

# Decision Trees II

Nikhil Sardana  
presented by Vinay Bhaip

October 2019

## 1 Introduction

Part 1 of this lecture covered the mathematics of decision trees. This lecture will transition you from graphs and mathematics to writing the actual code for the competition.

## 2 Review

Decision trees are created by splitting on a threshold of a feature which results in the most information gain, and recursively continuing this process for each of the children.

Let us review Information Gain with an example. Remember, the equation is:

$$IG(D_p, f) = I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right})$$

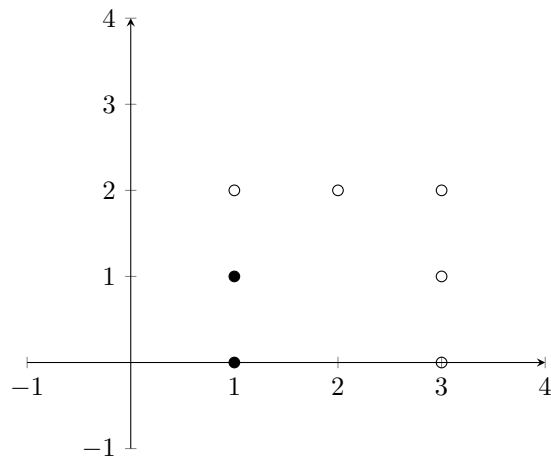
We will use the Gini Impurity measure, just like last time:

$$I(i) = 1 - \sum_{k=1}^c p(k|i)^2$$

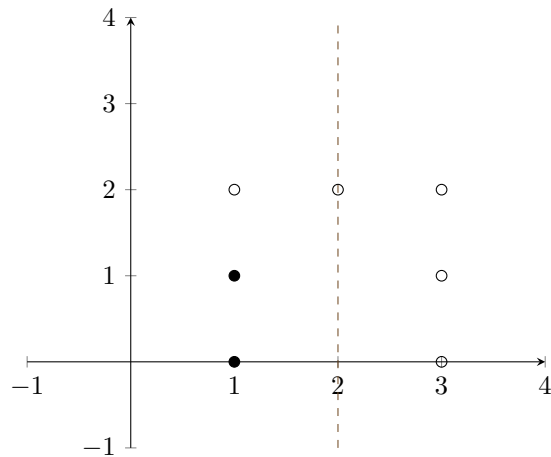
Consider the following data. We have two features,  $x1$  and  $x2$ . Our label is  $y$ .

x1	x2	y
1	0	0
1	1	0
1	2	1
2	2	1
3	0	1
3	1	1
3	2	1

Let's graph this data. Let the  $\bullet = 0$  and  $\circ = 1$ .

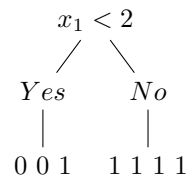


We want to draw a vertical or horizontal line that best separates the white dots from the black dots (ie has the greatest information gain). Since we are restricted to the coordinates of data points, a human will look at the graph and choose:



Since we are classifying things based on whether they are less than the threshold ( $x < 2$ ), the point in the middle goes with the right side.

The decision tree that goes with this is:



We can calculate the information gain.

$$\begin{aligned}
IG(D_p, f) &= I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right}) \\
N_{left} &= 3 \\
N_{right} &= 4 \\
N_p &= 7
\end{aligned}$$

Now we need the impurities.

$$I(i) = 1 - \sum_{k=1}^c p(k|i)^2$$

For the parent, we calculate the impurity of the entire dataset, which is:

$$I(D_p) = 1 - \left(\frac{2}{7}\right)^2 - \left(\frac{5}{7}\right)^2 = 1 - \frac{4}{49} - \frac{25}{49} = \frac{49 - 4 - 25}{49} = \frac{20}{49}$$

And for the children:

$$\begin{aligned}
I(D_{left}) &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 1 - \frac{4}{9} - \frac{1}{9} = \frac{4}{9} \\
I(D_{right}) &= 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 1 - 1 = 0
\end{aligned}$$

Note that the most pure data  $I(i) = 0$  is all one type.

