

TREE-BASED METHODS (ch 8)

12/4/2021

Non-linear methods

Main idea: Divide the feature space into (simple) regions (rectangles/boxes)

The predicted outcome (\hat{y} or class) is a function of the region each observation belongs to

→ stratification or segmentation of the feature space

The splitting process has a tree structure.

REGRESSION TREES: MOTIVATING EXAMPLE

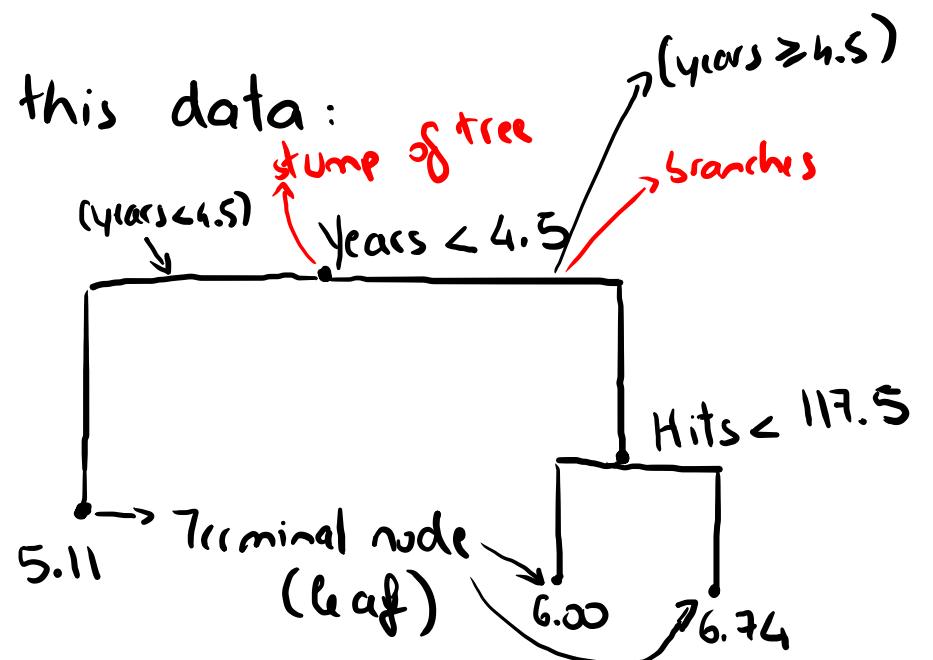
Predicting next year salary of baseball players using this year's statistics:

Predictor 1: # years played in major leagues

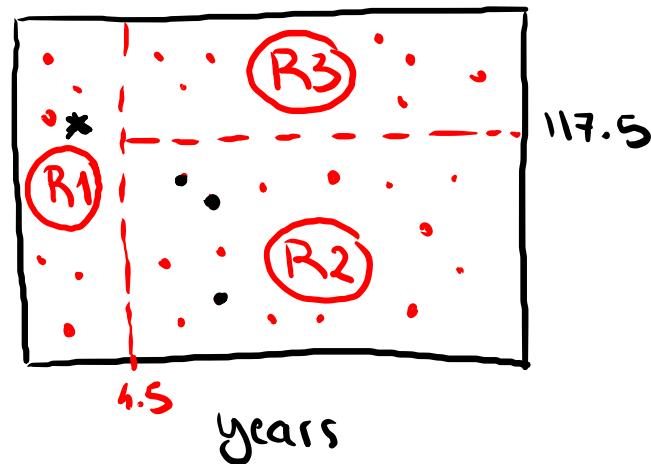
Predictor 2: # hits

Response: salary (next season)

Decision tree fitted to this data:



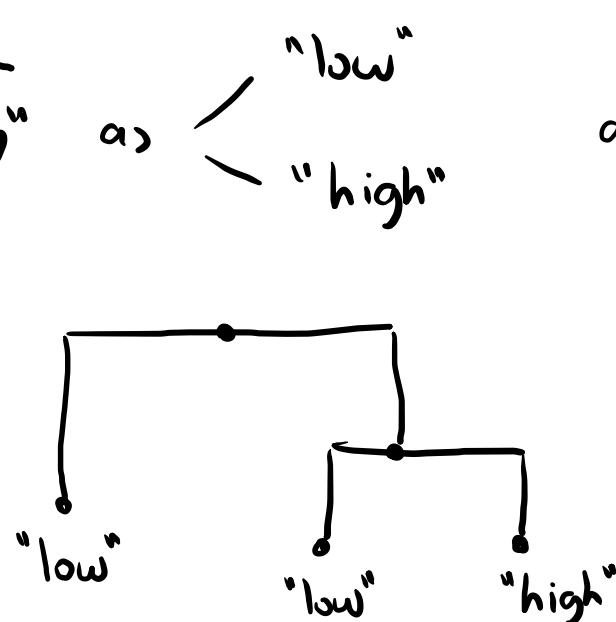
hits



3 terminal leaves \leftrightarrow 3 non-overlapping regions

CLASSIFICATION TREE: EXAMPLE

Predicting response variable "Salary" as
according to some threshold (eg mean 5.93)



