

$$P(\text{No} | \text{New Instance}) = P(\text{No}) * P(\text{color} = \text{Red} | \text{No}) * P(\text{Type} = \text{SUV} | \text{No}) * P(\text{origin} = \text{Domestic} | \text{No})$$

$$= 0.5 \times \frac{2}{5} \times \frac{3}{5} \times \frac{3}{5} = 0.072$$

$$P(\text{Yes} | \text{New Instance}) < P(\text{No} | \text{New Instance})$$

The new car is not getting stolen.

Q.9 The following table shows the dataset on hand. Our target is to predict whether to play golf or not. Draw a decision tree using the ID3 algorithm to find a feature best suitable as its root.

Outlook	Temp	Humidity	Wind	Play Tennis?
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rain	Mild	High	False	Yes
Rain	Cool	Normal	True	Yes
Rain	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rain	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rain	Mild	High	True	No

ID3 Algorithm :-

The attribute which gives the maximum information, it will be considered as root node.

Attribute 1:- Outlook

values(outlook) = Sunny, Overcast, Rain

$$S = \begin{matrix} 9+ & 5- \\ \uparrow & \uparrow \\ \text{yes} & \text{No} \end{matrix} \quad \text{Entropy}(S) = \frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

This is entropy for entire dataset, we have calculate entropy for sunny, overcast Rain.

$$S_{\text{sunny}} = [2+, 3-] \quad \text{Entropy}[S_{\text{sunny}}] = \frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$S_{\text{overcast}} = [4+, 0-] \quad \text{Entropy}[S_{\text{overcast}}] = \frac{4}{4} \log_2 \frac{4}{4} + \frac{0}{4} \log_2 \frac{0}{4} = 0$$

$$S_{\text{rain}} = [3+, 2-] \quad \text{Entropy}[S_{\text{rain}}] = \frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} = 0.971$$

$$\text{Information Gain}(S, \text{Outlook}) = \text{Entropy}(S) - \sum_S \frac{|S_v|}{|S|} \text{Entropy}(S_v) \quad \forall v \in (\text{sunny, overcast, Rain})$$

$$= 0.94 - \frac{5}{14} \times 0.971 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.971 = 0.2466$$



- Attribute 2 :- Temperature  
Values (Temperature) = Hot, Mild, Cool

$$S_{Hot} = [2+, 2-] \text{Entropy}[S_{Hot}] = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4}$$

$$S_{Mild} = [4+, 2-] \text{Entropy}[S_{Mild}] = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.9182$$

$$S_{Cool} = [3+, 1-] \text{Entropy}[S_{Cool}] = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.8112$$

$$\text{Information Gain}(S, \text{Temp}) = \text{Entropy}(S) - \sum_{v \in \{Hot, mild, cool\}} \frac{|S_v|}{S} \text{Entropy}(S_v)$$

$$= 0.94 - \frac{4}{14} \times 1 - \frac{6}{14} \times 0.9182 - \frac{4}{14} \times 0.8112 = 0.029$$

- Attribute 3 :- Humidity  
Values (Humidity) :- High, Normal

$$S_{High} = [3+, 4-] \text{Entropy}(S_{High}) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.9852$$

$$S_{Normal} = [6+, 4-] \text{Entropy}(S_{Normal}) = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7}$$

$$= 0.5916$$

Information Gain (S, Humidity)

$$= \text{Entropy}(S) - \sum \frac{|S_v|}{S} \cdot \text{Entropy}(S_v)$$

$v \in \{\text{High, Normal}\}$

$$= 0.94 - \frac{7}{14} \times 0.9852 - \frac{7}{14} \times 0.5916$$

$$= 0.1516$$

Attribute 4 = Wind

Values(Wind) = Weak, Strong

$$S_{\text{False}} = [6+, 2-] \text{Entropy}(S_{\text{False}}) = \frac{6}{8} \log_2 \frac{6}{8} + \frac{2}{8} \log_2 \frac{2}{8}$$

$$= 0.8112$$

$$S_{\text{True}} = [3+, 3-] \text{Entropy}(S_{\text{True}}) = \frac{3}{6} \log_2 \frac{3}{6} + \frac{3}{6} \log_2 \frac{3}{6}$$

$$\text{Information Gain}(S, \text{wind}) = \text{Entropy}(S) - \frac{|S_v|}{S} \text{Entropy}(S_v)$$

