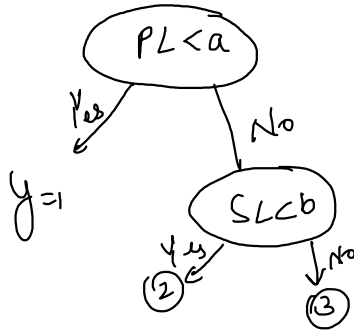
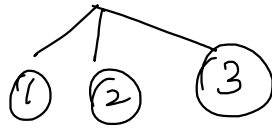


Decision Trees

IRIS Dataset = [SL, SW, PL, PW]



if $PL < a$:

$y_i = 1$

else:

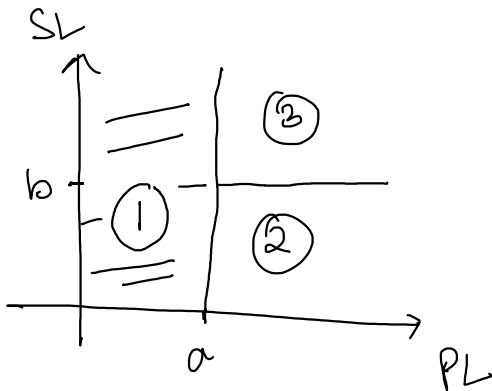
if $SL < b$:

$y_i = 2$

else:

$y_i = 3$

Recursive Partitioning: \rightarrow (Axis Parallel hyperplanes)



if $PL < a$:

$y_i = 1$

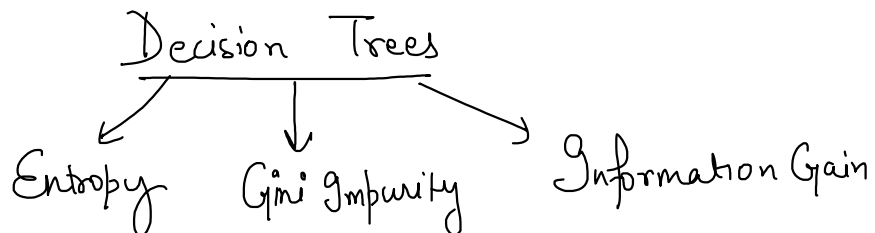
else:

if $SL < b$:

$y_i = 2$

else:

$y_i = 3$

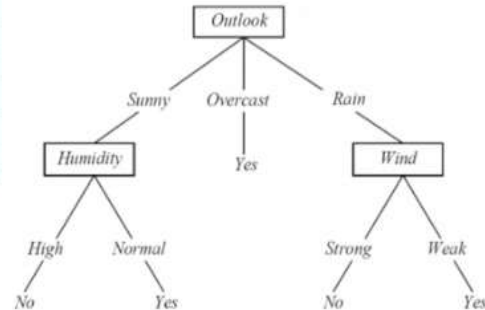


Entropy \rightarrow Randomness in dataset

$$H_D(Y) = - \sum_{i=1}^n p_i \lg(p_i)$$

$$\lg \Rightarrow \log_2$$

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No



Parent's Entropy

$$H_D(Y) = - \sum_{i=1}^n p_i \lg(p_i)$$

$$= - P(\text{Yes}) \lg P(\text{Yes}) - P(\text{No}) \lg P(\text{No})$$

$$P(\text{Yes}) = \frac{9}{14} \quad P(\text{No}) = \frac{5}{14}$$

$$H_D(Y) = - \frac{9}{14} \times \lg\left(\frac{9}{14}\right) - \frac{5}{14} \lg\left(\frac{5}{14}\right)$$

$$= 0.94$$

entropy of each column

0.94

Outlook

Sunny (5)	(2Y, 3N) →	$-\frac{2}{5} \lg \frac{2}{5} - \frac{3}{5} \lg \frac{3}{5} = 0.97$
Overcast (4)	(4Y, 0N) →	$-\frac{4}{4} \lg \frac{4}{4} - 0 \lg 0 = 0$
Rainy (3)	(3Y, 2N) →	$-\frac{3}{5} \lg \frac{3}{5} - \frac{2}{5} \lg \frac{2}{5} = 0.97$

weighted entropy = $\frac{5}{14} \times 0.97 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.97$

⇒ $0.97 \times \frac{5}{7} = 0.69$

Temperature

Hot (4)	(2Y, 2N) =	$-\frac{2}{4} \lg \frac{2}{4} - \frac{2}{4} \lg \frac{2}{4} = 1$
Mild (6)	(4Y, 2N) =	$-\frac{4}{6} \lg \frac{4}{6} - \frac{2}{6} \lg \frac{2}{6} = 0.91$
Cool (4)	(3Y, 1N) =	$-\frac{3}{4} \lg \frac{3}{4} - \frac{1}{4} \lg \frac{1}{4} = 0.81$

weighted entropy = $\frac{4}{14} \times 1 + \frac{6}{14} \times 0.91 + \frac{4}{14} \times 0.81$

