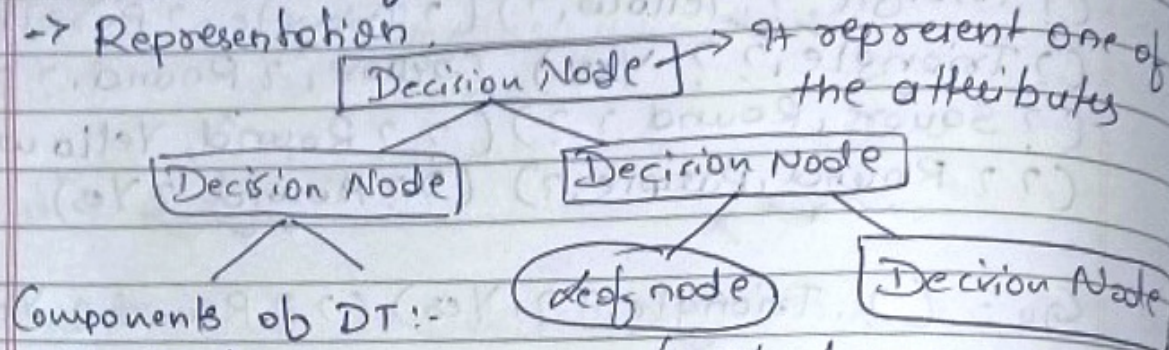


28/12/2023

## Decision Tree Learning.

- popular supervised classification Algorithm.
- inductive inference. (final conclusion)
- Representation.



Components of DT:-

- Root node
- Internal node
- leaf node
- Branches

→ It always represents final decision/target.

A DT is a tree in which each branch represents a choice b/w a number of alternatives and each leaf node represents a decision.

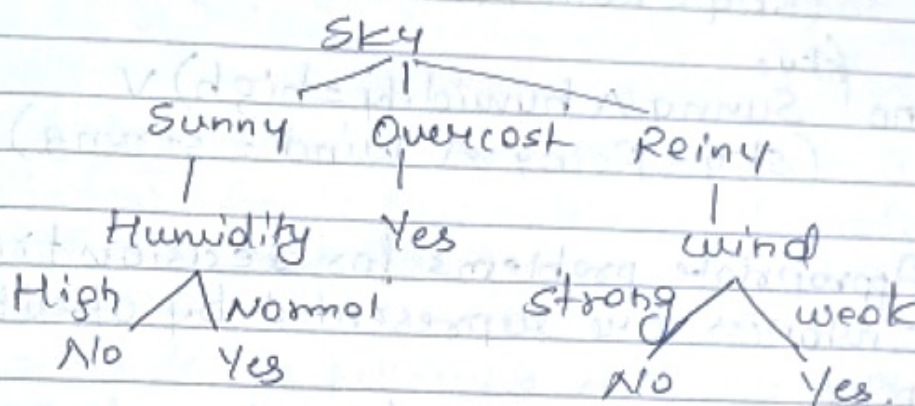
A DT is a tree in which each node represents a feature or attribute, each branch represents a decision or rule each leaf represents an outcome.

Components of a DT:-

- Root node refers to the start of the decision tree
- Node is a condition with multiple outcomes in the tree
- leaf is the final decision or the end point of a node
- In General, A DT represents a disjunction of conjunction of constraints on the attribute values of the instances.



- Each path from the root to a leaf corresponds to a conjunction of attribute value & the tree itself corresponds to a disjunction of these conjunction.



Rule/Expression:  $(\text{Sky} = \text{Sunny} \wedge \text{Humidity} = \text{Normal}) \vee (\text{Sky} = \text{Overcast}) \vee (\text{Sky} = \text{Rainy} \wedge \text{Wind} = \text{Weak})$

Appropriate p

- 1) Instances are represented by attribute value pairs
- 2) The target values are discrete
- 3) Disjunctive discrimination may be measured
- 4) Training data may contain missing attribute values



### Rule/Expression:-

yes (sky = Sunny  $\wedge$  Humidity = Normal)  $\vee$   
 (sky = overcast)  $\vee$   
 (sky = Rain  $\wedge$  wind = weak)  
 no (sky = Sunny  $\wedge$  humidity = high)  $\vee$   
 (sky = Rainy  $\wedge$  wind = strong)

- 1) Appropriate problems for decision tree learning
- 2) Instances are represented by attribute value pairs
- 3) The target function has discrete output values
- 4) Disjunctive description may be required
- 5) Training data may contain errors and may contain missing attribute values.

