

MATH 4322 - Lecture 16

Tree Based Methods: Classification Tree

Dr. Cathy Poliak, cpoliak@uh.edu

University of Houston

Tree Based Models

So far we have covered such parametric models as:

- Linear Regression
- Logistic Regression
- Linear Discriminant Analysis

and such resampling techniques as

- Cross-Validation
- Bootstrap

we are ready for another model class:

- **Tree-based models.**

Tree-based models can be applied to **both** regression and classification problems.

Recall Building a Tree

Below are the steps of **building a decision tree** (ISLR, page 309):

1. Use recursive binary splitting to grow a large tree on the training data, stopping only when each terminal node has fewer than some minimum number of observations.
2. Apply cost complexity pruning to the large tree in order to obtain a sequence of best subtrees, as a function of α .
3. Use K -fold cross validation to choose α . That is, divide the training observations into K folds. For each $k = 1, \dots, K$;
 - (a) Repeat steps 1 and 2 on all but the k th fold of the training data.
 - (b) Evaluate the mean squared prediction error on the data in the left-out k th fold, as a function of α .
4. Average the results for each value of α , and pick α to minimize the average error.
5. Return the subtree from Step 2 that corresponds to the chosen value of α .

For Regression: $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \rightarrow \text{as small as possible}$

Residual Sum of Squares $RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Example

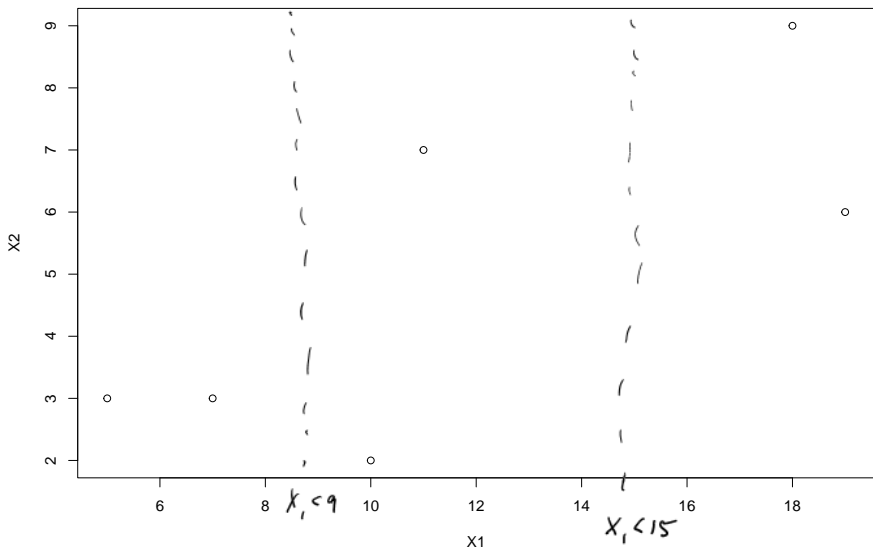
Let X_1 and X_2 be two predictors for the response Y . The following is a small data set.

| Y | 15 | 9 | 3 | 25 | 7 | 13 |
|-------|----|---|---|----|----|----|
| X_1 | 10 | 5 | 7 | 19 | 11 | 18 |
| X_2 | 2 | 3 | 3 | 6 | 7 | 9 |

$\bar{y} = 12$

Build a regression tree using the recursive binary splitting.

$$\begin{aligned}RSS &= (15-12)^2 + (9-12)^2 + \dots + (13-12)^2 \\&= 3^2 + 3^2 + 6^2 + 13^2 + 5^2 + 1^2 \\&= 294\end{aligned}$$



