

## Dataset

## Decision Tree Example: ID3

S.No	color	Attribute outline	Dot	shape
1	green	dashed	NO	triangle
2	green	dashed	Yes	triangle
3	yellow	dashed	NO	Square
4	red	dashed	NO	Square
5	red	Solid	NO	Square
6	red	Solid	Yes	triangle
7	green	Solid	NO	Square
8	green	dashed	NO	triangle
9	yellow	Solid	Yes	Square
10	red	Solid	NO	Square
11	green	Solid	Yes	Square
12	yellow	dashed	Yes	Square
13	yellow	Solid	NO	Square
14	red	dashed	Yes	triangle

Entropy formula :

$$E(S) = \sum_{i=1}^c -P_i \log_2(P_i)$$

Information gain formula:

$$\text{Gain}(T, x) = \text{Entropy}(T) - \text{Entropy}(T, x)$$

1. We need to get the attribute which is giving maximum information among all the attributes. Hence we need to calculate Information gain for every attribute. After this we need to identify which attribute has highest information gain. Then we will take that attribute for tree construction.

we are planning to calculate information gain ~~for~~ each attribute.

The attributes in the dataset are:

1. color
2. outline
3. dot

Attribute: color

possible values of attribute color: { green, yellow, red }

Attribute: outline

possible values of attribute outline: { dashed, solid }

Attribute: dot

possible values of attribute dot: { No, Yes }

We need to calculate entropy of entire dataset & entropy of individual attribute of dataset ~~the~~ ~~set~~ ~~to~~ find gain.

Entropy of the dataset:

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

In the dataset, we find 5 triangles & 9 squares.

$\Rightarrow S = [90, 54]$  out of 14 records in dataset.

$$\text{Entropy}(S) = - \frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

$$= 0.94$$

Similarly we need to compute entropy of green, yellow, and red for color attribute.

$\Rightarrow$  total records related to green in color attribute  
 $= 5$ .

out of them 5-records, 3 triangle & 2 squares:

$$\Rightarrow S_{\text{green}} \leftarrow [2\Delta, 3\triangle]$$

$$\Rightarrow \text{Entropy}(S_{\text{green}}) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$
$$= 0.971$$

$\Rightarrow$  total records related to yellow in color attribute.  
 $= 4$

out of them 4-records, 4 squares & 0-triangle.

$$S_{\text{yellow}} \leftarrow [4\Box, 0\triangle]$$

$$\text{Entropy}(S_{\text{yellow}}) = -\frac{4}{4} \log_2\left(\frac{4}{4}\right) - \frac{0}{4} \log_2\left(\frac{0}{4}\right)$$

~~no~~  
no need to calculate. It is zero only.

$$= 0$$

$\Rightarrow$  total records related to red in color attribute

$$S_{\text{red}} \leftarrow [3\Box, 2\triangle] \quad \text{total} = 5$$

$$\text{Entropy}(S_{\text{red}}) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$
$$= 0.971$$

We need to find:

Now, Information gain for color attribute.

$$\text{Gain}(S, \text{color}) = \text{Entropy}(S) - \sum \frac{|S_i|}{|S|} \text{Entropy}(S_i)$$

$$\text{Gain}(S, \text{Color}) = \text{Entropy}(S) - \sum_{v \in \{\text{green, red, yellow}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$S_v \rightarrow$  The NO. of possible values of 'v'

$$\text{Gain}(S, \text{color}) = \text{Entropy}(S) - \frac{5}{14} \text{Entropy}(S_{\text{green}}) - \frac{4}{14} \text{Entropy}(S_{\text{yellow}}) - \frac{5}{14} \text{Entropy}(S_{\text{red}})$$

$$\Rightarrow \text{Gain}(S, \text{color}) = 0.94 - \frac{5}{14} \times 0.971 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.971$$

$$= 0.2464.$$

Similarly, we need to calculate Information gain for the remaining attributes, also.

