# Ch 8.2.1, 8.2.2: Bagging and Random Forests
## Lecture 25 - CMSE 381

Prof. Guanqun Cao

Michigan State University
::
Dept of Computational Mathematics, Science & Engineering

Mon, Nov 3, 2025

**Last time:**
- 8.1 Decision Trees - regression

**This lecture:**
- 8.1 Decision Trees - classification
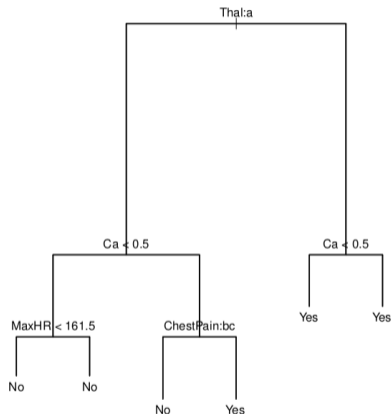- 8.2.1 Bagging
- 8.2.2 Random forest

**Announcements:**
- Homework 7 Due Sunday

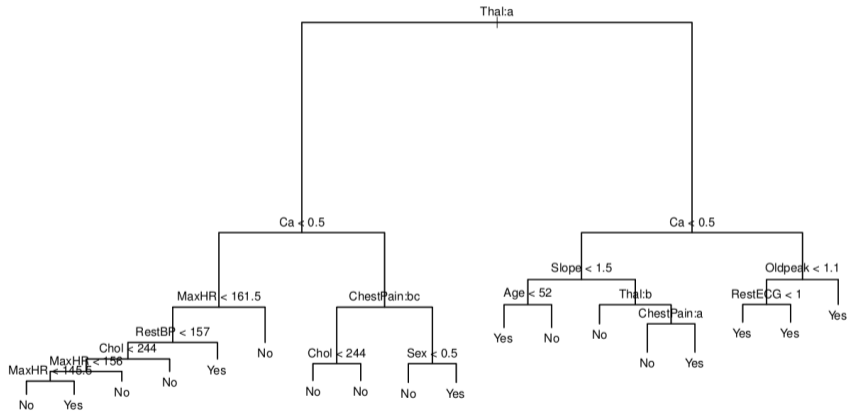| | | | | | |
|---|---|---|---|---|---|
| | F | 10/17 | *Review* | | |
| | M | 10/20 | Fall Break | | |
| | W | 10/22 | **Midterm #2** | | |
| 21 | F | 10/24 | Polynomial & Step Functions | 7.1-7.2 | HW #5 Due Sun 10/28 |
| 22 | M | 10/27 | Step Functions; Basis functions; Start Splines | 7.2-7.4 | |
| 23 | W | 10/29 | Regression Splines | 7.4 | |
| 24 | F | 10/31 | Decision Trees | 8.1 | HW #6 Due Sun 11/2 |
| 25 | M | 11/3 | Random Forests | 8.2.1, 8.2.2 | |
| 26 | W | 11/5 | Maximal Margin Classifier | 9.1 | |
| 27 | F | 11/7 | SVC | 9.2 | HW #7 Due Sun 11/9 |
| 28 | M | 11/10 | SVM | 9.3, 9.4 | |
| 29 | W | 11/12 | Single Layer NN | 10.1 | |
| 30 | F | 11/13 | Multi Layer NN | 10.2 | HW #8 Due Sun 11/16 |
| 31 | M | 11/17 | CNN | 10.3 | |
| 32 | W | 11/19 | Unsupervised learning / clustering | 12.1, 12.4 | |
| 33 | F | 11/21 | Virtual: Project Office Hours | | HW #9 Due Sun 11/23 |
| | M | 11/24 | *Review* | | |
| | W | 11/26 | **Midterm #3** | | |
| | F | 11/28 | Thanksgiving | | |
| | M | 12/1 | Virtual: Project Office Hours | | |
| | W | 12/3 | Virtual: Project Office Hours | | |
| | F | 12/5 | | | **Project Due** |

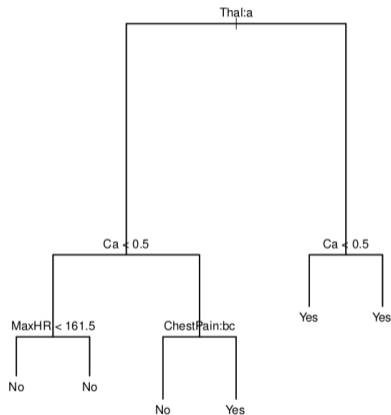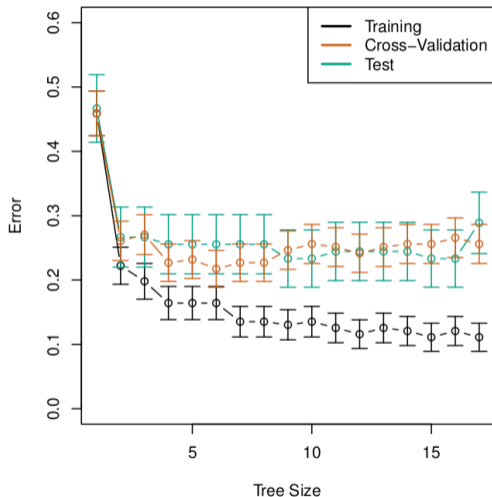Section 1