Split 1

$$x_i < 15 \quad y = C(15, 9, 3, 7) \quad \bar{y} = 8.5$$

$$RSS_1 = (15-8.5)^2 + (9-8.5)^2 + (3-8.5)^2 + (7-8.5)^2 = 75$$

$$x_i \geq 15 \quad y = C(25, 13) \quad \bar{y} = 19$$

$$RSS_2 = (25-19)^2 + (13-19)^2 = 72$$

$$RSS = RSS_1 + RSS_2 = 75 + 72 = 147$$

Split 2

$$x_i < 9 \quad y = C(9, 3) \quad \bar{y} = 6$$

$$RSS_1 = (9-6)^2 + (3-6)^2 = 18$$

$$9 \leq x_i < 15 \quad y = C(15, 7) \quad \bar{y} = 11$$

$$RSS_2 = (15-11)^2 + (7-11)^2 = 32$$

$$x_i \geq 15 \quad RSS_3 = 72$$

$$RSS = 18 + 32 + 72 = \underline{122}$$

R Result

```
library(tree)
tree.ex = tree(Y~X1 + X2, control = tree.control(6, mincut = 2, minn = 10))
tree.ex
```

node), split, n, deviance, yval

* denotes terminal node

1) root 6 294 12.0 \leftarrow No split

2) $X1 < 14.5$ 4 75 8.5

4) $X1 < 8.5$ 2 18 6.0 *

5) $X1 > 8.5$ 2 32 11.0 *

3) $X1 > 14.5$ 2 72 19.0 *

\in terminal node

Classification Trees

- Used to predict a qualitative (categorical) response.
- Recall regression trees - predicted response for an observation is given by the mean response of the training observations that belong to the same terminal node.
- Classification tree - predicts that each observation belongs to the **most commonly occurring class** of training observations in the region to which it belongs.
- Interpretation of classification tree - class prediction corresponding to a particular terminal node regions and the class proportions among the training observations that fall into that region.

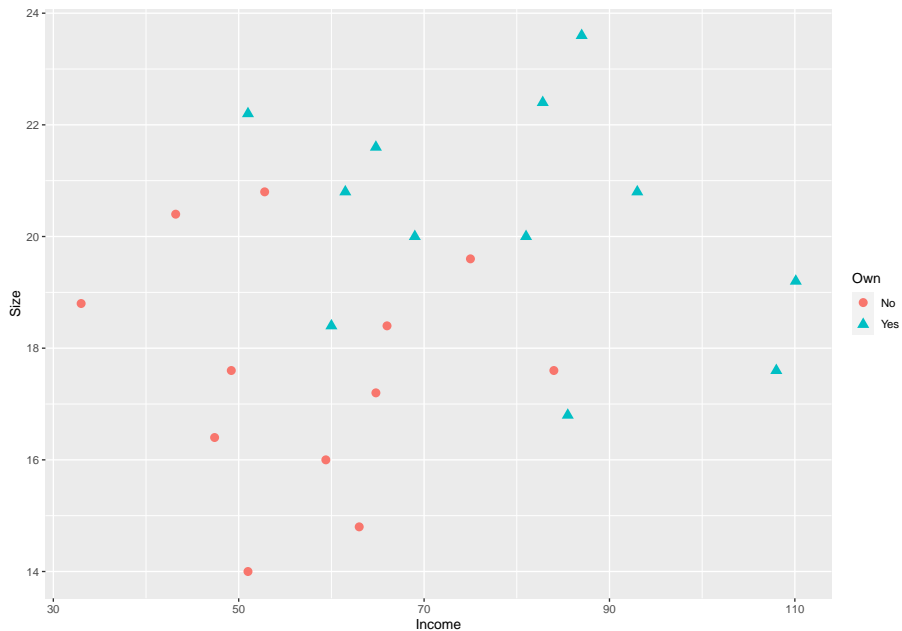
Growing a Classification Tree

- Task of growing a classification tree is similar to the regression trees, we use recursive binary splitting to grow a classification tree.
- Recall criterion for regression tree is **RSS**. Since the response is categorical we cannot calculate the residual standard error. Thus we use other criterion to grow a classification tree.
 - ▶ Classification error rate
 - ▶ Gini index
 - ▶ Entropy

Example

From *Applied Multivariate Statistical Analysis* by Johnson and Wichern:

A riding-mower manufacturer would like to find a way of classifying families in a city into those that are likely to purchase a riding mower and those who are not likely to buy one. A pilot random sample of 12 owners and 12 non-owners in the city is undertaken. The data are plotted below. The independent variables here are Income (x_1) and Lot Size (x_2). The categorical y variable has two classes: owners and non-owners.



