

→ 'Decision Tree using C4.5 Algorithm'

Note: ID3 favors attributes with large no. of values.

C4.5 Principle ⇒

It is a Decision Tree Classifier which can be employed to generate a decision, based on a certain sample of data.

Iterative Dichotomiser 3 (decision tree using Information Gain)

→ ID3 algorithm that we have seen earlier is biased towards multivalued attributes. C4.5 uses Gain Ratio.

$$\left[\begin{array}{l} \text{Gain Ratio: } \text{Gain Ratio}(A) = \frac{\text{Gain}(A)}{\text{SplitInfo}(A)} \\ \text{SplitInfo: } \text{SplitInfo}_A(D) = - \sum_{j=1}^V \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right) \end{array} \right]$$

Attribute with maximum Gain Ratio is selected as 'Splitting Attribute'

C4.5 Algorithm ⇒

1) For each attribute 'A', find the normalized information gain ratio from splitting on 'A'.

2) Let A-Best be the attribute with the highest normalized information gain.

3) Create a decision node that splits on 'A-Best'

Repeat these processes on the subsets obtained by splitting on A-Best.

4) Add these nodes as children of node.

P.T.O

"Dataset"

Work Period	Income	Married	Rank	Buy House
Short	High	No	High	No
Short	High	No	High	No
Average	Medium	Yes	High	Yes
Long	Low	Yes	High	Yes
Long	Low	Yes	Low	Yes
Long	Low	No	High	No
Average	Low	No	Low	Yes
Short	Medium	Yes	Low	No
Short	Low	No	Low	Yes
Long	Medium	Yes	High	Yes
Short	Low	No	Low	Yes
Average	Medium	Yes	Low	No
Average	Medium	Yes	High	Yes
Long	High	Yes	Low	No
Average	High	No	Low	No

We have to first compute the information Gain.

$$\text{Information Gain} = - \sum_{j=1}^V P_j \cdot \log_2(P_j)$$

$$\begin{aligned} \text{Information Gain}(D) &= \frac{-7}{15} \log_2\left(\frac{7}{15}\right) - \frac{8}{15} \log_2\left(\frac{8}{15}\right) \\ &= 0.9967 \end{aligned}$$

Now we need to compute information gain for each attribute :->

"Work-Period Attribute"

	Yes	No	Total
Short	2	3	5
Average	3	2	5
Long	3	2	5

"15"

$$\begin{aligned}
 \text{Info}_{\text{woorperiod}}^{(D)} &= \frac{5}{15} I(3,2) + \frac{5}{15} I(2,2) + \frac{5}{15} I(2,3) \\
 &= \frac{15}{15} I(3,2) \\
 &= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0.4421 + 0.5287 \\
 &= 0.9708
 \end{aligned}
 \tag{3}$$

Similarly let's try to compute Info for other attributes \Rightarrow

$$\text{Info}_{\text{Income}}^{(D)} = 0.5838$$

$$\text{Info}_{\text{married}}^{(D)} = 0.2923$$

$$\text{Info}_{\text{Rank}}^{(D)} = 0.9990$$

$$\begin{aligned}
 \text{Gain}_{\text{Income}}^{(D)} &= \text{Info}^{(D)} - \text{Info}_{\text{Income}}^{(D)} \\
 &= 0.9967 - 0.5838 \\
 &= \underline{0.4129}
 \end{aligned}$$

$$\text{Gain}_{\text{woorperiod}}^{(D)} = 0.0259$$

$$\text{Gain}_{\text{married}}^{(D)} = 0.1044$$

$$\text{Gain}_{\text{Rank}}^{(D)} = 0.0037$$

Now, we need to compute the SplitInfo for each attribute \Rightarrow

$$\text{SplitInfo}_{\text{Income}}^{(D)} = -\frac{4}{15} \log_2\left(\frac{4}{15}\right) - \frac{5}{15} \log_2\left(\frac{5}{15}\right) - \frac{6}{15} \log_2\left(\frac{6}{15}\right)$$

$$\begin{aligned}
 \text{SplitInfo}_{\text{Income}}^{(D)} &= 0.5085 + 0.5283 + 0.5287 \\
 &= 1.5655
 \end{aligned}$$

Similarly,

$$\left. \begin{aligned} \text{SplitInfo}_{\text{Work Period}}(D) &= 1.5849 \\ \text{SplitInfo}_{\text{Married}}(D) &= 0.9967 \\ \text{SplitInfo}_{\text{Ramp}}(D) &= 0.9967 \end{aligned} \right\}$$

