# Chapter 14: Introduction to Boosting

Modern Clinical Data Science
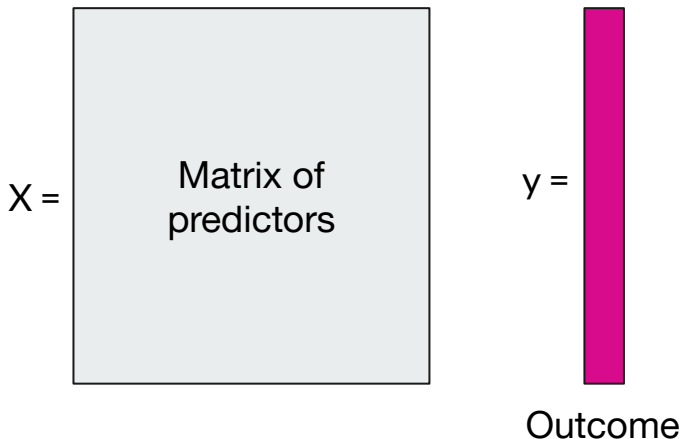Chapter Guides
Bethany Percha, Instructor

# How to Use this Guide

- Read the corresponding notes chapter first

- Try to answer the discussion questions on your own

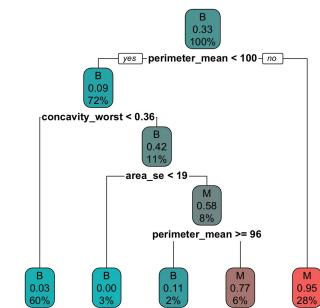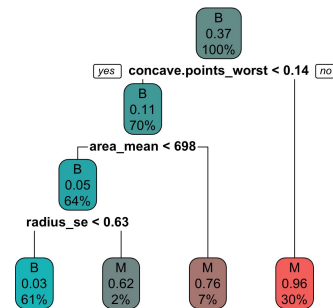- Listen to the chapter guide (should be 30 min, max) while following along in the notes

So far…

- Regression
- K-Nearest Neighbors (KNN)
- Decision trees
- Random forests

What are their stories?

X = | Matrix of predictors |

y = | Outcome |

**Ensemble learning:**
Combining multiple models that are only slightly better than random (weak learners) can produce a good model (strong learner).

**Boosting:**
A principled way of creating weak learners whose combined predictions get better over time.

# AdaBoost

First practical boosting algorithm. For a long time, no one knew why it worked.

1. Initialize the observation weights to $w_i^{(1)} = \frac{1}{N}$ for $i = 1, \ldots, N$.

2. For $m = 1, \ldots, M$:

   (a) Select a classifier, $G_m(x)$, that minimizes the weighted training error according to the current set of weights, $w_i^{(m)}$. Depending on the algorithm, it may be possible to train a single classifier on the weighted training set; in other cases, one may need to select the best-performing classifier from among a predefined set.

   (b) Compute

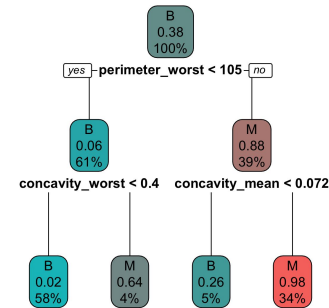   $$\text{err}_m = \frac{\sum_{i=1}^{N} w_i^{(m)} \cdot \mathcal{I}(y^{(i)} \neq G_m(x^{(i)}))}{\sum_{i=1}^{N} w_i^{(m)}}$$

   (c) Compute voting weight for classifier $m$:

   $$\alpha_m = \log\left(\frac{1 - \text{err}_m}{\text{err}_m}\right)$$

   (d) Set

   $$w_i^{(m+1)} := w_i^{(m)} \cdot \exp\left[\alpha_m \cdot \mathcal{I}(y^{(i)} \neq G_m(x^{(i)}))\right]$$

   for $i = 1, \ldots, N$.

