# Chapter 8: Tree-Based Methods

*→ non parametric supervised.*

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

*→ quantitative response Y* *→ categorical response Y.*

*These involve segmenting the predictor space, into a number of simple regions*
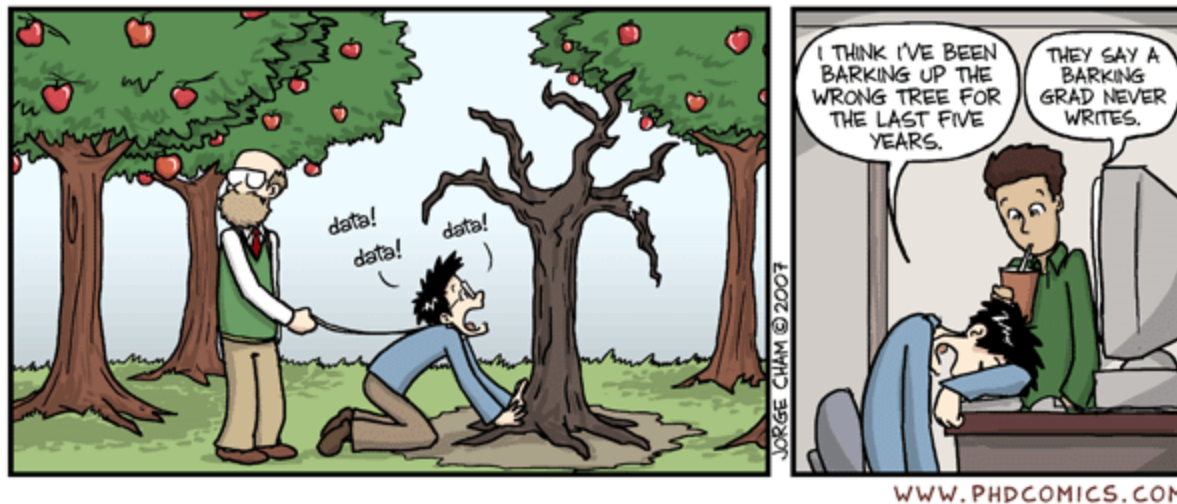*↳ all values $X_1, \ldots, X_p$ can take.*

*To make a prediction for an observation, we use the mean or mode of training observations in the region to which it belongs.*
*↓ regression* *↳ classification.*