Now Compute Gain Ratio for each attribute \Rightarrow

$$\text{Gain Ratio} = \frac{\text{Info. Gain}}{\text{Split Ratio.}}$$

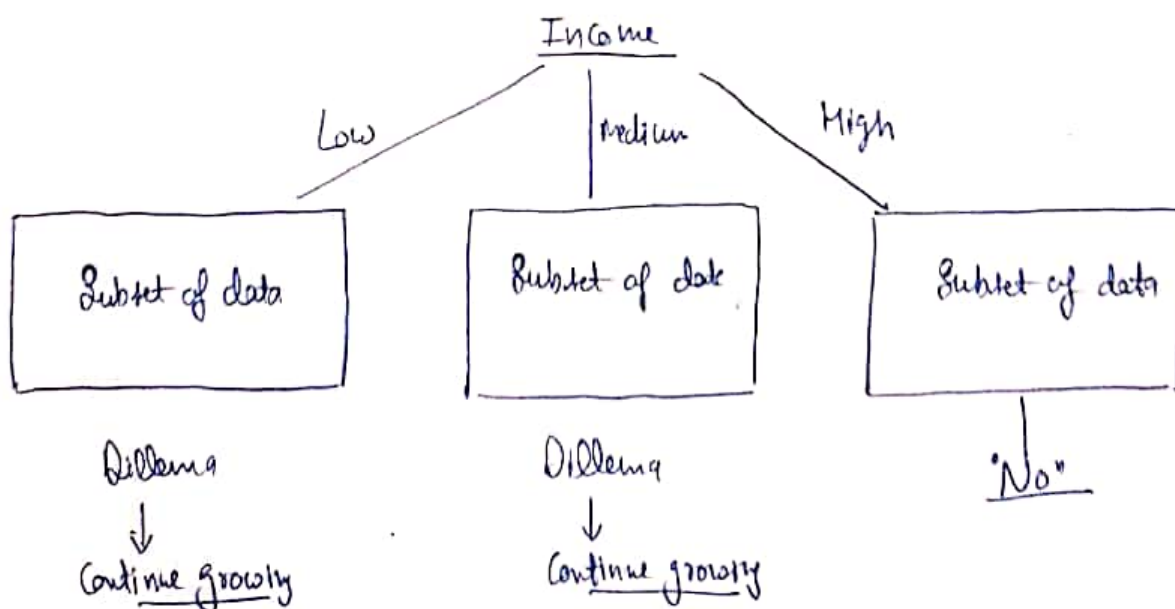
$$\text{Gain Ratio (Work Period)} = \frac{0.0259}{1.5849} / \frac{0.0259}{1.5849} = 0.0163$$

$$\text{Gain Ratio (Income)} = \frac{0.4129}{1.5655} = 0.2637$$

$$\text{Gain Ratio (Married)} = 0.1047 \quad \& \quad \text{Gain Ratio (Ramp)} = 0.0037$$

Since Gain Ratio of income is highest

\therefore Income will become the 'Root'



CART (Classification & Regression Trees)

⑤

↳ Another method to create Decision Tree

Here we use Gini-Index instead of 'Entropy'.

& Binary split

Gini-Index, is also known as Gini Impurity

↓

It calculates the amount of probability of a specific feature that it is classified incorrectly when selected randomly.

If all the elements are linked with a single class then it can be called pure.

Gini index varies b/w values 0 & 1 ∴

↳ Where '0' expresses the purity of classification.

i.e. All the elements belong to a specified class or only one class exists there.

↳ '1' indicates the random distribution of elements across various classes.

↳ The value of '0.5' of the Gini index shows an equal distribution of elements over some classes.

Let's now construct a Decision Tree using Gini Index

'P.T.O'

'Dataset'

Weekend	Weather	Parents	Money	Decision
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay in
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

off
Variable

While constructing a Decision tree, we will be selecting an attribute which has highest information & we will make it as a root node.

Here the attribute which is having the minimum 'Gini Index' will be the attribute that is having the maximum information, & we will select it as a root node.

Let's first compute the 'Gini index' of entire dataset.

In this dataset, we have four possible output variables \Rightarrow

'Cinema, Tennis, Stay in & Shopping'

In total we have 6 instances, where Cinema is present, 2 instances where Tennis is present, 1 instance of Stay in and 1 instance, where Shopping is present.