Finally, we get the information gain:

$$\begin{aligned}
IG(D_p, f) &= I(D_p) - \frac{N_{left}}{N_p} I(D_{left}) - \frac{N_{right}}{N_p} I(D_{right}) \\
IG(D_p, f) &= \frac{20}{49} - \frac{3}{7} * \frac{4}{9} - \frac{4}{7} * 0 = \frac{20}{49} - \frac{12}{63} = \frac{1}{7} * \left(\frac{20}{7} - \frac{12}{9}\right) = \frac{1}{7} * \left(\frac{180}{63} - \frac{84}{63}\right) \\
IG(D_p, f) &= \frac{1}{7} * \frac{96}{63} = \frac{96}{441}
\end{aligned}$$

### 3 Data Structure

Let's go back to the example from Section 2.1.

x1	x2	y
1	0	0
1	1	0
1	2	1
2	2	1
3	0	1
3	1	1
3	2	1

How will we represent this data in our code? It's simple! We'll use a list of lists, with each column a feature.

$$\begin{bmatrix} [1, 0, 0] \\ [1, 1, 0] \\ [1, 2, 1] \\ [2, 2, 1] \\ [3, 0, 1] \\ [3, 1, 1] \\ [3, 2, 1] \end{bmatrix}$$

Alternatively, you can make the labels their own list. Our sample I/O code has it this way:

$$Features = \begin{bmatrix} [1, 0] \\ [1, 1] \\ [1, 2] \\ [2, 2] \\ [3, 0] \\ [3, 1] \\ [3, 2] \end{bmatrix} \qquad Labels = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

To find the best information gain, we loop through each of the features (columns), and each of the thresholds, splitting the matrix and calculating the information gain. For example, in the previous example, we split on  $x_1 < 2$ , so our left child matrix would be:

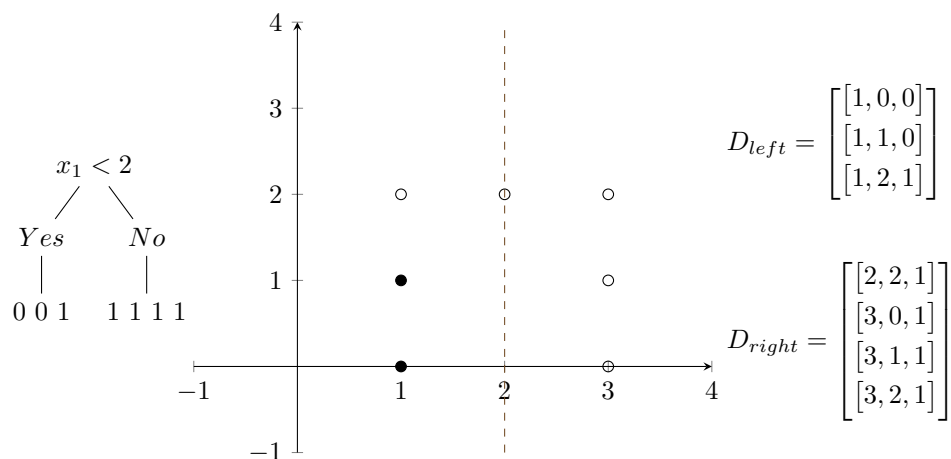
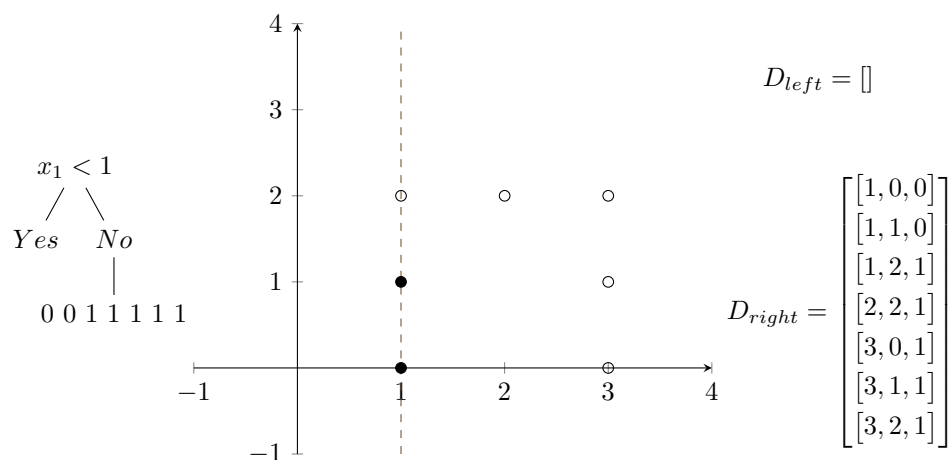
$$\begin{bmatrix} [1, 0, 0] \\ [1, 1, 0] \\ [1, 2, 1] \end{bmatrix}$$

