

# Building a Tree Model with the Retention Dataset

*Jerry Kiely*

*11 April 2017*

## Introduction

This analysis is concerned with data relating to student retention in the Engineering faculty of DIT. It is the same data considered in Assignment 3, but with more potential predictors available.

The purpose of the analysis will be to use rpart models to build a predictive model for student retention. The data includes risk risk factors regarding prior academic performance (e.g. leaving certificate results, leaving certificate maths grade), personal characteristics (gender, home address, CAO choices made etc).

## The Algorithm

rpart takes a data set and recursively partitions it based on a splitting criteria - choosing the split that maximizes the reduction in impurity for the node. For example in the case of the root, and splitting on lcpoints, deciding where to split within lcpoints involves:

- iterating over all values of lcpoints in turn
- calculating the reduction in impurity for splitting at each value
- keeping track of the value that results in the maximum impurity reduction
- splitting at the value that results in the maximum impurity reduction

this process is repeated at every level until other criteria such as maxdepth or minsplitlevel come in to play. The impurity of a node can be calculated using a gini score or an entropy score. An example of this process is given in the code - specifically finding the splitting point for the lcpoints predictor using a gini score.

## The Data

First we read the data in. For the sake of readability we will add a column called “result” that will hold a string value - “failed” or “passed”. Also we will drop the “id” column as unnecessary, and the “overall6” column as redundant. Next we split the data into a training and validation set. The training set will be used

to create our model, and the validation set will be used to see how well it performs, and to help with the process of improving it through pruning.

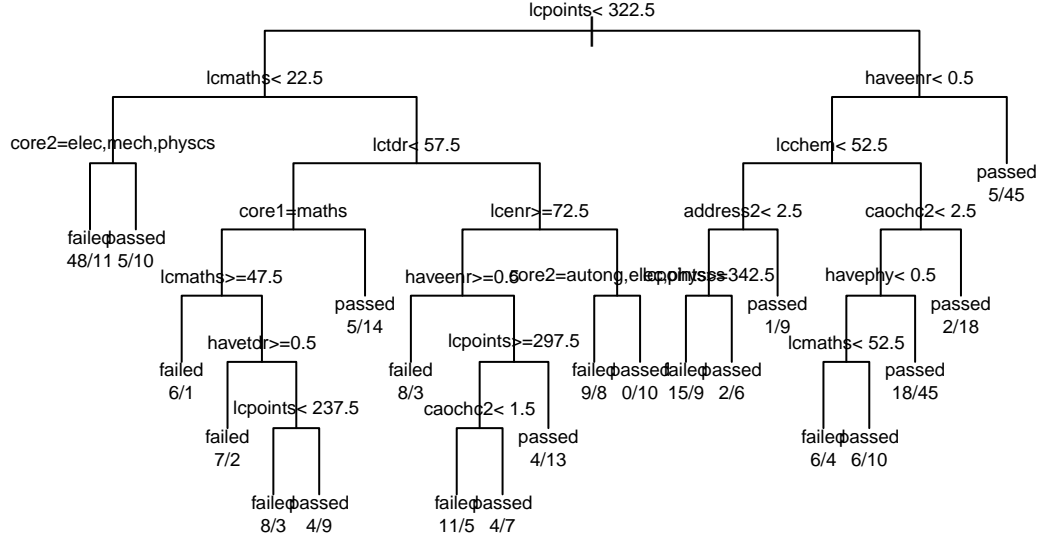


Figure 1: Tree built with default parameters and  $cp = 0$

## The Model

We fit a model with all variables included as potential predictors. The default values are used for all parameters except for the complexity parameter, which is set to 0. The plot of the tree for this model is included in Figure 1.

In Figure 2 we see a plot of the complexity parameters for the tree. The value of interest is the left-most value below the dotted line:

- below the dotted line because we are interested in the complexity parameter with the the lowest cross validation error
- left-most because we favour a smaller tree over a deeper, more complex, tree which could suffer from

