

# Entropy\_ID3\_Exercise.pdf

by Derek Dilger, for the team

	Color	Shape	Size	class
1	red	square	big	+
2	blue	square	big	+
3	red	round	small	-
4	green	square	small	-
5	red	round	big	+
6	green	round	big	-

each column header is an attribute

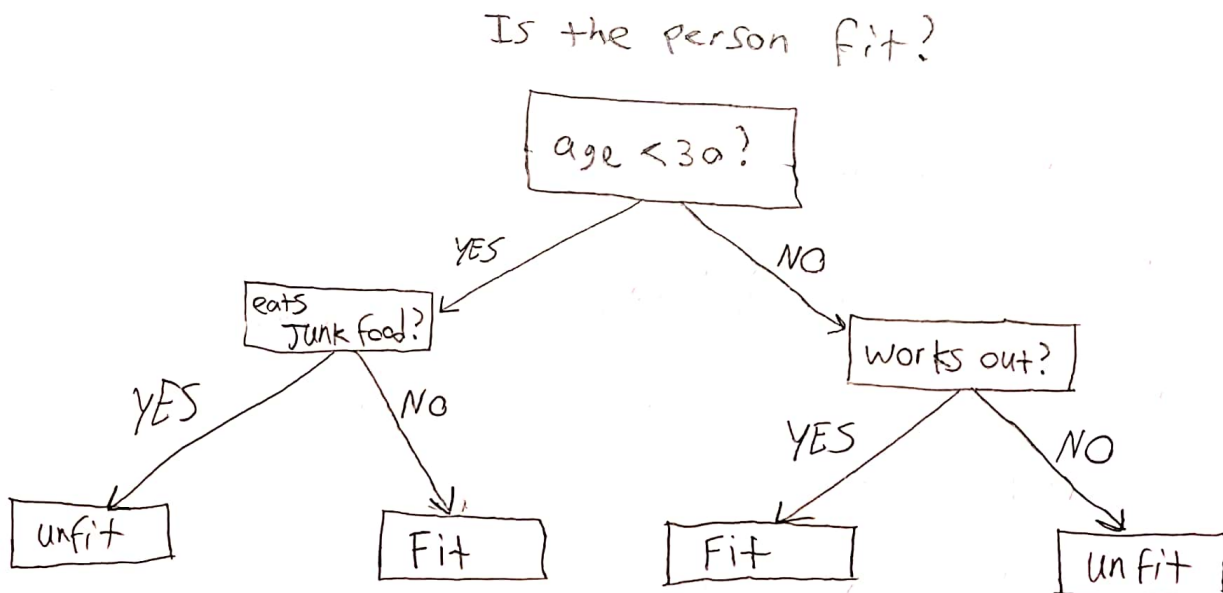
each row is like a tuple.  
here's tuple #1

each tuple here has 4 elements, but this generalizes to n elements.

↑ ↑ ↑  
"data set"

↑  
this column is the "target attribute".  
Our target attribute may take on two values, "+" and "-".

Here's an example of what we're working towards.  
This example does not use our data.



How was this tree made? By using the ID3 algorithm.  
This small fitness example can be done by hand. But with large data sets, the ID3 algorithm is used to make a large tree.

Here we will make the decision tree one iteration at a time. ID3 uses entropy to calculate which attribute is the most immediately discriminating. (as in, it's a greedy algorithm)

let our dataset be " $S$ "

let  $n$  be the total number of classes in the target attribute (here  $n=2$ . "+" or "-")

let  $p_i$  be the probability of being classed as " $i$ " in other words, 
$$\frac{\text{number of rows with class } i \text{ in the target column}}{\text{total number of rows}}$$

let entropy be calculated as

$$\text{Entropy}(S) = - \sum_{i=1}^n p_i * \log_2(p_i)$$

let Information Gain be calculated as

$$\text{IG}(S, A) = \text{Entropy}(S) - \sum ( (|S_v| / |S|) * \text{Entropy}(S_v) )$$

Where  $S_v$  is the set of rows in  $S$  for which the feature column  $A$  has value  $v$ ,  $|S_v|$  is the number of rows in  $S_v$ , and  $|S|$  is the number of rows in  $S$ .