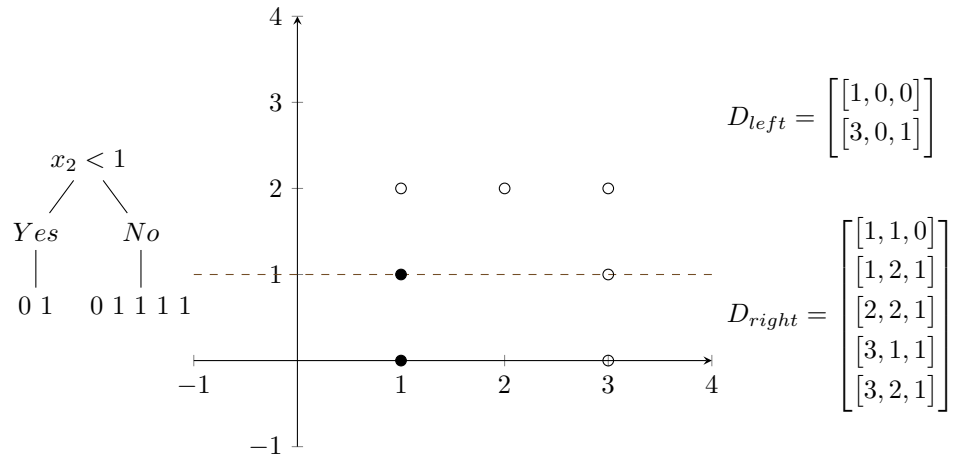
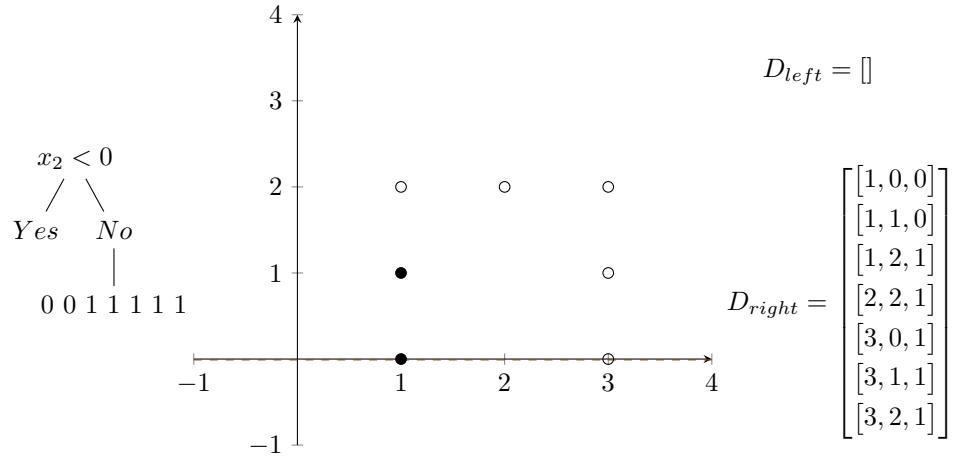
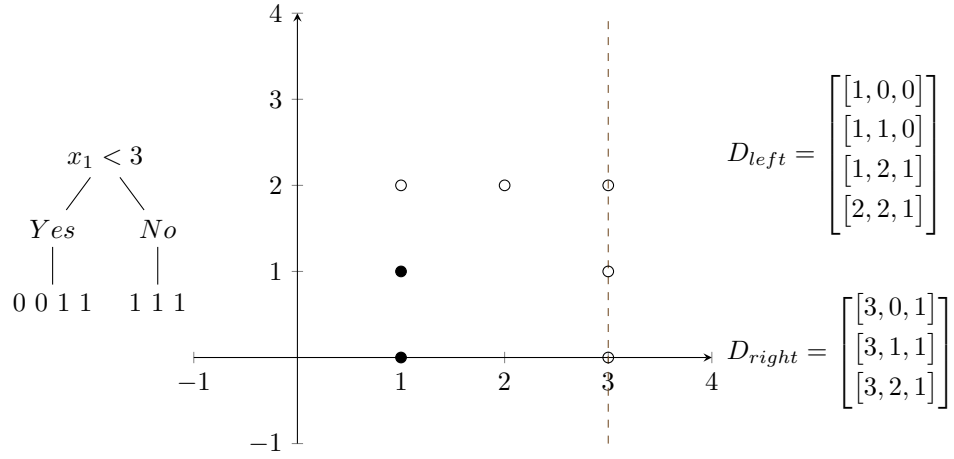
and the right child matrix:

