

(1)

ID3 Algorithm  $\rightarrow$ \* ~~Iterative~~ Iterative Dichotomiser 3 AlgorithmEntropy <sup>measures of</sup>  $\rightarrow$  disorder in a system\* In a particular Node, all the examples are ~~positive~~ or all are negative. that <sup>mean</sup> all the examples ~~are~~ belong to same class. then it is a Homogenous set of example and Entropy is low.

\* However if we have two classes all the examples half belong to one class and half belong to another class then Entropy is highest.

Information Gain

When we decide which attribute to split on we will use the principle of Information Gain.

\* if all the example have same target classification. then information gain is high.

\* ~~if~~ <sup>Suppose</sup> 90% belong to one class then information gain is quite high.

\* 50% belong to one class and 50% another class then information gain is low.

Gain is measure of how much we can reduce uncertainty. if the example belong to the same class there is no uncertainty.  
if the example spread among the classes almost uniformly there is high uncertainty.

What an Entropy does?

(2)

Entropy Controls how a Decision tree decides to split the data. It actually effects how a Decision tree draws its boundaries.

What is information gain and, why it is Matter in Decision Tree?

\* measures how much "information" a feature gives us about the class.

\* Matter → information gain is the main key that is used by Decision Tree Algorithms to Construct a Decision tree.

\* Decision Tree algorithm will always tries to Maximize information gain.

\* An attribute with highest information Gain will tested / split first.

\* ID3 is one of the first decision tree algorithm.

\* only categorical attributes supported by ID3.





buys - Computer		Entropy
yes	No	
9	5	

$$\text{info}(D) = (I 9,5) = -\frac{9}{14} \log_2 \left( \frac{9}{14} \right) - \frac{5}{14} \log_2 \left( \frac{5}{14} \right)$$

$$\Rightarrow 0.9403$$

(i) Calculate Entropy of class attribute:-

(ii) Calculate Gain Ratio of all other attributes.

		class		
		Yes	No	
age	Youth	2	3	5
	Middle-aged	4	0	4
	Senior	3	2	5
				14

$$\text{Gain Ratio}(A) = \frac{\text{Gain}(A)}{\text{Split Info}(A)}$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$\downarrow$  info of Dataset       $\downarrow$  info of Attribute

$$-\frac{2}{5} \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \log_2 \left( \frac{3}{5} \right)$$

$$\text{Info}_A(D) = \sum_{j=1}^n \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

$$\text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2)$$

$$\Rightarrow \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971$$

$$\Rightarrow 0.3467 + 0 + 0.3467$$

$$\Rightarrow 0.6934$$

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D)$$

$$\Rightarrow 0.9403 - 0.6934$$

$$\Rightarrow 0.2469$$

$$\text{Split Info}_{\text{age}}(D) = -\frac{5}{14} \log_2 \left( \frac{5}{14} \right) - \frac{4}{14} \log_2 \left( \frac{4}{14} \right) - \frac{5}{14} \log_2 \left( \frac{5}{14} \right)$$

$$\Rightarrow 1.5774$$

$$\text{Gain Ratio}_{(\text{age})} = \frac{0.246}{1.5774}$$

$$\Rightarrow 0.1559$$

(5)

		class		
		yes	no	
Income	low	3	1	4
	Medium	4	2	6
	high	2	2	4
				14

$$\text{Info}_{\text{income}}^{(1)} = \frac{4}{14} I(3,1) + \frac{6}{14} I(4,2) + \frac{4}{14} I(2,2)$$

$$\Rightarrow \frac{4}{14} * 0.8113 + \frac{6}{14} * 0.9183 + \frac{4}{14} * 1$$

$$\Rightarrow 0.2318 + 0.3935 + 0.2857$$

$$\Rightarrow 0.911$$

$$\text{Gain}(\text{income}) = 0.9403 - 0.911$$

$$\Rightarrow 0.0293$$

$$\text{Split Info}_{(\text{income})}^{(1)} \Rightarrow \frac{4}{14} * \log_2\left(\frac{4}{14}\right) + \frac{6}{14} * \log_2\left(\frac{6}{14}\right) + \frac{4}{14} * \log_2\left(\frac{4}{14}\right)$$

$$\Rightarrow 1.5566$$

$$\text{Gain Ratio}_{(\text{income})} = \frac{0.0293}{1.5566}$$

$$\Rightarrow 0.0188$$

		class		
		yes	no	
Student	yes	6	1	7
	no	3	4	7
				14

$$\text{Info}_{\text{student}}^{(1)} = \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4)$$

$$\Rightarrow \frac{7}{14} * 0.5917 + \frac{7}{14} * 0.9852$$

$$\Rightarrow 0.2958 + 0.4926$$

$$\Rightarrow 0.7884$$

$$\text{Gain (Student)} \Rightarrow 0.9403 - 0.7884$$

(6)

$$\Rightarrow 0.1519$$

$$\text{Split Info}_{\text{Student}}^{(D)} = -\frac{7}{14} \times \log_2\left(\frac{7}{14}\right) - \frac{7}{14} \times \log_2\left(\frac{7}{14}\right)$$