Classification Error Rate

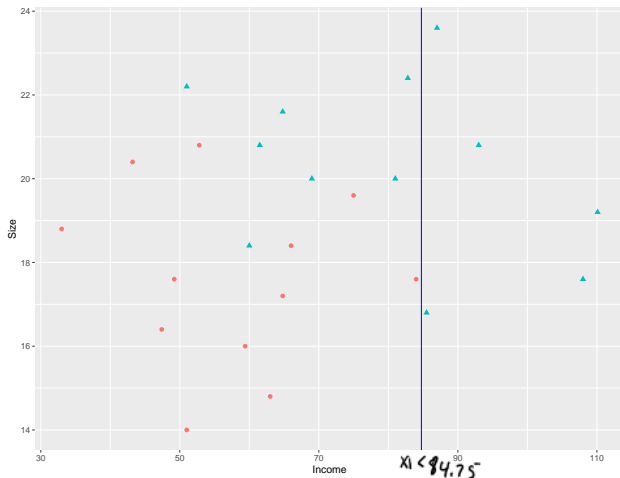
- A natural alternative to RSS
- The fraction of the training observations in that region that do not belong to the most common class:

$$E = 1 - \max_k(\hat{p}_{mk}).$$

Where \hat{p}_{mk} represents the portion of training observations in the m th region that are from the k th class.

- This is not sufficiently sensitive for tree-growing. This only finds the best split at that immediate place.