3. Output

$$G(x) = \text{sign}\left[\sum_{m=1}^{M} \alpha_m G_m(x)\right]$$

| Subject ID | friends ($X_1$) | money ($X_2$) | free time ($X_3$) | pet ($X_4$) | happy ($Y$) |
|:---:|:---:|:---:|:---:|:---:|:---:|
| 1 | 1 | 1 | 0 | 0 | -1 |
| 2 | 1 | 1 | 1 | 0 | -1 |
| 3 | 0 | 1 | 1 | 0 | -1 |
| 4 | 0 | 0 | 0 | 0 | -1 |
| 5 | 1 | 0 | 0 | 0 | -1 |
| 6 | 0 | 0 | 0 | 0 | -1 |
| 7 | 1 | 2 | 1 | 0 | 1 |
| 8 | 1 | 0 | 1 | 0 | 1 |
| 9 | 0 | 0 | 1 | 1 | 1 |
| 10 | 1 | 0 | 0 | 1 | 1 |

$$X_1 = \begin{cases} 0 & \text{no friends} \\ 1 & \text{friends} \end{cases} \qquad X_2 = \begin{cases} 0 & \text{poor} \\ 1 & \text{enough money} \\ 2 & \text{rich} \end{cases}$$

$$X_3 = \begin{cases} 0 & \text{no free time} \\ 1 & \text{some free time} \end{cases} \qquad X_4 = \begin{cases} 0 & \text{no pet} \\ 1 & \text{has a pet} \end{cases}$$

Now, let's go through the process of applying AdaBoost to this dataset, step by step.

(a) Initialize the observation weights for the training data. (Make them uniform.)

$$w_1^{(1)} =$$

$$w_2^{(1)} =$$

$$w_3^{(1)} =$$

$$w_4^{(1)} =$$

$$w_5^{(1)} =$$

$$w_6^{(1)} =$$

$$w_7^{(1)} =$$

$$w_8^{(1)} =$$

$$w_9^{(1)} =$$

$$w_{10}^{(1)} =$$

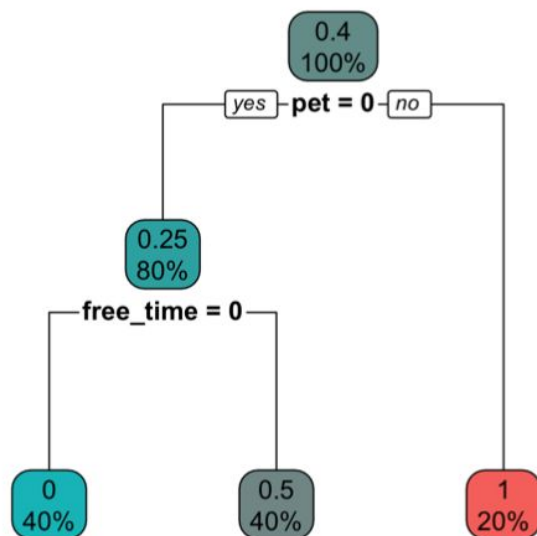Now, let's go through the process of applying AdaBoost to this dataset, step by step.

(a) Initialize the observation weights for the training data. (Make them uniform.)

$$w_1^{(1)} = 0.1 \qquad\qquad w_6^{(1)} = 0.1$$

$$w_2^{(1)} = 0.1 \qquad\qquad w_7^{(1)} = 0.1$$

$$w_3^{(1)} = 0.1 \qquad\qquad w_8^{(1)} = 0.1$$

$$w_4^{(1)} = 0.1 \qquad\qquad w_9^{(1)} = 0.1$$

$$w_5^{(1)} = 0.1 \qquad\qquad w_{10}^{(1)} = 0.1$$

(b) We will now grow a decision tree on this dataset, $G_1(x)$, using these weights. We'll use the `rpart` package in R, just as in Section 7.3. Here is the first tree.

| Datapoint ID | $w_i^{(1)}$ | happy ($Y$) | $G_1(x)$ |
|:---:|:---:|:---:|:---:|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | -1 |
| 3 | 0.1 | -1 | -1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 0.1 | -1 | -1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.1 | 1 | -1 |
| 8 | 0.1 | 1 | -1 |
| 9 | 0.1 | 1 | 1 |
| 10 | 0.1 | 1 | 1 |