Color	Shape	Size	class
red	square	big	+
blue	square	big	+
red	round	small	-
green	square	small	-
red	round	big	+
green	round	big	-

$$\text{Entropy}(S) = -\left(\frac{3}{6}\right) \log_2\left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log_2\left(\frac{3}{6}\right) = -\left(\frac{1}{2}\right)(-1) - \left(\frac{1}{2}\right)(-1) = 1.0$$

Now calculate IG for each feature (i.e. find which is most discriminating)

IG for Color:

$$|S| = 6 \quad (\text{total rows})$$

$$\left\{ \begin{array}{l} \text{For } v = \text{red}, |S_v| = 3 \leftarrow (\text{red rows}) \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{Entropy}(S_v) = -\left(\frac{2}{3}\right) \log_2\left(\frac{2}{3}\right) - \left(\frac{1}{3}\right) \log_2\left(\frac{1}{3}\right) = 0.9182958 \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{For } v = \text{blue}, |S_v| = 1 \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{Entropy}(S_v) = -\left(\frac{1}{1}\right) \log_2\left(\frac{1}{1}\right) - \underbrace{\left(\frac{0}{1}\right) \log_2\left(\frac{0}{1}\right)}_{\text{nonsense artifact of formula}} = 0 \end{array} \right.$$

note zero entropy because "blue" with it's 1 sample size, appears to predict "class" (perfectly).

$$\left\{ \begin{array}{l} \text{For } v = \text{green}, |S_v| = 2 \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{Entropy}(S_v) = 0 \leftarrow \text{zero, like blue because both greens predict "-" exactly and only.} \end{array} \right.$$

$$\begin{aligned} \text{IG}(S, \text{Color}) &= \text{Entropy}(S) - \left(\frac{|S_{\text{red}}|}{|S|}\right) \text{Entropy}(S_{\text{red}}) - \left(\frac{|S_{\text{blue}}|}{|S|}\right) \text{Entropy}(S_{\text{blue}}) - \left(\frac{|S_{\text{green}}|}{|S|}\right) \text{Entropy}(S_{\text{green}}) \\ &= 1.0 - \left(\frac{3}{6}\right)(0.9182958) - (0) - 0 = \boxed{0.54085} \end{aligned}$$

IG for Shape:

For  $V = \text{square}$ ,  $|S_V| = 3$

Coincidentally same as  
Entropy( $S_{\text{red}}$ )  
✓

$$\text{Entropy}(S_V) = -\left(\frac{2}{3}\right)\left(\log_2\left(\frac{2}{3}\right)\right) - \left(\frac{1}{3}\right)\left(\log_2\left(\frac{1}{3}\right)\right) = 0.9182958$$

The computations are laborious and excessive. The general formulas given describe fully how to calculate each value as needed for input into IG( $S$ , class). The rest are done on calculator. Edge cases have also been demonstrated, and how to deal with them.

Overall so far, we have

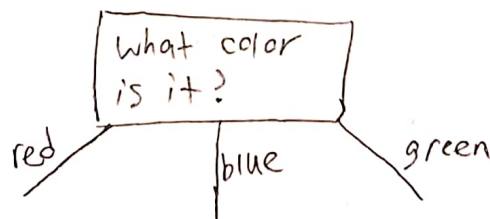
$$\text{IG}(S, \text{color}) = 0.54085$$

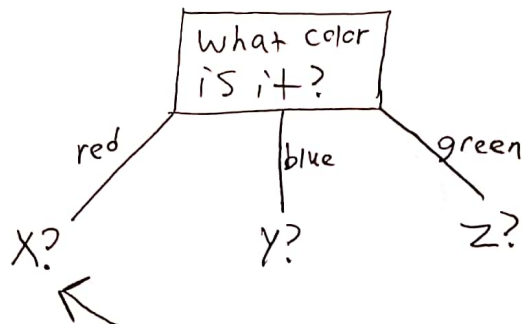
$$\text{IG}(S, \text{shape}) = 0.08170$$

$$\text{IG}(S, \text{size}) = 0.459148$$

So therefore we choose to discriminate based on color to begin with, as it has the highest information gain.

The tree so far is:





So what question do we put here?

We ask about the attribute with the highest Information gain, just as before.

However, now our dataset has changed to contain only tuples with color == red.

i.e.

Color	Shape	Size	Class
red	square	big	+
red	round	small	-
red	round	big	+

$S = \text{dataset for "X?"}$

$$IG(S, \text{color}) = 0.811278$$

$$IG(S, \text{Shape}) = 1.0$$

$$IG(S, \text{size}) = 1.5$$

So we ask about Size at "X?"