Attribute: outline = { dashed, solid }

$$S = [9 \square, 5 \Delta] \quad \text{Entropy}(S) = 0.94$$

$$S_{\text{dashed}} \leftarrow [3 \square, 4 \Delta] \quad \text{total} = 7 \text{ records.}$$

$$\text{Entropy}(S_{\text{dashed}}) = -\frac{3}{7} \log_2 \left( \frac{3}{7} \right) - \frac{4}{7} \log_2 \left( \frac{4}{7} \right)$$

$$= 0.9852$$

$$S_{\text{solid}} \leftarrow [6 \square, 1 \Delta] \quad \text{total} = 7 \text{ records.}$$

$$\text{Entropy}(S_{\text{solid}}) = -\frac{6}{7} \log_2 \left( \frac{6}{7} \right) - \frac{1}{7} \log_2 \left( \frac{1}{7} \right)$$

$$= 0.5916$$

$$\Rightarrow \text{Gain}(S, \text{outline}) = \text{Entropy}(S) - \sum_{v \in \{\text{dashed, solid}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, \text{outline}) = \text{Entropy}(S) - \frac{7}{14} \text{Entropy}(S_{\text{dashed}}) - \frac{7}{14}$$



$$\text{Gain}(S, \text{color}) = \text{Entropy}(S) - \sum_{v \in \left\{ \begin{smallmatrix} \text{green, red} \\ \text{yellow} \end{smallmatrix} \right\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$S_v \rightarrow$  The no. of possible values of 'v'

$$\text{Gain}(S, \text{color}) = \text{Entropy}(S) - \frac{5}{14} \text{Entropy}(S_{\text{green}}) - \frac{4}{14} \text{Entropy}(S_{\text{yellow}}) - \frac{5}{14} \text{Entropy}(S_{\text{red}})$$

$$\Rightarrow \text{Gain}(S, \text{color}) = 0.94 - \frac{5}{14} \times 0.971 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.971$$

$$= 0.2464$$

Similarly, we need to calculate Information Gain for the remaining attributes, also.

Attribute: outline = { dashed, solid }

$$S = [9 \square, 5 \Delta] \quad \text{Entropy}(S) = 0.94$$

$$S_{\text{dashed}} \leftarrow [3 \square, 4 \Delta] \quad \text{total} = 7 \text{ records.}$$

$$\text{Entropy}(S_{\text{dashed}}) = -\frac{3}{7} \log_2 \left( \frac{3}{7} \right) - \frac{4}{7} \log_2 \left( \frac{4}{7} \right)$$

$$= 0.9852$$

$$S_{\text{solid}} \leftarrow [6 \square, 1 \Delta] \quad \text{total} = 7 \text{ records.}$$

$$\text{Entropy}(S_{\text{solid}}) = -\frac{6}{7} \log_2 \left( \frac{6}{7} \right) - \frac{1}{7} \log_2 \left( \frac{1}{7} \right)$$

$$= 0.5916$$

$$\Rightarrow \text{Gain}(S, \text{outline}) = \text{Entropy}(S) - \sum_{v \in \left\{ \begin{smallmatrix} \text{dashed} \\ \text{solid} \end{smallmatrix} \right\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\text{Gain}(S, \text{outline}) = \text{Entropy}(S) - \frac{7}{14} \text{Entropy}(S_{\text{dashed}}) - \frac{7}{14} \text{Entropy}(S_{\text{solid}})$$

$$\text{gain}(S, \text{outline}) = 0.94 - \frac{7}{14} \times 0.9852 - \frac{7}{14} \times 0.5916$$

$$= 0.1516$$

Attribute: ~~out~~ = { No, Yes }

$$S = [9, 5]$$

$$\text{Entropy}(S) = 0.94$$

$$S_{\text{yes}} \leftarrow [3, 3] \quad \text{total} = 6 \text{ records.}$$

$$\text{Entropy}(S_{\text{yes}}) = -\frac{3}{6} \log_2 \left( \frac{3}{6} \right) - \frac{3}{6} \log_2 \left( \frac{3}{6} \right)$$

$$= 1$$

$$S_{\text{no}} \leftarrow [6, 2] \quad \text{total} = 8 \text{ records.}$$

$$\text{Entropy}(S_{\text{no}}) = -\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right)$$

$$= 0.8113$$

$$\text{gain}(S, \text{out}) = \text{Entropy}(S) - \sum_{v \in \{\text{yes}, \text{no}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= 0.94 - \frac{6}{14} \times 1.0 - \frac{8}{14} \times 0.8113$$

