

Assignment: Decision Tree

1. Difference between Entropy and Gini Impurity

Decision Tree

DATE 06/11/2022

* Assignment : 01

Entropy

1. Mathematical formula:

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

eg. If we have two class

$$E(s) = -P_Y \log_2(P_Y) - P_N \log_2(P_N)$$

where,

$E(s) \rightarrow$ Entropy

$P_Y \rightarrow$ Probability of 'yes'

$P_N \rightarrow$ Probability of 'no.'

Gini Impurity

1. Mathematical formula.

$$GI = 1 - (P_Y^2 + P_N^2)$$

where,

$GI =$ gini impurity

$P_Y \rightarrow$ Probability of 'yes'

$P_N \rightarrow$ Probability of 'no.'

2. For binary classification problem

$$E(s)_{\min} = 0 \text{ \& } E(s)_{\max} = 1$$

For multiclass classification

$$E(s)_{\min} = 0 \text{ \& } E(s)_{\max} = > 1$$

2. For binary classification

$$GI_{\min} = 0 \text{ \& } GI_{\max} = 0.5$$

3. Entropy is computationally slower due logarithmic calculation

3. Gini is computationally faster as compared to entropy due to squaring.

4. Entropy gives slightly better results and balanced trees.

4. Gini is less accurate as compared to entropy

5. Entropy is used for EDA

5. Gini is used to minimize misclassification.

2. Take a tennis dataset and build a decision tree from scratch with mathematical representation

Dataset:

Day	outlook	temperature	humidity	wind	Decision
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
3	overcast	hot	high	weak	Yes
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
7	overcast	cool	normal	strong	Yes
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
10	rainfall	mild	normal	weak	Yes
11	sunny	mild	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	rainfall	mild	high	strong	No

Step 1: Selecting root node.

Calculate gini index of each feature and select one which has less gini index.

a) Outlook	decision	T.V.	P(Y)	P(N)
sunny	$Y=2, N=3$	5	$2/5$	$3/5$
overcast	$Y=4, N=0$	4	$4/4=1$	0
rainfall	$Y=3, N=2$	5	$3/5$	$2/5$
		<u>14</u>		

$$G_{\text{sunny}} = 1 - P_Y^2 - P_N^2 = 1 - \frac{4}{25} - \frac{9}{25} = 0.48$$

$$G_{\text{overcast}} = 1 - 1 - 0 = 0$$

$$G_{\text{rainfall}} = 1 - \frac{9}{25} - \frac{4}{25} = 0.48$$

$$\text{Gini column} = \frac{\sum \text{no. of instance for class}}{\text{total no. of instance}} \times \text{Gini class}$$

$$\text{Gini outlook} = \frac{5}{14} \times 0.48 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.48$$

$$\text{Gini outlook} = 0.34$$

temperature	decision	T.V.	$P(Y)$	$P(N)$
hot	$y=2, N=2$	4	2/4	2/4
mild	$y=4, N=2$	6	4/6	2/6
cool	$y=3, N=1$	4	3/4	1/4
		14		

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

$$\text{Gini}_{\text{mild}} = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.44$$

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

$$\text{Gini}_{\text{temp}} = \frac{4}{14} \times 0.5 + \frac{6}{14} \times 0.44 + \frac{4}{14} \times 0.375 = 0.43$$

humidity	decision	T.V.	$P(Y)$	$P(N)$
high	$y=3, N=4$	7	3/7	4/7
normal	$y=6, N=1$	7	6/7	1/7
		14		

$$\text{Gini}_{\text{high}} = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.49$$

$$\text{Gini}_{\text{normal}} = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0.24$$

$$\text{Gini}_{\text{humidity}} = \frac{7}{14} \times 0.49 + \frac{7}{14} \times 0.24 = 0.36$$

d)	Wind	decision	T.V.	P(y)	P(N)
	weak	$y=6, N=2$	8	6/8	2/8
	strong	$y=3, N=3$	6	3/6	3/6
			14		

$$Gini_{weak} = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.37$$

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

$$Gini_{wind} = \frac{8}{14} \times 0.37 + \frac{6}{14} \times 0.5 = 0.42$$

Gini outlook	Gini temp	Gini humidity	Gini wind
0.34	0.43	0.36	0.42

Since, Gini outlook has less value compared to other three features selected as root node.

