

Ch 8.1: Decision Trees

Lecture 23 - CMSE 381

Prof. Elizabeth Munch

Michigan State University

::

Dept of Computational Mathematics, Science & Engineering

Fri, Nov 11, 2022

Last time:

- Cubic Splines

This lecture:

- 8.1 Decision Trees

Announcements:

- HW #7 Due tonight
-

Section 1

Decision Trees

Big idea

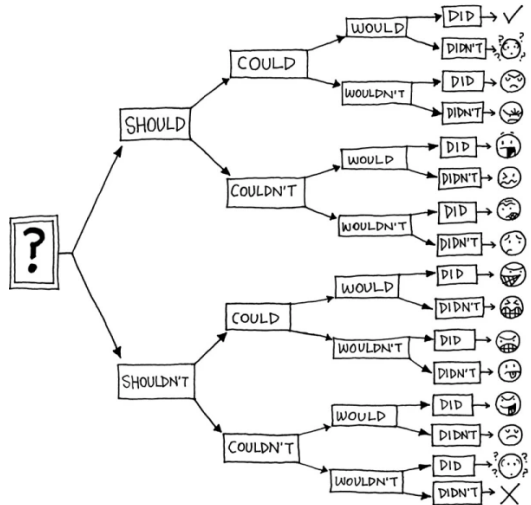


Image: <https://marekbennett.com/2014/02/14/decision-tree/>

Subset of Hitters data

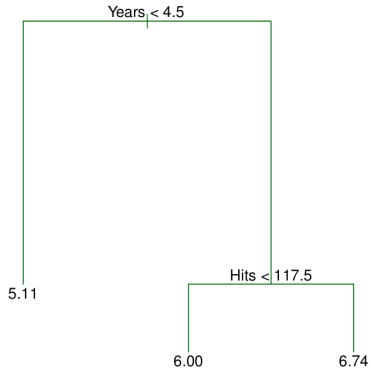
	Hits	Years	Salary	LogSalary
1	81	14	475.0	6.163315
2	130	3	480.0	6.173786
3	141	11	500.0	6.214608
4	87	2	91.5	4.516339
5	169	11	750.0	6.620073
...
317	127	5	700.0	6.551080
318	136	12	875.0	6.774224
319	126	6	385.0	5.953243
320	144	8	960.0	6.866933
321	170	11	1000.0	6.907755

First decision tree example

	Hits	Years	LogSalary
1	81	14	6.163315
2	130	3	6.173786
3	141	11	6.214608
4	87	2	4.516339
5	169	11	6.620073
...
317	127	5	6.551080
318	136	12	6.774224
319	126	6	5.953243
320	144	8	6.866933
321	170	11	6.907755



Interpretation of example

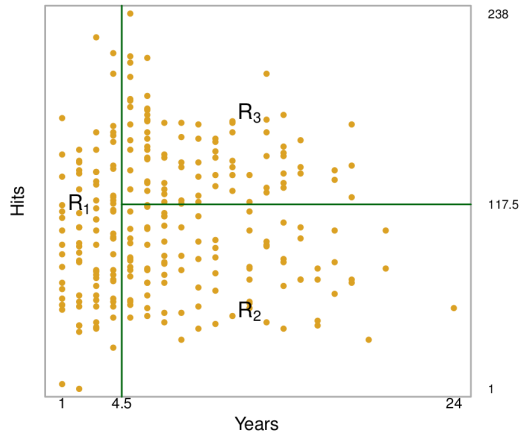
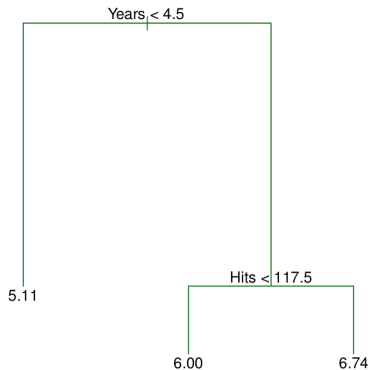