$$= 0.0478$$

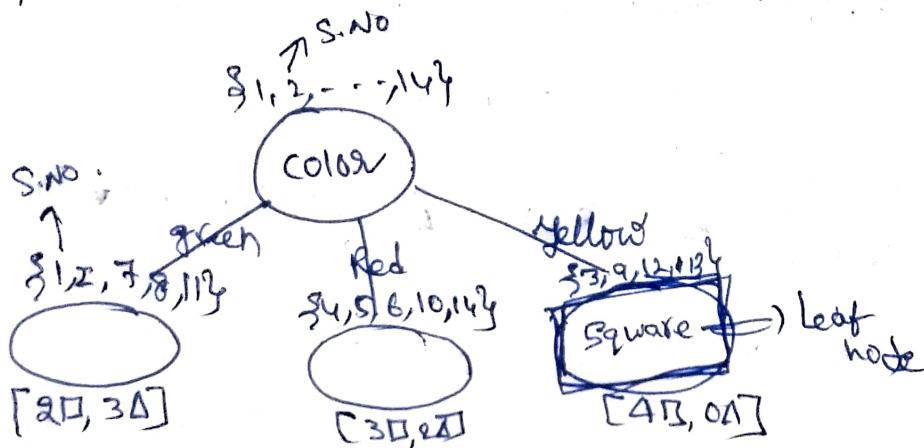
$$\text{gain}(S, \text{color}) = 0.2464$$

$$\text{gain}(S, \text{outline}) = 0.1516$$

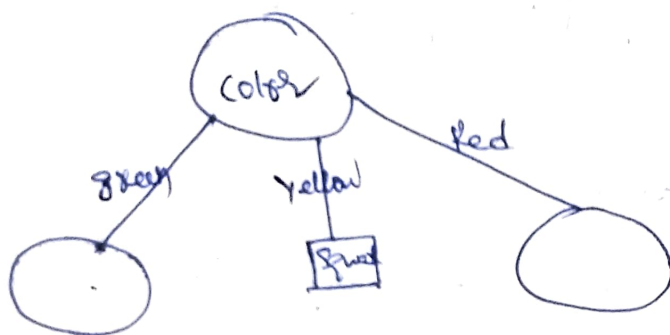
$$\text{gain}(S, \text{out}) = 0.0478$$

Here we can conclude that color attribute has highest gain. So, we consider color as root node.

The possibilities of color are green, red, and yellow.



11,



We need to compute true father. The tree at green & red. While as yellow become the leaf node.

Now we will start @ green branch side.

Now we need to take the examples of records which are related to green only. i.e. S.No {1, 2, 7, 8, 11}

S.No	outline	dot	shape
1	dashed	NO	triangle
2	dashed	yes	triangle
7	solid	NO	Square
8	dashed	NO	triangle
11	solid	yes	Square

$S_1(\text{color}) = [2□, 3△]$  total 5-records.

Choose attribute outline

$$\text{Entropy}(S_1, \text{green}) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) \\ = 0.97$$

Absolute: outline  
 $S_1(\text{dashed}) \leftarrow [00, 3A] \rightarrow \text{total} = 3 \text{ records}$

$$\text{Entropy}(S_1(\text{dashed})) = 0.0$$

$S_1(\text{solid}) \leftarrow [20, 0A] \rightarrow \text{total} = 2 \text{ records}$

$$\text{Entropy}(S_1(\text{solid})) = 0.$$

~~gain(S, green)~~

$$\text{gain}(S, \text{outline}) = \text{Entropy}(S_1) - \sum_{v \in \{ \text{dashed}, \text{solid} \}} \frac{|S_v|}{|S_1|} \text{Entropy}(S_v)$$

$$\text{gain}(S, \text{outline}) = 0.97 - \frac{3}{5} \times 0 - \frac{2}{5} \times 0 \\ = 0.97$$

Absolute: dot { NO, YES }

$S_1(\text{NO}) \leftarrow [10, 2A] \Rightarrow \text{total} = 3 \text{ records}$

$$\text{Entropy}(S_1(\text{NO})) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) \\ = 0.9183$$

$S_1(\text{Yes}) \leftarrow [10, 1A] \Rightarrow \text{total} = 2 \text{ records}$

$$\text{Entropy}(S_1(\text{Yes})) = -\frac{1}{1} \log_2\left(\frac{1}{1}\right) - \frac{1}{1} \log_2\left(\frac{1}{1}\right) \\ = 1.$$



$$\begin{aligned} \text{gain}(S_1, \text{dot}) &= \text{Entropy}(S_1) - \sum_{S \in S_1} \frac{|S|}{|S_1|} \cdot \text{Entropy}(S) \\ &= 0.97 - \frac{2}{5} \times 1.0 - \frac{3}{5} \times 0.918 \\ &= 0.0192 \end{aligned}$$