Now let's compute the Gini of entire dataset \Rightarrow

$$\begin{aligned} \text{Gini}(S) &= 1 - \left[\left(\frac{6}{10}\right)^2 + \left(\frac{2}{10}\right)^2 + \left(\frac{1}{10}\right)^2 + \left(\frac{1}{10}\right)^2 \right] \\ &= 0.58. \end{aligned}$$

Now we try to compute the Gini of every attribute.

(7)

Computation of Gini index for Money Attribute.

It has two possible values (Rich & Poor)

[7 instances] [3 instances]

for Money = Poor, there are three instances with Cinema

$$\text{Gini}(s)_{\text{Money}=\text{Poor}} = 1 - \left[\left(\frac{3}{3} \right)^2 \right] = \boxed{0}$$

For Money = Rich there are two instances with Tennis
3 instances with Cinema & 1 instance with stay in & Shopping each

$$\begin{aligned} \text{Gini}(s)_{\text{Money}=\text{Rich}} &= 1 - \left[\left(\frac{2}{7} \right)^2 + \left(\frac{3}{7} \right)^2 + \left(\frac{1}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right] \\ &= \boxed{0.694} \end{aligned}$$

Now, we need to compute weighted average for the money attribute.

$$\text{Weighted Average (Money)} = \underset{\substack{\downarrow \\ \text{Gini value}}}{0} * \underset{\substack{\downarrow \\ \text{Proportion}}}{\left(\frac{3}{10} \right)} + 0.694 * \frac{7}{10} = \boxed{0.486}$$

Similarly, we will compute the Gini for remaining attributes
of Parents & Weather.

Computation of Gini Index for Parents Attribute.

It has two possible values of Yes (5 instances) & No (5 instances)

for Parents = Yes, there are 5 instances all with Cinema

$$\therefore \text{Gini}(s)_{\text{Parents}=\text{Yes}} = 1 - \left[\left(\frac{5}{5} \right)^2 \right] = 0$$

For $\text{Parents} = \text{No}$, there are two instances with "Tennis" one instance with "Stay in" shopping & Cinema each.

∴ $\text{Gini}(S)_{\text{Parents}=\text{No}}$

$$= 1 - \left[\left(\frac{2}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 \right] = \boxed{0.72}$$

Now let's compute Weighted Average (Parents)

$$= 0 * \left(\frac{5}{10} \right) + 0.72 * \left(\frac{5}{10} \right) = \boxed{0.36}$$

Computation of Gini index for Weather Attribute.

It has three possible values of Sunny (3 instances), Rainy (3 instances) & Windy (4 instances).

for Weather = Sunny, there are 2 instances with "Cinema" & 1 with "Tennis"

$$\text{Gini}(S)_{\text{Weather}=\text{Sunny}} = 1 - \left[\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right] = \boxed{0.444}$$

for Weather = Rainy, there are 2 instances with "Cinema" & 1 with "stay in"

$$\text{Gini}(S)_{\text{Weather}=\text{Rainy}} = 1 - \left[\left(\frac{2}{3} \right)^2 + \left(\frac{1}{3} \right)^2 \right] = \boxed{0.444}$$

for Weather = Windy, there are 3 instances with "Cinema" & 1 with "shopping"

$$\text{Gini}(S)_{\text{Weather}=\text{Windy}} = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = \boxed{0.375}$$

Now Weighted Average (weather)

(9)

$$= 0.444 * \left(\frac{3}{10}\right) + 0.444 * \left(\frac{3}{10}\right) + 0.375 * \left(\frac{4}{10}\right)$$
$$= \boxed{0.416}$$