29/12/2023

### ID3 (Quadratic Dichotomiser 3)

CHS (Successes of ID3)

CART (Classification & Regression Tree)

### ID3

Based on the statistical value of each attribute we decide which attribute should select as a root node. To do this we use <sup>the</sup> Algorithm.

### \* ID3

→ learns the DT by constructing from top down. Beginning with the question "which attribute should be at the root of the tree". To answer this question each attribute is tested using a statistical test. To determine how well it classifies training examples ID3 (examples, target attribute, attributes)

- examples and the training examples
- target attributes is the att. whose values has to be predicted by the tree
- att is the list of other attributes that may be tested by the tree.

→ Given a dataset how we build a tree using ID3

- The initial step is to create a root node.
- The att that best classifies the training example is taken as the root node. In order to choose the best attribute the following steps are followed by ID3.

Step 1:- Compute the entropy for the dataset i.e.  $ENTROPY(S)$   $S \rightarrow$  Dataset

Step 2:- For every att or feature

→ Calculate the entropy for all other values i.e. entropy of A.

→ Calculate the average information entropy for the current attribute.

→ Calculate Gain for the current attribute

Step 3:- Pick the att with the highest gain value

Step 4:- Repeat until we get the desired output.

Entropy:- Given a collection S containing +ve & -ve examples the entropy of (S) is given as:

$$Entropy(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

S is a sample of training examples

P<sub>+</sub> is the proportion of +ve

P<sub>-</sub> " " " -ve



Entropy is defined as the entropy amount of uncertainty in the dataset.

It can also be written as

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

$$Entropy = - \frac{p}{p+n} \log_2 \left( \frac{p}{p+n} \right) - \frac{n}{p+n} \log_2 \left( \frac{n}{p+n} \right)$$

Average Info Entropy is calculated as

$$I(Attributes) = - \sum \frac{P_i + n_i}{p+n} Entropy(n_i)$$

Info Gain / Gain

It is the difference in the entropy before & after splitting the dataset on attribute A.

$$Gain = Entropy(S) - I(Attribute)$$

D

S1 Outlook Temp Humidity Wind PlayTennis

1	Sunny	hot	high	weak	No	p=9
2	Sunny	hot	high	strong	No	n=5
3	Overcast	hot	high	weak	Yes	p+n=14
4	Rainy	mild	high	weak	Yes	
5	Rainy	cool	normal	weak	Yes	
6	Rainy	cool	normal	strong	No	
7	Overcast	cool	normal	strong	Yes	
8	Sunny	mild	high	weak	No	
9	Sunny	cool	normal	weak	Yes	
10	Rainy	mild	normal	weak	Yes	
11	Sunny	mild	normal	strong	Yes	
12	Overcast	mild	high	strong	Yes	
13	"	hot	normal	weak	Yes	
14	Rainy	mild	high	strong	No	

classmate

classmate

$$\log_2 \frac{9}{14} = \frac{\ln(9/14)}{\ln(2)} = \frac{-0.1518}{0.3610}$$

Step 1:- Calculated entropy for entire dataset

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

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

$$= (-0.643) - (-0.637) - (0.357) - (-1.486)$$

$$= 0.4096 + 0.5305$$

0.940

→ Entropy is 0 when all the numbers belongs to class.

→ Entropy is 1 " the collection contains equal no. of +ve & -ve examples.

→ If the collection contains unequal no. of +ve & -ve examples then the entropy will be b/w 0 & 1.

