
Input 1 for model fitting, 2 for importance or DIF scoring, 3 for
data conversion ([1:3], <cr>=1): 2
  Option 2 yields importance scores.
Name of batch output file: bbimp.out
Input 1 for classification tree, 2 for regression tree ([1:2], <cr>=1): 2
Choose type of regression model:
  1=linear, 2=quantile, 3=Poisson, 4=proportional hazards,
  5=multiresponse or itemresponse, 6=longitudinal data (with T variables).
Input choice ([1:6], <cr>=1):
Input 1 for least squares, 2 least median of squares ([1:2], <cr>=1):
```

Input 1 for default options, 2 otherwise ([1:2], <cr>=1):

Input name of data description file (max 100 characters);

enclose with matching quotes if it has spaces: bbdsc.txt

Reading data description file ...

Training sample file: bbdatt.txt

Missing value code: NA

Warning: N variables changed to S

Warning: B variables changed to C

Dependent variable is Logsalary

Length of longest data entry = 17

Total number of cases: 263

Col. no.	Categorical variable	#levels	#missing values
16	Leag86	2	0
17	Div86	2	0
18	Team86	24	0
19	Pos86	23	0
24	Leag87	2	0
25	Team87	24	0

Checking data ...

Total #cases	#cases w/ miss. D	#missing ord. vals	#X-var	#N-var	#F-var	#S-var	#B-var	#C-var
263	0	0	3	0	0	16	0	6

No weight variable in data file

No. cases used for training: 263

Input expected fraction of noise variables erroneously selected

([0.00:0.99], <cr>=0.01):

This sets the ‘alpha’ value such that, under the null hypothesis that all variables are noise, the proportion erroneously selected is alpha.

You can create a description file with the selected variables included or excluded

Input 2 to create such a file, 1 otherwise ([1:2], <cr>=1): 2

This option lets GUIDE automatically write a new description file with the unimportant variables given the X designation.

Input 1 to keep only selected variables, 2 to exclude selected variables ([1:2], <cr>=1):

Input file name: bbsub.dsc

You can also output the importance scores and variable names to a file

Input 1 to create such a file, 2 otherwise ([1:2], <cr>=1):

Input file name: bbimp.scr

A file by that name already exists

Input 1 to overwrite it, 2 to choose another name ([1:2], <cr>=1):

Input file is created!

Results The importance scores are given at the end of the output file bbimp.out.

Predictor variables sorted by importance scores
Importance Scores

Scaled	Unscaled	Rank	Variable
100.0	1.93217E+01	1	Hitcr
95.4	1.84423E+01	2	Batcr
85.8	1.65721E+01	3	Runcr
82.0	1.58486E+01	4	Rbcr
74.7	1.44374E+01	5	Yrs
66.6	1.28632E+01	6	Wlkcr
41.8	8.06758E+00	7	Hit86
41.0	7.91880E+00	8	Hrcr
32.0	6.17371E+00	9	Run86
31.0	5.98389E+00	10	Bat86
30.8	5.95530E+00	11	Rb86
27.6	5.32659E+00	12	Wlk86
17.6	3.39316E+00	13	Hr86
8.8	1.70322E+00	14	Pos86
7.3	1.40858E+00	15	Puto86
5.3	1.02078E+00	16	Team87
----- cut-off -----			
3.7	7.07834E-01	17	Err86
3.4	6.53864E-01	18	Asst86
2.3	4.43732E-01	19	Team86
1.9	3.61595E-01	20	Leag87
1.4	2.77317E-01	21	Leag86
1.3	2.51640E-01	22	Div86

Variables with unscaled scores above 1 (the cut-off line) are deemed important.