$H_b(Y, temp)$

= 0.91

Humidity

High (7)	(3Y, 4N) →	$-\frac{3}{7} \lg \frac{3}{7} - \frac{4}{7} \lg \frac{4}{7} = 0.98$
Normal (7)	(6Y, 1N) →	$-\frac{6}{7} \lg \frac{6}{7} - \frac{1}{7} \lg \frac{1}{7} = 0.59$

$$\text{Normal} \rightarrow -\frac{6}{7} \lg \frac{6}{7} - \frac{1}{7} \lg \frac{1}{7} = 0.59$$

$$\begin{aligned} H_b(Y, \text{humidity}) \\ \text{weighted entropy} &= \frac{7}{14} \times 0.98 + \frac{7}{14} \times 0.59 \\ &= \frac{7}{14} \times 1.57 = 0.78 \end{aligned}$$

$$\begin{aligned} \text{Windy} \begin{cases} \xrightarrow{6 (3Y, 3N)} \text{True} & \rightarrow -\frac{3}{6} \lg \frac{3}{6} - \frac{3}{6} \lg \frac{3}{6} = 1 \\ \xrightarrow{8 (6Y, 2N)} \text{False} & \rightarrow -\frac{6}{8} \lg \frac{6}{8} - \frac{2}{8} \lg \frac{2}{8} = 0.81 \end{cases} \end{aligned}$$

$$\begin{aligned} \text{weighted entropy} \\ H_b(Y, \text{Windy}) &= \frac{6}{14} \times 1 + \frac{8}{14} \times 0.81 \\ &= 0.89 \end{aligned}$$

Choosing the column for split

a) Compare weighted entropies of column & choose column with least entropy

outlook	temp	humidity	windy
0.69	0.91	0.78	0.89

0.69

0.91

0.94

0.94

(b) Information Gain (IG) = Parent's entropy - weighted entropy of each column

highest IG \Rightarrow $IG(Y, outlook) = 0.94 - 0.69 = 0.25$

$$IG(Y, temp) = 0.94 - 0.91 = 0.03$$

$$IG(Y, humidity) = 0.94 - 0.78 = 0.16$$

$$IG(Y, windy) = 0.94 - 0.89 = 0.05$$

outlook is chosen because it has highest information gain.

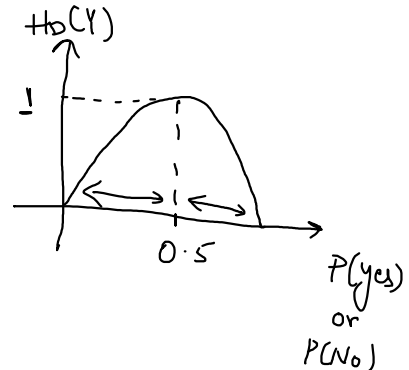
RR
Properties

(1) $P(\text{yes}) = \frac{1}{2}$ $P(\text{No}) = \frac{1}{2}$

$$H_D(Y) = -\frac{1}{2} \lg \frac{1}{2} - \frac{1}{2} \lg \frac{1}{2} = 1$$

(2) $P(\text{yes}) = 1$ $P(\text{No}) = 0$

$$H_D(Y) = -\frac{1}{1} \lg \frac{1}{1} - 0 \lg 0 = 0$$



Gini Impurity $I_G \neq I_G$ (I_G)

$$I_G(Y) = 1 - \sum_{i=1}^n p_i^2$$

→ binary classification ⇒

$$I_G(Y) = 1 - [p(+)^2 + p(-)^2]$$

→ multiclass classification ⇒

$$I_G(Y) = 1 - [p(y_1)^2 + p(y_2)^2 + p(y_3)^2 + \dots + p(y_n)^2]$$

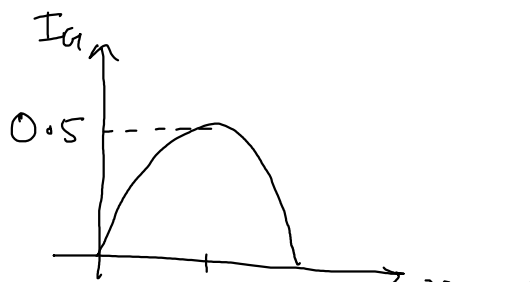
Properties:

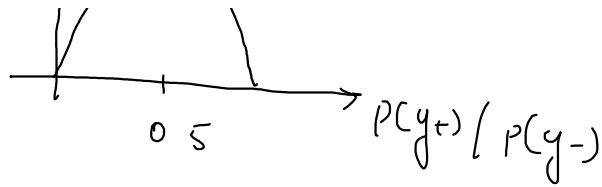
① $p(y_+) = 0.5$ $p(y_-) = 0.5$

$$I_G = 1 - [0.5^2 + 0.5^2] = 1 - [0.5] = 0.5$$

② $p(y_+) = 1$ $p(y_-) = 0$

$$I_G = 1 - [1^2 + 0^2] = 0$$





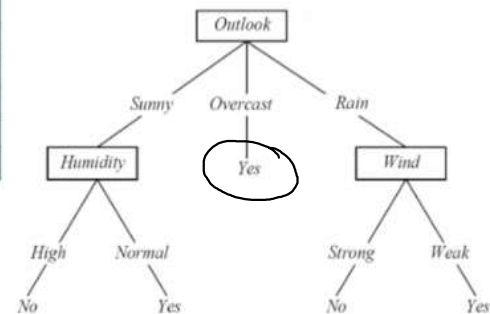
Comparison of Gini Impurity & Entropy:

$$1) \quad H_D(Y) = - \sum_{i=1}^n p_i \lg p_i \quad I_G = 1 - \sum_{i=1}^n p_i^2$$

Since, I_G is computationally efficient, it can be used for larger datasets

Entropy shouldn't be used for large dataset.

Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No



$$\begin{aligned}
 \text{Parent's } I_G &= 1 - [P(\text{yes})^2 + P(\text{No})^2] \\
 &= 1 - \left[\left(\frac{9}{14} \right)^2 + \left(\frac{5}{14} \right)^2 \right] \\
 &= 0.45
 \end{aligned}$$

I_{C_j} for each column:

0.45

Outlook

$$\begin{cases} 5 \text{ (2Y, 3N)} \text{ Sunny} \Rightarrow 1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right] = 0.48 \\ 4 \text{ (4Y, 0N)} \text{ overcast} \Rightarrow 1 - \left[\left(\frac{4}{4} \right)^2 + 0^2 \right] = 0 \\ 5 \text{ (3Y, 2N)} \text{ rainy} \Rightarrow 1 - \left[\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right] = 0.48 \end{cases}$$

weighted $I_{C_j} \Rightarrow \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48$

$\Rightarrow \frac{5}{14} \times 0.96 \Rightarrow 0.342$

Temperature

$$\begin{cases} 4 \text{ (2Y, 2N)} \text{ Hot} \Rightarrow 1 - \left[\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right] = 0.5 \\ 6 \text{ (4Y, 2N)} \text{ Mild} \Rightarrow 1 - \left[\left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right] = 0.444 \\ 4 \text{ (3Y, 1N)} \text{ Cool} \Rightarrow 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.375 \end{cases}$$

weighted $I_{C_j} = \frac{4}{14} \times 0.5 + \frac{6}{14} \times 0.444 + \frac{4}{14} \times 0.375 = 0.44$

Humidity

$$\begin{cases} 7 \text{ (3Y, 4N)} \text{ High} \Rightarrow 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 0.48 \\ 7 \text{ (6Y, 1N)} \text{ Normal} \Rightarrow 1 - \left[\left(\frac{6}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right] = 0.24 \end{cases}$$

weighted $I_{C_j} = \frac{7}{14} \times 0.48 + \frac{7}{14} \times 0.24 = 0.36$

$\frac{6 \text{ (3Y, 3N)}}{\text{Windy}} \Rightarrow 1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right] = 0.5$

$$\text{Windy} \begin{cases} 6 (3Y, 3N) \text{ True} \Rightarrow 1 - \left[\left(\frac{3}{6} \right)^2 + \left(\frac{3}{6} \right)^2 \right] = 0.5 \\ 8 (6Y, 2N) \text{ False} \Rightarrow 1 - \left[\left(\frac{6}{8} \right)^2 + \left(\frac{2}{8} \right)^2 \right] = 0.375 \end{cases}$$

$$\text{Weighted } I_G \Rightarrow \frac{6}{14} \times 0.5 + \frac{8}{14} \times 0.375$$

$$\Rightarrow 0.42$$

Choose column for split.

