# Classification Trees

# Quarto

Quarto enables you to weave together content and executable code into a finished presentation. To learn more about Quarto presentations see https://quarto.org/docs/presentations/.

# Bullets

When you click the **Render** button a document will be generated that includes:

- Content authored with markdown
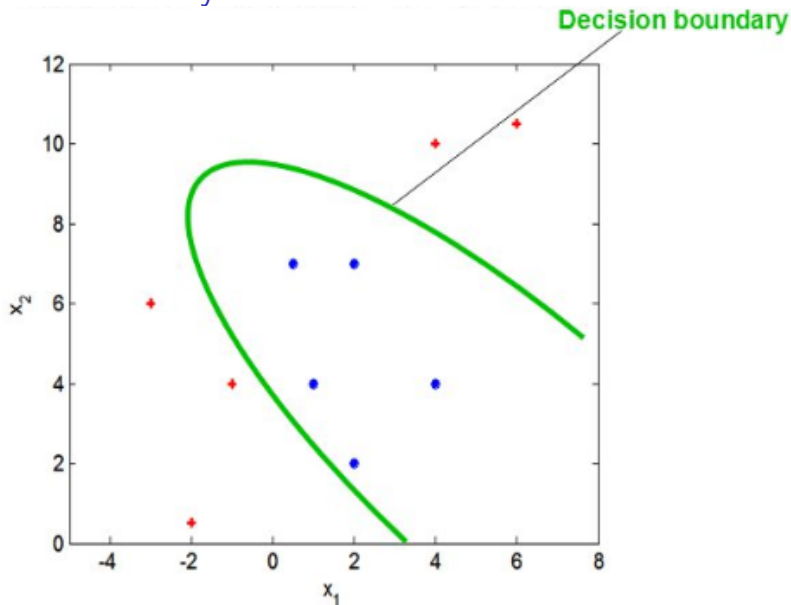- Output from executable code

# Code

When you click the **Render** button a presentation will be generated that includes both content and the output of embedded code. You can embed code like this:

```
[1] 2
```

# Reading Materials
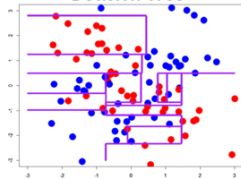
- Max Kuhn. Chapter 14. Section 14.1
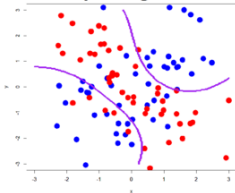
# Decision Boundary in Classification



Classification is a process of finding the **decision boundary** that
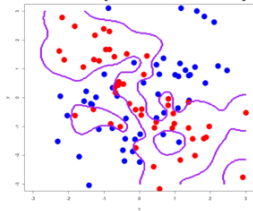
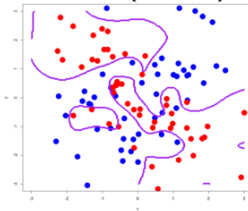# Decision Boundary in Classification

# Decision Tree

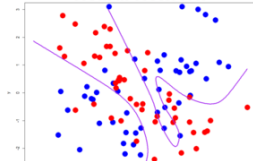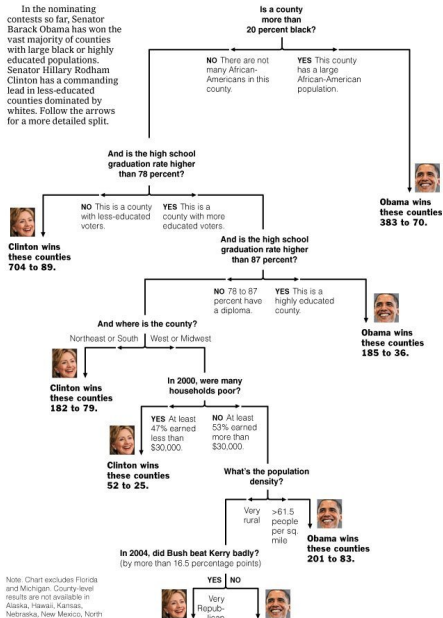▶ Decision Tree for classification is **Classification Tree**
▶ Decision Tree for Regression is **Regression Tree**

