

1D3

- ① Compute the **ENTROPY** for data-set **ENTROPY(S)**
- ② For every attribute/feature
 - Ⓐ Calculate entropy for all other values **ENTROPY(A)**
 - Ⓑ Take **AVERAGE INFORMATION ENTROPY** for the current attribute
 - Ⓒ Calculate **GAIN** for the current attribute
- ③ Pick the **HIGHEST GAIN ATTRIBUTE**
- ④ **REPEAT** until we get the tree we desired.

$$\begin{aligned}
 \text{Entropy}(S) &= \sum -P(I) \cdot \log_2 P(I) \\
 &= -\frac{p}{p+n} \log_2 \left(\frac{p}{p+n} \right) - \frac{n}{p+n} \log_2 \left(\frac{n}{p+n} \right) \quad \text{--- (I)}
 \end{aligned}$$

$S = p+n$ (Total sample space)
 $\begin{matrix} \text{Yes} & \text{No} \\ \downarrow & \downarrow \\ p & n \end{matrix}$

$$\text{Entropy}(A) = \sum_{i=1}^x \frac{P_i + n_i}{P+n} I(P_i, n_i) \quad \text{--- (II)}$$

\downarrow
 Attribute

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum [P(S/A) \cdot \text{Entropy}(S/A)]$$

$$\left\{ \log_2 x = \frac{\log_{10} x}{\log_{10} 2} \right\}$$

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

$$I(p, n) = - \frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$\text{Entropy}(\text{decision}) = -P(\text{yes}) \times \log_2 P(\text{yes}) - P(\text{no}) \times \log_2 P(\text{no})$$

Yes = 9
No = 5
S = Yes + No = 9 + 5 = 14

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

$$= 0.94$$

Attributes = { Outlook, Temperature, Humidity, Wind }

For each attribute (let say wind)

- Calculate Entropy for each wind i.e. strong & weak

Wind	Play Tennis	Wind	Play Tennis
Weak	No	Strong	No
Weak	Yes	Strong	No
Weak	Yes	Strong	Yes
Weak	Yes	Strong	Yes
Weak	No	Strong	Yes
Weak	Yes	Strong	No
Weak	Yes		
Weak	Yes		

Wind	P	N	Entropy
strong	3	3	1
weak	6	2	0.811

$$\text{Entropy}(S, w = \text{strong}) = - \frac{3}{8} \log_2 \left(\frac{3}{8} \right) - \frac{3}{8} \log_2 \left(\frac{3}{8} \right) = ?$$

$$\text{Entropy}(S, w = \text{weak}) = - \frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) = ?$$

Average Information Entropy,

$$I(\text{Wind}) = P_{\text{strong}} + N_{\text{strong}} \cdot \text{Entropy}(\text{Wind} = \text{strong}) +$$

$$\frac{p_{weak} + n_{weak}}{p+n} \cdot \text{Entropy}(\text{wind} = \text{weak})$$

$$= \frac{3+3}{9+5} \times 1 + \frac{6+2}{9+5} \cdot 0.811$$

$$= 0.892$$

$$\text{Gain}^{(S)} = \text{Entropy}(S) - I(\text{Attribute}) = 0.94 - 0.892$$

$$= 0.048$$

Humidity	Yes	No	Entropy
High	3	4	0.985
Normal	6	1	0.591

$$\text{Entropy}(\text{humidity}) = 0.940$$

$$I(\text{Humidity}) = 0.788, \quad \underline{\underline{\text{Gain} = 0.152}}$$

Temperature	Yes	No	Entropy
Cool	3	1	0.811
Hot	2	2	

Mild	4	2	0.918	Entropy (temp) = 0.940
------	---	---	-------	------------------------