Number of important splitting variables = 16

Number of unimportant splitting variables = 6

Here are the contents of the file `bbimp.scr`:

Rank	Score	Variable
1	19.322	Hitcr
2	18.442	Batcr
3	16.572	Runcr
4	15.849	Rbcr
5	14.437	Yrs
6	12.863	Wlkcr
7	8.068	Hit86
8	7.919	Hrcr
9	6.174	Run86
10	5.984	Bat86
11	5.955	Rb86
12	5.327	Wlk86
13	3.393	Hr86

14	1.703	Pos86
15	1.409	Puto86
16	1.021	Team87
17	0.708	Err86
18	0.654	Asst86
19	0.444	Team86
20	0.362	Leag87
21	0.277	Leag86
22	0.252	Div86

And here are the contents of the file `bbsub.dsc`:

```
"bbdat.txt"
"NA"
colnumber  varname    vartype
1 Id x
2 Name x
3 Bat86 n
4 Hit86 n
5 Hr86 n
6 Run86 n
7 Rb86 n
8 Wlk86 n
9 Yrs n
10 Batcr n
11 Hitcr n
12 Hrcr n
13 Runcr n
14 Rbcr n
15 Wlkcr n
16 Leag86 x
17 Div86 x
18 Team86 x
19 Pos86 c
20 Puto86 n
21 Asst86 x
22 Err86 x
23 Salary x
24 Leag87 x
25 Team87 c
26 Logsalary d
```

8 Other features

8.1 Pruning with test samples

GUIDE typically has three pruning options for deciding the size of the final tree: (i) cross-validation, (ii) test sample, and (iii) no pruning. Test-sample pruning is available only when there are no derived variables, such as creation of dummy indicator variables when ‘b’ variables are present. If test-sample pruning is chosen, the program will ask for the name of the file containing the test samples. This file must have the same column format as the training sample file. Pruning with test-samples or no pruning are non-default options.

8.2 Prediction of test samples

GUIDE can produce R code to predict future observations from all except kernel and nearest neighbor classification and ensemble models. This is also a non-default option.

Predictions of the training data for all models can be obtained, however, at the time of tree construction. This feature can be used to obtain predictions on “test samples” (i.e., observations that are not used in tree construction) by adding them to the training sample file. There are two ways to distinguish the test observations from the training observations:

1. Use a *weight* variable (designated as W in the description file) that takes value 1 for each training observation and 0 for each test observation.
2. Replace the D values of the test observations with the missing value code.

For tree construction, GUIDE does not use observations in the training sample file that have zero weight.

8.3 GUIDE in R and in simulations

GUIDE can be used in simulations or used repeatedly on bootstrap samples to produce an ensemble of tree models. For the latter,

1. Create a file (with name `data.txt`, say) containing one set of bootstrapped data.
2. Create a data description file (with name `desc.txt`, say) that refers to `data.txt`.

3. Create an input file (with name `input.txt`, say) that refers to `desc.txt`.
4. Write a batch program (Windows) or a shell script (Linux or Macintosh) that repeatedly:
 - (a) replaces the file `data.txt` with new bootstrapped samples;
 - (b) calls GUIDE with the command: `guide < input.txt`; and
 - (c) reads and processes the results from each GUIDE run.

In R, the command in step 4b depends on the operating system. If the GUIDE program and the files `data.txt` and `input.txt` are in the same folder as the working R directory, the command is:

Linux/Macintosh: `system("guide < input.txt > log.txt")`

Windows: `shell("guide < input.txt > log.txt")`

If the files are not all in the same folder, full path names must be given. Here `log.txt` is a text file that stores messages during execution. If GUIDE does not run successfully, errors are also written to `log.txt`.

8.4 Generation of powers and products

GUIDE allows the creation of certain powers and products of regressor variables on the fly. Specifically, variables of the form $X_1^p X_2^q$, where X_1 and X_2 are numerical predictor variables and p and q are integers, can be created by adding one or more lines of the form

```
0 i p j q a
```

at the end of the data description file. Here i and j are integers giving the column numbers of variables X_1 and X_2 , respectively, in the data file and a is one of the letters **n**, **s**, or **f** (corresponding to a numerical variable used for both splitting and fitting, splitting only, or fitting only).