Split 1



$$\text{Split}$$

$$X_1 < 84.75$$

$$\hat{P}_{\text{Yes}} = \frac{7}{19} = 0.3684$$

$$\hat{P}_{\text{No}} = \frac{12}{19} = 0.6316$$

Own
 ● No
 ▲ Yes

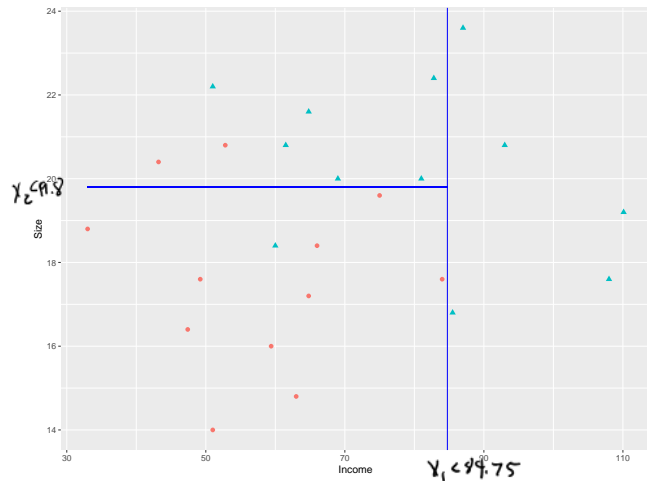
$$X_1 > 84.75$$

$$\hat{P}_{\text{Yes}} = 1$$

$$\hat{P}_{\text{No}} = 0 \quad \text{Node purity}$$

$$E = 1 - 1 = 0$$

Split 2



$$X_1 \geq 84.75$$

$$\hat{P}_{\text{Yes}} = 1, \hat{P}_{\text{No}} = 0$$

$$X_1 < 84.75 \text{ \& } X_2 < 19.8$$

$$\hat{P}_{\text{Yes}} = \frac{1}{11} = 0.0909$$

$$\hat{P}_{\text{No}} = \frac{10}{11} = 0.9091$$

$$X_1 < 84.75 \text{ \& } X_2 \geq 19.8$$

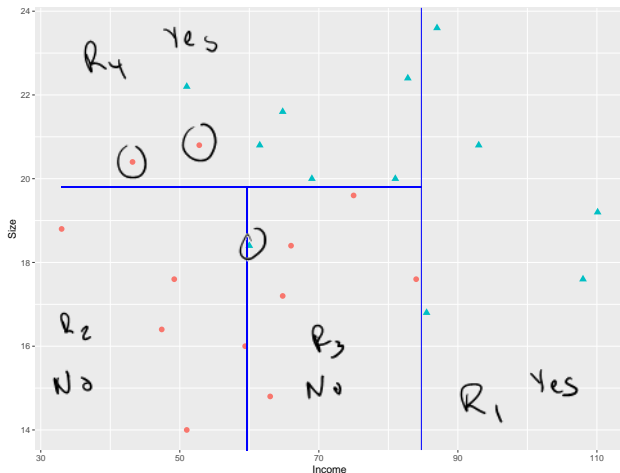
$$\hat{P}_{\text{Yes}} = \frac{6}{8} = 0.75$$

$$\hat{P}_{\text{No}} = \frac{2}{8} = 0.25$$

$$E = 0$$

$$= 1 - \max(\hat{P}_i)$$

Split 3



3 miss classified
out of 24

R1 Income > 84.75

R2

Income ≤ 59.74
and Size ≤ 19.8

R3

$59.74 < \text{Income} \leq 84.75$
and size ≤ 19.8

R4

Income ≤ 59.74
and size > 19.8

Fit A Classification Tree In R

```
library(tree)
mower$Own = as.factor(mower$Own)
tree.mower = tree(Own ~ Income + Size, data = mower)
summary(tree.mower)
```

Classification tree:

```
tree(formula = Own ~ Income + Size, data = mower)
```

Number of terminal nodes: 4

Residual mean deviance: 0.7202 = 14.4 / 20

Misclassification error rate: 0.125 = 3 / 24

```
tree.mower
```

node), split, n, deviance, yval, (yprob)

* denotes terminal node

- ```
1) root 24 33.270 No (0.50000 0.50000)
 2) Income < 84.75 19 25.010 No (0.63158 0.36842)
 4) Size < 19.8 11 6.702 No (0.90909 0.09091)
 8) Income < 59.7 5 0.000 No (1.00000 0.00000) *
 9) Income > 59.7 6 5.407 No (0.83333 0.16667) *
 5) Size > 19.8 8 8.997 Yes (0.25000 0.75000) *
 3) Income > 84.75 5 0.000 Yes (0.00000 1.00000) *
```

The calculation for deviance for a classification tree is  $-2 \sum_{i=1}^k n_i \log(p_i)$ .