$I(\text{temperature}) = 0.911, \text{Gain} = 0.029$

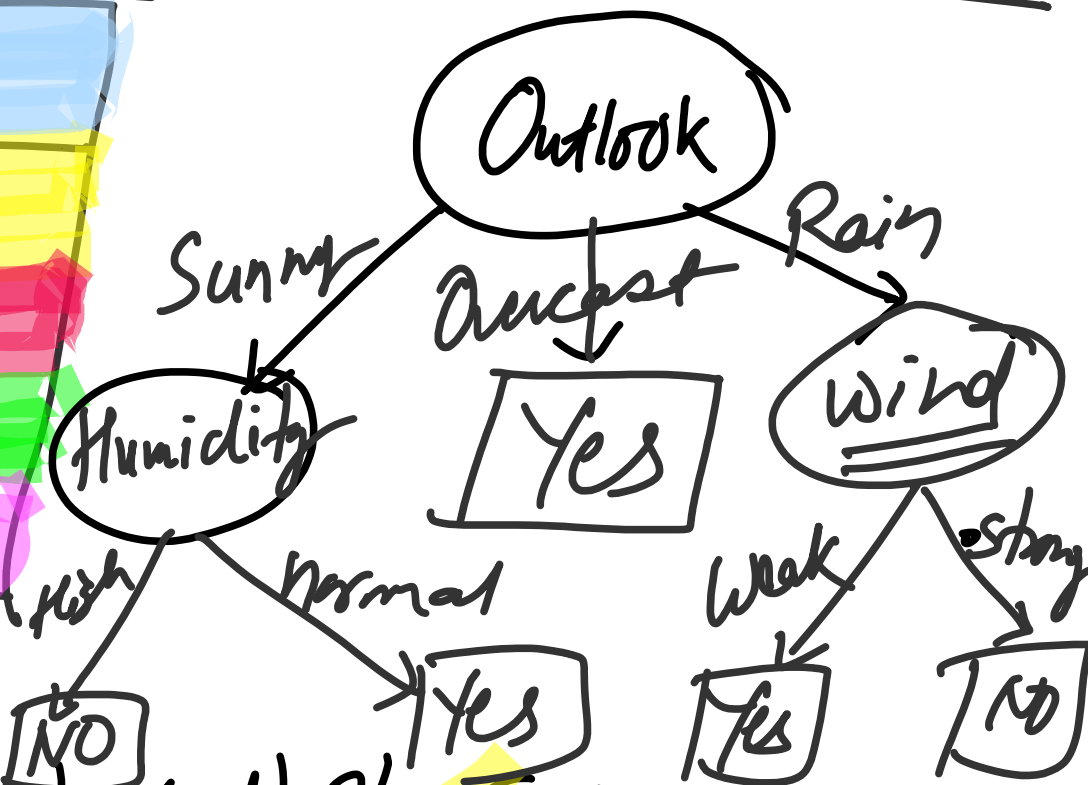
Outlook

Entropy = 0.941

$I(\text{outlook}) = 0.693$
 $\text{Gain}(\text{outlook}) = 0.247$

Outlook	Yes	No	Entropy
Sunny	2	3	0.971
Overcast	4	0	0
Rain	3	2	0.971

Attribute	Gain
Outlook	0.247
Temperature	0.029
Humidity	0.152
Wind	0.048



Outlook	Temperature	Humidity	Wind	Play Tennis
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Sunny	Mild	High	Weak	No
Sunny	Cool	High	Weak	No

No	No	No	Yes	No	Allergy
No	No	Yes	No	No	Strep Throat
Yes	No	No	Yes	Yes	Allergy
No	Yes	No	Yes	Yes	Cold
Yes	Yes	No	Yes	Yes	Cold

Total sample space (S) = Strep Throat + Allergy + Cold = 10

Info Gain(S) = $-\left[\frac{3}{10} \log_2\left(\frac{3}{10}\right) + \frac{3}{10} \log_2\left(\frac{3}{10}\right) + \frac{4}{10} \log_2\left(\frac{4}{10}\right)\right] = 1.57$