Compute the misclassification error of this tree, err$_1$. Compare this tree to the one we constructed by hand in Chapter 7.

$$\text{err}_1 =$$

(b) We will now grow a decision tree on this dataset, $G_1(x)$, using these weights. We'll use the `rpart` package in R, just as in Section 7.3. Here is the first tree.

| Datapoint ID | $w_i^{(1)}$ | happy $(Y)$ | $G_1(x)$ |
|---|---|---|---|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | -1 |
| 3 | 0.1 | -1 | -1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 0.1 | -1 | -1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.1 | 1 | -1 |
| 8 | 0.1 | 1 | -1 |
| 9 | 0.1 | 1 | 1 |
| 10 | 0.1 | 1 | 1 |



Compute the misclassification error of this tree, $\text{err}_1$. Compare this tree to the one we constructed by hand in Chapter 7.

$$\text{err}_1 = \frac{0.1(2)}{0.1(10)} = 0.2$$

(c) Based on how $G_1(x)$ performs, calculate $\alpha_1$, its voting weight.

$$\alpha_1 =$$

(d) Re-weight the observation weights for the training examples.

$w_1^{(2)} =$

$w_2^{(2)} =$

$w_3^{(2)} =$

$w_4^{(2)} =$

$w_5^{(2)} =$

$w_6^{(2)} =$

$w_7^{(2)} =$

$w_8^{(2)} =$

$w_9^{(2)} =$

$w_{10}^{(2)} =$

(c) Based on how $G_1(x)$ performs, calculate $\alpha_1$, its voting weight.

$$\alpha_1 = \log\left(\frac{1 - \text{err}_1}{\text{err}_1}\right) = 1.386$$

(d) Re-weight the observation weights for the training examples.

$w_1^{(2)} = 0.1$

$w_2^{(2)} = 0.1$

$w_3^{(2)} = 0.1$

$w_4^{(2)} = 0.1$

$w_5^{(2)} = 0.1$

$w_6^{(2)} = 0.1$

$w_7^{(2)} = 0.1 * \exp(1.386) = 0.4$

$w_8^{(2)} = 0.1 * \exp(1.386) = 0.4$

$w_9^{(2)} = 0.1$

$w_{10}^{(2)} = 0.1$

(e) Now we grow another decision tree, $G_2(x)$, using these new weights as inputs to the `rpart` package.



| Datapoint ID | $w_i^{(2)}$ | happy $(Y)$ | $G_2(x)$ |
|---|---|---|---|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | -1 |
| 3 | 0.1 | -1 | -1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 0.1 | -1 | -1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.4 | 1 | 1 |
| 8 | 0.4 | 1 | 1 |
| 9 | 0.1 | 1 | 1 |
| 10 | 0.1 | 1 | -1 |

Compute the misclassification error of this tree, err$_2$.

err$_2$ =

(e) Now we grow another decision tree, $G_2(x)$, using these new weights as inputs to the `rpart` package.



| Datapoint ID | $w_i^{(2)}$ | happy $(Y)$ | $G_2(x)$ |
|:---:|:---:|:---:|:---:|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | -1 |
| 3 | 0.1 | -1 | -1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 0.1 | -1 | -1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.4 | 1 | 1 |
| 8 | 0.4 | 1 | 1 |
| 9 | 0.1 | 1 | 1 |
| 10 | 0.1 | 1 | -1 |

Compute the misclassification error of this tree, err$_2$.

$$\text{err}_2 = \frac{0.1}{0.1(8) + 0.4(2)} = 0.0625$$

(f) Based on how $G_2(x)$ performs, calculate $\alpha_2$, its voting weight.

$$\alpha_2 = \log\left(\frac{1 - \text{err}_2}{\text{err}_2}\right) = 2.708$$

(g) Re-weight the observation weights for the training examples.