Classification Decision Tree

# Basic idea



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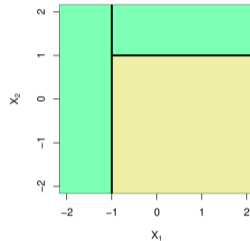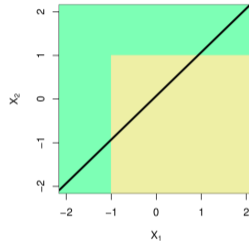- $E = 1 - \max_k(\hat{p}_{mk})$

# Example

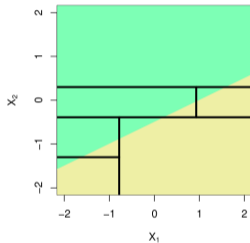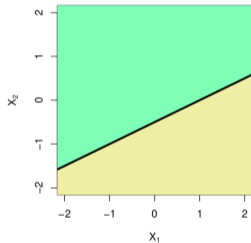# Pruning the example

# Coding!

Second part of day 24's jupyter notebook.
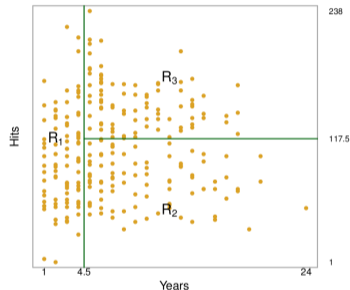
# Linear models vs trees

**Pros:**                                    **Cons:**

- Split into regions by greedily decreasing RSS (or error rate)
- Prune tree by using cost complexity
- Not robust - Next, figure out how to aggregate trees

Section 2

## 8.2.1 Bagging

## The bootstrap

**Want to do (but can't):**
Build separate models from independent
training sets, and average resulting
predictions:

- $\hat{f}^1(x), \cdots, \hat{f}^B(x)$ for $B$ separate
  training sets
- Return the average

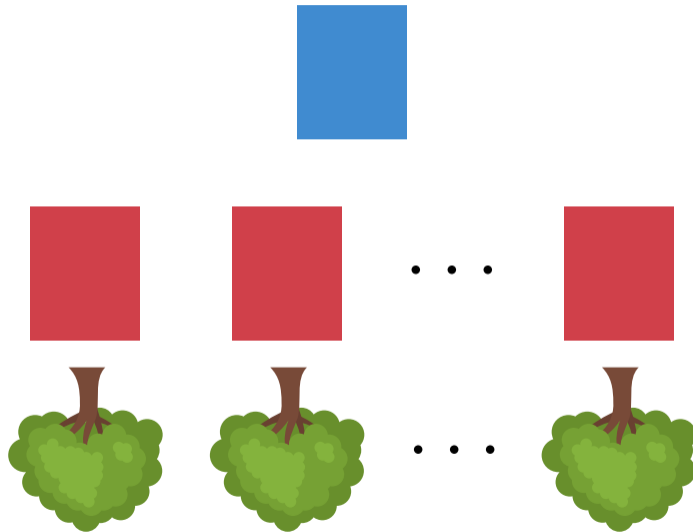$$\hat{f}_{avg}(x) = \frac{1}{B} \sum_{b=1}^{B} \hat{f}^b(x)$$

**Boostrap modification:**

- Work with fixed data set
- Take $B$ samples from this data set
  (with replacement)
- Train method on $b$th sample to get
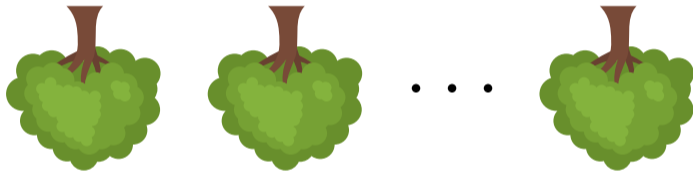  $\hat{f}^{*b}(x)$
- Return average of predictions
  (regression)

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

or majority vote (classification)