Let's look at the details of how one fits a tree to data :

- ① Dividing the space into J non-overlapping regions
- ② Predicting a value of the response y in each region

Step 2 : Assume we have a partition with J leaves / regions: R_1, \dots, R_J

Then:

Regression : Given a new observation \underline{x} , consider the region R_j it falls into and predict \hat{y} by the average y for all training observations in region R_j

Classification Given a new observation \underline{x} , consider the region R_j it falls into and use for classification

$$\hat{p}(c | R_j) = \text{proportion of training observations in region } R_j \text{ that belong to class } c \quad (\text{similar to K-NN})$$

Use these probabilities within a Bayes classifier, e.g. assign \underline{x} to class with highest probability (threshold 0.5)

Step 1 HOW TO BUILD A TREE

Aim : Divide the feature space into boxes

Idea : Choose regions so as to minimize some optimality criterion

Two Key aspects :

- ① Choice of optimality criterion
- ② Find a computationally efficient approach

CHOICE OF OPTIMALITY CRITERION

In regression , the most natural choice is RSS :

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

↑
mean response in region R_j

In classification, a natural choice is

$$E = \sum_{j=1}^J E_j \quad \text{where } E_j = 1 - \max_c \hat{p}(c|R_j)$$

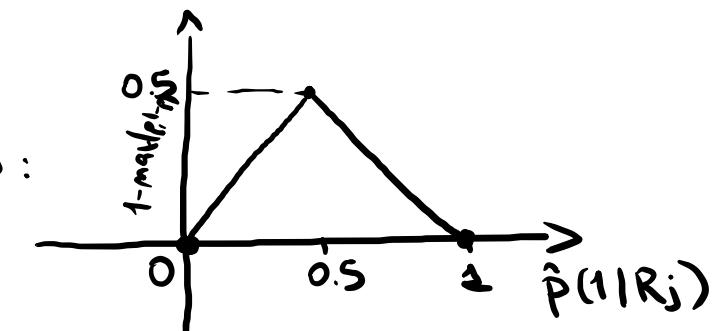
$\frac{\downarrow}{\text{fraction of observations in region } R_j \text{ that are misclassified}}$ $\underbrace{\text{probability associated to the predicted class (0.5 threshold)}}$

This function is not very "smooth" for optimisation purposes:

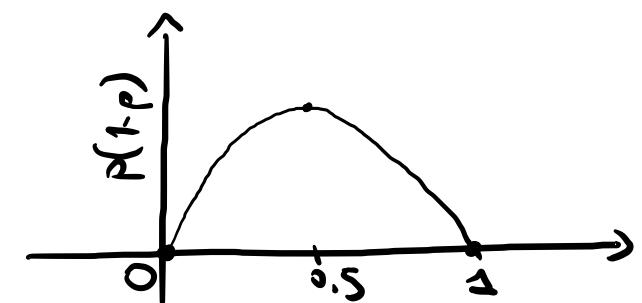
There are two alternatives:

Sini Index : $G_j = \sum_{c=1}^k \hat{p}(c|R_j) (1 - \hat{p}(c|R_j))$

variance of Bernoulli

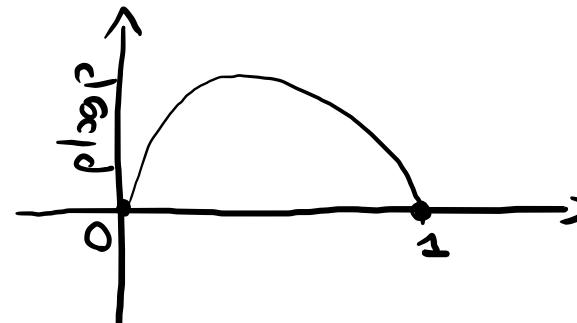


$$E_j = 1 - \max(\hat{p}(1|R_j), 1 - \hat{p}(1|R_j))$$



Entropy (or cross-entropy)

$$D_j = - \sum_{c=1}^k \hat{p}(c|R_j) \log \hat{p}(c|R_j)$$



RECURSIVE BINARY SPLITTING

- top-down approach : start with the whole space as a single region
- binary : at each step, one region is selected and split into two regions
- greedy : at each step, the best split is made without looking ahead (local rather than global optimum)