# Example of Classification Tree

Decision Tree: The Obama-Clinton Divide

In the nominating contests so far, Senator Barack Obama has won the vast majority of counties with large black or highly educated populations. Senator Hillary Rodham Clinton has a commanding lead in less-educated counties dominated by whites. Follow the arrows for a more detailed split.
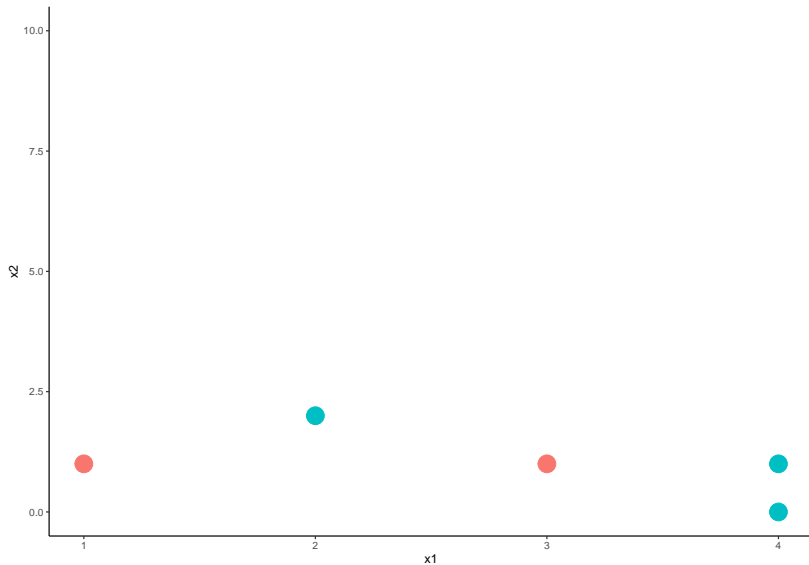
**Is a county more than 20 percent black?**

**NO** There are not many African-Americans in this county.

**YES** This county has a large African-American population.

**Obama wins these counties 383 to 70.**

**And is the high school graduation rate higher than 78 percent?**

**NO** This is a county with less-educated voters.

**YES** This is a county with more educated voters.

**Clinton wins these counties 704 to 89.**

**And is the high school graduation rate higher than 87 percent?**

**NO** 78 to 87 percent have a diploma.

**YES** This is a highly educated county.

**Obama wins these counties 185 to 36.**

**And where is the county?**

Northeast or South | West or Midwest

**Clinton wins these counties 182 to 79.**

**In 2000, were many households poor?**

**YES** At least 47% earned less than $30,000.

**NO** At least 53% earned more than $30,000.

**Clinton wins these counties 52 to 25.**

**What's the population density?**

Very rural | >61.5 people per sq. mile

**Obama wins these counties 201 to 83.**

**In 2004, did Bush beat Kerry badly?**
(by more than 16.5 percentage points)

**YES | NO**

Very Republican

Note. Chart excludes Florida and Michigan. County-level results are not available in Alaska, Hawaii, Kansas, Nebraska, New Mexico, North Dakota or Maine, Texas
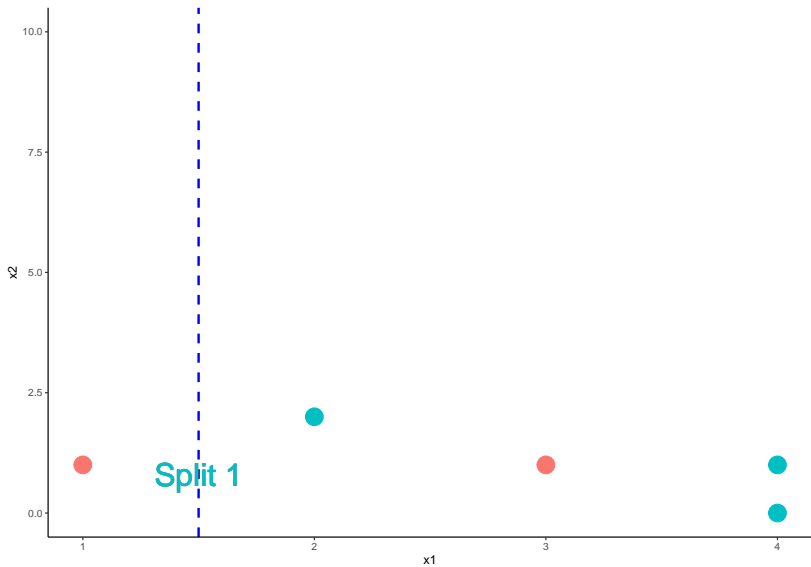
# Classification Tree