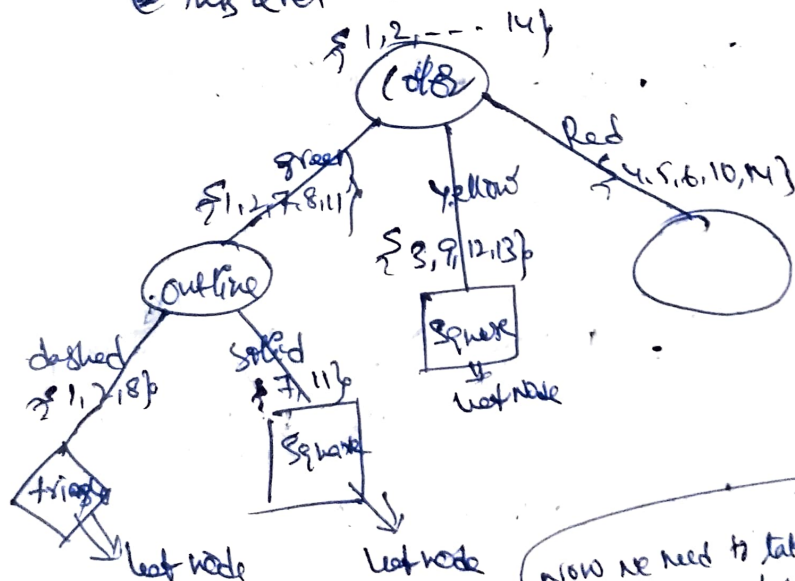
finally

$$\text{gain}(S_1, \text{outline}) = 0.97$$

$$\text{gain}(S_1, \text{dot}) = 0.0192$$

outline has highest Information gain, hence it can be taken as node at this level.

@ this level tree will be



Now we need to take examples which are ~~not~~ scaled to red on

S.NO	outline	dot	shape
4	dashed	NO	square
5	solid	NO	square
6	solid	YES	Triangle
10	solid	NO	square
14	dashed	YES	triangle

treat this as  
82

$$S_2 = [3\bar{D}, 2\Delta] \Rightarrow \text{total } 5\text{-records}$$

$$\begin{aligned} \text{Entropy}(S_2) &= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \\ &= 0.97 \end{aligned}$$

Attribute: outline  $\Rightarrow \{ \text{dashed, solid} \}$

$$S_2(\text{dashed}) \leftarrow [1\bar{D}, 1\Delta] \Rightarrow 2\text{-records}$$

$$\begin{aligned} \text{Entropy}(S_2(\text{dashed})) &= -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \\ &= 1 \end{aligned}$$

$$S_2(\text{solid}) \leftarrow [2\bar{D}, 1\Delta] \Rightarrow 3\text{-records}$$

$$\begin{aligned} \text{Entropy}(S_2(\text{solid})) &= -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) \\ &= 0.9183 \end{aligned}$$

$$\begin{aligned} \text{gain}(S_2, \text{outline}) &= \text{Entropy}(S_1) - \sum_{V \in \{ \text{dashed, solid} \}} \frac{|S_V|}{|S|} \text{Entropy}(S_V) \\ &= 0.97 - \frac{2}{5} \times 1.0 - \frac{3}{5} \times 0.918 \\ &= 0.0192 \end{aligned}$$

Attribute: out  $\Rightarrow \{ \text{no, yes} \}$

$$S_2(\text{no}) \leftarrow [3\bar{D}, 0\Delta] \Rightarrow 3\text{-records}$$

$$\begin{aligned} \text{Entropy}(S_2(\text{no})) &= -\frac{3}{3} \log_2\left(\frac{3}{3}\right) - \frac{0}{3} \log_2\left(\frac{0}{3}\right) \\ &= 0 \end{aligned}$$

$$S_2(\text{yes}) \leftarrow [0\bar{D}, 2\Delta]$$

~~$S_2(yes) \leftarrow [3]$~~

$$S_2(yes) \leftarrow [0], 2A \quad \text{total} = 2 \text{ records}$$

$$\text{Entropy}(S_2(yes)) = 0.$$

$$\text{gain}(S_2, \text{dot}) = \text{Entropy}(S_2) - \sum_{v \in \{yes, no\}} \frac{|S_v|}{|S_2|} \text{Entropy}(S_v)$$

~~$= 0.97$~~

$$= 0.97 - \frac{2}{5} \times 0.0 - \frac{3}{5} \times 0. = 0.97$$

finally

$$\text{gain}(S_2, \text{outline}) = 0.0192$$

$$\text{gain}(S_2, \text{dot}) = 0.97$$

As "dot" attribute has highest information gain, it can be taken as node at this stage.



final tree