a) choose column with least I_G

outlook	temp	humidity	windy
0.342	0.44	0.36	0.42

b) Info Gain: Parent's I_G - weighted I_G for each column

highest I_G
chosen for split

$$I_G(Y, \text{outlook}) = 0.45 - 0.342 = 0.108$$

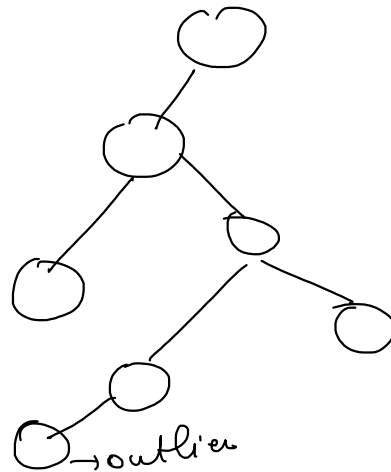
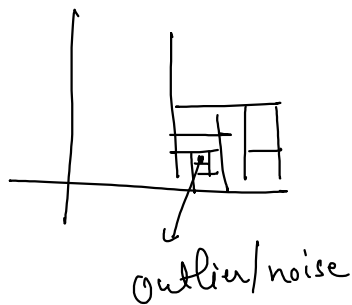
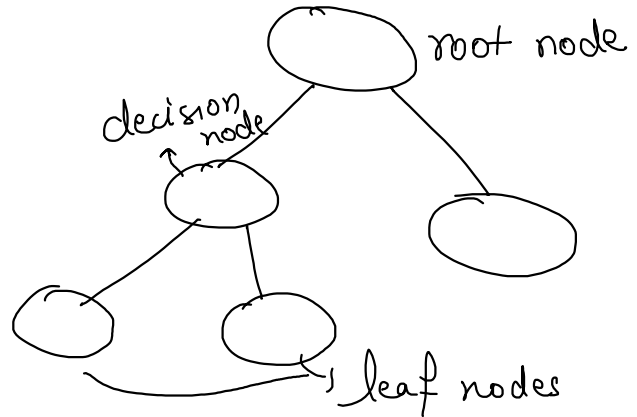
$$I_G(Y, \text{temp}) = 0.45 - 0.44 = 0.01$$

$$I_G(Y, \text{humidity}) = 0.45 - 0.36 = 0.09$$

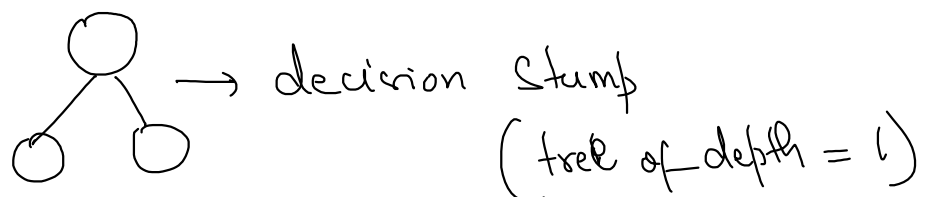
$$I_G(Y, \text{windy}) = 0.45 - 0.42 = 0.03$$

When to stop a tree:

- pure node
- If you have very few samples in terminal nodes



Hyperparameter: $\text{max_depth} \uparrow \Rightarrow \text{height (depth)} \uparrow \Rightarrow \text{overfitting} \uparrow$
(max value = 32)



Splitting Numerical

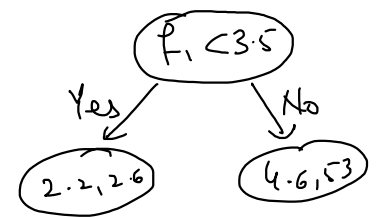
→ features

f_1	y
2.2	1
2.6	1
3.5	0
4.6	1
5.3	0

① sort the f_1

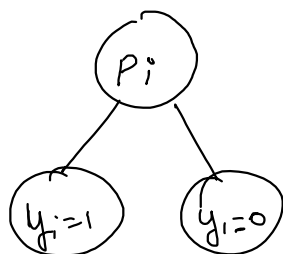
② $f_1 < 2.6$ $f_1 < 5.3$

$f_1 < 3.5 \rightarrow$ lowest I_H/H
 $f_1 < 4.6$



Feature Engineering: column with lot of categories

PINCODE

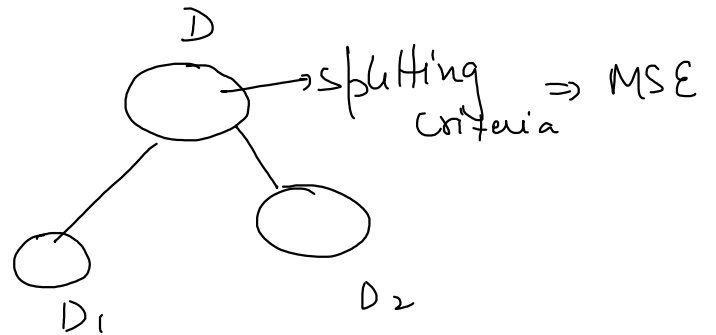


P	P _{pool}	y_i
P_j	C.P	1
P_j		0
		0
		0
		0
		0

$$P(y=0/p_j) = \frac{P_j \cap y=0}{P_j}$$

$$P(y=1/p_j) = \frac{P_j \cap y}{P_j}$$

Regression in DT



Advantages:

- easy to interpret
- important features can be extracted
- No need to standardize