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

*— Simple and useful for interpretation.*

*— not competitive w/ other supervised approaches (e.g. lasso) for prediction.*

*↗ boosting, bagging, random forests (later)*

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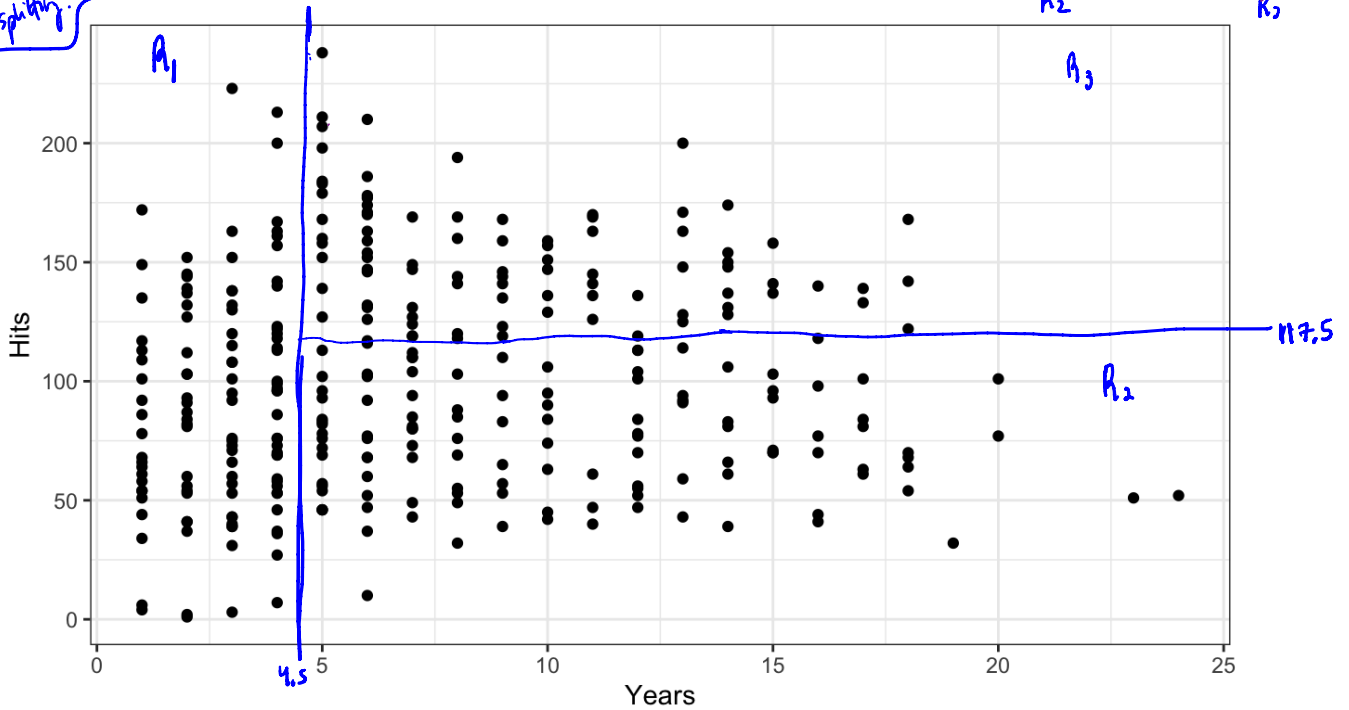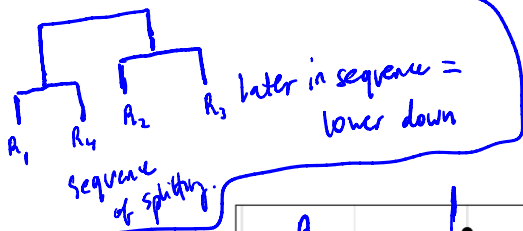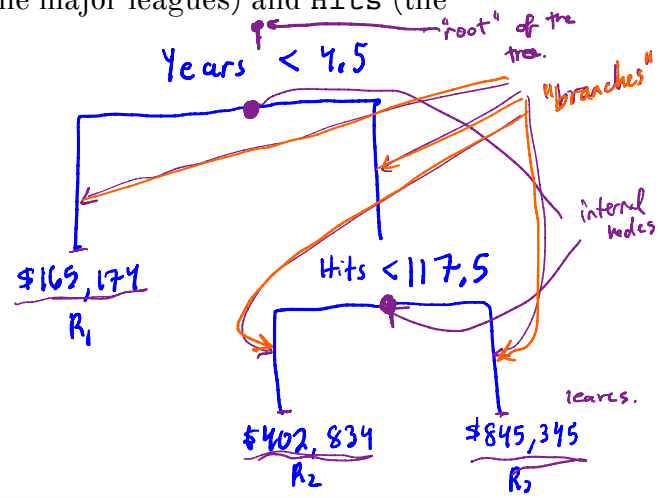


Credit: http://phdcomics.com/comics.php?f=852

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

# 1 Regression Trees

*start with*

**Example:** We want to predict baseball salaries using the `Hittlers` 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).

we can make a series of splitting rules
(according to a tree fit to this data)
to create regions and predict salary as
the mean in each region.

later in sequence =
lower down

sequence of splitting

$R_1$ $R_4$ $R_2$ $R_3$

Years < 4.5 — "root" of the tree.

"branches"

internal nodes

$165,174
$R_1$

Hits < 117.5

leaves.

$402,834
$R_2$

$845,345
$R_3$



$R_1$

$R_3$

$R_2$

117.5

4.5

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.

**Terminology** $R_1, R_2, R_3 = $ "terminal nodes" or "leaves" of the tree.

points along tree where predictor space is split = "internal nodes"
segments of tree that connect nodes = "branches"

**Interpretation** Years is most important factor in determining salary.
↳ given that a player has less experience (< 4.5 years), # hits plays very little role in his salary
↳ among player who had been in the league 5+ years, # hits does affect salary: ↑hits$_2$ ↑salary.

probably an oversimplification but it is easy to interpret & has nice graphical representation.

We now discuss the process of building a regression tree. There are $\overset{w}{\text{to}}$ steps:

$\nearrow$ quantitative response.

1. Divide <u>predictor space</u>

   $\nearrow$ set of possible values for $X_1, \ldots, X_p$

   into $J$ distinct and non-overlapping regions $R_1, \ldots, R_J$

2. Predict

   For every observation that falls into the region $R_j$ we make the same prediction

   = mean of response $Y$ for every training observation in $R_j$.