To illustrate, suppose we wish to fit a piecewise quadratic model in the variable **Yrs** for the baseball data. This is easily done by adding one line to the file `bbdsc.txt`. First we assign the **s** (for splitting only) designator to every numerical predictor except **Yrs**. This will prevent all variables other than **Yrs** from acting as regressors in the piecewise quadratic models. To create the variable Yrs^2 , add the line

```
0 9 2 9 0 f
```

to the end of `bbdsc.txt`. The 9's in the above line refers to the column number of the variables `Yrs` in the data file, and the `f` tells the program to use the variable Yrs^2 for fitting terminal node models only. Note: The line defines Yrs^2 as $\text{Yrs}^2 \times \text{Yrs}^0$. Since we can equivalently define the variable by $\text{Yrs}^2 = \text{Yrs}^1 \times \text{Yrs}^1$, we could also have used the line: "0 9 1 9 1 f".

The resulting description file now looks like this:

```
bbdat.txt
NA
column, varname, vartype
1 Id x
2 Name x
3 Bat86 s
4 Hit86 s
5 Hr86 s
6 Run86 s
7 Rb86 s
8 Wlk86 s
9 Yrs n
10 Batcr s
11 Hitcr s
12 Hrcr s
13 Runcr s
14 Rbcr s
15 Wlkcr s
16 Leag86 c
17 Div86 c
18 Team86 c
19 Pos86 c
20 Puto86 s
21 Asst86 s
22 Err86 s
23 Salary x
24 Leag87 c
25 Team87 c
26 Logsalary d
0 9 2 9 0 f
```

When the program is given this description file, the output will show the regression coefficients of `Yrs` and Yrs^2 in each terminal node of the tree.

8.5 Data formatting functions

The program includes a utility function for reformatting data files into forms required by some statistical software packages:

1. R/Splus: Fields are space delimited. Missing values are coded as NA. Each record is written on one line. Variable names are given on the first line.
2. SAS: Fields are space delimited. Missing values are coded with periods. Character strings are truncated to eight characters. Spaces within character strings are replaced with underscores (`_`).
3. TEXT: Fields are comma delimited. Empty fields denote missing values. Character strings longer than eight characters are truncated. Each record is written on one line. Variable names are given on the first line.
4. STATISTICA: Fields are comma delimited. Commas in character strings are stripped. Empty fields denote missing values. Each record occupies one line.
5. SYSTAT: Fields are comma delimited. Strings are truncated to eight characters. Missing character values are replaced with spaces, missing numerical values with periods. Each record occupies one line.
6. BMDP: Fields are space delimited. Categorical values are sorted in alphabetic order and then assigned integer codes. Missing values are indicated by asterisks. Variable names longer than eight characters are truncated.
7. DataDesk: Fields are space delimited. Missing categorical values are coded with question marks. Missing numerical values are coded with asterisks. Each record is written on one line. Spaces within categorical values are replaced with underscores. Variable names are given on the first line of the file.
8. MINITAB: Fields are space delimited. Categorical values are sorted in alphabetic order and then assigned integer codes. Missing values are coded with asterisks. Variable names longer than eight characters are truncated.
9. NUMBERS: Same as **TEXT** option except that categorical values are converted to integer codes.
10. C4.5: This is the format required by the C4.5 (Quinlan, 1993) program.
11. ARFF: This is the format required by the WEKA (Witten and Frank, 2000) programs.

Following is a sample session where the iris data are reformatted for R or Splus.

```

0. Read the warranty disclaimer
1. Create an input file for batch run
2. Fit a model without creating input file
3. Convert data to other formats
4. Variable importance scoring and differential item functioning
Input your choice: 3
Input name of log file: log.txt
Input 1 if D variable is categorical, 2 if real, 0 if none ([0:2], <cr>=1):

Input name of data description file (maximum 100 characters; enclose within quotes
if it contains spaces or non alphanumeric characters): irisdesc.txt
Reading data description file ...
Training sample file: irisdata.txt
Missing value code: ?