$$\Rightarrow 1$$

$$\text{Gain Ratio}_{\text{Student}} \Rightarrow \frac{0.1519}{1} \Rightarrow 0.1519$$

		class		
		Yes	No	
Credit rating	Fair	6	2	8
	Excellent	3	3	6
				14

$$\text{Info}_{\text{Credit rating}}^{(D)} = \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$$

$$\Rightarrow \frac{8}{14} \times 0.8113 + \frac{6}{14} \times 1$$

$$\Rightarrow 0.4636 + 0.4286$$

$$\Rightarrow 0.8922$$

$$\text{Gain (Credit rating)} \Rightarrow 0.9403 - 0.8922$$

$$\Rightarrow 0.0481$$

$$\text{Split Info}_{\text{Credit-rating}}^{(D)} \Rightarrow -\frac{8}{14} \times \log_2\left(\frac{8}{14}\right) - \frac{6}{14} \log_2\left(\frac{6}{14}\right)$$

$$\Rightarrow 0.9852$$

$$\text{Gain Ratio}_{\text{Credit-rating}} \Rightarrow \frac{0.0481}{0.9852}$$

$$\Rightarrow 0.0488$$

\* As the Gain Ratio of 'age' is highest.

\* Age is the best attributes & become the Root Node.

Age

Junior

Middle-aged

Senior

Income	Student	Credit rating	Class
high	NO	fair	NO
high	NO	Excellent	NO
Medium	NO	fair	NO
low	YES	fair	YES
Medium	YES	Excellent	YES

Income	Student	Credit rating	Class
Medium	NO	fair	YES
low	YES	fair	YES
low	YES	Excellent	NO
Medium	YES	fair	YES
Medium	NO	Excellent	NO

uncertainty

Income	Student	Credit rating	Class
high	NO	fair	YES
low	YES	Excellent	YES
Medium	NO	Excellent	YES
high	YES	fair	YES

same class

replace it by yes

buys Computer

YES NO  
2 3

$$\text{info}(D) = I(2,3) \Rightarrow -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$\Rightarrow 0.971$$

Calculate Gain Ratio of all other attributes

		Class		
		yes	NO	
Income	low	1	0	1
	Medium	1	1	2
	high	0	2	2
				5

$$\text{Info}_{\text{income}}(D) = \frac{1}{5} I(1,0) + \frac{2}{5} I(1,1) + \frac{2}{5} I(0,2)$$

$$\Rightarrow \frac{1}{5} \times 0 + \frac{2}{5} \times 1 + \frac{2}{5} \times 0$$

$$\Rightarrow 0 + 0.4 + 0$$

$$\Rightarrow 0.4$$

$$\text{Gain}(\text{income}) = 0.971 - 0.4$$

$$\Rightarrow 0.571$$



$$\text{Split Info}_{\text{income}}^{(1)} = -\frac{1}{5} \times \log_2\left(\frac{1}{5}\right) - \frac{2}{5} \times \log_2\left(\frac{2}{5}\right) - \frac{2}{5} \times \log_2\left(\frac{2}{5}\right)$$

$$\Rightarrow 1.5219$$

$$\text{Gain Ratio (Income)} \Rightarrow \frac{0.571}{1.5219}$$

$$\Rightarrow 0.3751$$

		class		
		yes	no	
Credit rating	fair	1	2	3
	Excellent	1	1	2
				5

$$\text{Info}^{(1)}_{\text{Credit rating}} \Rightarrow \frac{3}{5} I(1,2) + \frac{2}{5} I(1,1)$$

$$\Rightarrow \frac{3}{5} \times 0.9183 + \frac{2}{5} \times 1$$

$$\Rightarrow 0.3443 + 0.4$$

$$\Rightarrow 0.7443$$

$$\text{Gain (Credit rating)} \Rightarrow 0.971 - 0.7443$$

$$\Rightarrow 0.2267$$

$$\text{Split Info}_{\text{(Credit rating)}} = -\frac{3}{5} \times \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \times \log_2\left(\frac{2}{5}\right)$$

$$\Rightarrow 0.9709$$

$$\text{Gain Ratio}_{\text{(Credit Rating)}} \Rightarrow \frac{0.2267}{0.9709}$$

$$\Rightarrow 0.2335$$

		class		
		yes	no	
Student	yes	2	0	2
	no	0	3	3
				5

$$\text{info}^{(1)}_{\text{Student}} = \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3)$$