How do we construct the regions $R_1, \ldots, R_J$? How to divide the predictor space?

regions could have any shape, but that's too hard (to do & to interpret).

$\Rightarrow$ divide predictor space into high dimensional rectangles ("boxes").

The goal is to find boxes $R_1, \ldots, R_J$ that minimize the RSS. $= \sum\limits_{j=1}^{J} \sum\limits_{i \in R_j} \left( y_i - \hat{y}_{R_j} \right)^2$ where

Unfortunately it is computationally infeasible to consider every possible partition.

$\hat{y}_{R_j}$ = mean response of training data in $R_j^{\text{th}}$ box.

$\Rightarrow$ take <u>top-down</u>, <u>greedy</u> approach called <u>recursive binary splitting</u>.

The approach is *top-down* because

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

each split is indicated via two new branches down the tree.

The approach is *greedy* because

at each step in the building process, the best split is made at that particular step.

$\hookrightarrow$ not looking ahead to make a split that will lead to a better tree later.

In order to perform <u>recursive binary splitting</u>,

① Select the predictor and cutpoint $s$ s.t. splitting the predictor space into regions $\{X \mid X_j < s\}$ and $\{X \mid X_j \geq s\}$ leads to greatest possible reduction in RSS,

  ↳ region of predictor space where $X_j$ takes values $< s$.

  ↳ consider all possible $X_1, \dots, X_p$ and cutpoints $s$ (based on training data), then choose predictor & cutpoint that result in lowest RSS.

  i.e. consider all possible half planes $R_1(j,s) = \{X \mid X_j < s\}$ and $R_2(j,s) = \{X \mid X_j \geq s\}$
  we seek $j,s$ that minimize

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

  ← finding $j,s$ can be quickly done when $p$ not too large.

② Repeat process, looking for next best $j,s$ combo but instead of splitting entire space, split $R_1(j,s)$ or $R_2(j,s)$ to minimize RSS.

③ Continue until stopping criteria is met, i.e. no region contains more than <u>5</u> obs.

④ predict using mean of training observations in the region to which test observation falls.

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

because the resulting tree may be too complex.

  ↳ less regions $R_1, \dots, R_J$

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

<u>Idea</u>: only split the tree if it resulted in "large enough" drop in RSS.
  ↑

bad idea because a seemingly "worthless" split early in the tree might be followed by a good split later (large drop in RSS).

A strategy is to grow a very large tree $T_0$ and then *prune* it back to obtain a *subtree*.

**Better idea**

How to prune the tree?
<u>goal</u>: select a subtree that leads to lowest test error rate. → could use CV to estimate the error for every possible subtree
<u>solution</u>: "cost complexity pruning" aka "weakest link pruning". This is expensive! (large # of potential subtrees).

Consider a sequence of trees indexed by a nonnegative tuning parameter $\alpha$
For each value of $\alpha$, ∃ a corresponding subtree $T \subset T_0$ s.t.

  penalizing complexity in the tree

$$\underset{\substack{\text{# of} \\ \text{terminal} \\ \text{nodes in} \\ \text{tree } T}}{|T|} \quad \sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \hat{y}_{R_m})^2 + \underline{\alpha |T|} \text{ is as small as possible.}$$

  $R_m = m^{th}$ terminal node region
  $\hat{y}_{R_m} =$ predicted response in $R_m$

when $\alpha = 0$
$T = T_0$

$\alpha \uparrow \Rightarrow$ pay for having many terminal nodes ⇒ smaller subtree.

$\alpha$ controls trade off between subtree complexity & fit to training data

Select $\alpha$ via CV, then use full dataset & chosen $\alpha$ to get subtree $T$.

Algorithm for building a regression tree:

① Use recursive binary splitting to grow a large tree on training data, stopping only when each terminal node has fewer than some Min. # of observations.

② Use cost complexity pruning to get a sequence of best trees as a function of $\alpha$.

③ Use k-fold CV to choose $\alpha$

    Divide training data into K folds, for each $k = 1, --, K$

    (a) Repeat ① and ② on all data but $k^{th}$ fold.

    (b) evaluate the predicted MSE on $k^{th}$ fold as a function of $\alpha$.

    Average results for each value of $\alpha$ and pick $\alpha$ that minimizes CV error.

④ Return the subtree from ② that corresponds to chosen $\alpha$ from ③.

---

Example: Fit regression tree to ⌄subset of Hitters using 9 features. predicting salary.