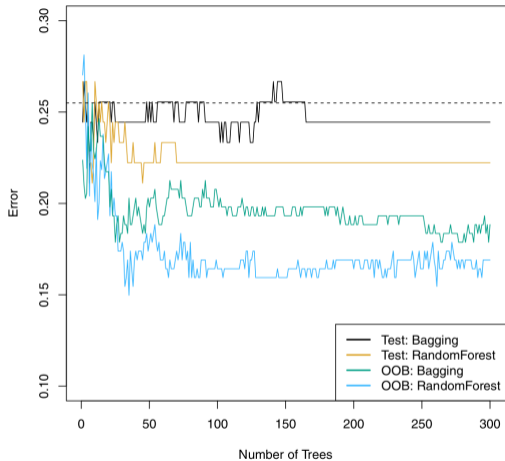# Tree version

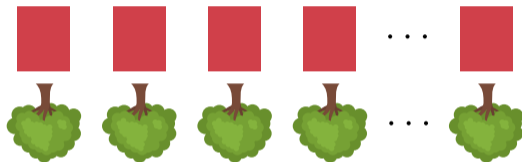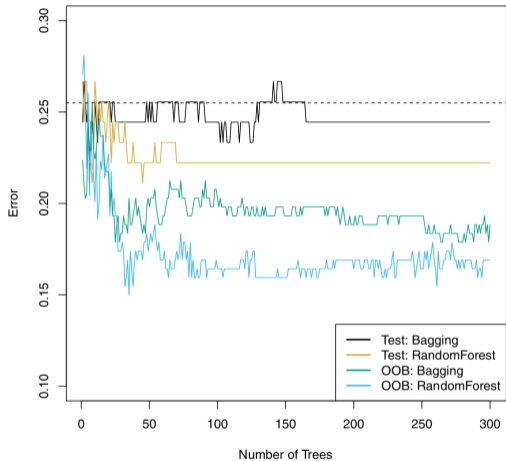# Example: Heart classification data

# Out of Bag Error Estimation

- On average, bootstrap sample uses about 2/3 of the data
- Remaining observations not used are called *out-of-bag* (OOB) observations
- For each observation, run through all the trees where it wasn't used for building
- Return the average (or majority vote) of those as test prediction

# Error using OOB

# Section 3

## Random Forests

# The idea

- Goal is to decorrelate the bagged trees:
    - ▶ If there is a strong predictor, the first split of most trees will be the same
    - ▶ Most or all trees will be highly correlated
    - ▶ Averaging highly correlated quantities doesn't decrease variance as much as uncorrelated

- The random forest fix:
    - ▶ Each time a split is considered, only use a random subset of $m$ the predictors
    - ▶ Fresh sample taken every time
    - ▶ Typically $m \approx \sqrt{p}$
    - ▶ On average, $(p - m)/p$ of splits won't consider strong predictor
    - ▶ $m = p$ gives back bagging

# Example on gene expression

# Coding time!

# TL:DR

- Bagging: trees grown independently on random samples. Trees tend to be similar to each other, can result in getting caught in local optima
- Random forest: trees independently on samples, but split is done using random subset of features

# Next time

| | | | | | |
|---|---|---|---|---|---|
| | F | 10/17 | *Review* | | |
| | M | 10/20 | Fall Break | | |
| | W | 10/22 | **Midterm #2** | | |
| 21 | F | 10/24 | Polynomial & Step Functions | 7.1-7.2 | HW #5 Due Sun 10/28 |
| 22 | M | 10/27 | Step Functions; Basis functions; Start Splines | 7.2-7.4 | |
| 23 | W | 10/29 | Regression Splines | 7.4 | |
| 24 | F | 10/31 | Decision Trees | 8.1 | HW #6 Due Sun 11/2 |
| 25 | M | 11/3 | Random Forests | 8.2.1, 8.2.2 | |
| 26 | W | 11/5 | Maximal Margin Classifier | 9.1 | |
| 27 | F | 11/7 | SVC | 9.2 | HW #7 Due Sun 11/9 |
| 28 | M | 11/10 | SVM | 9.3, 9.4 | |
| 29 | W | 11/12 | Single Layer NN | 10.1 | |
| 30 | F | 11/13 | Multi Layer NN | 10.2 | HW #8 Due Sun 11/16 |
| 31 | M | 11/17 | CNN | 10.3 | |
| 32 | W | 11/19 | Unsupervised learning / clustering | 12.1, 12.4 | |
| 33 | F | 11/21 | Virtual: Project Office Hours | | HW #9 Due Sun 11/23 |
| | M | 11/24 | *Review* | | |
| | W | 11/26 | **Midterm #3** | | |
| | F | 11/28 | Thanksgiving | | |
| | M | 12/1 | Virtual: Project Office Hours | | |
| | W | 12/3 | Virtual: Project Office Hours | | |
| | F | 12/5 | | | **Project Due** |

Q of the day:
You have two very different datasets to create two very different models.
You have to use random forest on one and bagging on the other.
Which one would benefit more from random forest? what criteria would you use for the making the decision?