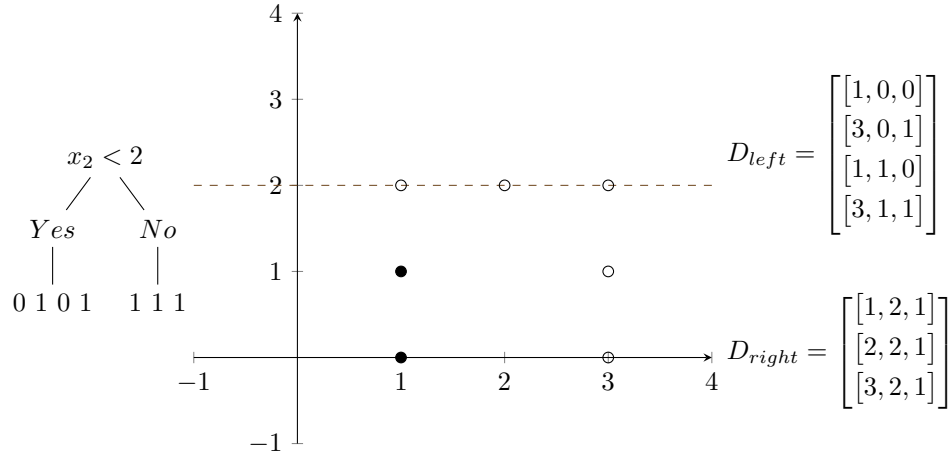
$$\begin{bmatrix} [2, 2, 1] \\ [3, 0, 1] \\ [3, 1, 1] \\ [3, 2, 1] \end{bmatrix}$$

## 4 Summary

To ensure complete understanding and summarize the topics we've covered, let's loop through all the features and each of the possible thresholds, splitting the matrix for each one.







## 5 The Full Tree

Let's build an entire decision tree. First, consider the following data.

x1	x2	y
0	0	0
1	0	0
1	1	0
0	1	0
1	3	1
0	3	1
2	2	1
4	1	1
3	1	1
3	2	1

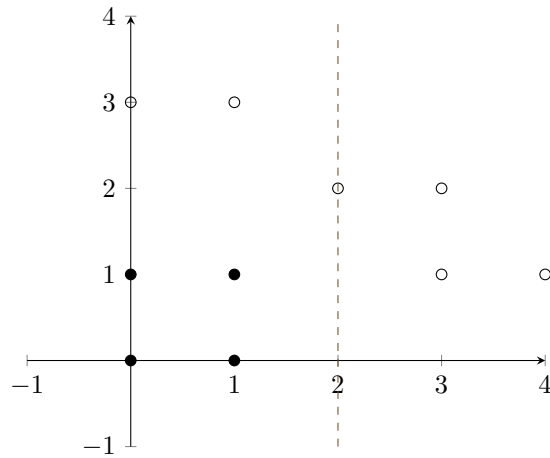
We graph it and calculate that  $x_1 < 2$  gives us the greatest information gain.

$$I(D_p) = 1 - \left(\frac{4}{10}\right)^2 - \left(\frac{6}{10}\right)^2 = 1 - \frac{16}{100} - \frac{36}{100} = \frac{100 - 52}{100} = \frac{48}{100} = \frac{12}{25}$$

$$I(D_{left}) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 1 - \frac{16}{36} - \frac{4}{36} = \frac{16}{36} = \frac{4}{9}$$

$$I(D_{right}) = 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 = 1 - 1 = 0$$