$w_1^{(3)} = 0.1$

$w_2^{(3)} = 0.1$

$w_3^{(3)} = 0.1$
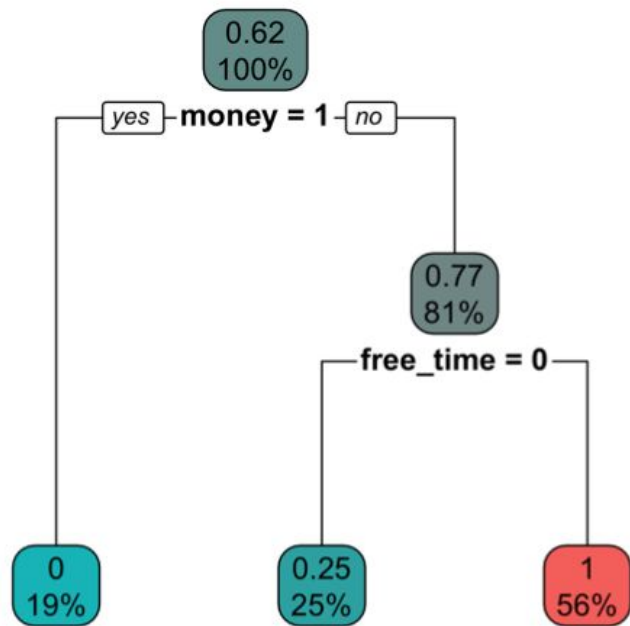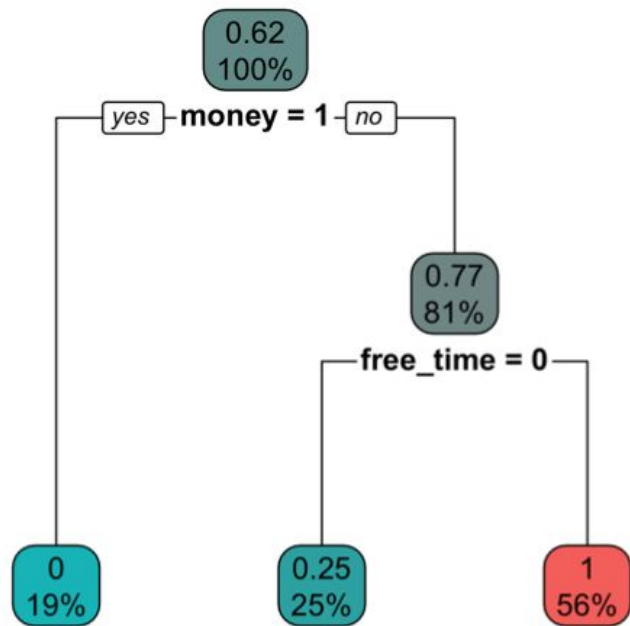
$w_4^{(3)} = 0.1$

$w_5^{(3)} = 0.1$

$w_6^{(3)} = 0.1$
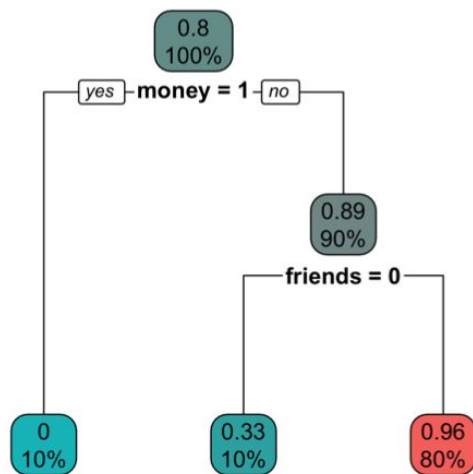
$w_7^{(3)} = 0.4$

$w_8^{(3)} = 0.4$

$w_9^{(3)} = 0.1$

$w_{10}^{(3)} = 0.1 * \exp(2.708) = 1.5$

(h) Now we grow another decision tree, $G_3(x)$, using these new weights as inputs to the `rpart` package.