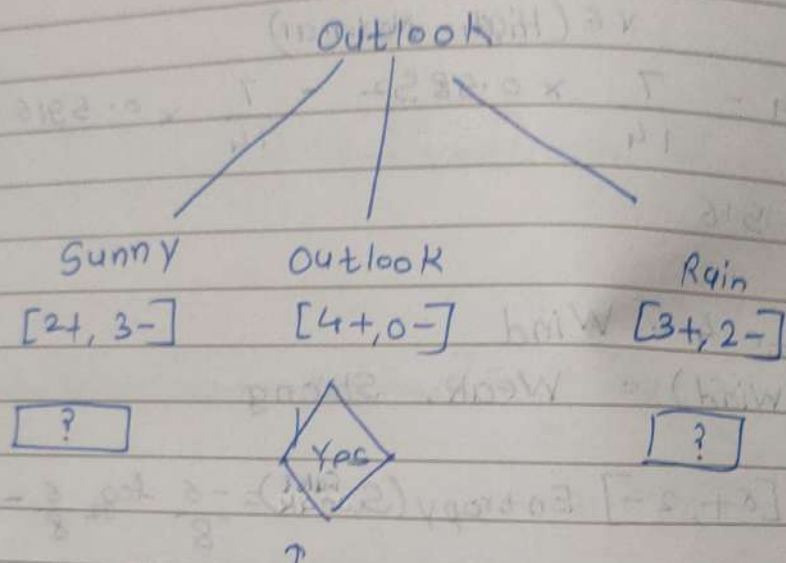
$v \in \{\text{False, True}\}$

$$= 0.94 - \frac{8}{14} \times 0.8112 - \frac{6}{14} \times 1.1119$$

$$= 0.04788$$



$IG(S, outlook) = 0.2466 \Rightarrow \text{maximum}$   
 $IG(S, Temperature) = 0.029$   
 $IG(S, Humidity) = 0.1516$   
 $IG(S, wind) = 0.04788$



It is Yes because there are no values

Only consider Sunny's rows for further calc.

Day	Temp.	Humidity	Wind	Play tennis
D <sub>1</sub>	Hot	High	False	No
D <sub>2</sub>	Hot	High	True	No
D <sub>3</sub>	Mild	High	False	No
D <sub>4</sub>	Cool	Normal	False	Yes
D <sub>5</sub>	Mild	Normal	True	Yes

Attribute :- Temperature

$$S_{\text{Sunny}} = [2+, 3-] \text{ Entropy}[S_{\text{Sunny}}] = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.97$$

$$S_{\text{Hot}} = [0+, 2-] \text{ Entropy}[S_{\text{Hot}}] = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} = 0$$

$$S_{\text{Mild}} = [1+, 1-] \text{ Entropy}[S_{\text{Mild}}] = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$S_{\text{Cool}} = [1+, 0-] \text{ Entropy}[S_{\text{Cool}}] = -\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} = 0$$

Information Gain ( $S_{\text{Sunny}}, \text{Temp}$ ) = Entropy(S) -  $\sum \frac{S_v}{S} \text{Ent}(S_v)$

$$1 \times \frac{5}{5} - \frac{2}{5} \times 0.97 - \frac{2}{5} \times 0 - \frac{1}{5} \times 1 - \frac{1}{5} \times 0 = 0.57$$

$$= 0.97 - \frac{2}{5} \times 0.97 - \frac{2}{5} \times 0 - \frac{1}{5} \times 1 - \frac{1}{5} \times 0$$

$$= 0.97 - 0.388 = 0.57$$

Attribute :- Humidity

$$S_{\text{High}} = [0+, 3-] \text{ Entropy}(S_{\text{High}}) = -\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} = 0$$

$$S_{\text{Normal}} = [2+, 0-] \text{ Entropy}(S_{\text{Normal}}) = -\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} = 0$$



$$\text{Information Gain}(S_{\text{Sunny}}, \text{Humidity}) = \text{Entropy}(S) - \sum \frac{|S_i|}{S} \text{Entropy}(S_i)$$

$$= 0.97 - 0 - 0 = 0.97$$

Attribute - Wind

$$S_{\text{False}} = [1+, 2-] \text{Entropy}(S_{\text{False}}) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.9182$$

$$S_{\text{True}} = [1+, 1-] \text{Entropy}(S_{\text{True}}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