Warning: N variables changed to S
Dependent variable is class
Length of longest data entry = 11
Total number of cases =      150
Number of classes =      3
Choose one of the following data formats:
      Field  Miss.val.codes
No. Name    Separ  char.   numer. Remarks
-----
 1 R/Splus   space  NA      NA      1 line/case, var names on 1st line
 2 SAS       space  .       .       strings trunc., spaces -> '_'
 3 TEXT      comma  empty   empty   1 line/case, var names on 1st line
 4 STATISTICA comma  empty   empty   1 line/case, commas stripped
                                var names on 1st line
 5 SYSTAT    comma  space   .       1 line/case, var names on 1st line
                                strings trunc. to 8 chars
 6 BMDP      space          *       strings trunc. to 8 chars
                                cat values -> integers (alph. order)
 7 DATADESK space  ?       *       1 line/case, var names on 1st line
                                spaces -> '_'
 8 MINITAB   space          *       cat values -> integers (alph. order)
                                var names trunc. to 8 chars
 9 NUMBERS   comma  NA      NA      1 line/case, var names on 1st line
                                cat values -> integers (alph. order)
10 C4.5      comma  ?       ?       1 line/case, dependent variable last
11 ARFF      comma  ?       ?       1 line/case
-----
0                                abort this job
Input your choice ([0:11], <cr>=1):
Input name of new data file: iris.rdata
Follow the commented lines in "iris.rdata" to read the data into R or Splus

```

References

- Andrews, D. and Herzberg, A. (1985). *Data: a collection of problems from many fields for the student and research worker*. Springer.
- Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24:123–140.
- Breiman, L. (2001). Random forests. *Machine Learning*, 45:5–32.
- Breiman, L., Friedman, J. H., Olshen, R. A., and Stone, C. J. (1984). *Classification and Regression Trees*. Wadsworth, Belmont.
- Broekman, B. F. P., Niti, M., Nyunt, M. S. Z., Ko, S. M., Kumar, R., and Ng, T. P. (2011). Validation of a brief seven-item response bias-free Geriatric Depression Scale. *American Journal of Geriatric Psychiatry*, 19:589–596.
- Broekman, B. F. P., Nyunt, S. Z., Niti, M., Jin, A. Z., Ko, S. M., Kumar, R., Fones, C. S. L., and Ng, T. P. (2008). Differential item functioning of the Geriatric Depression Scale in an Asian population. *Journal of Affective Disorders*, 108:285–290.
- Chan, K.-Y. and Loh, W.-Y. (2004). LOTUS: An algorithm for building accurate and comprehensible logistic regression trees. *Journal of Computational and Graphical Statistics*, 13:826–852.
- Chaudhuri, P. and Loh, W.-Y. (2002). Nonparametric estimation of conditional quantiles using quantile regression trees. *Bernoulli*, 8:561–576. www.stat.wisc.edu/~loh/treeprogs/guide/quantile.pdf.
- Hallin, M. and Ingenbleek, J.-F. (1983). The Swedish automobile portfolio in 1977. *Scandinavian Actuarial Journal*, 83:49–64.
- Hoaglin, D. C. and Velleman, P. F. (1995). A critical look at some analyses of major league baseball salaries. *American Statistician*, 49:277–285.
- Hosmer, D. W., Lemeshow, S., and May, S. (2008). *Applied Survival Analysis*. Wiley, 2nd edition.
- Kim, H. and Loh, W.-Y. (2001). Classification trees with unbiased multiway splits. *Journal of the American Statistical Association*, 96:589–604. www.stat.wisc.edu/~loh/treeprogs/cruise/cruise.pdf.

- Kim, H. and Loh, W.-Y. (2003). Classification trees with bivariate linear discriminant node models. *Journal of Computational and Graphical Statistics*, 12:512–530. www.stat.wisc.edu/~loh/treeprogs/cruise/jcgs.pdf.
- Kim, H., Loh, W.-Y., Shih, Y.-S., and Chaudhuri, P. (2007). Visualizable and interpretable regression models with good prediction power. *IIE Transactions*, 39:565–579. www.stat.wisc.edu/~loh/treeprogs/guide/iie.pdf.
- Koenker, R. W. and Bassett, G. (1978). Regression quantiles. *Econometrica*, 46:33–50.
- Loh, W.-Y. (2002). Regression trees with unbiased variable selection and interaction detection. *Statistica Sinica*, 12:361–386. www3.stat.sinica.edu.tw/statistica/j12n2/j12n21/j12n21.htm.
- Loh, W.-Y. (2006a). Logistic regression tree analysis. In Pham, H., editor, *Handbook of Engineering Statistics*, pages 537–549. Springer.