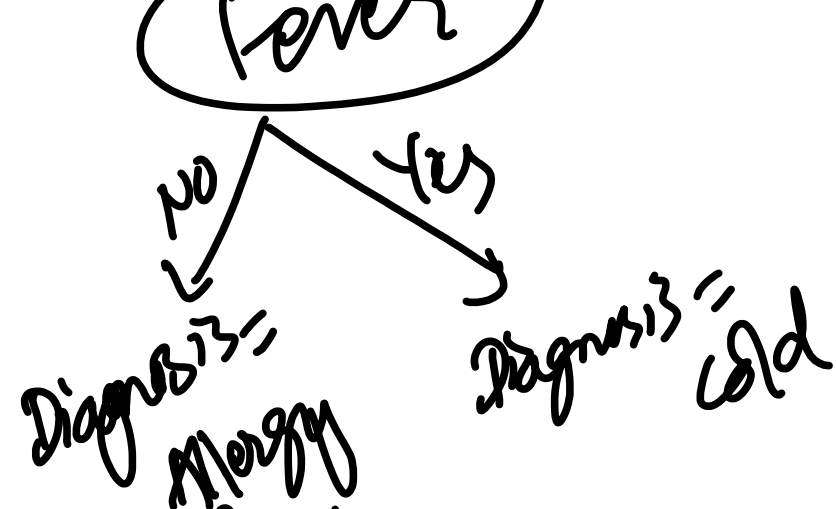
Finding splitting attribute: (Select attribute with highest gain)

(1) Strep Throat:-

Decision	S.T.	A	C
Yes	2	1	2
No	1	2	2

Gain of each attribute	
Attribute	Gain
Strep Throat	0.05
Fever	0.72
Swollen Glands	0.88
Congestion	0.46
Headache	0.05



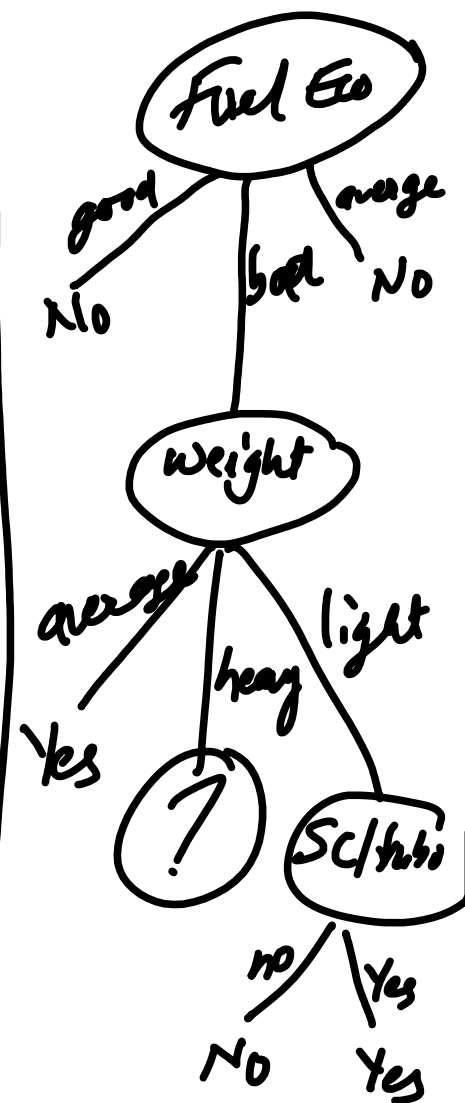


Diagnosis: cold

ID3 Question

from the following data, find which car is fast using decision tree with ID3 algorithm.

Model	Engine	SC/Murbo	Weight	Fuel Eco.	Fast
Prius	small	no	average	good	no
Civic	small	no	light	average	no
WRX STI	small	yes	average	bad	yes
M3	medium	no	heavy	bad	yes
RS4	large	no	average	bad	yes
GTI	medium	no	light	bad	no
XJR	large	yes	heavy	bad	no
S500	large	no	heavy	bad	no
911	medium	yes	light	bad	yes
Corvette	large	no	average	bad	yes
Insight	small	no	light	good	no
RSX	small	no	average	average	no
IS350	medium	no	heavy	bad	no
MR2	small	yes	average	average	no
E320	medium	no	heavy	bad	no



CART

- CART (for classification and regression tree) is often used as generic acronym for decision tree, although it is a specific implementation.
- Similar to C4.5, CART can handle continuous attributes

- CART uses the Gini diversity index which is defined in equation as

$$Gini_x = 1 - \sum_{x \in X} p(x)^2$$

- CART constructs a sequence of subtrees, uses cross-validation to estimate the misclassification cost of each subtree and chooses the one with the lowest cost.