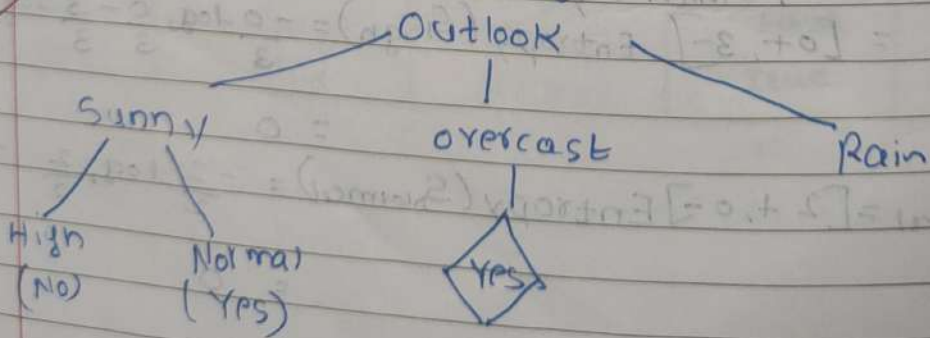
$$\text{Gain}(S_{\text{Sunny}}, \text{Wind}) = 0.97 - \frac{3}{5} \times 0.9182 - \frac{2}{5} \times 1$$

$$= 0.0192$$

$$\text{IG}(S_{\text{Sunny}}, \text{Temp}) = 0.57$$

$$\text{IG}(S_{\text{Sunny}}, \text{Humidity}) = 0.97 \Rightarrow \text{maximum}$$

$$\text{IG}(S_{\text{Sunny}}, \text{Wind}) = 0.0192$$



Now only consider Rain's rows for further calc.

Temperature	Humidity	Wind	Play Tennis
Mild	High	Weak	Yes
Cool	Normal	Weak	Yes
Cool	Normal	Strong	No
Mild	Normal	Weak	Yes
Mild	High	Strong	No

Attribute: Temperature

$$S_{\text{mild}} = [2+, 1-] \text{ Entropy}(S_{\text{mild}}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ = 0.9183$$

$$S_{\text{cool}} = [1+, 1-] \text{ Entropy}(S_{\text{cool}}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \\ = 1.0$$

$$IG_{(\text{mild})} = \text{Entropy}(S_{\text{rain}}) - \sum \frac{|S_v|}{S} \text{Entropy}(S_v) \\ \forall v \in (\text{mild, cool}) \\ = 0.971 - \frac{3}{5} \times 0.918 - \frac{2}{5} \times 1 \\ = 0.8912$$



Attribute: Humidity  
 $S_{high} = [1+, 1-]$  Entropy( $S_{high}$ ) =  $-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$   
 $= 1$

$S_{normal} = [2+, 1-]$  Entropy( $S_{normal}$ ) =  $-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$   
 $= 0.9182$

$IG(S_{rain}, Humidity) = 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.9182$

$= 0.0912$

Attribute: Wind

$S_{false} = [3+, 0-]$  Entropy( $S_{false}$ ) =  $-\frac{3}{3} \log_2 \frac{3}{3} - 0$   
 $= 0$

$S_{true} = [0+, 2-]$  Entropy( $S_{true}$ ) =  $-\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2}$   
 $= 0$