Step 2:-

calculate entropy for each attribute

outlook	+ve	-ve
Sunny	2	3
Rainy	3	2
Overcast	4	0

calculate entropy(outlook=sunny)

$$= - \frac{p}{p+n} \log_2 \left( \frac{p}{p+n} \right) - \frac{n}{p+n} \log_2 \left( \frac{n}{p+n} \right)$$

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

$$= -0.4 (-1.3219) - (0.6) (-0.2218) / \log_2 2 = (-0.737)$$

$$0.5287 + 0.133888$$

$$= 0.6626$$

$$entropy(outlook=rainy)$$

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

$$= (-0.6) (-0.7219) - (0.4) (-1.3219)$$

$$= 0.4422 + 0.5287 = 0.971$$



$$\text{entropy (outlook} = \text{overcast})$$

$$= \frac{4}{9} \log_2 \left( \frac{4}{9} \right) - \frac{5}{9} \log_2 \left( \frac{5}{9} \right)$$

$$= -1(0) - 0$$

$$= 0$$

Calculate average information entropy.

$$I(\text{outlook}) = \frac{P_{\text{rainy}} + P_{\text{sunny}}}{P_{\text{th}}} \text{Entropy (outlook} = \text{sunny)}$$

$$+ \frac{P_{\text{rainy}} + P_{\text{rainy}}}{P_{\text{th}}} (\text{Entropy (outlook} = \text{rainy}) + \frac{P_{\text{overcast}} + P_{\text{overcast}}}{P_{\text{th}}})$$

$$(\text{Entropy (outlook} = \text{overcast})$$

$$= \frac{2+3}{14} \log_2 \left( \frac{5}{14} \right) + \frac{3+2}{14} \log_2 \left( \frac{9}{14} \right) + \frac{4+0}{14} \log_2 \left( \frac{0}{14} \right)$$

$$= 0.357(0.97) + 0.357(0.971) + 0$$

$$= 0.3466 + 0.3466$$

$$= 0.693$$

Calculate gain for the outlook

$$\text{Gain} = \text{Entropy}(S) - I(\text{outlook})$$

$$= 0.940 - 0.693$$

$$[\text{Gain} = 0.247]$$

Calculate entropy for Temp

Temp P n

hot 3 2

mild 4 2

cool 3 1

$$\text{En (Temp} = \text{hot}) = -\frac{2}{5} \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \log_2 \left( \frac{3}{5} \right)$$

$$= -\frac{2}{5} (1.609) - \frac{3}{5} (0.916)$$

$$= -1.109 + 0.549 = -0.56$$

$$\text{En (Temp} = \text{mild}) = -\frac{4}{6} \log_2 \left( \frac{4}{6} \right) - \frac{2}{6} \log_2 \left( \frac{2}{6} \right)$$

$$= -\frac{2}{3} (0.66) (0.353) - 0.33(0)$$

$$= 0.232$$

$$= 0.923$$

$$\text{En (Temp} = \text{cold}) = -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right)$$

$$= -0.75 ($$

$$= 0.811$$

$$I(\text{Temp}) = \frac{2+2}{14} (1) + \frac{4+2}{14} (0.92) + \frac{3+1}{14} (0.811)$$

$$= 0.285 + 0.394 + 0.231$$

$$= 0.912$$

$$\text{Gain} = 0.940 - 0.912$$

$$= 0.028$$

humidity +ve -ve +2 -1

high 3 4

normal 6 1

$$- \frac{3}{7} \log_2 \left( \frac{3}{7} \right) - \frac{4}{7} \log_2 \left( \frac{4}{7} \right)$$

$$= -0.428 (0.222) - 0.571 (-0.807)$$

$$= 0.985$$

$$- \frac{6}{7} \log_2 \left( \frac{6}{7} \right) - \frac{1}{7} \log_2 \left( \frac{1}{7} \right)$$

$$= -0.853 (-2.564) - 0.142 (-2.807)$$

$$= -2.197 + 0.398$$

$$= 0.591$$



$$\pm (\text{humidity}) = \frac{3+4}{14} (0.985) + \frac{6+1}{14} (0.591)$$

$$= 0.788$$

$$0.940 - 0.788$$

$$\text{Gain} = 0.152$$

Wind +ve -ve  
weak 6 2  
strong 3 3

$$-\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right)$$

$$= -(0.75) (-0.415) - (0.25) (-2)$$

$$= 0.812 + 0.5$$

$$-\frac{3}{8} \log_2 \left( \frac{3}{8} \right) - \frac{3}{8} \log_2 \left( \frac{3}{8} \right)$$

