

Q1)

If color is green?

Level 1: Root

$$\text{Entropy} = \sum_{i=1}^c -P_i \log_2 P_i$$

Color	Diam	label
green	3	Apple
yellow	3	Apple
Red	1	Grape
Yellow	3	Lemon

$$\begin{aligned} \text{Entropy} &= -P(A) \log_2 P(A) - P(G) \log_2 P(G) - P(L) \log_2 P(L) \\ &= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{2}{5} \log_2 \frac{2}{5} - \frac{1}{5} \log_2 \frac{1}{5} \\ &= -0.4 \log_2 0.4 - 0.4 \log_2 0.4 - 0.2 \log_2 0.2 \\ &\approx 0.52 + 0.52 + 0.46 \approx 1.52 \end{aligned}$$

Avg Imp = $\frac{5}{5} \times 1.52 = \underline{\underline{1.52}}$

Level 2:

left subtree: $-P(A) \log_2 P(A) = -(\frac{1}{2}) \log_2 (\frac{1}{2}) = -1 \log_2 1 \approx \underline{\underline{0}}$

Right Subtree: $-P(A) \log_2 P(A) - P(G) \log_2 P(G) - P(L) \log_2 P(L)$

$$= -\frac{1}{4} \log_2 \frac{1}{4} - \frac{2}{4} \log_2 \frac{2}{4} = \frac{1}{4} \log_2 \frac{1}{4}$$

$$= -0.25 \log_2 0.25 - 0.5 \log_2 0.5 - 0.25 \log_2 0.25$$

$$= 0.5 + 0.5 + 0.5 \approx \underline{\underline{1.5}}$$

Average impurity or weighted Average = $\frac{4}{5} \times 1.5 + \frac{1}{5} \times 0 = \underline{\underline{1.2}}$

Information gain = $1.52 - 1.2 = \underline{\underline{0.32}}$

If the feature color is green Information Gain's from 2 methods are:

Gini = 0.14

Entropy = 0.32 Info gain is Entropy is higher than Gini

IF diameter $>= 3$

Level 1: Root

Entropy = $-P(A) \log_2 P(A) - P(G) \log_2 P(G) - P(L) \log_2 P(L)$

$$\approx \underline{\underline{1.52}}$$

Avg. Imp = $\frac{5}{5} \times 1.52 = \underline{\underline{1.52}}$

Level 2

Entropy left Subtree: $-P(G) \log_2 P(G) = -\frac{2}{2} \log_2 \frac{2}{2} = \underline{\underline{0}}$

Entropy right Subtree: $-P(A) \log_2 P(A) - P(L) \log_2 P(L)$

$$= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= -0.67 \log_2 0.67 - 0.33 \log_2 0.33 \approx 0.38 + 0.52$$

$$\approx \underline{\underline{0.90}}$$

Avg Impurity = $\frac{2}{5} \times 0 + \frac{3}{5} \times 0.9 = \underline{\underline{0.54}}$

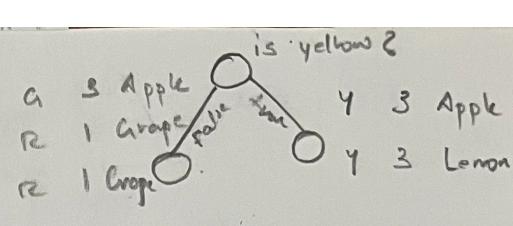
Information gain = $1.52 - 0.54 = \underline{\underline{0.98}}$

If the feature "diameter is $>= 3$ ", then information gain is 2 methods:

Gini = 0.37

Entropy = 0.98 Information Gain is Entropy is higher than Gini

If color is Yellow ?