~~$IG(S_{rain}, Wind) = 0.97 - \frac{3}{5} \times 0 - \frac{2}{5} \times 0$~~

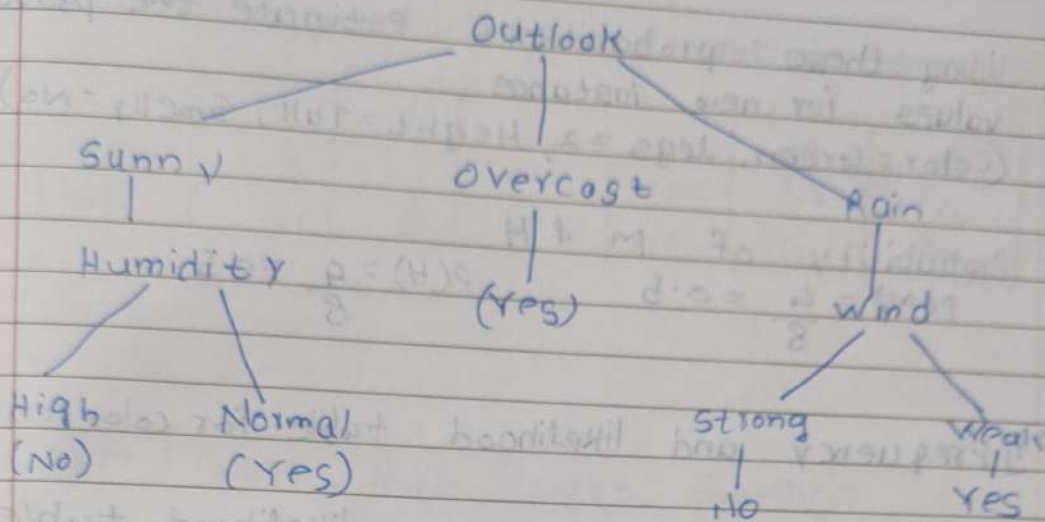
~~$= 0.97$~~

$IG(Temp) = 0.0912$

$IG(Humidity) = 0.0912$

$IG(Wind) = 0.97$

$\Rightarrow$  maximum



2.3. Estimate conditional probabilities (Naive Bayes Classifier) of each attribute {color, legs, height, smelly} for species classes: {M, H} using data given in table.

Color	Legs	Height	Smelly	Species
White	3	Short	Yes	M
Green	2	Tall	No	M
Green	3	Short	Yes	M
White	3	Short	Yes	M
Green	2	Short	No	H
White	2	Tall	No	H
White	2	Tall	No	H
White	2	Short	Yes	H



Using these probabilities estimate the probability values for new instance  
(Color = Green, legs = 2, Height = Tall, Smelly = No)

→ Probability of M & H  
 $P(M) = \frac{4}{8} = 0.5$        $P(H) = \frac{4}{8} = 0.5$

1] Frequency and likelihood table for color feature

Frequency table

	Species	
	M	H
Color White	2	3
Green	2	2

likelihood table

	P(M)	P(H)
White	$\frac{2}{4}$	$\frac{3}{4}$
Green	$\frac{2}{4}$	$\frac{1}{4}$

2] Frequency & likelihood table for legs

	Species	
	M	H
3	3	0
2	1	4

3	
2	

Species	
M	H
$\frac{3}{4}$	0
$\frac{1}{2}$	1

3] Frequency & likelihood table for height

	Species	
	M	H
Short	3	2
Tall	1	2

Species	
M	H
$\frac{3}{4}$	$\frac{1}{2}$
$\frac{1}{4}$	$\frac{1}{2}$

4] Frequency & likelihood table for smelly

	M	H		M	H
Yes	3	1	Yes	3/4	1/4
No	1	3	No	1/4	3/4

Probability:-

$$P(M/x) =$$

$$P(M / \text{color, legs, Height, Smelly}) =$$

$$= P(\text{color} | M) * P(\text{legs} | M) * P(\text{Height} | M) * P(\text{Smelly} | M)$$

$$= \frac{2}{4} * \frac{1}{2} * \frac{1}{4} * \frac{1}{4} * \frac{4}{8}$$

$$= \frac{1}{256} = 0.0039$$

~~Probability -~~

~~$$P(H/x) = P(H / \text{color, legs, Height, Smelly})$$~~

~~$$= \frac{1}{4} * 1 * \frac{1}{2} * \frac{3}{4} * \frac{4}{8} = \frac{3}{64}$$~~

~~$$= 0.046$$~~

$$P(H) = \frac{25+12}{38} = 1 - \left[ \frac{2}{38} + \frac{12}{38} \right] = 1 - \frac{14}{38} = \frac{24}{38} = \frac{12}{19}$$



Q.2 Draw a decision tree for a given dataset using the Gini Index Algorithm. You can develop tree until level 2.

Age	Income	Age	Student	Credit rating	Buys Computer
<=30	High		No	Fair	No
<=30	High		No	Excellent	No
31-40	High		No	Fair	Yes
>40	Medium		Yes	Fair	Yes
>40	Low		Yes	Excellent	No
>40	Low		Yes	Excellent	Yes
<=30	Medium		No	Fair	No
<=30	Low		Yes	Fair	Yes
>40	Medium		Yes	Fair	Yes
<=30	Medium		Yes	Excellent	Yes
31-40	Medium		No	Excellent	Yes
31-40	High		Yes	Fair	Yes
>40	Medium		No	Excellent	No

→ There are two possible output values - Yes & No.  
The data has 9 instances of Yes & 5 of No.

