# Supervised Learning:
# Trees, Bagging, Random Forests

Noah Simon & Daniela Witten

July 15-17, 2015
Summer Institute for Statistics of Big Data
University of Washington

---

## Decision Trees

- A completely different approach to non-linear modeling!
- Partition the range of $X$'s into boxes.
- Then make a prediction for $Y$ within each box.
- Leads to very easily interpretable models . . . though with some loss of prediction accuracy compared to other approaches.

---

## Baseball Data

Want to predict a baseball player's Salary based on various characteristics of his past performance.

---
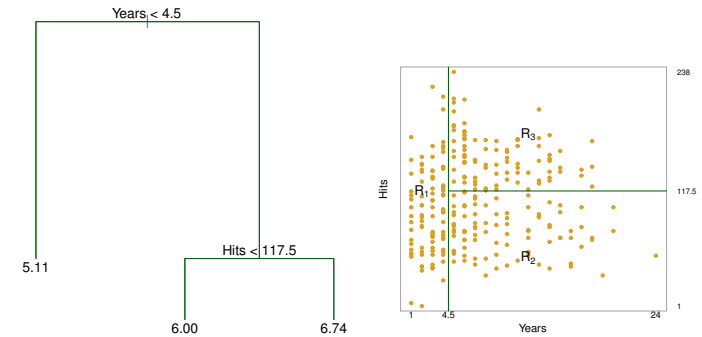
## A Regression Tree for Baseball Data

# Regression Trees

1. Partition the range of $X$ into $M$ disjoint regions, $R_1, \ldots, R_M$.

2. Within the $m$th region, predict $Y$ as

$$\hat{f}(X) = \sum_{m=1}^{M} c_m I(X \in R_m)$$

where $c_m$ is the mean response value among the training observations in $R_m$.
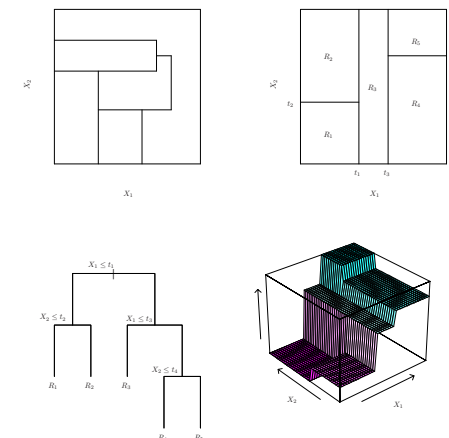
# A Regression Tree for Baseball Data

# Recursive Binary Splitting

Q: How do we determine the regions $R_1, \ldots, R_M$?
A: Recursive binary splitting.

# Recusive Binary Splitting

# How Do We Get the Regions?

- Seek the variable $j$ and cut-point $s$ that minimizes

$$\min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2,$$

  where $R_1(j,s) = \{X \mid X_j \leq s\}$ and $R_2(j,s) = \{X \mid X_j > s\}$.

- For any choice of $j$ and $s$, we get

$$\hat{c}_1 = \mathrm{Ave}(y_i \mid x_i \in R_1(j,s)), \qquad \hat{c}_2 = \mathrm{Ave}(y_i \mid x_i \in R_2(j,s)).$$

- For any variable $j$, computing the split point $s$ is very fast. So scanning through the variables to find the best $(j,s)$ is fast.
- Once we find the best $(j,s)$, we partition the data into the two regions, and repeat splitting in each region.
- Continue along these lines.

# How Do We Make Predictions?

- Consider a partition into $M$ regions, $R_1, \ldots, R_M$.
- We model the response as a constant within each region,

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

- To minimize sum of squared errors, the best $\hat{c}_m$ is

$$\hat{c}_m = \mathrm{Ave}(y_i \mid x_i \in R_m).$$

# How Big Should We Grow the Tree?

- A large tree might overfit the data . . . high variance.
- A small tree might miss important structure . . . high bias.
- Tree size is a tuning parameter that governs model complexity.
- Optimal tree size should be chosen adaptively from the data.
- Tempting to grow the tree until the decrease in sum-of-squares due to further splits exceeds some threshold. But this strategy is too short-sighted.
- Instead, we build a very big tree, and then prune it down.