We do the same for blue & green, Y? and Z?.

dataset for "Y?"

Color	Shape	Size	Class
blue	square	big	+

note this degenerate dataset with one tuple. We don't ask a question here, but rather, conclude that the class is "+". The tuple correlates exactly with class's outcome.

dataset for "Z?"

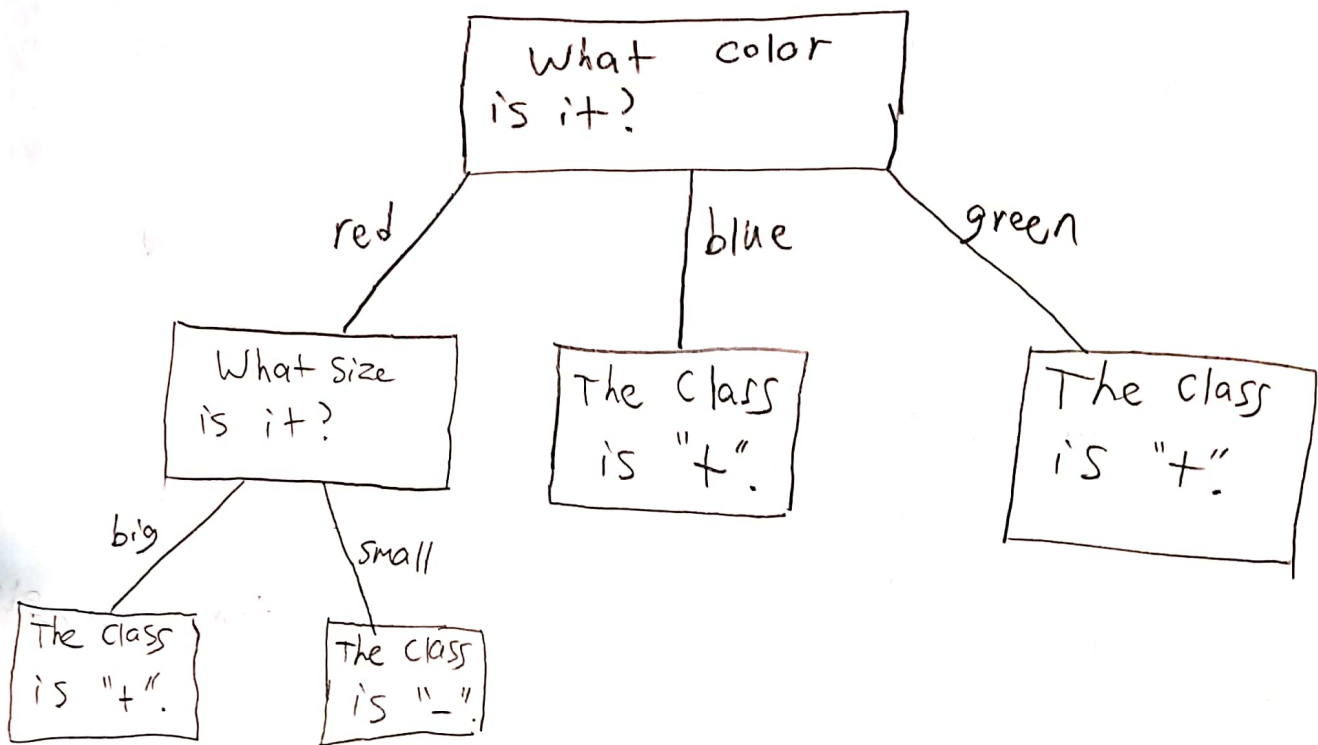
Color	Shape	Size	Class
green	square	small	-
green	round	big	-

degenerate.

Perfect correlation between color and class.



Now the decision tree is:



Note that

(red and big) = class "+"

(red and small) = class "-"

for all tuples.

(from the previous table. It would be computed as before.)

So we have classified everything.

<end of part (i) in project>

(ii) If "pattern of shirt" attribute is added, only a few possibilities may occur. Either it offers more information gain at a certain step of classification, or it does not. In the case it does, it may or may not allow us to classify using fewer steps. In the case it does not, we will simply use the tree as shown, because

the shirt's didn't help.

With a data size of millions of shirts, almost certainly the shirt color attribute would help classify at some step in the tree.\* This decision would influence millions of dollars of revenue.

If the manager and CEO are adequately paying attention (doing their job), then certainly the data scientist would make an impact.

\* (unlike in our dwarf tree with small dataset)