$$Gini(S) = 1 - \left[ \left( \frac{9}{14} \right)^2 + \left( \frac{5}{14} \right)^2 \right] = 1 - \left[ \frac{81 + 25}{196} \right] = \frac{116}{196} = 0.592$$

Gini Index for Age attribute

It has 3 possible attributes =  $\leq 30$ ,  $31-40$ ,  $>40$

For age  $\leq 30$

$$\text{Gini}(S) = 1 - \left[ \left( \frac{2}{5} \right)^2 + \left( \frac{3}{5} \right)^2 \right] = 1 - \left[ \frac{4+9}{25} \right] = \frac{12}{25} = 0.48$$

For age  $31-40$

$$\text{Gini}(S) = 1 - \left[ \left( \frac{4}{4} \right)^2 + (0)^2 \right] = 0$$

Age  $>40$

$$\text{Gini}(S) = 1 - \left[ \left( \frac{3}{5} \right)^2 + \left( \frac{2}{5} \right)^2 \right] = \frac{12}{25} = 0.48$$

$$\text{weighted avg(Age)} = \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48$$

$$= 0.3428$$

Gini Index for IncomeAge

It has 3 possible attr. = High, Low, Medium

For IncomeAge = Low High

$$\text{Gini}(S) = 1 - \left[ \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right] = 1 - \left[ \frac{8}{16} \right] = 0.5$$

For IncomeAge = Low

$$\text{Gini}(S) = 1 - \left[ \left( \frac{3}{4} \right)^2 + \left( \frac{1}{4} \right)^2 \right] = 1 - \left[ \frac{10}{16} \right] = 0.375$$

For IncomeAge = Medium

$$\text{Gini}(S) = 1 - \left[ \left( \frac{4}{8} \right)^2 + \left( \frac{2}{8} \right)^2 \right] = 1 - \left[ \frac{20}{32} \right] = 0.444$$



$$\begin{aligned} \text{weighted avg (IncomeAge)} \\ &= \frac{4}{14} \times 0.5 + \frac{4}{14} \times 0.375 + \frac{6}{14} \times 0.444 \\ &= 0.4401 \end{aligned}$$

Gini Index for student Attribute  
for student = No

$$\text{Gini}(S) = 1 - \left[ \left( \frac{3}{7} \right)^2 + \left( \frac{4}{7} \right)^2 \right] = \left[ \frac{49 - 25}{49} \right] = 0.4897$$

For student = Yes

$$\text{Gini}(S) = 1 - \left[ \left( \frac{6}{7} \right)^2 + \left( \frac{1}{7} \right)^2 \right] = 1 - \left[ \frac{37}{49} \right] = 0.2448$$

weighted age (Student)

$$= \frac{7}{14} \times 0.4897 + \frac{7}{14} \times 0.2448 = 0.3672$$

Gini Index for credit - ranking

for credit - ranking = Fair

$$\text{Gini}(S) = 1 - \left[ \left( \frac{6}{8} \right)^2 + \left( \frac{2}{8} \right)^2 \right] = 1 - \frac{40}{64} = 0.375$$

for credit - ranking = Excellent

$$\text{Gini}(S) = 1 - \left[ \left( \frac{3}{6} \right)^2 + \left( \frac{3}{6} \right)^2 \right] = 1 - \frac{18}{36} = 0.5$$

weighted avg (credit - ranking)

$$= \frac{8}{14} \times 0.375 + \frac{6}{14} \times 0.5$$

$$= 0.2142 + 0.2142 = 0.4284$$

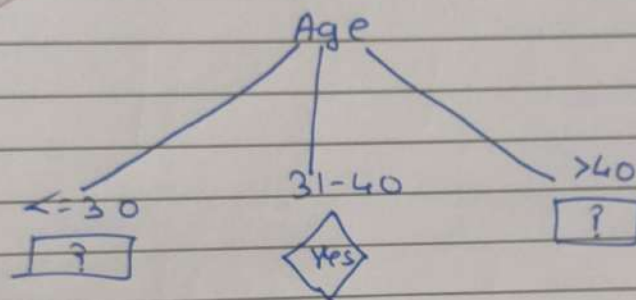
$$\text{Gini}(\text{Age}) = 0.3428 \Rightarrow \text{Minimum}$$

$$\text{Gini}(\text{IncomeAge}) = 0.4401$$

$$\text{Gini}(\text{Student}) = 0.3872$$

$$\text{Gini}(\text{credit-ranking}) = 0.4284$$

Age attr. has minimum Gini Index



8/10/11/23