# Cost-Complexity Pruning, Part I

- Define a subtree $T \subset T_0$ to be a tree obtained by pruning $T_0$.
- The $m$th terminal node represents region $R_m$.
- Let $|T|$ denote number of terminal nodes in $T$.
- Pruning the tree is complicated because there are so many subtrees!
- Need a way to navigate tree spacee: cost-complexity pruning.
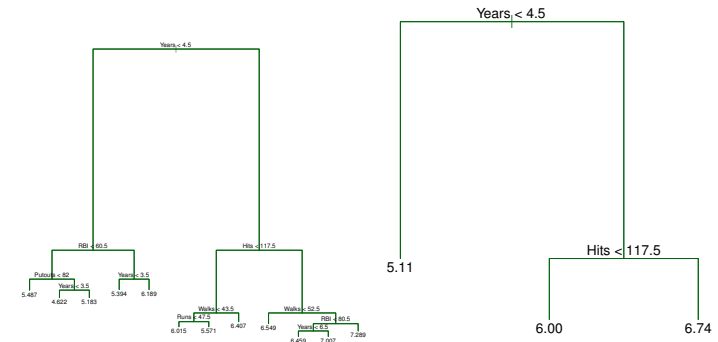
# Cost-Complexity Pruning, Part II

- Let $N_M = \#\{x_i \in R_m\}$, $\hat{c}_m = \frac{1}{N_m} \sum_{x_i \in R_m} y_i$, and $Q_m(T) = \frac{1}{N_m} \sum_{x_i \in R_m} (y_i - \hat{c}_m)^2$.
- Define the cost complexity criterion as

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|.$$

- For each $\alpha$, find the subtree $T_\alpha \subset T_0$ that minimizes $C_\alpha(T)$.
- The tuning parameter $\alpha \geq 0$ governs the trade-off between tree size and its goodness of fit to the data.
  - When $\alpha = 0$ we get the full tree.
  - When $\alpha$ is large we get a pruned tree.
- Choose $\alpha$ by cross-validation.

---

# Pruned Tree for Baseball Data

---

# Classification Trees, Part I

- Suppose $Y$ is qualitative with values in $1, 2, \ldots, K$.
- Need to adjust criteria for splitting nodes and pruning tree.
- For the $m$th node, representing the region $R_m$ containing $N_m$ observations, define

$$\hat{p}_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k),$$

the proportion of class $k$ observations in node $m$.
- Classify the observations in node $m$ to the majority class,
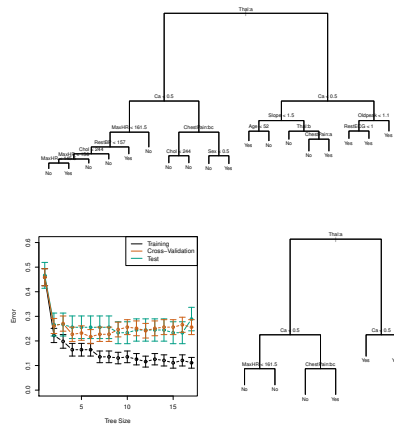
$$k(m) = \operatorname{argmax}_k \hat{p}_{mk}.$$

---

# Classification Trees, Part II

- Three possible measures of node impurity:
  - Misclassification Error: $1 - \hat{p}_{mk(m)}$.
  - Gini Index: $\sum_{k=1}^{K} \hat{p}_{mk}(1 - \hat{p}_{mk})$.
  - Cross-Entropy or Deviance: $-\sum_{k=1}^{K} \hat{p}_{mk} \log \hat{p}_{mk}$.
- CE and Gini are more sensitive to changes in node probabilities, and hence are a better choice than MCE for growing the tree.
- Misclassification error is typically used for tree pruning.