Step 2: Further splitting of root node to get lowest label. (Sunny, overcast, rainfall)

* Splitting Sunny.

a) Gini temp \Rightarrow

temp	decision	T.V.	P(y)	P(N)
hot	$y=0, N=2$	2	0	2/2=1
cool	$y=1, N=0$	1	1/1=1	0
mild	$y=1, N=1$	2	1/2	1/2
		5		

classmate

(3)

D	outlook	temperatu	humidit	wind	Decisio
1	sunny	hot	high	weak	No
2	sunny	hot	high	strong	No
8	sunny	mild	high	weak	No
9	sunny	cool	normal	weak	Yes
11	sunny	mild	normal	strong	Yes

$$Gini_{hot} = 1 - (0)^2 - (1)^2 = 0$$

$$Gini_{cool} = 1 - (1)^2 - (0)^2 = 0$$

$$Gini_{mild} = 1 - \left(\frac{1}{4}\right) - \left(\frac{1}{4}\right) = 0.5$$

$$Gini_{temp} = \frac{2}{5} \times 0 + \frac{1}{5} \times 0 + \frac{2}{5} \times 0.5 = 0.2$$

b) $Gini_{humidity} =$

humidity	decision	TV	$P(y)$	$P(n)$
high	$y=0, N=3$	3	$0/3=0$	$3/3=1$
normal	$y=2, N=0$	$\frac{2}{5}$	$2/2=1$	$0/2=0$

$$Gini_{high} = 1 - (0)^2 - (1)^2 = 0$$

$$Gini_{normal} = 1 - (1)^2 - (0)^2 = 0$$

$$Gini_{humidity} = \frac{3}{5} \times 0 + \frac{2}{5} = 0$$

c) Gini wind

wind	decision	TN.	P(Y)	P(N)
strong	$y=1, N=1$	2	$1/2$	$1/2$
weak	$y=1, N=2$	3	$1/3$	$2/3$
		05		

$$\text{Gini}_{\text{strong}} = 1 - \left(\frac{1}{4}\right) - \left(\frac{1}{4}\right) = 0.5$$

$$\text{Gini}_{\text{weak}} = 1 - \left(\frac{1}{9}\right) - \left(\frac{4}{9}\right) = 0.44$$

$$\text{Gini}_{\text{wind}} = \frac{2}{5} \times 0.5 + \frac{3}{5} \times 0.44 = 0.46$$

Since, $\text{Gini}_{\text{temp}} = 0.2$, $\text{Gini}_{\text{humidity}} = 0$, $\text{Gini}_{\text{wind}} = 0.46$

$\text{Gini}_{\text{humidity}}$ is selected as leaf node with no further splitting

* Splitting Overcast:

Since overcast has gini value 0, no further splitting required.

ID	outlook	temperature	humidity	wind	Decision
3	overcast	hot	high	weak	Yes
7	overcast	cool	normal	strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes

* Splitting of rainfall.

a) Gini temp :

temp.	decision	TV	P(Y)	P(N)
mild	y=2, N=1	3	2/3	1/3
cool	y=1, N=1	2	1/2	1/2
		<u>0.5</u>		

D	outlook	temperatu	humidit	wind	Decisio
4	rainfall	mild	high	weak	Yes
5	rainfall	cool	normal	weak	Yes
6	rainfall	cool	normal	strong	No
10	rainfall	mild	normal	weak	Yes
14	rainfall	mild	high	strong	No

$$\text{Gini}_{\text{mild}} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.44$$

$$\text{Gini}_{\text{cool}} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Gini}_{\text{temp}} = \frac{3}{5} \times 0.44 + \frac{2}{5} \times 0.5 = 0.46$$

b) Gini humidity

humidity	decision	TV	P(Y)	P(N)
high	y=1, N=1	2	1/2	1/2
normal	y=2, N=1	3	2/3	1/3
		<u>0.5</u>		

$$\text{Gini}_{\text{high}} = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Gini}_{\text{normal}} = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 0.30$$

$$\text{Gini}_{\text{humidity}} = \frac{2}{5} \times 0.5 + \frac{3}{5} \times 0.30 = 0.38$$

c) Gini wind

wind	decision	TV	P(Y)	P(N)
weak	y=3, N=0	3	3/3=1	0/3=0
strong	y=0, N=2	2	0/2=0	2/2=1
		<u>0.5</u>		

$$\text{Gini}_{\text{weak}} = 1 - (1)^2 - (0)^2 = 0$$

$$\text{Gini}_{\text{strong}} = 1 - (0)^2 - (1)^2 = 0$$

$$\text{Gini}_{\text{wind}} = \frac{3}{5} \times 0 + \frac{2}{5} \times 0 = 0$$

$Gini_{temp} = 0.46$, $Gini_{humidity} = 0.38$, $Gini_{wind} = 0$

Hence, $Gini_{wind}$ is selected as child node with no further splitting