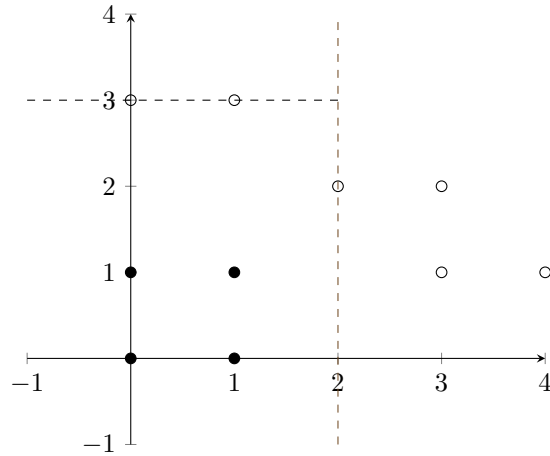
$$IG(D_p, f) = \frac{12}{25} - \frac{6}{10} * \frac{4}{9} - \frac{4}{10} * 0 = \frac{12}{25} - \frac{4}{15} = \frac{36}{75} - \frac{20}{75} = \frac{16}{75}$$



The dataset is split:

$$D_{left} = \begin{bmatrix} [0, 0, 0] \\ [1, 0, 0] \\ [1, 1, 0] \\ [0, 1, 0] \\ [0, 3, 1] \\ [1, 3, 1] \end{bmatrix} \quad D_{right} = \begin{bmatrix} [2, 2, 1] \\ [3, 1, 1] \\ [3, 2, 1] \\ [4, 1, 1] \end{bmatrix}$$

Now we make our next decision. We start with the left child, and it is obvious that  $x_2 < 3$  generates the greatest information gain. (Again, our thresholds are restricted to coordinates in our dataset).





The data is now split

$$D_{left} = \begin{bmatrix} [0, 0, 0] \\ [1, 0, 0] \\ [1, 1, 0] \\ [0, 1, 0] \end{bmatrix}$$
$$D_{right} = \begin{bmatrix} [0, 3, 1] \\ [1, 3, 1] \end{bmatrix}$$

We can now calculate the information gain. We know the impurity of the parent dataset, which was the impurity of the left child on the last split.

$$I(D_p) = \frac{4}{9}$$

Now, we calculate the children impurities to find our information gain:

$$I(D_{left}) = 0$$

$$I(D_{right}) = 0$$

$$IG(D_p, f) = \frac{4}{9} - \frac{4}{6} * 0 - \frac{2}{6} * 0 = \frac{4}{9}$$

All our data has been classified correctly, so the decision tree is done. We will leave it as an exercise to the reader to draw the decision tree.

## 6 Code

Here's how we recommend you structure your decision tree code for the competition. It should be noted: do not split recursively until all data is pure. You will overfit the testing data. Overfitting means that we train a model that memorizes the training data, rather than finding the generalized patterns embedded. When you overfit your data, the training accuracy will be high, but the testing accuracy will be low. Instead, set a baseline for your information gain. If the maximum information gain is less than, say, 0.1, then stop the recursion and make it a leaf of the tree. Play with this baseline to see how your training and testing accuracy change. Remember: your goal is the greatest testing accuracy.

```
def calculateImpurity(matrix):  
    #do calculations  
    return impurity  
  
def splitmatrix(matrix, feature, threshold):  
  
    #split matrix into leftchild, rightchild  
    return leftchild, rightchild
```

```

def informationGain(parent, leftchild, rightchild):
    #do calculations
    return infoGain

def bestsplit(parent, depth):
    for each feature in parent:
        for each threshold in feature:
            lchild, rchild = splitmatrix(parent, feature, threshold)
            igan = informationGain(parent, lchild, rchild)

            #if igan is the greatest:
            #    save the leftchild, rightchild, feature, threshold

    #print this node (feature, threshold) of the decision tree
    # hint: use depth to indent
    #recur on leftchild
    #recur on rightchild

```

## 7 Competition

- Do the Kaggle competition at (UPDATE LINK STUFF HERE). It ends in a week on (UPDATE DATE HERE). Refer to the lecture on how to set up Kaggle if you have questions, or feel free to message us.