▶ In two dimension, classification Tree's decision boundary is a collection of horizontal and vertical line
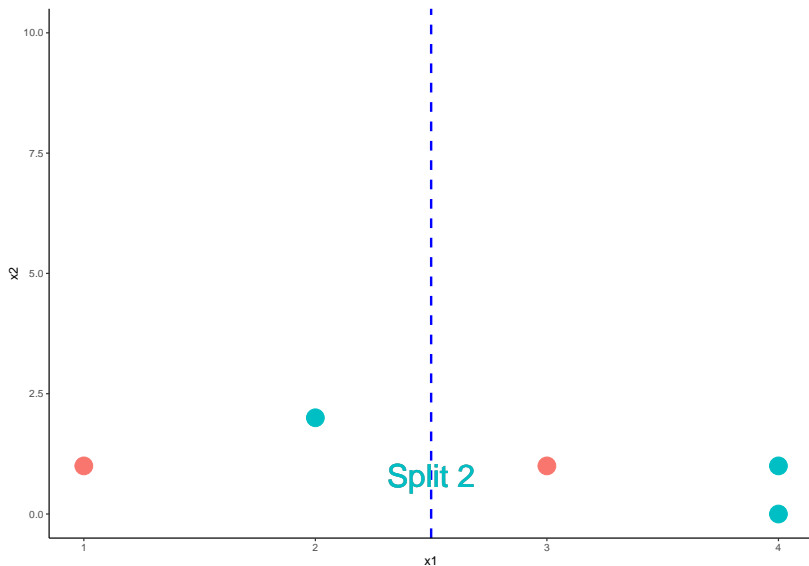
# Data



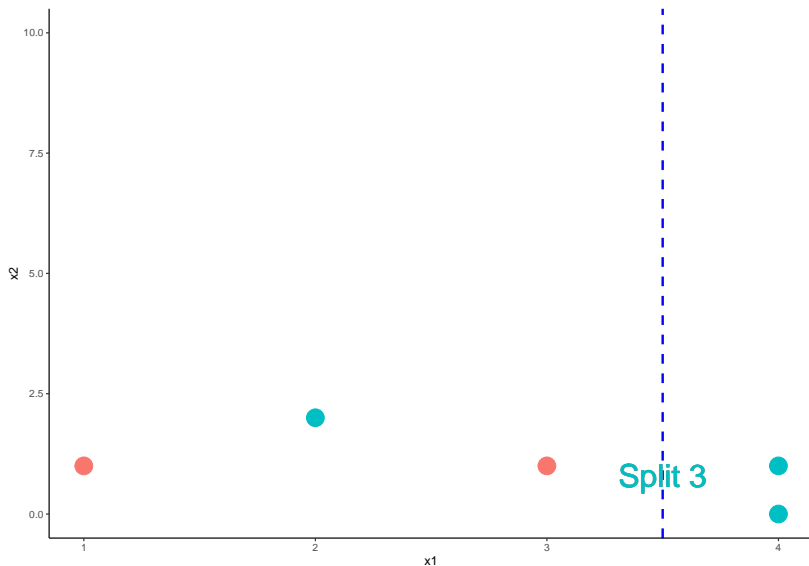▶ The tree starts by a vertical or horizontal line that **best** seperate the data

# One way to seperate the reds and greens

# One way to seperate the reds and greens

# One way to seperate the reds and greens

▶ **Question**: Which is the best split?

# Partial Answer

- It looks like Split 1 and 3 are better than Split 2 since it misclassifies less
- Which is the better split between Split 1 and Split 3?
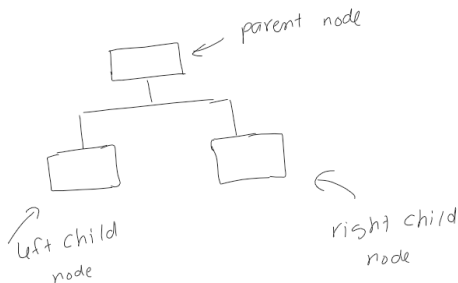- We need to find a way to measure *how good a split is*

# Impurity Measure

▶ The impurity of a node (**a node = a subset of the data or the original data**) measure how uncertain the node is.

▶ For example, node A with 50% reds and 50% greens would be more uncertained than node B with 90% reds and 10% greens. Thus, node A has greater impurity than node B.

▶ More uncertained = Greater impurity

# Impurity Measure

▶ A split that *gains* more impurity is the **better split**!