① is the large tree

② CV error to estimate test MSE as a function of $\alpha$

③ subtree selected.



most important variable

deviance

apply cost-complexity pruning

size of tree (one-to-one relationship w/ $\alpha$)

2 inner node subtree.

# 2 Classification Trees

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

Recall for regression tree, the predicted response for an observation is given by the mean response of the training observations that belong to the same terminal node.

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 mode

We are also often interested in the class prediction proportions that fall into each terminal node.

$\hookrightarrow$ this can give us some idea of how reliable the prediction is:

e.g. terminal node | training data in terminal node is 100% Class 1 | vs. | 55% Class 1 45% Class 2 | "node purity"

both terminal nodes will predict as "Class 1"

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

Use recursive binary splitting to grow a classification tree.

But RSS cannot be used as criterion for splitting.

Instead, natural alternative is classification error rate.

= fraction of training observations that do not belong to the most common class.

$$= 1 - \max_{k} \{ \hat{p}_{mk} \}$$

proportion of training observations in the $m^{th}$ region from $k^{th}$ class.

It turns out that classification error is not sensitive enough to use for growing the tree.

preferred measures: $\qquad$ as splitting criteria

more sensitive to node purity than classification error rate.

① Gini index $\quad G = \sum_{k=1}^{K} \hat{p}_{mk} (1 - \hat{p}_{mk})$ measure of total variance across K classes.

$\hookrightarrow$ will take small values if all $\hat{p}_{mk}$'s are close to zero or one $\Rightarrow$ good measure of node purity, $\downarrow G \Rightarrow$ nodes contain primarily obs. from 1 class.

② Entropy $\quad D = -\sum_{k=1}^{K} \hat{p}_{mk} \log \hat{p}_{mk}$

$\hookrightarrow$ will take small values if $\hat{p}_{mk}$ close to 0 or 1 $\Rightarrow$ $\downarrow D$ when nodes are more pure.

note: neither Gini nor entropy work well w/ unbalanced classes in data.

There are other options out there to split on.

Gini and entropy are actually quite similar.

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

minimize. $\left( \begin{array}{c} \text{measure} \\ \text{of distance} \\ \text{in data} \end{array} + \alpha |T| \right)$

Any of 3 methods ( classification error, Gini, or entropy ) can also be used for pruning

But if prediction accuracy of final pruned tree is the goal, classification error rate should be used.
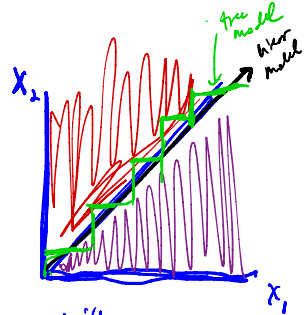
# 3 Trees vs. Linear Models

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

eg. linear regression: $f(x) = \beta_0 + \sum_{j=1}^{p} x_j \beta_j$

regression tree: $f(x) = \sum_{m=1}^{M} c_m \mathbb{I}(x \in R_m)$.

where $R_1, \ldots, R_M$ are partitions of the feature space.

Which method is better? It depends on the problem.

- If the true relationship between features and response is approximately linear, will out perform a regression tree.

- If highly nonlinear and complex relationships, decision trees may be better.

Also, trees may be preferred because of interpretation and visualization.

## 3.1 Advantages and Disadvantages of Trees

Advantages

- easy to explain, even easier than linear regression.

- Some people think decision trees more closely mirror human decision making.

- Can be displayed graphically, easy to interpret for non expert (especially if small).

- Can handle categorical predictors without creating dummy variables.

Disadvantages

- do not have same level of predictive performance as other methods we've seen.

- Not robust: small changes in data can have large changes in estimated tree (high variability)

$\downarrow$

We can aggregate many trees to try and improve this! (Next).

7

# 4 Bagging

Decision trees suffer from *high variance.*

i.e. if we split data in half (randomly) and fit decision trees to each half, results could be quite different

vs. low variance will yield similar results if applied repeatedly to distinct datasets.

$\hookrightarrow$ linear regression is low variance $n >> p$.

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

Recall: for a given set of $n$ independent observations $Z_1, ..., Z_n$ each w/ variance $\sigma^2$,

$$Var\left(\bar{Z_n}\right) = Var\left(\frac{1}{n}\sum_{i=1}^{n}Z_i\right) \stackrel{indep}{=} \frac{1}{n^2}\sum_{i=1}^{n}Var Z_i = \frac{1}{n^2}\cdot n \cdot \sigma^2 = \frac{\sigma^2}{n}$$

i.e. averaging a set of observations reduces variance.

