

(6.)

A Detailed Explanation of the Algorithm used to fit a Regression Tree

↳ Based loosely on ISLR Algorithm 8.1.

Step #1: Grow a large tree on the training data.

- 1: Grow a large tree on the training data.
- We do this by, for each node (including the first node), we check if the stopping criterion has been met. Often this stopping condition is to stop if a node has fewer than a given minimum number of observations.
- * If criterion has been met, do not split that terminal node any further.
 - * If criterion has NOT yet been met, then make the best binary split at that node. The best split is the one that ~~minimizes~~ minimizes, ~~the error~~ for 2 regions (of the node in question) that form a partition over the whole node in question space, call those 2 regions R_1 and R_2 , the quantity:

$$RSS = \sum_{i: X_i \in R_1(j, s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i: X_i \in R_2(j, s)} (y_i - \hat{y}_{R_2})^2 \quad \textcircled{A}$$

where: y_i is the value of the response value for obs i ,

2 regions
of node,
regions are defined by
j and s

$$\left\{ \begin{array}{l} R_1(j, s) = \{x \mid x_j < s\}; \text{ for given predictor } x_j \text{ and} \\ R_2(j, s) = \{x \mid x_j \geq s\}; \text{ cutpoint } s \end{array} \right.$$

The above ~~SS~~ sums over all obs in each of the 2 regions

\hat{y}_{R_1} and \hat{y}_{R_2} are the predictions (often the averages) of response var of the obs in their region.

↓ BACK ↓

we search over all values of j and s , and b/c decision trees are greedy algorithms, we make the split corresponding to the lowest value of RSS .

Step #2: Apply cost complexity pruning to the above large tree (call it T_0), to obtain a sequence of best trees as a function of α :

$$\{T_0, T_{\alpha_1}, T_{\alpha_2}, \dots, T_{\alpha_Q}\}$$

$\alpha=0 \Rightarrow$ large tree

we do this by..... We create an ^{increasing} ordered sequence of values for α as:

$$\alpha = \{\alpha_0, \alpha_1, \alpha_2, \dots, \alpha_Q\} \text{ with } 0 < \alpha_1 < \alpha_2 < \dots < \alpha_Q.$$

For each α (non-negative value of our tuning param), we compute

$$\text{penalized } RSS_{\alpha} = \sum_{m=1}^{|T|} \sum_{X_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

where \hat{y}_{R_m} = prediction in m th region/rectangle

m = indexes the ~~nodes~~ nodes/rectangles.

$|T|$ = number of nodes/rectangles in tree

for each α , there corresponds a subtree T_{α} that minimizes RSS_{α} . We select as T_{α_i}

This looks a bit like LASSO, where α is penalty tuning param to control complexity of algorithm.

~~The output of this is a sequence of length $(Q+1)$ of pen RSS_{α}~~
 ~~$RSS_{\alpha} = \{RSS_0, RSS_{\alpha_1}, \dots, RSS_{\alpha_Q}\}$~~

Step #3: Use K -fold cross validation to choose α . Namely, divide obs into K folds, for each $k=1, 2, \dots, K$, ~~leave~~ leave fold k out, train on all folds except fold k , evaluate ~~SSE~~ ^{SSE} on left-out k th fold, and average over all K vals of ~~SSE~~ ^{SSE} to produce a MSE value for each value of α . Select the ~~value~~ ^{value} of α that minimizes MSE.

Step #4: Return/select the subtree from Step 2 that corresponds to the value of α selected above in Step #3.