Coding a regression decision tree

Regions defined by the tree

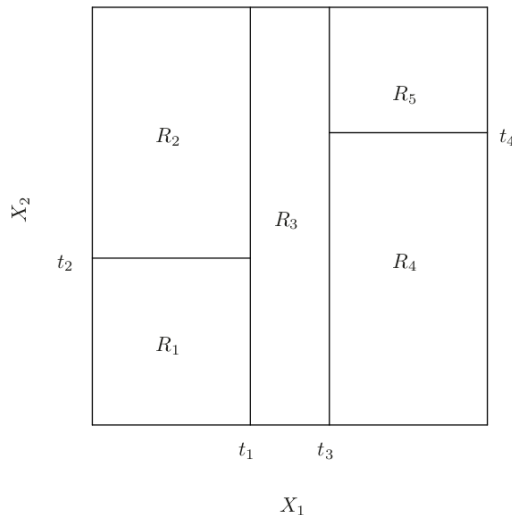


Viewing Regions Defined by Tree



How do we actually get the tree? Two steps

- 1 We divide the predictor space — that is, the set of possible values for X_1, X_2, \dots, X_p — into J distinct and non-overlapping regions, R_1, R_2, \dots, R_J .
- 2 For every observation that falls into the region R_j , we make the same prediction = the mean of the response values for the training observations in R_j .



Step 1: How do we decide on R_j s?

Goal:

Find boxes R_1, \dots, R_J that minimize

$$\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$$

\hat{y}_{R_j} = mean response for training
observations in j th box

Recursive Binary Splitting

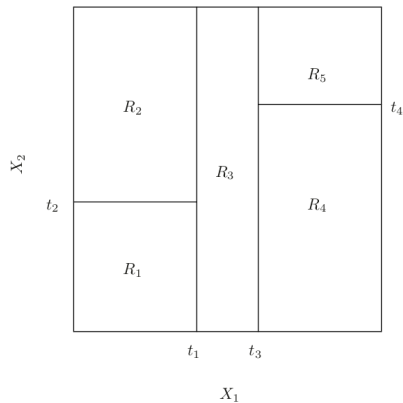
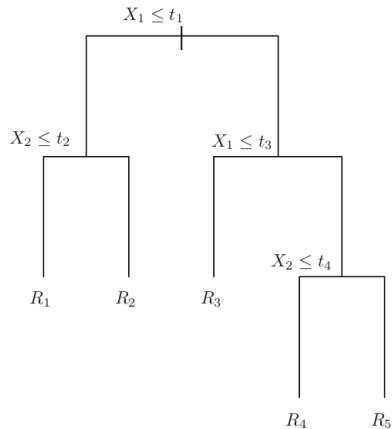
- Pick X_j
- Pick s so that splitting into $\{X \mid X_j < s\}$ and $\{X \mid X_j \geq s\}$ results in largest possible reduction in RSS

$$R_1(j, s) = \{X \mid X_j < s\}$$

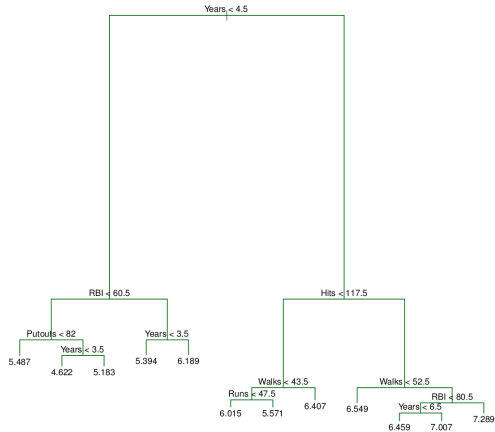
$$R_2(j, s) = \{X \mid X_j \geq s\}$$

$$\sum_{i \mid x_i \in R_1(j, s)} (y_i - \hat{y}_{R_1})^2 + \sum_{i \mid x_i \in R_2(j, s)} (y_i - \hat{y}_{R_2})^2$$

