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

#### Prof. Mengsen Zhang

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

Fri, Mar 28, 2025

1/23

#### 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

|    | M | 3/17 | Midterm #2  |              | Sun 3/16              |
|----|---|------|---|--------------|-----------------------|
| 21 | W | 3/19 | Polynomial & Step Functions                       | 7.1-7.2      |                       |
| 22 | F | 3/21 | Step Functions; Basis<br>functions; Start Splines | 7.2-7.4      |                       |
| 23 | M | 3/24 | Regression Splines                                | 7.4          |                       |
| 24 | w | 3/26 | Decision Trees                                    | 8.1          | HW #6 Due<br>Wed 3/26 |
| 25 | F | 3/28 | Random Forests                                    | 8.2.1, 8.2.2 | HW #7 Due<br>Sun 3/30 |
| 26 | M | 3/31 | Maximal Margin Classifier                         | 9.1          |                       |
| 27 | W | 4/2  | SVC   | 9.2          |                       |
| 28 | F | 4/4  | SVM   | 9.3, 9.4     | HW #8 Due<br>Sun 4/6  |
| 29 | M | 4/7  | Single Layer NN                                   | 10.1         |                       |
| 30 | W | 4/9  | Multi Layer NN                                    | 10.2         |                       |
| 31 | F | 4/11 | CNN   | 10.3         | HW #9 Due<br>Sun 4/13 |
| 32 | м | 4/14 | Unsupervised learning /<br>clustering             | 12.1, 12.4   |                       |
| 33 | W | 4/16 | Virtual: Project Office Hours                     |              |                       |
|    | F | 4/18 | Review  |              |                       |
|    | M | 4/21 | Midterm #3  |              |                       |
|    | W | 4/23 |   |              |                       |
|    | F | 4/25 |   |              | Project Due           |

# Section 1

## Classification Decision Tree

### Basic idea



•  $\hat{p}_{mk} =$  proportion of training observations in  $R_m$  from the kth 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





 $\mathsf{Pros}/\mathsf{Cons}$ 

Pros:



TL;DR

- 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:

- *f*<sup>1</sup>(x), ..., *f*<sup>B</sup>(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



Dr. Zhang (MSU-CMSE)

13/23

### Prediction on new data point



### 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



Number of Trees

# Coding time!

- 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

|    | М   | 3/17 | Midterm #2  |              | Sun 3/16              |
|----|-----|------|---|--------------|-----------------------|
| 21 | W   | 3/19 | Polynomial & Step Functions                       | 7.1-7.2      |                       |
| 22 | F   | 3/21 | Step Functions; Basis<br>functions; Start Splines | 7.2-7.4      |                       |
| 23 | М   | 3/24 | Regression Splines                                | 7.4          |                       |
| 24 | w   | 3/26 | Decision Trees                                    | 8.1          | HW #6 Due<br>Wed 3/26 |
| 25 | ΓF. | 3/28 | Random Forests                                    | 8.2.1, 8.2.2 | HW #7 Due<br>Sun 3/30 |
| 26 | М   | 3/31 | Maximal Margin Classifier                         | 9.1          |                       |
| 27 | W   | 4/2  | SVC   | 9.2          |                       |
| 28 | F   | 4/4  | SVM   | 9.3, 9.4     | HW #8 Due<br>Sun 4/6  |
| 29 | М   | 4/7  | Single Layer NN                                   | 10.1         |                       |
| 30 | W   | 4/9  | Multi Layer NN                                    | 10.2         |                       |
| 31 | F   | 4/11 | CNN   | 10.3         | HW #9 Due<br>Sun 4/13 |
| 32 | М   | 4/14 | Unsupervised learning /<br>clustering             | 12.1, 12.4   |                       |
| 33 | W   | 4/16 | Virtual: Project Office Hours                     |              |                       |
|    | F   | 4/18 | Review  |              |                       |
|    | М   | 4/21 | Midterm #3  |              |                       |
|    | W   | 4/23 |   |              |                       |
|    | F   | 4/25 |   |              | 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?