# Impurity Gain



$$IG = I_{parent} - \frac{N_{left}}{N}I_{left} - \frac{N_{right}}{N}I_{right}$$

▶ IG is Impurity Gain of the split
▶ $N_{left}$ and $N_{right}$ are the number of points in the left child node and right child node, respectively.
▶ $N_{left} + N_{right} = N$

# Impurity Measure

▶ Impurity can be measured by: classification error, Gini Index, and Entropy.

# Impurity Measure

▶ Let $p_0$ and $p_1$ be the proportion of class 0 and class 1 in a node.

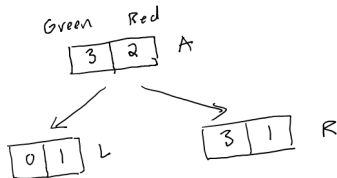By Classification Error: $I = min\{p_0, p_1\}$

By Gini Index: $I = 1 - p_0^2 - p_1^2$

By Entropy: $I = -p_0 \log_2(p_0) - p_1 \log_2(p_1)$

# Calculation

▶ Let's calculate the impurity gain of the three splits to decide which split is the best
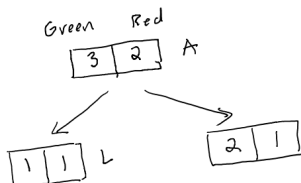
# IG By Classification Error



▶ Let **green** and **red** be class 0 and class 1, respectively.

For Split 1: $N = 5, N_{left} = 1, N_{right} = 4$

▶ Node *parent*, A: $p_0 = \frac{2}{5}, p_1 = \frac{3}{5}$. Thus, $I_A = \min(\frac{2}{5}, \frac{3}{5}) = \frac{2}{5}$
▶ Node *child left*, L: $p_0 = \frac{0}{1} = 0, p_1 = \frac{1}{1} = 1$. Thus,
  $I_L = \min(0, 1) = 0$
▶ Node *child right*, R: $p_0 = \frac{3}{4}, p_1 = \frac{1}{4}$. Thus,
  $I_R = \min(\frac{3}{4}, \frac{1}{4}) = \frac{1}{4}$
▶ Impurity Gain of Split 1:

$$IG = \frac{2}{5} - \frac{1}{5} \cdot 0 - \frac{4}{5} \cdot \frac{1}{4} = 0.2$$
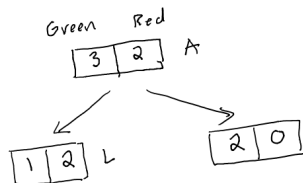
# IG By Classification Error



For Split 2: $N = 5, N_{left} = 2, N_{right} = 3$

▶ Node *parent*, A: $p_0 = \frac{2}{5}, p_1 = \frac{3}{5}$. Thus, $I_A = \min(\frac{2}{5}, \frac{3}{5}) = \frac{2}{5}$

▶ Node *child left*, L: $p_0 = \frac{1}{2}, p_1 = \frac{1}{2}$. Thus, $I_L = \frac{1}{2}$

▶ Node *child right*, R: $p_0 = \frac{2}{3}, p_1 = \frac{1}{3}$. Thus,
$I_R = \min(\frac{2}{3}, \frac{1}{3}) = \frac{1}{3}$

▶ Impurity Gain of Split 2:

$$IG = \frac{2}{5} - \frac{2}{5} \cdot \frac{1}{2} - \frac{3}{5} \cdot \frac{1}{3} = 0$$

# IG By Classification Error

For Split 3: $N = 5, N_{left} = 3, N_{right} = 2$



- ▶ Node *parent*, A: $p_0 = \frac{2}{5}, p_1 = \frac{3}{5}$. Thus, $I_A = \min(\frac{2}{5}, \frac{3}{5}) = \frac{2}{5}$
- ▶ Node *child left*, L: $p_0 = \frac{1}{3}, p_1 = \frac{2}{3}$. Thus, $I_A = \min(\frac{1}{3}, \frac{2}{3}) = \frac{1}{3}$
- ▶ Node *child right*, R: $p_0 = \frac{2}{2}, p_1 = \frac{0}{2}$. Thus, $I_R = \min(1, 0) = 0$
- ▶ Impurity Gain of Split 3:

$$IG = \frac{2}{5} - \frac{3}{5} \cdot \frac{1}{3} - \frac{2}{5} \cdot 0 = 0.2$$

# Comparing IG By Classification Error

| Split | IG |
|-------|-----|
| 1 | 0.2 |
| 2 | 0 |
| 3 | 0.2 |

▶ By classification error, Split 1 and Split 3 are tie as the best because they have the same impurity gain.