# Classification Tree

Predict presence of heart disease for 303 patients with chest pain.
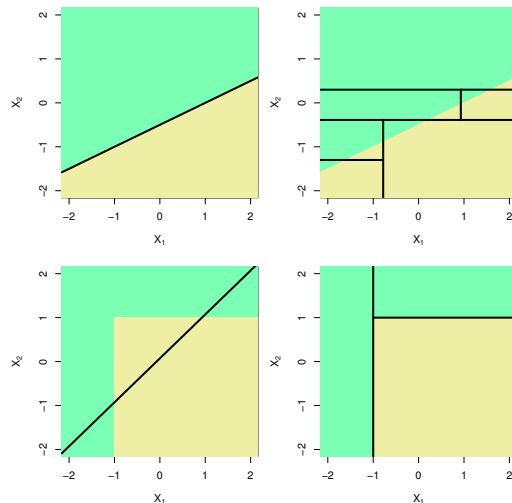
---

# Advantages & Disadvantages of Trees

- Easy to interpret — nice visual!
- Thought to model human decision-making.
- Tends to have lower accuracy than other approaches, like generalized additive models.
- Can suffer from extremely high variance — a small change in data can result in a very different series of splits.
- Difficulty in capturing additive structure: e.g.

$$Y = c_1 I(X_1 < t_1) + c_2 I(X_2 < t_2) + \epsilon.$$

- Accuracy can be improved using bagging and boosting.

---

# Trees Versus Linear Models

---

# Bootstrap Aggregation

- Suppose we fit a model to training data $\mathbf{Z} = \{(x_1, y_1), \ldots, (x_N, y_N)\}$, to obtain prediction $\hat{f}(x)$ at $x$.
- Bootstrap aggregation, or bagging, averages this prediction over many bootstrap samples in order to reduce its variance.
- The bagging estimate is defined as

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

where the prediction $\hat{f}^{*b}(x)$ was obtained using the bootstrap sample $\mathbf{Z}^{*b}$, $b = 1, \ldots, B$.

# More on Bagging

- If $\hat{f}(x)$ is a linear function of the data, then you don't accomplish much . . . for $B$ sufficiently large, $\hat{f}_{\mathrm{bag}}(x)$ and $\hat{f}(x)$ will be very similar.
- However, bagging very non-linear fits — like those from a regression tree — can lead to a substantial reduction in variance, and very improved predictions.
- If $Y$ is a qualitative response, then the bagged prediction can be obtained by 1 of 2 options:
  - majority vote: $\mathrm{argmax}_k \hat{f}_{\mathrm{bag}}(x)$; or
  - averaging probabilities across the $B$ bootstrap fits.

# Bagged Trees

Elements of Statistical Learning (2nd Ed.) ©Hastie, Tibshirani & Friedman 2009 Chap 8

Elements of Statistical Learning (2nd Ed.) ©Hastie, Tibshirani & Friedman 2009 Chap 8
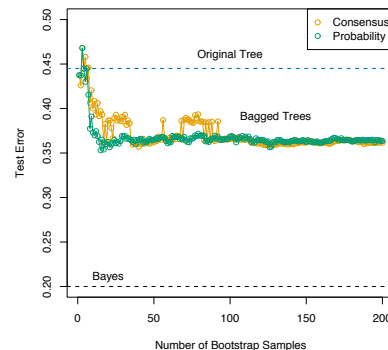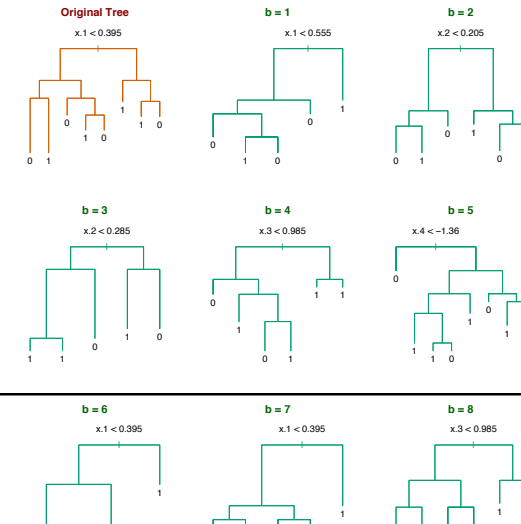
FIGURE 8.10. *Error curves for the bagging example of Figure 8.9. Shown is the test error of the original tree and bagged trees as a function of the number of bootstrap samples. The orange points correspond to the consensus vote, while the green points average the probabilities.*

# Bagging

- In regression, bagging can only help you!
- Not so for classification: bagging a good classifier can help, but bagging a bad classifier can hurt.
- Suppose $Y = 1$ for all $X$, and

$$\hat{Y}(X) = \begin{cases} 1 \text{ with probability } 0.4 \\ 0 \text{ with probability } 0.6 \end{cases}.$$

- Then the misclassification error of $\hat{Y}(X) = 0.6$, and that of the bagged classifier is 1.