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.

i.e. take $B$ training sets

calculate $\hat{f}^1(x), \hat{f}^{(2)}(x), ..., \hat{f}^{(B)}(x)$

obtain a low variance statistical learning model

$$\hat{f}_{AVG}(x) = \frac{1}{B}\sum_{b=1}^{B}\hat{f}^b(x).$$

Of course, this is not practical because we generally do not have access to multiple training sets. Collecting training data can be expensive.

Instead we could take repeated samples (w/ replacement) from the training data set (these are called "bootstrapped training data sets" because we are bootstrapping samples from the population using only one training data set, i.e. "pulling ourselves up from our bootstraps")

$\hookrightarrow$ assuming empirical distribution in sample is similar to population dsn, i.e. we have a representative sample.

Then we could train our method on $b^{th}$ bootstrapped training data set to get $\hat{f}^{*b}(x)$ and avg:

$$\hat{f}_{bag}(x) = \frac{1}{B}\sum_{i=1}^{B}\hat{f}^{*b}(x).$$

8

this is called bagging, short for bootstrap aggregation.

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

To apply bagging to regression trees, ① construct B regression trees using B bootstrapped data sets

② average resulting predictors

These trees are grown deep and not pruned.

⟹ each tree have low bias & ↑ variance.

averaging trees reduces variance by combining hundreds or thousands of trees!

↳ won't lead to overfitting, but can be slow.

How can bagging be extended to a classification problem?

For a given test observation, record class prediction from each tree and take majority vote: overall prediction is the class that occurs most often.

# 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.

Key: trees are fit to bootstrapped subsets of observations.

⟹ on average each tree uses ≈ $\frac{2}{3}$ of the data to fit the tree

— has to do w/ probability of being selected in a bootstrap sample as $B \uparrow \infty$

i.e. ≈ $\frac{1}{3}$ of observations on not used to fit the tree (out-of bag OOB observations).

idea: predict the response for $i^{th}$ observation using all the trees in which that observation was OOB.

This will lead to ≈ B/3 predictions for $i^{th}$ observation.

Then average (or majority vote) these predictions to get single OOB prediction for $i^{th}$ observation.

We can then get OOB prediction for each training observation to get OOB MSE (OOB classification error), which is an estimate of test error!

because we only use predictions from trees that didn't use those data points in the fitting.

## 4.2 Interpretation

Bagging typically result in improved accuracy in predictions over a single tree.

But it can be difficult to <u>interpret</u> the resulting model!

⤷ one of the biggest advantages of trees $\frac{1}{4}$.

⤷ no longer represent model using a single tree.

⟹ no longer clear which variables are the <u>most important</u> to predict the response!

Bagging improves prediction at the cost of interpretability.

What can we do?

We can obtain an overall summary of the importance of each predictor using RSS (or Gini index)

— record total amount RSS (or Gini) is decreased due to splits for a given predictor, averaged over B trees.

— large values indicates an important predictors.

# 5 Random Forests

*Random forests* provide an improvement over bagged trees by a small tweak that decorre-lates the trees.

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

But when building trees, a random sample of $m$ predictors is chosen a split candidates from the full set of predictors.

$\quad \hookrightarrow$ each split is allowed only to use those chosen predictors

$\quad \hookrightarrow$ fresh sample of predictors taken at each split.

$\quad \hookrightarrow$ typically $m \approx \sqrt{p}$

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. Why?

Suppose there is one strong predictor in the data set and a number of moderately strong predictors. In the collection of trees, most or all will use the strong predictors as the top split.

$\quad \Longrightarrow$ all of of the bagged trees will look quite similar.

$\quad \Longrightarrow$ predictions will be highly correlated.

(bagging) and averaging highly correlated values does not lead to much variance reduction!

Random forests overcome this by forcing each split to consider a subset of predictors.

$\quad \Longrightarrow$ on average $\frac{(p-m)}{p}$ of the splits will not even consider the strong predictor $\Longrightarrow$ other predictors will have higher chance of being split on.

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

If $m = p \Longrightarrow$ random forest = bagging.

Using small $m$ will typically help when we have a lot of correlated predictors.

$\quad -$ As with bagging, we will not have overfitting w/ large B

$\quad -$ And we can examine variable importance in the same way.

11

# 6 Boosting *✻ very popular right now ( Adaboost and XG Boost).*

*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 boostrap sampling, instead each tree is fit on a modified version of the original data set.

Boosting has three tuning parameters:

1.

2.

3.