- Loh, W.-Y. (2006b). Regression tree models for designed experiments. In Rojo, J., editor, *The Second Erich L. Lehmann Symposium—Optimality*, volume 49, pages 210–228. Institute of Mathematical Statistics Lecture Notes-Monograph Series. arxiv.org/abs/math.ST/0611192.
- Loh, W.-Y. (2008a). Classification and regression tree methods. In Ruggeri, F., Kenett, R., and Faltin, F. W., editors, *Encyclopedia of Statistics in Quality and Reliability*, pages 315–323. Wiley, Chichester, UK. www.stat.wisc.edu/~loh/treeprogs/guide/eqr.pdf.
- Loh, W.-Y. (2008b). Regression by parts: Fitting visually interpretable models with GUIDE. In Chen, C., Härdle, W., and Unwin, A., editors, *Handbook of Computational Statistics*, pages 447–469. Springer. www.stat.wisc.edu/~loh/treeprogs/guide/handbk.pdf.
- Loh, W.-Y. (2009). Improving the precision of classification trees. *Annals of Applied Statistics*, 3:1710–1737. www.stat.wisc.edu/~loh/treeprogs/guide/aoas260.pdf.
- Loh, W.-Y. (2011). Classification and regression trees. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 1:14–23.

- Loh, W.-Y. (2012). Variable selection for classification and regression in large p , small n problems. In Barbour, A., Chan, H. P., and Siegmund, D., editors, *Probability Approximations and Beyond*, volume 205 of *Lecture Notes in Statistics—Proceedings*, pages 133–157, New York. Springer.
- Loh, W.-Y. (2014). Fifty years of classification and regression trees (with discussion). *International Statistical Review*, 34:329–370.
- Loh, W.-Y., Chen, C.-W., and Zheng, W. (2007). Extrapolation errors in linear model trees. *ACM Trans. Knowl. Discov. Data*, 1(2):6. www.stat.wisc.edu/~loh/treeprogs/guide/acm.pdf.
- Loh, W.-Y., He, X., and Man, M. (2015). A regression tree approach to identifying subgroups with differential treatment effects. *Statistics in Medicine*. DOI: [10.1002/sim.6454](https://doi.org/10.1002/sim.6454).
- Loh, W.-Y. and Shih, Y.-S. (1997). Split selection methods for classification trees. *Statistica Sinica*, 7:815–840. www3.stat.sinica.edu.tw/statistica/j7n4/j7n41/j7n41.htm.
- Loh, W.-Y. and Zheng, W. (2013). Regression trees for longitudinal and multiresponse data. *Annals of Applied Statistics*, 7:495–522.
- Marc, L. G., Raue, P. J., and Bruce, M. L. (2008). Screening performance of the 15-item Geriatric Depression Scale in a diverse elderly home care population. *American Journal of Geriatric Psychiatry*, 16:914–921.
- Murnane, R. J., Boudett, K. P., and Willett, J. B. (1999). Do male dropouts benefit from obtaining a GED, postsecondary education, and training? *Evaluation Reviews*, 23:475–502.
- Peters, A. and Hothorn, T. (2012). Improved predictors. R package version 0.8-13.
- Quinlan, J. R. (1992). Learning with continuous classes. In *5th Australian Joint Conference on Artificial Intelligence*, pages 343–348.
- Quinlan, J. R. (1993). *C4.5: Programs for Machine Learning*. Morgan Kaufmann.
- Rousseeuw, P. J. and Leroy, A. M. (1987). *Robust Regression and Outlier Detection*. Wiley.

- Schmoor, C., Olschewski, M., and Schumacher, M. (1996). Randomized and non-randomized patients in clinical trials: experiences with comprehensive cohort studies. *Statistics in Medicine*, 15:263–271.
- Singer, J. D. and Willett, J. B. (2003). *Applied Longitudinal Data Analysis*. Oxford University Press, New York, NY.
- Witten, I. and Frank, E. (2000). *Data Mining: Practical Machine Learning Tools and Techniques with JAVA Implementations*. Morgan Kaufmann, San Fransico, CA. www.cs.waikato.ac.nz/ml/weka.
- Yeh, I.-C. (2007). Modeling slump flow of concrete using second-order regressions and artificial neural networks. *Cement and Concrete Composites*, 29:474–480.