## Wisdom of Crowds

- If the classifiers are not bad, then bagging can help.
- To see this, a simplified example: the wisdom of crowds.
- Suppose that in a two-class example, $Y(x) = 1$.
- Suppose we have $B$ independent weak learners, $Y_b^*(x)$, each with error rate $e_b = e < 0.5$.
- Let $S_1(x) = \sum_{b=1}^B I(Y_b^*(x) = 1)$.
- Weak learners are independent, so $S_1(x) \sim \text{Bin}(B, 1 - e)$, and therefore $\Pr(S_1(x) > B/2) \to 1$ as $B$ gets large.
- Bagged predictions are dependent, but the idea still holds.

---

## Wisdom of Crowds



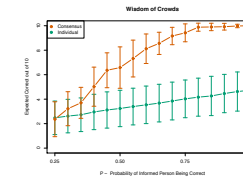Elements of Statistical Learning (2nd Ed.) ©Hastie, Tibshirani & Friedman 2009 Chap 8

**FIGURE 8.11.** *Simulated academy awards voting. 50 members vote in 10 categories, each with 4 nominations. For any category, only 15 voters have some knowledge, represented by their probability of selecting the "correct" candidate in that category (so $P = 0.25$ means they have no knowledge). For each category, the 15 experts are chosen at random from the 50. Results show the expected correct (based on 50 simulations) for the consensus, as well as for the individuals. The error bars indicate one standard deviation. We see, for example, that if the 15 informed for a category have a 50% chance of selecting the correct candidate, the consensus doubles the expected performance of an individual.*

---

## The Price of Bagging

Bagging trees can lead to improved predictions, but interpretability is lost.

---

## Out Of Bag Samples

- No need to do cross-validation to estimate prediction error for bagging.
- Instead, use out of bag (OOB) samples to compute the error.
- That is, for each observation $z_i = (x_i, y_i)$, predict the response *by taking the average of the predictions corresponding to bootstrap samples in which $z_i$ didn't appear.*
- Gives similar results to cross-validation.

# Random Forests

- Modification of bagging that builds a large collection of de-correlated trees.
- Very good predictive performance, though limited interpretability.

# Variance of an Average

- The average of $B$ uncorrelated r.v.'s, each with variance $\sigma^2$, has variance $\sigma^2/B$.
- But what about $B$ r.v.'s with pairwise correlation $\rho$?
  - The average has variance

$$\rho\sigma^2 + \frac{1-\rho}{B}\sigma^2.$$

  - As $B$ increases, the variance of the average approaches $\rho\sigma^2$.
- Therefore, bagging can only improve the variance so much if the bootstrapped predictions are correlated.
- Idea behind random forests is to reduce correlation among bootstrapped predictions.

# Random Forests

- Perform bagging on trees, but when building each tree, limit yourself to $m \leq p$ variables for each split.
- Then,
  - average the resulting $B$ trees for regression
  - take a majority vote of the resulting $B$ trees for classification.

# More On Random Forests

- Excellent prediction; very little tuning required.
- Great "out-of-box" method.
- Loss of interpretability relative to trees.
- Recommendations of the inventors:
  - For regression, take $m = p/3$, and minimum node size 5.
  - For classification, take $m = \sqrt{p}$, and minimum node size 1.

  However, better to treat these as tuning parameters. Select using OOB error estimate.
- Using more trees usually won't overfit (though this can depend on the size of the trees used). Often a few hundred to a few thousand trees are used.

**FIGURE 15.4.** OOB *error computed on the* spam *training data, compared to the test error computed on the test set.*

## Why Don't We Perform Random Forests for Other Types of Estimators?

- We typically perform random forests for trees, but the same idea could be applied more generally.
- But not all estimators can be improved by shaking up the data like this.
- Highly non-linear estimators, like trees, benefit the most.
- Doesn't really help linear estimates.

## Let's Try It Out in R!

Chapter 8 R Lab
www.statlearning.com