CART 2 Formulas

Gini index is a metric for classification task in CART

$$(1) \text{ Gini index(attribute=value)} = 1 - \sum_{i=1}^N (P_i)^2$$

$$(2) \text{ Gini index(attribute)} = \sum_{v \in \text{value}} P_v \times GI(v)$$

Outlook	Yes	No	No. of metrics
Sunny	2	3	5
Overcast	4	0	4
Rain	3	2	5

$$\begin{aligned} \text{Gini(outlook=sunny)} &= 1 - \sum_{i=1}^N (P_i)^2 \\ &= 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48 \end{aligned}$$

$$\begin{aligned} \text{Gini(outlook=overcast)} &= 1 - \left(\frac{4}{4}\right)^2 - \left(\frac{0}{4}\right)^2 \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Gini(outlook=rain)} &= 1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 \\ &= 0.48 \end{aligned}$$

Now calculate weighted sum of Gini Index for outlook

Calculate weighted sum of Gini index for outlook,

$$\text{Gini (outlook)} = \sum_{v=\text{value}} P_v \times G(v) = \underbrace{\frac{5}{14} \times 0.48}_{v=\text{sunny}} + \underbrace{\frac{4}{14} \times 0}_{v=\text{overcast}} + \underbrace{\frac{5}{14} \times 0.48}_{v=\text{rain}} = 0.342$$

Temperature	Yes	No	No. of metrics
Hot	2	2	4
Mild	4	2	6
Cool	3	1	4

$$\text{Gini (temperature=hot)} = 1 - \left(\frac{2}{4}\right)^2 - \left(\frac{2}{4}\right)^2 = 0.5$$

$$\text{Gini (temperature=mild)} = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.444$$

$$\text{Gini (temperature=cool)} = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

Weighted sum for temperature, $\text{Gini (temperature)} = \frac{4}{14} \times 0.5 + \frac{6}{14} \times 0.444 + \frac{4}{14} \times 0.375$
 $= 0.440$

Humidity	Yes	No	No. of metrics
High	3	4	7
Normal	6	1	7

$$\text{Gini (humidity=high)} = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.489$$

$$\text{Gini (humidity=normal)} = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0.244$$

Weighted sum of Gini index for humidity

Gini Humidity) = $\frac{7}{14} \times 0.489 + \frac{7}{14} \times 0.244 = 0.367$

Wind	Yes	No	No. of metrics
Weak	6	2	8
Strong	3	3	6

Gini (wind=weak) = $1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.375$

Gini (wind=strong) = $1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 0.5$

Weighted sum of Gini Indices for wind,

Gini (wind) = $\frac{8}{14} \times 0.375 + \frac{6}{14} \times 0.5 = 0.428$

Decision for node value

Features	Gini Index
Outlook	0.342
Temperature	0.440
Humidity	0.367
Wind	0.428

lowest value is the decision factor in CART

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

Outlook

Sunny

overcast

Rain

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
1	Sunny	-	-	-	-
2	Sunny	-	-	-	-
8	Sunny	-	-	-	-
9	Sunny	-	-	-	-

[3, 7, 12, 13]

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
3	Overcast	Hot	High	Weak	Yes
7	Overcast	Cool	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes

[4, 5, 6, 10, 14]

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
10	Rain	Mild	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

11	1				
----	---	--	--	--	--

Sunny	Yes	No	No. of instances
Temperature			
Hot	0	2	2
Mild	1	0	1
Cool	1	1	2

Gini (outlook=sunny and temperature=hot)

$$= 1 - \left(\frac{0}{2}\right)^2 - \left(\frac{2}{2}\right)^2 = 0$$

Weighted sum

$$\text{Gini (outlook=sunny)} = \frac{2}{5} \times 0 + \frac{1}{5} \times 0 + \frac{2}{5} \times 0.5 = 0.2$$

Outlook=Sunny	Gini Index
Features	
Temperature	
Humidity	
Wind	