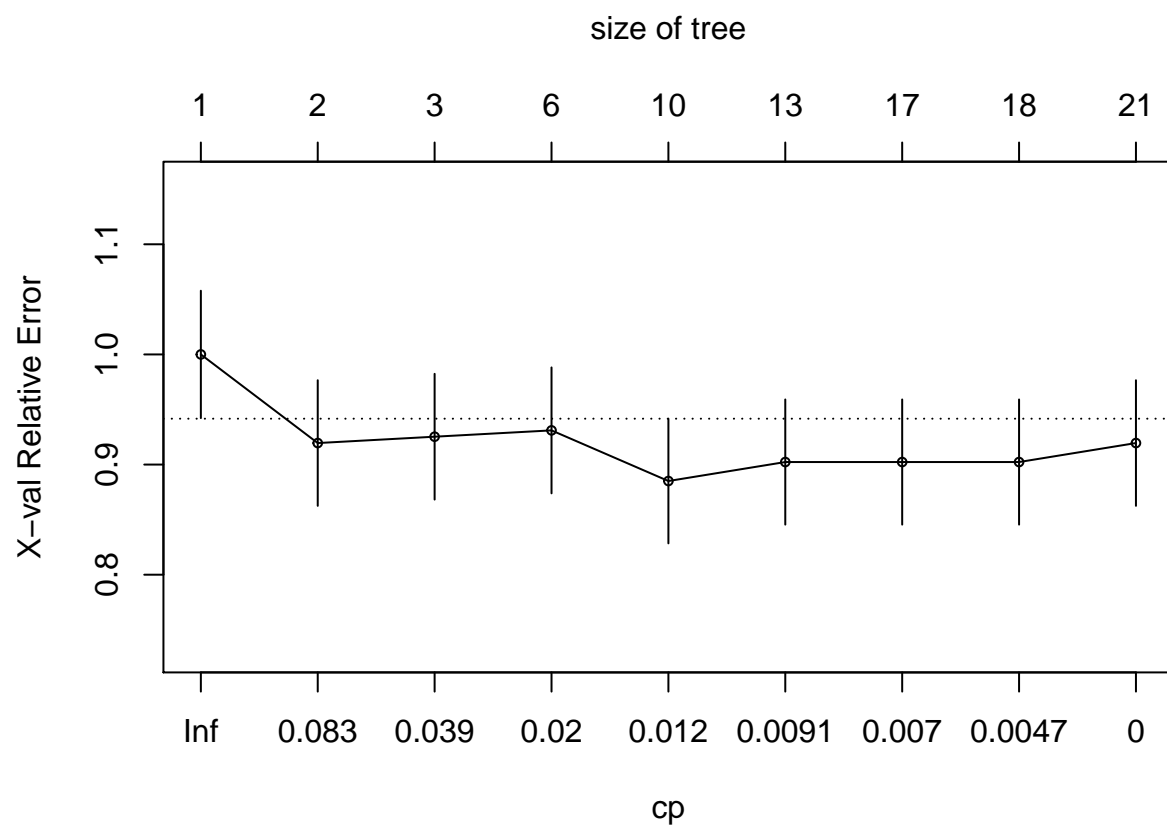


Figure 2: Complexity parameters versus cross validation results versus size of tree

overfitting

so we choose a value of 0.0095785 for our complexity parameter, and prune our tree based on that.

Table 1: table detailing complexity parameters

CP	nsplit	rel error	xerror	xstd
0.1321839	0	1.0000000	1.0000000	0.0578211
0.0517241	1	0.8678161	0.9195402	0.0570275
0.0287356	2	0.8160920	0.9252874	0.0570936
0.0143678	5	0.7298851	0.9310345	0.0571582
0.0095785	9	0.6666667	0.8850575	0.0565998
0.0086207	12	0.6379310	0.9022989	0.0568203
0.0057471	16	0.6034483	0.9022989	0.0568203
0.0038314	17	0.5977011	0.9022989	0.0568203
0.0000000	20	0.5862069	0.9195402	0.0570275

In Figure 3 we see our pruned tree. As can be seen it is a lot less complex than the earlier version.

## The Analysis

Now we can analyze the results of our pruning. We will now predict based on our validation set to see how well our model performs.

Table 2: Confusion matrix for pruned tree

	failed	passed
failed	0.4230769	0.5769231
passed	0.1946903	0.8053097

Looking at the confusion matrix in Table 2 we seem to do a better job of predicting passing than predicting failing - i.e. the true failed is below 0.5 and the true passed well above 0.5. An option might be to consider treating the model as a predictor of successfully passing rather than failing.