$$-2[0 + 5.407 + 8.997 + 0] / (24 - 4)$$

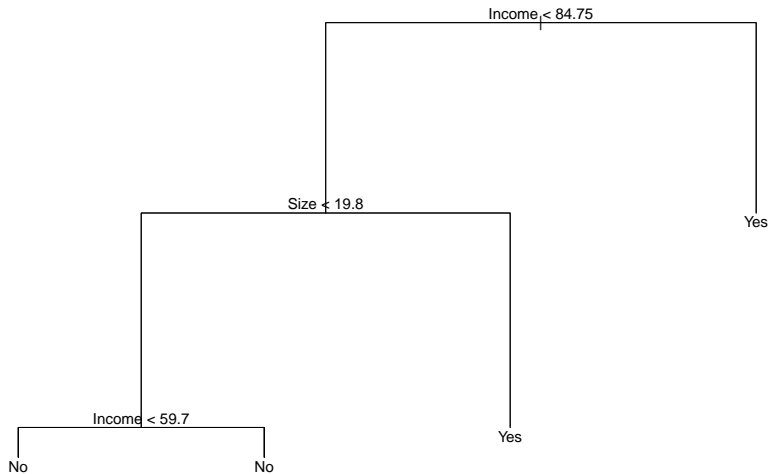
$$-2[12 \ln(0.5) + 12 \ln(0.5)]$$

$$-2[12 \ln(0.63158) + 7 \ln(0.36842)]$$

$$\rightarrow (\tilde{p}_{No}, \tilde{p}_{Yes})$$



# Plot of Tree



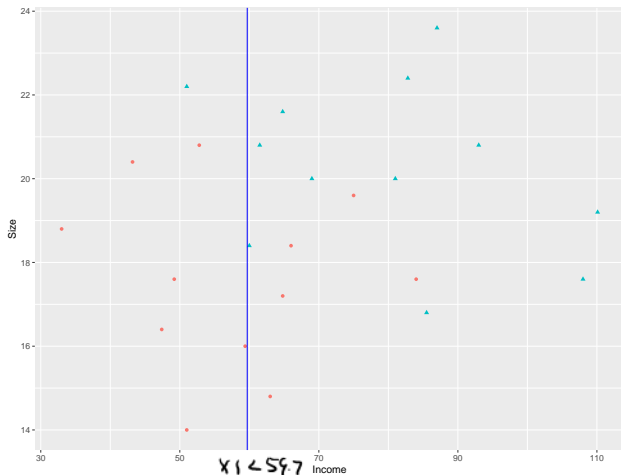
# Gini Index

- A measure of total variance across the  $K$  classes.

$$G = \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$$

- This is measure of node purity - a small value indicates that a node contains predominantly observations from a single class.
- Measures the degree or probability of a particular variable being wrongly classified when it is randomly chosen.
- For more information and more examples see: [Gini Example](#).

# Split 1



$$X_1 < 59.7$$

$$\hat{P}_{Yes} = \frac{1}{8} = 0.125$$

$$\hat{P}_{No} = \frac{7}{8} = 0.875$$

$$X_1 \geq 59.7$$

Own

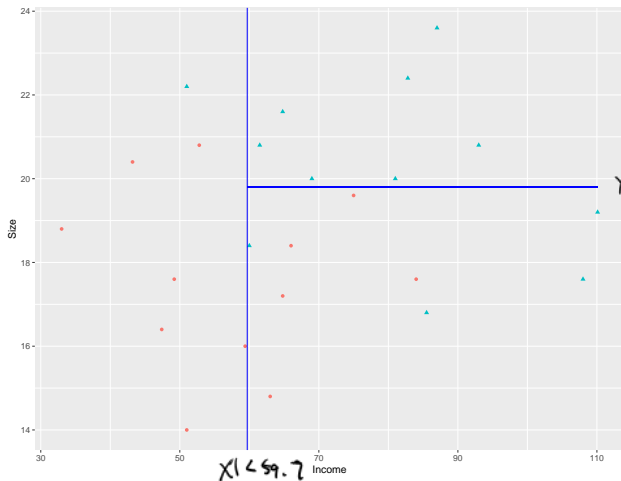
- No
- Yes

$$\hat{P}_{Yes} = \frac{11}{16} = 0.6875$$

$$\hat{P}_{No} = \frac{5}{16} = 0.3125$$

$$G = (0.125)(0.875) + (0.875)(0.125) + (0.6875)(0.3125) + (0.3125)(0.6875) = 0.6484$$

# Split 2



$$X_1 < 59.7$$

$$\hat{P}_{Yes} = 0.125$$

$$\hat{P}_{No} = 0.875$$

$$X_1 \geq 59.7 \text{ and } X_2 < 19.8$$

$$X_2 < 19.8$$

$$\hat{P}_{Yes} = 0.444$$

$$\hat{P}_{No} = 0.5556$$

$$X_1 \geq 59.7 \text{ and } X_2 \geq 19.8$$

$$\hat{P}_{Yes} = 1$$

$$\hat{P}_{No} = 0$$

$$G = 0.7126$$

# Results in R Based on Gini Index

```
tree.mower.gini = tree(Own ~ Income + Size, data = mower, split = "gini")
summary(tree.mower.gini)
```

Classification tree:

```
tree(formula = Own ~ Income + Size, data = mower, split = "gini")
```