$$= -(0.4) (-1.322) - (0.4) (-0.769)$$

$$= 0.5(-1) - 0.5(-1)$$

$$I = \frac{6+2}{14} (0.811) + \frac{3+3}{14} (1)$$

$$= 0.571 (0.811) + 0.428$$

$$0.463 + 0.428$$

$$\text{Gain} = 0.940 - 0.891$$

$$= 0.049$$

Gain Attribute ✓ root node (high gain value)

0.247

Outlook

0.028

temp

0.152

humidity

0.049

wind

Sunny

Outlook  
outlook  
rainy

Yes

2

01/01/2024

with table when outlook = sunny.

Outlook Temp Humidity wind Play tennis

Sunny hot high weak No

Sunny hot high strong No

Sunny mild high weak Yes

Sunny cool normal weak Yes

Sunny mild normal strong Yes

$p = 2$   $n = 3$

Entropy (Sunny)

$$= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= -0.4 (-1.3219) - 0.6 (-0.7369)$$

$$= 0.971$$

Temp +ve -ve

hot 0 2

mild 0 1

temp Cool 1 0

$$\text{temp not} = \frac{0}{0+2} \log_2 \left( \frac{0}{2} \right) - \frac{2}{2} \left( \log_2 \left( \frac{2}{2} \right) \right)$$

$$= 0 (0) - 1 (0)$$

$$= 0$$

$$\text{temp mild} = \frac{1}{2} \log_2 \left( \frac{1}{2} \right) - \frac{1}{2} \log_2 \left( \frac{1}{2} \right)$$

$$= 1$$



$$k_{cool} = \frac{1}{5} \log_2\left(\frac{1}{5}\right) - \frac{1}{5} \log_2\left(\frac{2}{5}\right)$$

$$= 0 \times 0 -$$

$$= 0$$

Calculate  $I(temp)$

$$= \frac{2+2}{2+3} \log_2(0) + \frac{1+1}{2+3} (1) + \frac{0+0}{5} (0)$$

$$= 0 + 0.4 + 0$$

$$I(temp) = 0.4$$

Calculate gain

$$Gain = Entropy(S_{temp}) - I(temp)$$

$$= 0.971 - 0.4$$

$$Gain = 0.571$$

Humidity +ve -ve

high 0 3

normal 2 0

$$H_{w=high} = \frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$= 0 \times ( ) = 2(0)$$

$$= 0$$

$$H_{w=normal} = \frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{0}{5} \log_2\left(\frac{0}{5}\right)$$

$$= 2(0) - 0 \times 1$$

$$= 0$$

Calculate

$$I(Humidity) = \frac{2+3}{5} (0) + \frac{2+0}{5} (0) = 0$$

$$Gain = 0.971 - 0$$

$$Gain = 0.971$$

wind +ve -ve

weak 1 2

strong 1 1

$$w_{weak} = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right)$$

$$= 0.33(-1.584) - 0.67(-0.584)$$

$$= 0.528 + 0.392$$

$$= 0.92$$

$$I_{strong} = 1$$

$$I(wind) = \frac{1+2}{5} (0.918) + \frac{1+1}{5} (1)$$

$$= \frac{3}{5} (0.918) + \frac{2}{5} (1)$$

$$= 0.6 (0.918)$$

$$= 0.55 + 0.4$$

$$= 0.950$$

$$Gain = 0.971 - 0.950$$

$$Gain = 0.02$$

Outlook

Sunny Overcast Rainy

Humidity Yes

high Normal

No Yes



Outlook	temp	humidity	wind	Play tennis
rainy	mild	high	weak	Yes
rainy	cool	normal	weak	Yes
rainy	cool	normal	strong	No
rainy	mild	normal	weak	Yes
rainy	mild	high	strong	No