Rinse and repeat



Pruning



Weakest Link Pruning

Also called Cost complexity pruning

For every α , there is a subtree T that minimizes:

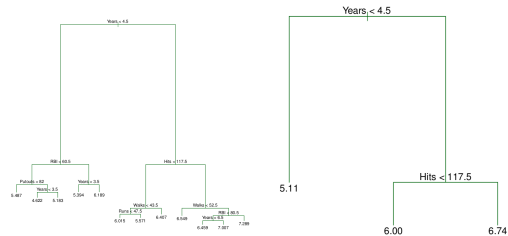
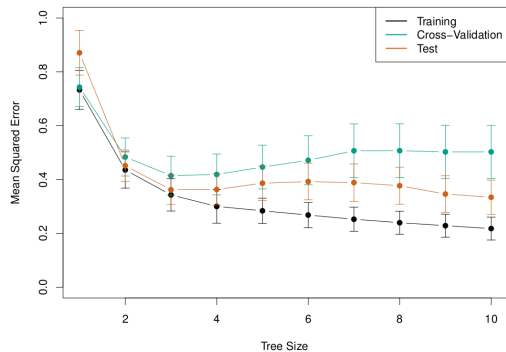
$$\sum_{m=1}^{|T|} \sum_{i|x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \alpha |T|$$

- $|T|$ = number of terminal nodes of T
- R_m is rectangle for m th terminal node
- \hat{y}_{R_m} is mean of training observations in R_m

Algorithm 8.1 *Building a Regression Tree*

1. Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
 2. Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of α .
 3. Use K-fold cross-validation to choose α . That is, divide the training observations into K folds. For each $k = 1, \dots, K$:
 - (a) Repeat Steps 1 and 2 on all but the k th fold of the training data.
 - (b) Evaluate the mean squared prediction error on the data in the left-out k th fold, as a function of α .Average the results for each value of α , and pick α to minimize the average error.
 4. Return the subtree from Step 2 that corresponds to the chosen value of α .
-

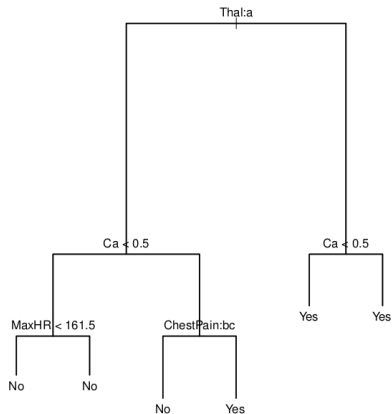
Messing with α



Section 2

Classification Decision Tree

Basic idea



- \hat{p}_{mk} = proportion of training observations in R_m from the k th class
- $E = 1 - \max_k(\hat{p}_{mk})$