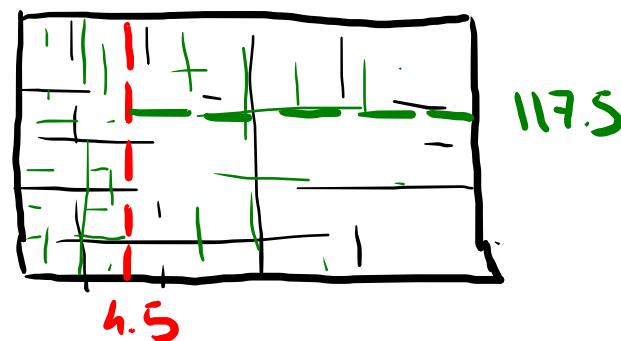
ALGORITHM

- ① For each predictor X_k and all its cut points s , define two regions :
- $$R_-(k, s) = \{x : x_k < s\}, R_+(k, s) = \{x : x_k \geq s\}$$

Compute the total optimality by summing up the optimality of R_- and R_+ .

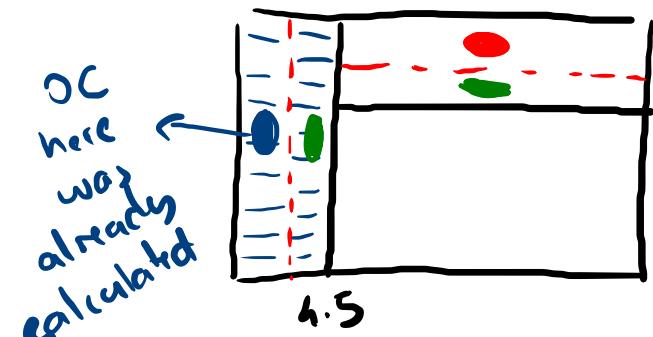
- Remarks :
- For categorical variables, the splitting divides the categories into 2 groups
 - For continuous variables, a selection of cutpoints is made (efficiently)

- ② choose (K, s) that minimizes the optimality criterion
(RSS, ζ_{ini} / Entropy)



- ③ Repeat the procedure inside each of the two regions, getting a new region
- ④ Continue until some stopping criterion is reached, e.g. maximal number of observations in each region ...

Computational tricks



Big issue with trees : Very high variance , predict very well training data but not as well the test data (overfitting)

Some solutions :

- ① Fit a smaller size tree (lower variance but typically too high bias)
- ② Stop dividing when reduction in RSS / Gini/ Entropy is below some (high) threshold
(problems because of the greediness of the algorithm → an "unimportant" split early or could result in good splits later on)
- ③ Build a big tree and then prune it back , ie remove branches

HOW TO PRUNE A TREE

OPTIMAL WAY : consider all subtrees

- calculate optimality criterion on a validation set or via CV for each subtree
- choose the best one

Computationally too expensive.

Solution : COST-COMPLEXITY PRUNING

Consider a sequence of subtrees indexed by a tuning parameter $\alpha (\geq 0)$

ALGORITHM

① Use recursive binary splitting to build a big tree (eg 5 observation maximum in each region)
Denote this with T_0

② For each value of α (in a pre-selected sequence) choose subtree $T_\alpha \subset T_0$ such that

$$T_\alpha = \underset{T \subset T_0}{\operatorname{argmin}}$$

$$\sum_{j=1}^{|T|} \text{OC}(R_j) + \alpha |T|$$

↑ number of terminal nodes in T

penalise too complex subtrees

objective function

Choice of α :

- $\alpha = 0 \rightarrow T_0$
- As α increases, the tree is pruned in a nested way
- choice of α corresponds to finding optimal bias-variance trade-off

CV for choice of α (K-fold)

Assume we have a sequence of α values = $(\alpha_1, \dots, \alpha_n)$

- ① Split the data in K folds : choose $K-1$ folds as training data and the remaining one as test data.
- ② For each fold and for each α : use recursive binary splitting on training data to build T_α and then select subtree T_α . Measure error on test data.
- ③ Compute average error across folds for each α :
 $CV(\alpha_1), CV(\alpha_2), \dots, CV(\alpha_n)$
- ④ Choose α_{opt} as the one that minimises CV error.
- ⑤ Typically, one would refit the optimal model on full data.