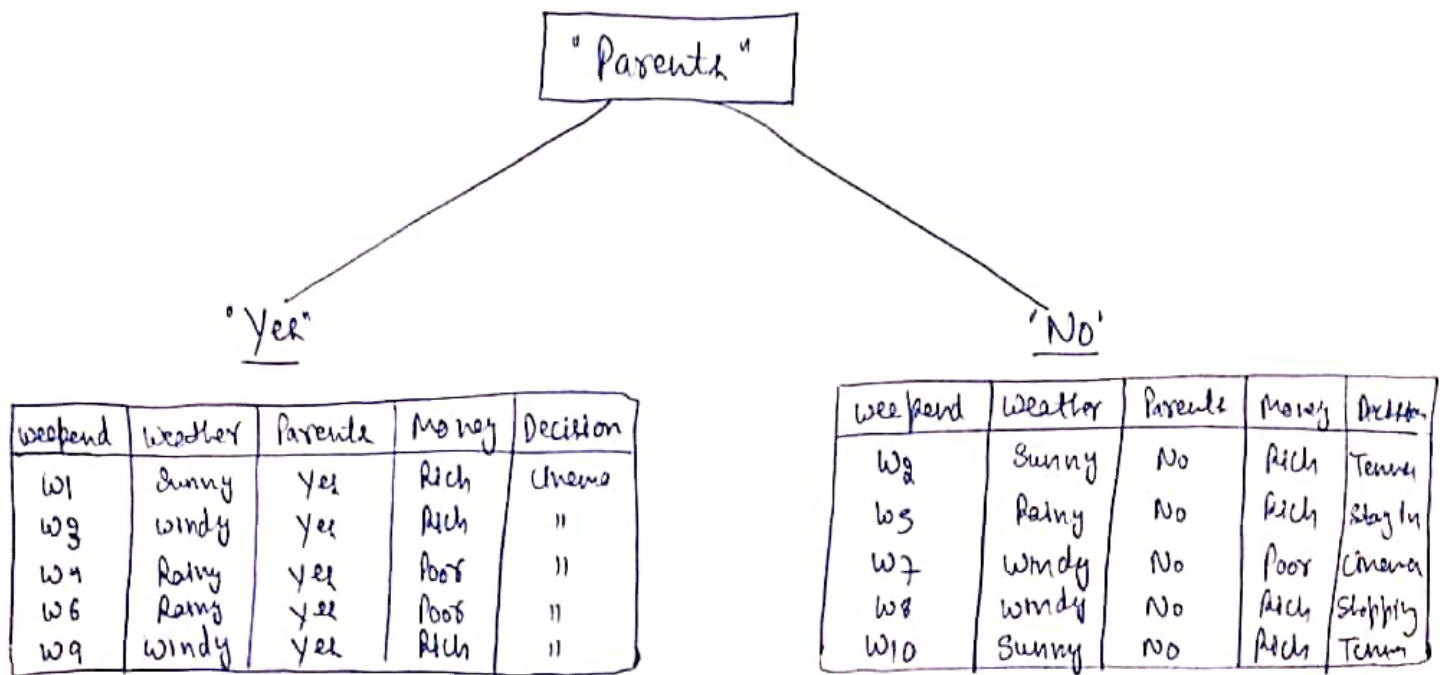
We have computed Gini Index for all the three attributes \Rightarrow

$$\left[\begin{array}{l} \text{For weather - Gini index} = 0.416 \\ \text{For Parents - Gini index} = 0.36 \\ \text{For Money - Gini index} = 0.486 \end{array} \right]$$

Now out of these three, Gini index for Parents is minimum which means, it is having highest information.

\therefore we select Parents as root node. (as it has smallest Gini index)

Now we have to divide the data, based on possible values of Parents [Yes, No]



Whenever Parents = "Yes", we have a decision at "Cinema" always.

No need to do computation any further from the end

P.T.O

When our Parents = No, we don't have a specific decision to take we have four decisions [Tennis, Stay in Cinema, Shopping]

We will grow our tree using the subset.

Weekend	Weather	Parents	Money	Decision
w2	Sunny	No	Rich	Tennis
w5	Rainy	No	Rich	Stay in
w7	windy	No	Poor	Cinema
w8	windy	No	Rich	Shopping
w10	Sunny	No	Rich	Tennis

Computation of Gini Index for Parents = No | Weather Attribute

We have two instances of Sunny,

for Parents = No | Weather = Sunny, there are two examples/instances with 'Tennis'.

$$\text{Gini}(S) = 1 - \left[\frac{2}{2} \right]^2 = \boxed{0}$$

We have one instance of Rainy

for Parents = No | Weather = Rainy, there is 1 example with "Stay in".

$$\text{Gini}(S) = 1 - \left[\left(\frac{1}{1} \right)^2 \right] = \boxed{0}$$

We have two instances of Windy

for Parents = No | Weather = Windy, there is 1 example with 'Cinema' & 1 example with 'Shopping'.