Level 1: Root

$$\text{Entropy} = -P(A) \log_2 P(A) - P(G) \log_2 P(G) - P(L) \log_2 P(L)$$

$$\approx 1.52$$

$$\text{Avg Impurity} = \frac{5}{5} \times 1.52 = 1.52$$

Level 2:

$$\begin{aligned} \text{Left Subtree : } -P(A) \log_2 P(A) - P(G) \log_2 P(G) &= -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \\ &\approx 0.33 \log_2 0.33 - 0.66 \log_2 0.66 \\ &= 0.52 + 0.38 \approx 0.90 \end{aligned}$$

$$\begin{aligned} \text{Right Subtree : } -P(A) \log_2 P(A) - P(L) \log_2 P(L) &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \\ &= 0.5 \log_2 0.5 - 0.5 \log_2 0.5 \\ &= 0.5 + 0.5 \approx 1 \end{aligned}$$

$$\text{Avg Imp} = \frac{3}{5} * 0.90 + \frac{2}{5} * 1$$

$$= 0.54 + 0.4 = 0.94$$

$$\text{Information Gain} = \text{Avg Imp}(Root) - \text{Avg Imp}(level 2)$$

$$= 1.52 - 0.94$$

$$= 0.58$$

If the feature color is "Yellow", then information gain by 2 methods are:

$$Gini = 0.17$$

$$\text{Entropy} = 0.54$$

Information gain by entropy is higher than Gini

Q2)

Jupyter Notebook:

https://github.com/abeednabith/CS483_AI/blob/main/assignment4/Q2_GINI_DT.ipynb

Dataset:

https://github.com/abeednabith/CS483_AI/blob/main/assignment4/Q2_GINI_DT.csv

Q2)

Dataset :

Age	Competition	Type	Profit
Old	Yes	Software	Down
Old	No	Software	Down
Old	No	Hardware	Down
Mid	Yes	Software	Down
Mid	Yes	Hardware	Down
Mid	No	Hardware	Up
Mid	No	Software	Up
New	Yes	Software	Up
New	No	Hardware	Up
New	No	Software	Up
Mid	No	Hardware	?

If Competition is "Yes" ?

Level 1: Root

$$\text{Imp} = P(D) * (1 - P(D)) + P(V) * (1 - P(V))$$

$$= \frac{5}{10} * \left(1 - \frac{5}{10}\right) + \frac{5}{10} * \left(1 - \frac{5}{10}\right)$$

$$= \frac{5}{10} * \frac{5}{10} + \frac{5}{10} * \frac{5}{10}$$

$$= 0.5 * 0.5 + 0.5 * 0.5 \approx 0.5$$

$$\text{Avg Imp} = \frac{10}{10} * 0.5 = 0.5$$

Level 2:

$$\text{Left Subtree Impurity} = P(D) * (1 - P(D)) + P(V) * (1 - P(V)) = \frac{2}{6} * \left(1 - \frac{2}{6}\right) + \frac{4}{6} * \left(1 - \frac{4}{6}\right)$$

$$= \frac{2}{6} * \frac{4}{6} + \frac{4}{6} * \frac{2}{6} = 0.33 * 0.66 + 0.66 * 0.33$$

$$= 0.43$$

$$\text{Right Subtree Impurity} = P(D) * (1 - P(D)) + P(V) * (1 - P(V))$$

$$= \frac{3}{4} * \left(1 - \frac{3}{4}\right) + \frac{1}{4} * \left(1 - \frac{1}{4}\right) = \frac{3}{4} * \frac{1}{4} + \frac{1}{3} * \frac{3}{4}$$

$$= 0.75 * 0.25 + 0.25 * 0.75$$

$$= 0.258 + 0.15$$

$$= 0.40$$

$$= 0.375$$

$$\text{Information Gain} = \text{Avg. Imp}(Root) - \text{Avg. Imp}(Level 1)$$

$$2 \text{ levels} = 0.5 - 0.40$$

$$= \underline{\underline{0.1}}$$

If "Type is Software"

Level 1: Root

$$\text{Impurity} = \underline{\underline{0.5}}$$

$$\text{Avg. Imp} = \underline{\underline{0.5}}$$

Level 2:

Left Subtree

$$\text{Impurity} = P(D) \times (1 - P(D)) + P(V) \times (1 - P(V))$$

$$= \frac{2}{4} \times \left(1 - \frac{2}{4}\right) + \frac{3}{6} \times \left(1 - \frac{3}{6}\right) = \frac{2}{4} \times \frac{2}{4} + \frac{3}{6} \times \frac{3}{6} \approx \underline{\underline{0.5}}$$

Right Subtree

$$\text{Impurity} = P(D) \times (1 - P(D)) + P(V) \times (1 - P(V))$$

$$= \frac{3}{6} \times \left(1 - \frac{3}{6}\right) + \frac{3}{6} \times \left(1 - \frac{3}{6}\right) \approx \underline{\underline{0.5}}$$

$$\text{Avg. Imp} = \frac{4}{10} \times 0.5 + \frac{6}{10} \times 0.5 = 0.4 \times 0.5 + 0.6 \times 0.5 = \underline{\underline{0.06}}$$

$$\text{Information Gain} = 0.5 - 0.06 = \underline{\underline{0.44}}$$

If "Age" feature

Level Root

$$\text{Impurity} = \underline{\underline{0.5}} ; \text{ Avg Imp} = \underline{\underline{0.5}}$$

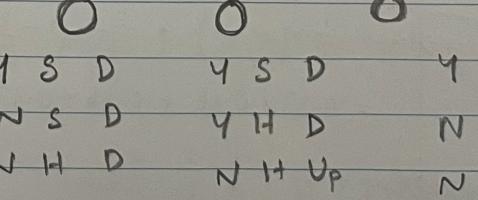
Level 2

$$\text{Left} \Rightarrow P(D) \times (1 - P(D)) + P(V) \times (1 - P(V))$$

$$= \frac{3}{3} \times \left(1 - \frac{3}{3}\right) + 0 = \underline{\underline{0}}$$

$$\text{Mid} \Rightarrow \frac{2}{4} \times \left(1 - \frac{2}{4}\right) + \frac{2}{4} \times \left(1 - \frac{2}{4}\right) = \underline{\underline{0.5}}$$

$$\text{Right} \Rightarrow 0 + \frac{3}{3} \times 1 - \frac{3}{3} = \underline{\underline{0}}$$



$$\text{Avg Imp} = \frac{3}{10} \times 0 + \frac{4}{10} \times 0.5 + \frac{3}{10} \times 0 = \underline{\underline{0.2}}$$

$$\text{Information Gain} = \underline{\underline{0.5 - 0.2}} = \underline{\underline{0.30}}$$

With these 3 criteria's Information gain is higher with the criteria "Type is Software". So we choose this decision tree for further split.

Level 1

$$\text{Avg Imp} = 0.5$$

Level 2

~~Age~~

Old N H D

Mid Y H D

Mid N H Up

New N H Up

Old Y S D

Old N S D

Mid Y S D

Mid N S U

New Y S U

New N S U

$$\text{Avg Imp} = 0.06$$

Level 3

For Level 3, If "Completion is Yes"

Left Subtree

Old N H D

Mid N H Up

New N H Up

Right Subtree

Old Y S D

Mid Y S D

New Y S Up

$$\text{Left Subtree} \Rightarrow \text{Imp} = \frac{1}{3} \times \left(1 - \frac{1}{3}\right) + \frac{2}{3} \times \left(1 - \frac{2}{3}\right)$$

$$= 0.33 \times 0.66 + 0.66 \times 0.33 \approx \underline{\underline{0.43}}$$

$$\text{Right Subtree} \Rightarrow \text{Imp} = \frac{1}{1} \times 1 - \frac{1}{1} = \underline{\underline{0}}$$

$$\text{Avg Imp} = \frac{3}{10} \times 0.43 = \underline{\underline{0.129}}$$

$$\text{Left Subtree} \Rightarrow \frac{1}{3} \times \left(1 - \frac{1}{3}\right) + \frac{2}{3} \times \left(1 - \frac{2}{3}\right) \approx \underline{\underline{0.43}}$$