| Datapoint ID | $w_i^{(3)}$ | happy ($Y$) | $G_3(x)$ |
|---|---|---|---|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | -1 |
| 3 | 0.1 | -1 | -1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 0.1 | -1 | 1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.4 | 1 | 1 |
| 8 | 0.4 | 1 | 1 |
| 9 | 0.1 | 1 | -1 |
| 10 | 1.5 | 1 | 1 |

(i) Compute the misclassification error of this tree, err₃.

$$\text{err}_3 =$$

$$\frac{0.2}{0.1(7) + 0.4(2) + 1.5} = 0.067$$

(j) Based on how $G_3(x)$ performs, calculate $\alpha_3$, its voting weight.

$$\alpha_3 = \log\left(\frac{1 - \mathrm{err}_3}{\mathrm{err}_3}\right) = 2.639$$

(k) Re-weight the observation weights for the training examples.

$w_1^{(4)} = 0.1$

$w_2^{(4)} = 0.1$

$w_3^{(4)} = 0.1$

$w_4^{(4)} = 0.1$

$w_5^{(4)} = $ 0.1*exp(2.639) = 1.4

$w_6^{(4)} = 0.1$
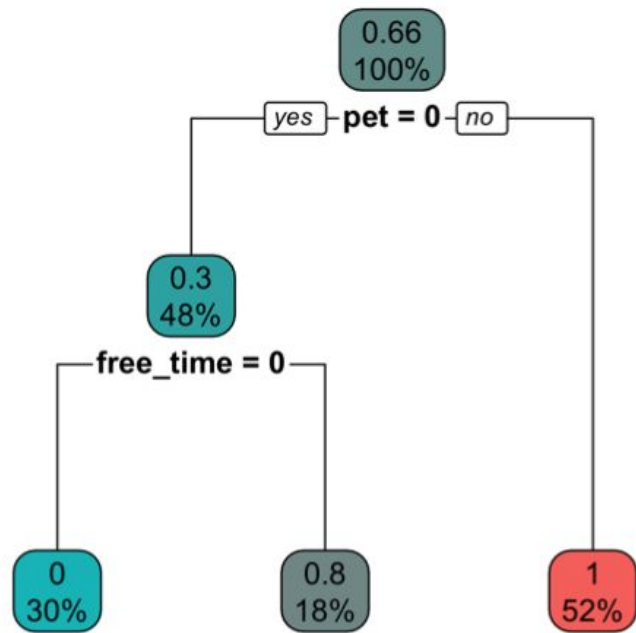
$w_7^{(4)} = 0.4$

$w_8^{(4)} = 0.4$

$w_9^{(4)} = $ 0.1*exp(2.639) = 1.4

$w_{10}^{(4)} = 1.5$

(l) Now we grow another decision tree, $G_4(x)$, using these new weights as inputs to the `rpart` package.

| Datapoint ID | $w_i^{(4)}$ | happy ($Y$) | $G_4(x)$ |
|:---:|:---:|:---:|:---:|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | 1 |
| 3 | 0.1 | -1 | 1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 1.4 | -1 | -1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.4 | 1 | 1 |
| 8 | 0.4 | 1 | 1 |
| 9 | 1.4 | 1 | 1 |
| 10 | 1.5 | 1 | 1 |



(m) Compute the misclassification error of this tree, $\text{err}_4$.

$$\text{err}_4 = \frac{0.2}{0.1(5) + 0.4(2) + 1.4(2) + 1.5} = 0.036$$

(n) Based on how $G_4(x)$ performs, calculate $\alpha_4$, its voting weight.

$$\alpha_4 = \log\left(\frac{1 - \text{err}_4}{\text{err}_4}\right) = 3.296$$

(o) Output the weighted average of the four classifiers' votes for each training example:

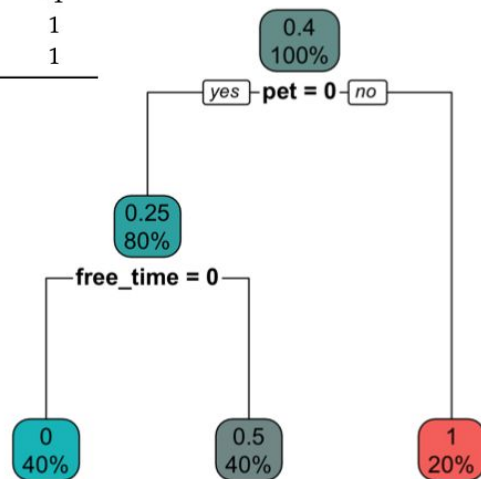$$G(x) = \text{sign}\left[\alpha_1 G_1(x) + \alpha_2 G_2(x) + \alpha_3 G_3(x) + \alpha_4 G_4(x)\right]$$

What is the final training error?