$$\text{Gini}(S) = 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right] = \boxed{0.5}$$

$$\text{Weighted Average (Parents = No | Weather)} = 0 * \left(\frac{2}{3} \right) + 0 * \left(\frac{1}{3} \right) + 0.5 * \left(\frac{2}{3} \right) = \boxed{0.2}$$

Computation of Gini Index for Parents = No | Money Attribute) (11)

We have four "Rich instances"

For Parents = No | Money = Rich, there is 1 example with "stay in" & "Shopping" each & 2 examples of "Tennis".

$$\therefore \text{Gini}(S) = 1 - \left[\left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 + \left(\frac{2}{4}\right)^2 \right] = \boxed{0.625}$$

We have 1 "Poor instance"

For Parents = No | Money = Poor, there is 1 example with "Cinema".

$$\therefore \text{Gini}(S) = 1 - \left[\left(\frac{1}{1}\right)^2 \right] = 0$$

$$\therefore \text{Weighted Average (Parents = No | Money)} = 0.625 * \left(\frac{4}{5}\right) + 0 * \left(\frac{1}{5}\right) = \boxed{0.5}$$

For Parents = No | weather \rightarrow Gini Index = 0.2
For Parents = No | Money \rightarrow Gini Index = 0.5

Now weather is selected as it has the smallest "Gini index"

Now for weather we have three categories of Sunny, Rainy & Windy
 \therefore we will get three branches

Now subset of data when Parent = No & weather = Sunny we have all the instances at Tennis.

Weekend	Weather	Parents	Money	Decision
w2	Sunny	No	Rich	Tennis
w10	Sunny	No	Rich	Tennis

No decision

Now for Parents = No & weather = Rainy, we have all the instances at 'Stay In'. (12)

Weekend	Weather	Parents	Money	Decision
WS	Rainy	No	Rich	Stay In

No dilemma

Now for Parents = No & weather = Windy, we need to split.

Weekend	Weather	Parents	Money	Decision
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping

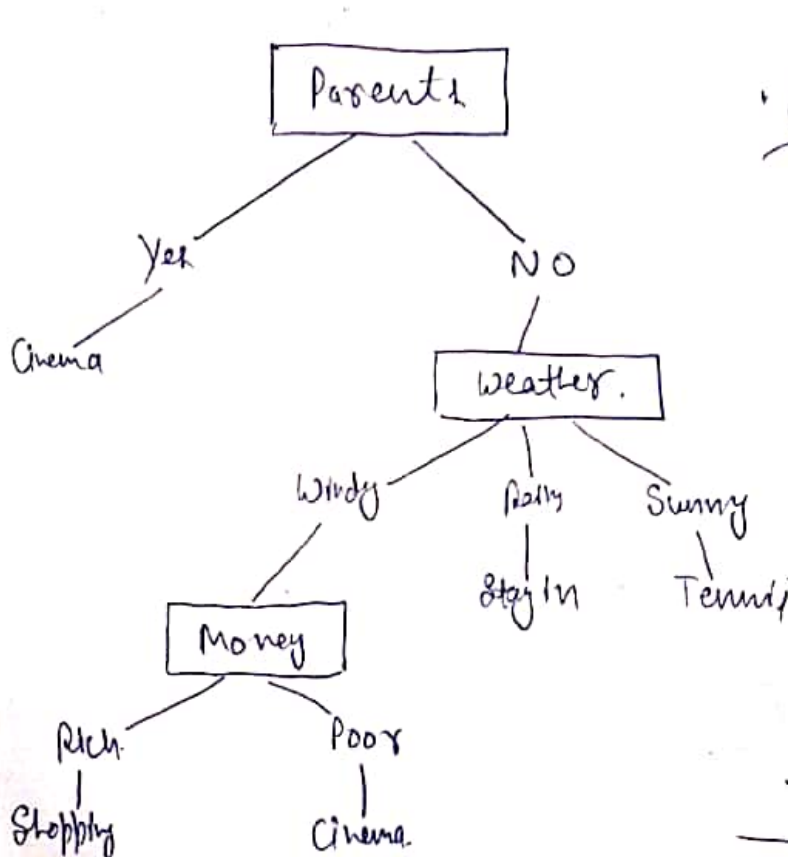
We need to split

Now we need to consider remaining attribute.

When Parents = No, weather = Windy, ^{money = poor} Decision is 'Cinema'

When Parents = No, weather = Windy, Money = Rich, Decision is 'Shopping'.

∴ Decision Tree looks like ⇒



'Child 3pm'
Graph CL10