Number of terminal nodes: 3

Residual mean deviance: 0.8759 = 18.39 / 21

Misclassification error rate: 0.2083 = 5 / 24

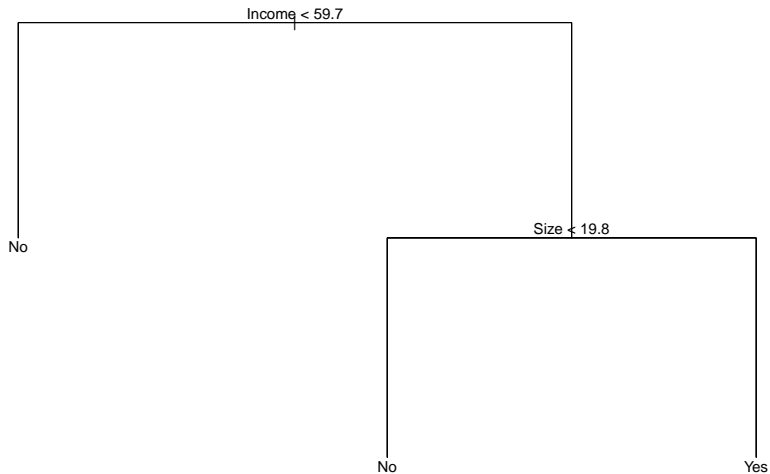
```
tree.mower.gini
```

```
node), split, n, deviance, yval, (yprob)
```

\* denotes terminal node

- 1) root 24 33.270 No ( 0.5000 0.5000 )
- 2) Income < 59.7 8 6.028 No ( 0.8750 0.1250 ) \*
- 3) Income > 59.7 16 19.870 Yes ( 0.3125 0.6875 )
  - 6) Size < 19.8 9 12.370 No ( 0.5556 0.4444 ) \*
  - 7) Size > 19.8 7 0.000 Yes ( 0.0000 1.0000 ) \*

# Plot of Tree



# Entropy

Given by

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

- Since  $0 \leq \hat{p}_{mk} \leq 1$ , then  $0 \leq -\hat{p}_{mk} \log \hat{p}_{mk}$
- The entropy will take on a value near zero if the  $\hat{p}_{mk}$  are all near zero or one.
- The entropy will be a small value if the  $m$ th node is pure.
- Also called **Information gain** measurement.
- This gives us the maximum information about a class.
- Harder to compute because of the  $\log$ , thus Gini index is preferred.

# What Method To Use

- When building a classification tree, either the Gini index or the entropy are typically used to evaluate the quality of a particular split. This leads to increased **node purity**.
- Any of the three approaches might be used when **pruning** the tree.
- The classification error rate is preferable if prediction accuracy of the final pruned tree is the goal.
- The `tree` function uses the classification error rate as the default but can also use the Gini Index.



## Example 2

- We will use the *Heart* data. See Canvas to get the data.
- This data contains info on patients with chest pains, and we'd like to classify if a patient has a heart disease (*AHD*) depending on multiple factors.
- Import the data into R type and run the following:

```
set.seed(100)
train = sample(1:nrow(Heart),nrow(Heart)/2+0.5)
Heart$AHD = as.factor(Heart$AHD)
Heart$ChestPain = as.factor(Heart$ChestPain)
Heart$Thal = as.factor(Heart$Thal)
Heart$Sex = as.factor(Heart$Sex)
tree.heart = tree(AHD ~ . -X, Heart,subset = train)
```

# Lab Questions

1. Type and run `summary(tree.heart)`. What are the number of terminal nodes?  

a) 9                      b) 10                      c) 11                      d) 13
2. What is the training error rate?  

a) 38.5%                      b) 10.96%                      c) 14%                      d) 51.6%

3. Type and run the following in R

```
plot(tree.heart)
```

```
text(tree.heart)
```

Is a person predicted to have heart disease if the  $Ca > 0.5$ , slope  $< 1.5$  and sex = a.

a) Yes

b) No

# Test Error Rate

In order to properly evaluate the performance of a classification tree on these data, we must estimate the test error rather than simply computing the training error. Type and run the following in

```
Heart.test = Heart[-train,]
tree.pred = predict(tree.heart, Heart.test, type = "class")
table(tree.pred, Heart.test$AHD)
```

4. What is the test error rate?

a) 18.34%

c) 21%

b) 81.66%

d) 17.5%

# Pruning

- Next, we consider whether pruning the tree might lead to improved results.
- We use the argument `FUN=prune.misclass` in the `cv.tree()` function in order to indicate that we want the classification error rate to guide the cross-validation and pruning process, rather than the default for the `cv.tree()` function, which is deviance.
- The `cv.tree()` function reports the number of terminal nodes of each tree considered (`size`) as well as the corresponding error rate and the value of the cost-complexity parameter used (`k`, which corresponds to  $\alpha$ ).