temp +ve -ve

mild 2 1

cool 1 1

mild  $-\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3}$

$0.66(-0.5849) - 0.33 \log(-1.5849)$

$0.386 + 0.522$

$= 0.918$

cool = 1

$I(temp) = \frac{3}{5} (0.971) + \frac{2}{5} (1)$

$0.6(0.971) + 0.4$

$0.5806 + 0.4$

$= 0.95$

Gain =  $0.971 - 0.95$

$= 0.02$

humidity +ve -ve

high 1 1

normal 2 1

normal  $\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3}$

$= 0.918$

high = 1

$I(humidity) = \frac{1+1}{5} (1) + \frac{2+1}{5} (0.918)$

$= 0.95$

Gain = 0.02

wind	+ve	-ve	
weak	3	0	= 0
strong	0	2	= 0

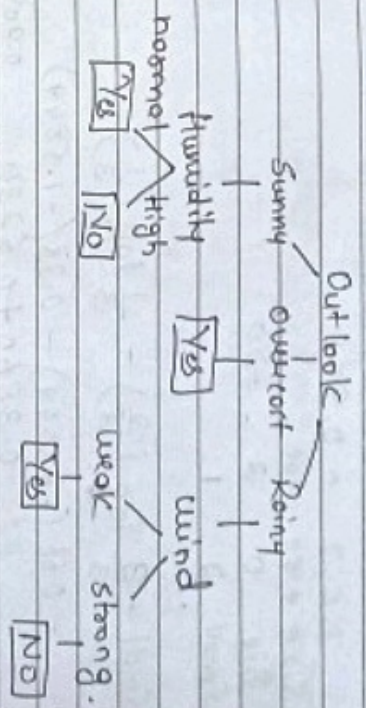
$I = 0$

Gain = 0.971

Temp 0.02

humidity 0.02

wind 0.971



1st Internals.

05/01/2023

Issues in Decision Tree Learning.

- \* Avoiding overfitting the data.
- \* Incorporating Continuous-valued attributes
- \* Alternate measures for selecting the attribute
- \* Handling training examples with missing attribute value
- \* " Attributes with differ in cost.



Ex	Size	Color	Shape	Class/Label	
	Big	Red	Circle	No	
	Small	Red	Triangle	No	+ve
	Small	Red	Circle	Yes	-ve
	Big	Blue	Circle	No	
	Small	Blue	Circle	Yes	

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

$$= -0.4(-1.3219) - 0.6(-0.7369) = 0.58996 + 0.4421 = 0.971$$

Size +ve -ve  
Big 0 2  $\Rightarrow 0$

Small 2 1

$$Small = \frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = -0.66(-0.5849) - 0.33(-1.5849) = 0.3860 + 0.5230 = 0.909$$

$$I(Size) = \frac{0+2}{5}(0) + \frac{2+1}{5}(0.909)$$

$$= 0 + 0.6(0.918) = 0.5508$$

$$Gain = 0.971 - 0.5508 = 0.42$$

Color +ve -ve

Red 2 1  $\Rightarrow 0.918$

Blue 1 1  $\Rightarrow 1$

$$\frac{2+1}{5}(0.918) + \frac{1+1}{5}(1)$$

$$= 0.5508 + 0.4$$

$$I(Color) = 0.9508 = 0.951$$

$$Gain = 0.971 - 0.9508 = 0.02$$

Shape +ve -ve

Circle 1 3  $\Rightarrow 0.918$

Triangle 0 1  $\Rightarrow 0$

$$= -\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right) = 0.811$$

$$I(Shape) = 0.8$$

$$Gain = 0.171 \quad (0.971 - 0.8)$$

Size Color Shape Label

Small Red Triangle No +ve -ve

Small Red Circle Yes +ve

Small Blue Circle Yes

No

Small

Color +ve -ve

Red 1 1  $\Rightarrow 1$

Blue 1 0  $\Rightarrow 0$

$$I(Color) = \frac{1+1}{3}(1) + \frac{1+0}{3}(0) = 0.918$$

$$= 0.66$$

$$Gain = 0.918 - 0.66$$

$$Gain = 0.258$$



10

### Talmon's Law



9

21

10

1

7

1

203

11

1

1

56

1

1

12



4

5

3

1

11

28

10

1

1

1

2

1

1

1

14

10

1