$$\Rightarrow \frac{2}{5} \times 0 + \frac{3}{5} \times 0$$

$$\Rightarrow 0$$

$$\text{Gain (Student)} \Rightarrow 0.971 - 0 \Rightarrow 0.971$$



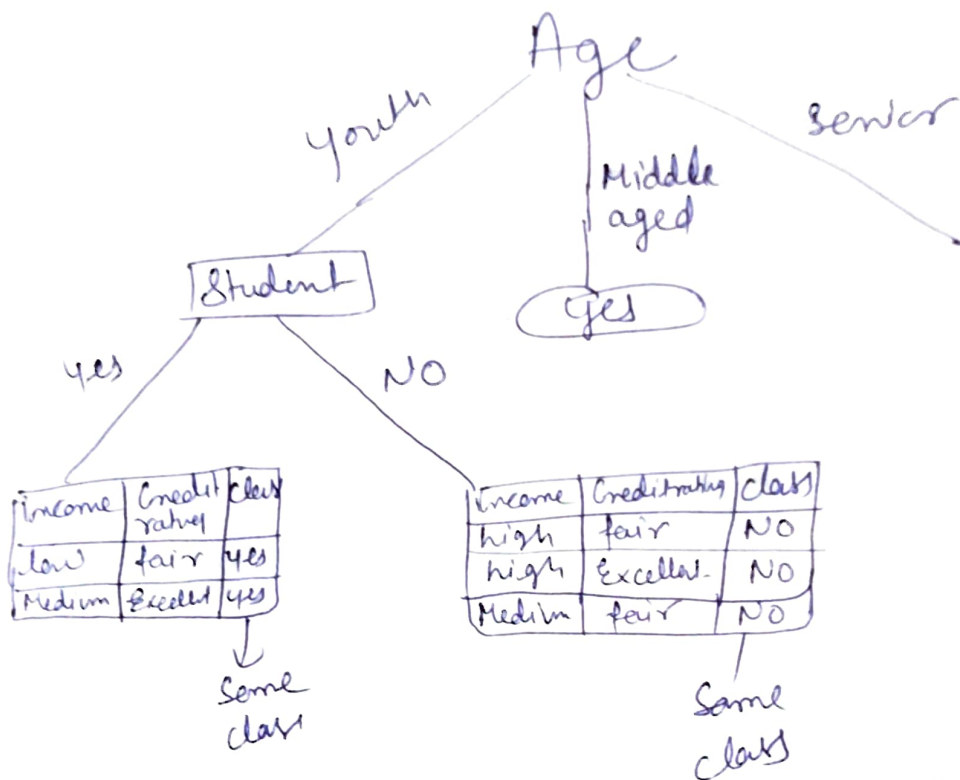
$$\text{Split Info}_{\text{Student}} \Rightarrow -\frac{2}{5} \times \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

(9)

$$\Rightarrow 0.9709$$

$$\text{Gain Ratio}_{(\text{Student})} \Rightarrow \frac{0.971}{0.9709}$$

$$\Rightarrow 1$$



Same we will do for Senior

buys Computer

yes NO  
3 2

$$\text{Info}(\text{D}) \Rightarrow I(3,2) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$

$$\Rightarrow 0.971$$

		class		
		yes	NO	
Income	low	1	1	2
	Medium	2	1	3
	high	0	0	0
				5

$$\text{Info}_{\text{Income}}(\text{D}) = \frac{2}{5} I(1,1) + \frac{3}{5} I(2,1)$$

$$\Rightarrow \frac{2}{5} \times 1 + \frac{3}{5} \times 0.9183$$

$$\Rightarrow 0.4 + 0.551$$

$$\Rightarrow 0.951$$

$$\text{Gain}(\text{income}) \Rightarrow 0.971 - 0.951 \Rightarrow 0.02$$

$$\text{Split Info}_{\text{income}}^{(D)} \Rightarrow -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$\Rightarrow 0.9709$$

$$\text{Gain Ratio}(\text{income}) \Rightarrow \frac{0.02}{0.9709}$$

$$\Rightarrow 0.0205$$

		class		
		yes	no	
Credit rating	fair	3	0	3
	Excellent	0	2	2
				5

$$\text{Info}_{\text{Credit rating}}^{(D)} = \frac{3}{5} I(3/5) + \frac{2}{5} I(0/2)$$

$$\Rightarrow \frac{3}{5} \times 0 + \frac{2}{5} \times 0$$

$$\Rightarrow 0$$

$$\text{Gain}_{\text{Credit rating}} \Rightarrow 0.971 - 0 \Rightarrow 0.971$$

$$\text{Split Info}_{\text{Credit rating}}^{(D)} \Rightarrow -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$

$$\Rightarrow 0.9709$$

$$\text{Gain Ratio}_{\text{Credit rating}} \Rightarrow \frac{0.971}{0.9709} \Rightarrow 0.1$$

		class		
		yes	no	
Student	yes	2	1	3
	No	1	1	2
				5

$$\text{Info}_{\text{Student}}^{(D)} \Rightarrow \frac{3}{5} I(3/5) + \frac{2}{5} I(2/2)$$

$$\Rightarrow \frac{3}{5} \times 0.9183 + \frac{2}{5} \times 1$$

$$\Rightarrow 0.551 + 0.4$$

$$\Rightarrow 0.951$$

$$\text{Gain}_{\text{Student}} \Rightarrow 0.971 - 0.951$$

$$\text{Split Info}_{\text{Student}} \Rightarrow -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right)$$

$$\Rightarrow 0.9702$$

$$\text{Gain Ratio} = \frac{0.02}{0.9702} = 0.02$$

