Chapter 8: Tree-Based Methods

We will introduce tree-based methods for regression and classification.

These involve segmenting the predictor space into a number of simple regions. To make a prediction for an observation, we use the mean or mode of training observations in the region to which it belongs.

The set of splitting rules can be summarized in a tree \Rightarrow "decision trees".

- simple and useful for interpretation - not competitive u/ other supervised approaches (e.g. lasso) for prediction.



Combining a large number of trees can often result in dramatic improvements in prediction accuracy at the expense of interpretation.



WWW. PHDCOMICS. COM

Credit: <u>http://phdcomics.com/comics.php?f=852</u>

Decision trees can be applied to both regression and classification problems. We will start with regression.

1 Regression Trees



Example: We want to predict baseball salaries using the Hittsers data set based on Years (the number of years that a player has been in the major leagues) and Hits (the number of hits he made the previous year).



The predicted salary for players is given by the mean response value for the players in that box. Overall, the tree segments the players into 3 regions of predictor space.



Ris Riz, Rig = terminal modes or leaves of the tree points along the tree where predictor space is split = internal modes sequents of tree that connect modes = branches

Years is the most important factor in deformining Salary Lygium that a player has less experience, # hits in previous year play little role in his salary. Ly among players who have been in the league 5+ years, # hits 1, Tsalary.



We now discuss the process of building a regression tree. There are to steps:

2. Predict For arey observation that fulls in into region R; we make the same prediction, the mean of the response Y for training values in R;

How do we construct the regions R_1, \ldots, R_J ? How to divide the predictor space? regions could have any shape, but that is too hard (to do, to to interpret)

>> divide predictor space into high dimensional rectangles (or " boxes)

The goal is to find boxes R_1, \ldots, R_J that minimize the RSS. $= \sum_{j=1}^{J} \sum_{i \in R_j} (q_i - \hat{y}_{R_j})^2$ where Uufortunately, it is computationally infeasible to consider every possile possile possile <math>possile $possile possile <math>\hat{y}_{R_j} = men response of framing data in R_j^{m}$ => take <u>top-down</u>, greedy approach alled recursive binary splitbly.

The approach is *top-down* because

We start at the top of the tree (all observations belong the a single region) and successively splitthe predictor space.

The approach is *greedy* because

In order to perform recursive binary splitting,

The process described above may produce good predictions on the training set, but is likely to overfit the data.

A smaller tree, with less splits might lead to lower variance and better interpretation at the cost of a little bias.

I dea: only split free it it results in large enough drop MRSS. Bad idea because seemingly worthern splits early in the tree might be followed by a "good" split.

A strategy is to grow a very large tree T_0 and then <u>prune</u> it back to obtain a subtree. How to prove the hee? gral: silent a subtrar but lends the lonest fest error rate. -> could use CV for every subtrar but the expensive (large # of possible subtras). Solution : " Cost complexity pruning" consider a sequence of trees indexed by a nonregative tuning parameter of. For each value of d, there exists a corresponding subtre TCTo st. $\begin{array}{l} |T| \\ \sum \left[(\gamma_i - \gamma_{R_m})^2 + \alpha_i |T| \right] \text{ is as small as possible.} \\ m = (\chi_i \in R_m) \\ \qquad \# \text{ terminal nodes of the tree.} \\ R_m = m^{n} \text{ terminal node rigion} \end{array}$ gran = predicted response for Rm. of controls the trade-off between subtree's complexity of fit the training data. T = To. d = 7 When - T = To. d = 7 when - T = To. d = 7 price to per dir having many tornhol hodes $t \Rightarrow$ smaller to many tornhol gikm = predicted response for Rm. Then on back and use full data I chosen of to get subtree.

Algorithm for building a regression tree:

Use recursive binary splitting to grow a lage tree on traching data, stopping only when each terminel mode has ferrer than some # (5?) observations.
Apply cost complexity praning to the large tree on traching data, stopping only when each terminel mode description of the large tree of the sequence of last trees as a function of d.
Use K-fild CV to down d
Divide training data Mo K filds, for each k=1,..., K
(a) repeate () and (2) on all but kth fild
(b) Evaluate the MSE on data to kth fold as a function of d.
Average usualts for each veloc of d to get CV error. Choose d which unbimizes CV k

(4) hadren to subtree from (2) pot corresponds to a from (3).

Example: Fit regression tree to littlers use 9 features



2 Classification Trees

A *classification tree* is very similar to a regression tree, except that it is used to predict a categorical response.

For a classification tree, we predict that each observation belongs to the *most commonly occurring class* of training observation in the region to which it belongs.

The task of growing a classification tree is quite similar to the task of growing a regression tree.

It turns out that classification error is not sensitive enough.

When building a classification tree, either the Gini index or the entropy are typically used to evaluate the quality of a particular split.

3 Trees vs. Linear Models

Regression and classification trees have a very different feel from the more classical approaches for regression and classification.

Which method is better?

3.1 Advantages and Disadvantages of Trees

4 Bagging

Decision trees suffer from high variance.

Bootstrap aggregation or *bagging* is a general-purpose procedure for reducing the variance of a statistical learning method, particularly useful for trees.

So a natural way to reduce the variance is to take many training sets from the population, build a separate prediction model using each training set, and average the resulting predictions.

Of course, this is not practical because we generally do not have access to multiple training sets.

While bagging can improve predictions for many regression methods, it's particularly useful for decision trees.

These trees are grown deep and not pruned.

How can bagging be extended to a classification problem?

4.1 Out-of-Bag Error

There is a very straightforward way to estimate the test error of a bagged model, without the need to perform cross-validation.

4.2 Interpretation

5 Random Forests

Random forests provide an improvement over bagged trees by a small tweak that decorrelates the trees.

As with bagged trees, we build a number of decision trees on bootstrapped training samples.

In other words, in building a random forest, at each split in the tree, the algorithm is not allowed to consider a majority of the predictors.

The main difference between bagging and random forests is the choice of predictor subset size m.

6 Boosting

Boosting is another approach for improving the prediction results from a decision tree.

While bagging involves creating multiple copies of the original training data set using the bootstrap and fitting a separate decision tree on each copy,

Boosting does not involve bootstrap sampling, instead each tree is fit on a modified version of the original data set.

Boosting has three tuning parameters:

1.

2.

3.