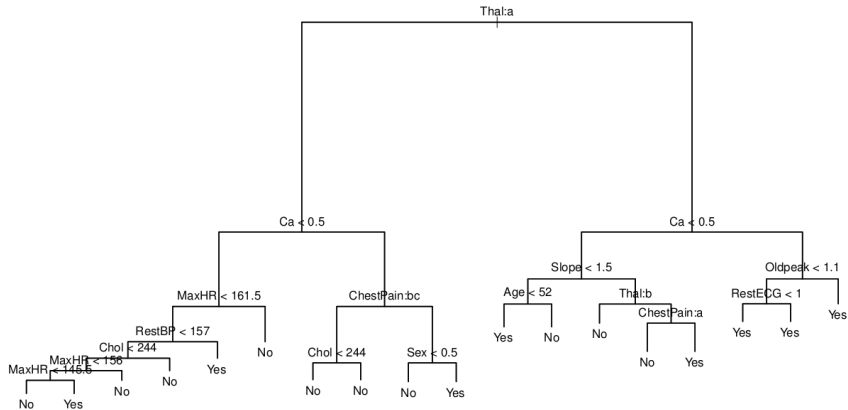
Gini index

$$G = \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$$

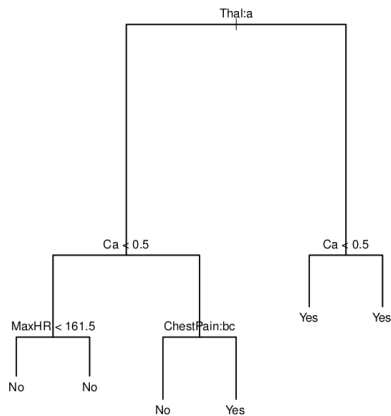
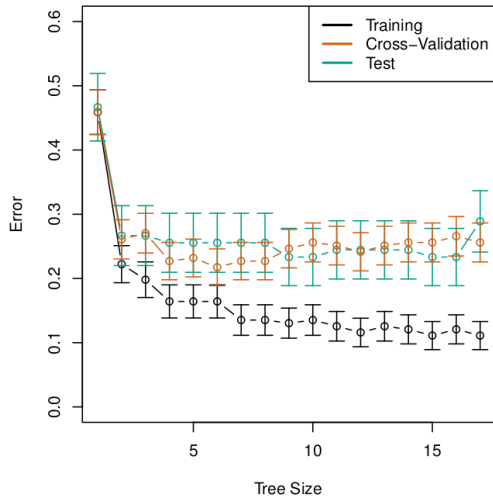
Entropy

$$D = - \sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$$

Example

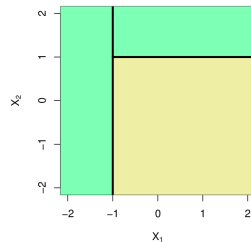
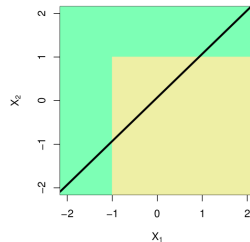
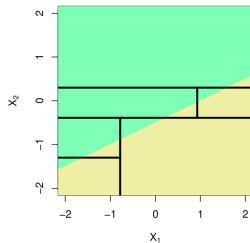
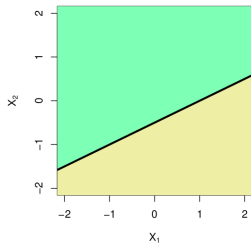


Pruning the example



More coding!

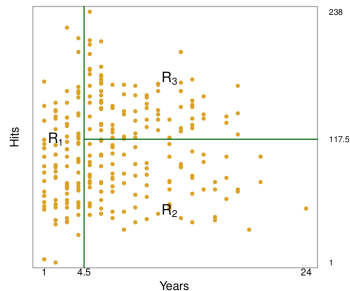
Linear models vs trees



Pros:

Cons:

- Split into regions by greedily decreasing RSS
- Prune tree by using cost complexity
- Not robust - Next time, figure out how to aggregate trees



Next time

20	F	Nov 4	Polynomial & Step Functions.	7.1,7.2	
21	M	Nov 7	Step Functions	7.2	
22	W	Nov 9	Basis functions, Regression Splines	7.3,7.4	
23	F	Nov 11	Decision Trees	8.1	HW #7 Due
24	M	Nov 14	Ensemble methods	8.2	
25	W	Nov 16	Maximal Margin Classifier	9.1	
26	F	Nov 18	SVC	9.2	HW #8 Due
27	M	Nov 21	SVM	9.3, 9.4, 9.5	
28	W	Nov 23	Single layer NN	10.1	
	F	Nov 25	No class - Thanksgiving		
29	M	Nov 28	Multi Layer NN	10.2	HW #9 Due
30	W	Nov 30	CNN	10.3	
31	F	Dec 2	Unsupervised Learning & Clustering	12.1, 12.4	
32	M	Dec 5	More Clustering	12.4	HW #10 Due
	W	Dec 7	Review		
	F	Dec 9	Midterm #3	Bring your cheat sheet and a non-internet-connected calculator	