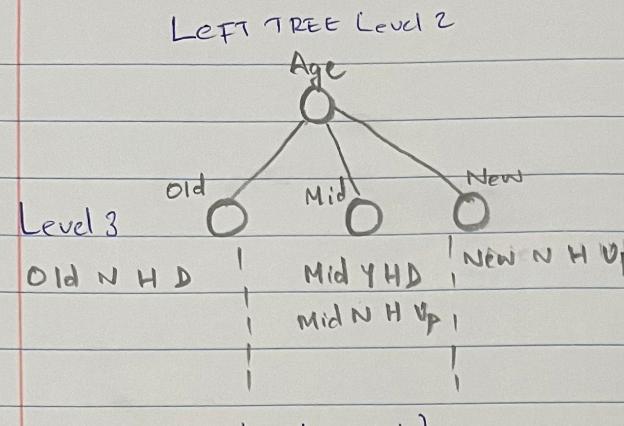
$$\text{Right Subtree} \Rightarrow \frac{2}{3} \times \left(1 - \frac{2}{3}\right) + \frac{1}{3} \times \left(1 - \frac{1}{3}\right) \approx \underline{\underline{0.43}}$$

$$\text{Avg Imp} = \frac{3}{10} \times 0.43 + \frac{3}{10} \times 0.43 = \underline{\underline{0.129}}$$

$$= \underline{\underline{0.258}}$$

$$\text{Information gain for 3 levels} = 0.5 - 0.06 - 0.258 = \underline{\underline{0.053}}$$

level 3, if "Age" feature



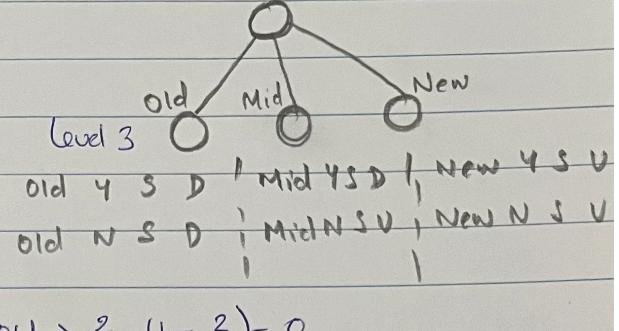
$$\text{Old} \Rightarrow P(D) \times (1 - P(D)) = \frac{1}{2} \times \left(1 - \frac{1}{2}\right) = 0$$

$$\text{Mid} \Rightarrow P(D) \times (1 - P(D)) + P(V) \times (1 - P(V)) = 0$$

$$\text{New} \Rightarrow P(D) \times (1 - P(V)) = \frac{1}{2} \times 1 - \frac{1}{2} = 0$$

$$\text{Avg Imp} = \underline{\underline{0}} + \underline{\underline{0}} + \underline{\underline{0}} = \underline{\underline{0}}$$

Right tree level 2



$$\text{Old} \Rightarrow \frac{2}{2} \times \left(1 - \frac{2}{2}\right) = 0$$

$$\text{Mid} \Rightarrow \frac{1}{2} \times \left(1 - \frac{1}{2}\right) = 0.5 \times 0.5 = 0.25$$

$$\text{New} \Rightarrow \frac{2}{2} \times \left(1 - \frac{2}{2}\right) = 0$$

$$\text{Avg Imp} = 0 + 0.25 + 0 = \underline{\underline{0.25}}$$

$$\text{Information gain for 3 levels} = 0.5 - 0.06 - 0.25$$

$$= \underline{\underline{0.19}}$$

From this calculation the following decisions are best suitable for tree:

Level 1 Root + Level 2 = Type is "Software"

Level 3 = Use feature "Age"; ie Old, Mid, New

And/Or

So, the profit value for : Mid NO Software Up