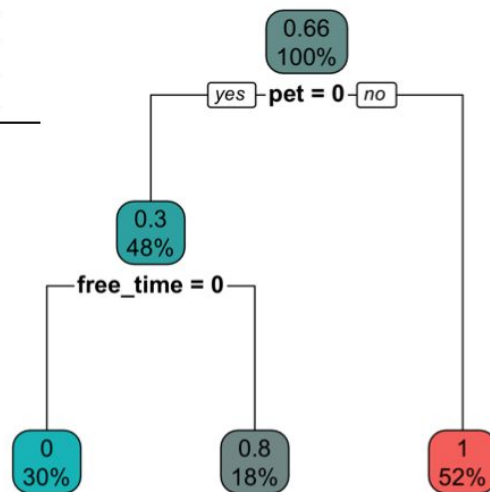| | | 1.386 | 2.708 | 2.639 | 3.296 | |
| Datapoint ID | happy $(Y)$ | $G_1(x)$ | $G_2(x)$ | $G_3(x)$ | $G_4(x)$ | $G(x)$ |
|---|---|---|---|---|---|---|
| 1 | -1 | -1 | -1 | -1 | -1 | |
| 2 | -1 | -1 | -1 | -1 | 1 | |
| 3 | -1 | -1 | -1 | -1 | 1 | |
| 4 | -1 | -1 | -1 | -1 | -1 | |
| 5 | -1 | -1 | -1 | 1 | -1 | |
| 6 | -1 | -1 | -1 | -1 | -1 | |
| 7 | 1 | -1 | 1 | 1 | 1 | |
| 8 | 1 | -1 | 1 | 1 | 1 | |
| 9 | 1 | 1 | 1 | -1 | 1 | |
| 10 | 1 | 1 | -1 | 1 | 1 | |

## Question 14.1

The most difficult thing to understand in all of this is how the updated weights play into the construction of subsequent trees. A clue comes from the dataset percentages shown in the nodes of each tree. For example, trees 1 and 3 actually

| Datapoint ID | $w_i^{(1)}$ | happy ($Y$) | $G_1(x)$ |
|---|---|---|---|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | -1 |
| 3 | 0.1 | -1 | -1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 0.1 | -1 | -1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.1 | 1 | -1 |
| 8 | 0.1 | 1 | -1 |
| 9 | 0.1 | 1 | 1 |
| 10 | 0.1 | 1 | 1 |