## Type and run the following in R

```
set.seed(3)
cv.heart = cv.tree(tree.heart,FUN = prune.misclass)
par(mfrow = c(1,2))
plot(cv.heart$size,cv.heart$dev,type = "b")
plot(cv.heart$k,cv.heart$dev,type = "b")

par(mfrow = c(1,1))
```

5. How many nodes do we really desire?

a) 14

c) 6

b) 10

d) 2

# Pruned Tree

Type and run the following:

```
prune.heart = prune.misclass(tree.heart,best = 6)
tree.pred = predict(prune.heart,Heart.test,type = "class")
table(tree.pred,Heart.test$AHD)
```

What is the test error rate?

# Trees vs Linear Models

Linear regression assumes a model of the form

$$f(X) = \beta_0 + \sum_{j=1}^p X_j \beta_j,$$

Regression trees assume a model of the form

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

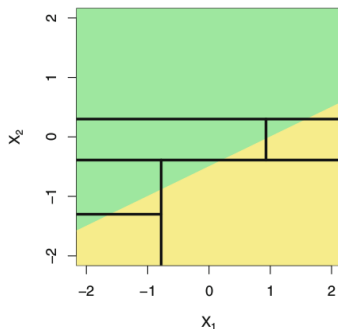
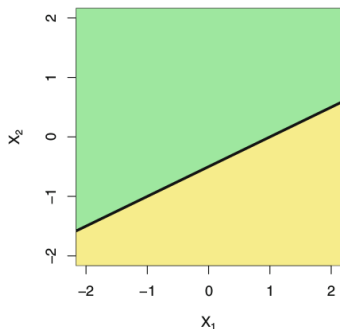
**Which is better?** **Data Scientist answer:** it depends on the problem and data at hand,

- If relationship between predictors  $X_1, \dots, X_p$  and  $Y$  is approximately linear  $\implies$  linear model it is.
- Otherwise, if that relationship is highly non-linear and complex  $\implies$  decision trees may have an edge.



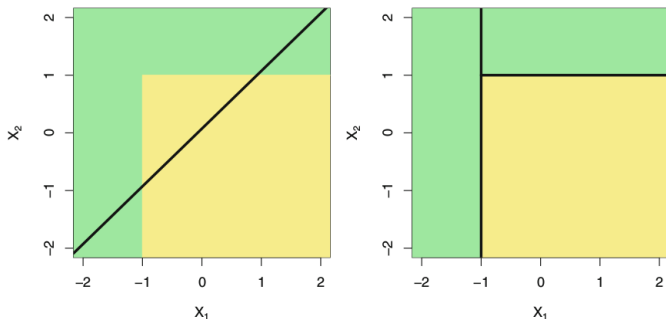
# Classification Example: Linear preferred over Trees

**Example.** In the case of a **classification problem** with a **linear boundary** - linear approach is equipped to perform better. Below you see results of fitting a linear model (left) and a decision tree (right).



# Classification Example: Trees preferred or Linear

**Example.** In classification problem with a more complex, non-linear boundary - decision trees have better chances. Below you see results of fitting a linear model (left) and a decision tree (right).



# Decision Trees: Advantages and Disadvantages

Several **advantages** of decision trees as a model:

1. **Easy** to **explain**, **visualize** and **interpret**.
2. More closely **mirror human decision-making** than regression and certain other methods.
3. **Easily** handle both **quantitative** and **qualitative** predictors (and responses).

The biggest downside:

1. Trees generally **do not have the same level of predictive accuracy** as some other regression and classification approaches.

However, by **aggregating many decision trees** (e.g. bagging, random forests), **the predictive performance of trees can be vastly improved**.