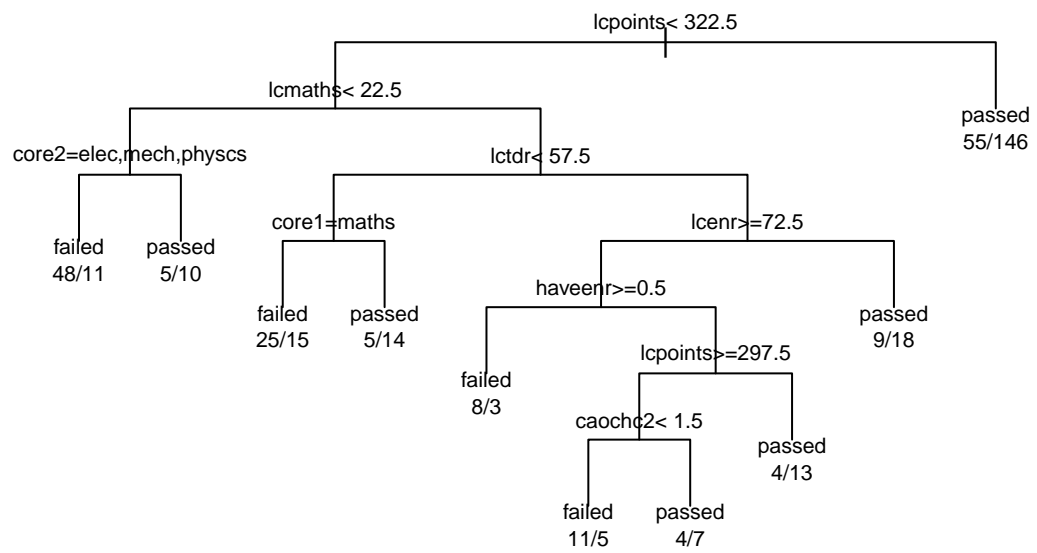


Figure 3: Tree pruned with found complexity parameter

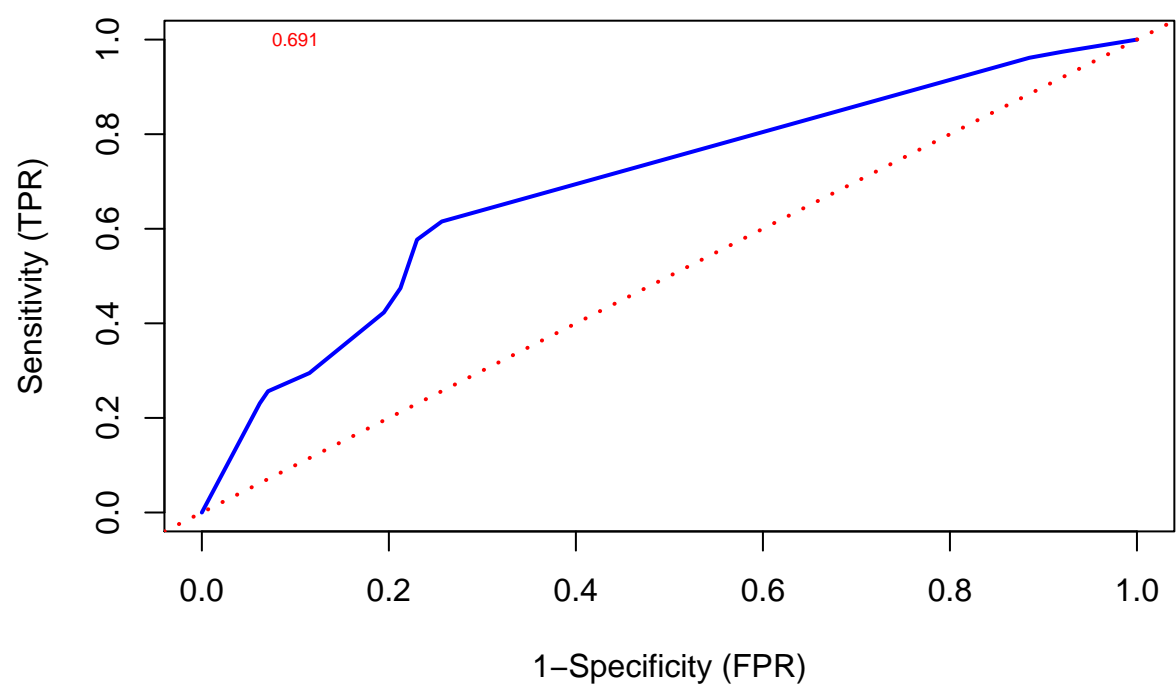


Figure 4: ROC curve with AUC value

We will next look at an ROC curve for the model in Figure 4. The blue curve gives an indication of how much better than random guessing our model is, with the red dotted line representing the random guess. The value in red is the AUC, or the area under the curve, for the ROC curve. A better than random model would have an AUC of  $> 0.5$  (the area under the red dotted line).