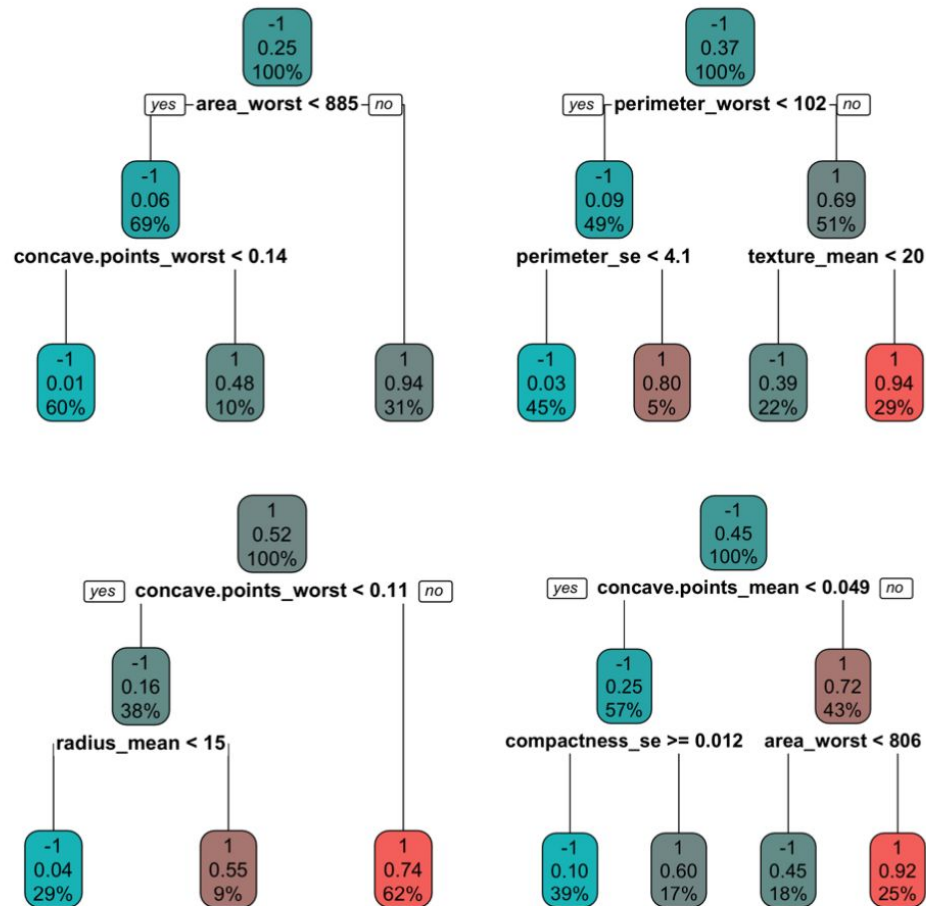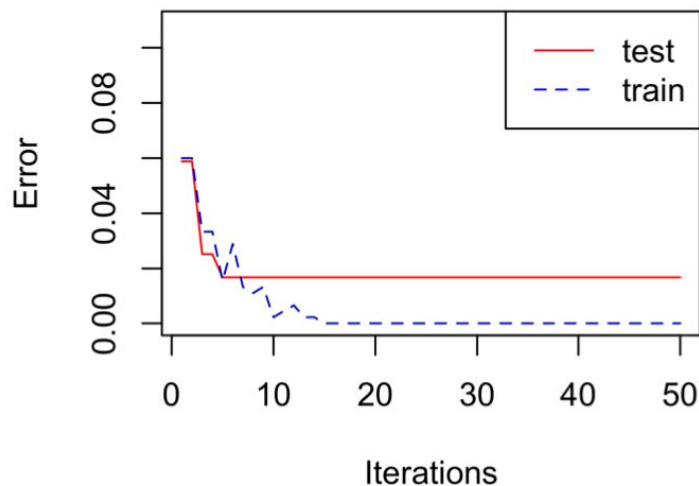| Datapoint ID | $w_i^{(4)}$ | happy ($Y$) | $G_4(x)$ |
|---|---|---|---|
| 1 | 0.1 | -1 | -1 |
| 2 | 0.1 | -1 | 1 |
| 3 | 0.1 | -1 | 1 |
| 4 | 0.1 | -1 | -1 |
| 5 | 1.4 | -1 | -1 |
| 6 | 0.1 | -1 | -1 |
| 7 | 0.4 | 1 | 1 |
| 8 | 0.4 | 1 | 1 |
| 9 | 1.4 | 1 | 1 |
| 10 | 1.5 | 1 | 1 |

**Gradient boosting:**
A more general conception of boosting in which later models are fit to the errors (pseudo-residuals) of earlier models.

How does this guarantee diversity?

# Wisconsin Breast Cancer dataset revisited

**Question 14.2**

Compare and contrast boosting and random forests on the basis of:

- Whether each tree uses all or part of the dataset

- Whether they consider a subset of the predictors at each split or all the predictors

- Whether you can parallelize the construction of different trees (i.e., build them at the same time on different processors)

- Whether the votes of different classifiers are independent

- Ease of use and interpretability