

Page No.
 (Date: / /)

Decision Tree Based clustering

ID3 (Iterative Dichotomiser)

It is a decision tree learning algorithm invented by Ross Quinlan used to generate a decision tree from a dataset.

ID3 is the precursor to C4.5 algorithm and is typically used in machine learning and natural language processing domains.

Example 1 - Play Tennis Database

<u>Day</u>	<u>Outlook</u>	<u>Temp</u>	<u>Humidity</u>	<u>wind</u>	<u>play</u>
1	Sunny	80	High	Weak	No
2	Sunny	81	High	Strong	No
3	Overcast	82	High	Weak	Yes
4	Rain	65	High	Weak	Yes
5	Rain	40	Normal	Weak	Yes
6	Rain	41	Normal	Strong	No
7	Overcast	42	Normal	Strong	Yes
8	Sunny	66	High	Weak	No
9	Sunny	43	Normal	Weak	Yes
10	Rain	67	Normal	Weak	Yes
11	Sunny	68	Normal	Strong	Yes
12	Overcast	69	High	Strong	Yes
13	Overcast	83	Normal	Weak	Yes
14	Rain	70	High	Strong	No

Steps

- ① Start with a training data set, which we call as 'S'. It should have attributes & classification
- ② Determine the best attribute in the data set S.
- ③ Split S into subsets that correspond to the possible values of the best attribute.
- ④ Make a decision tree node that contains the best attribute.
- ⑤ Recursively make new decision tree nodes, with the subsets of data created in step 3). Attributes can't be reused.
If there is no more attributes to split on, choose the most popular classification.

⇒ In the above example database, attributes are outlook, Temperature, humidity and wind

And The classification is → YES
→ NO

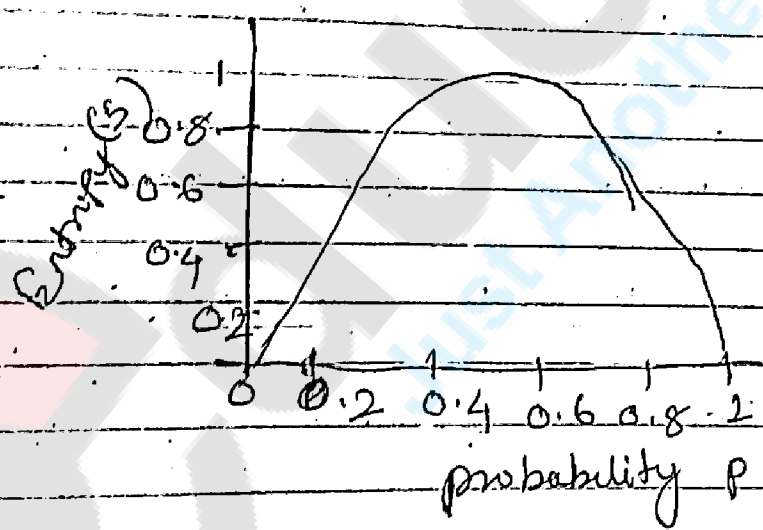
To find out the best attribute, we have to calculate (A) Entropy (B) Information gain

Entropy is an information theory and machine learning sense, measures the homogeneity of data set's distribution.

$$Entropy(S) = \sum_{i=1}^C -p_i \log_2 p_i$$

where C → corresponds to the number of different classifications.

p_i → corresponds to the proportion of the data with the classification i .



Entropy calculation

In the above example, two classes @ Yes & NO

$$Entropy(S) = \sum_{i=1}^2 -p_i \log_2 p_i$$

$$= -P_{Yes} \log_2 P_{Yes} - P_{No} \log_2 P_{No}$$

No. of Yes = 9 out of 14

$$P_{Yes} = P_{Yes} = \frac{9}{14}$$

No. of No = 5

$$P_{No} = \frac{5}{14}$$

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

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

Information Gain measures the reduction in entropy that results from partitioning the data on an attribute A, that is it represents how effective an attribute is at classifying the data.

Given a set of training data S and attribute A the formula for information gain is

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Calculation of Information Gain for all Attribute

① Outlook

values of Outlook = Sunny, Overcast, Rain

Page No.
 Date: / /

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

$$\left[\frac{|S_{\text{sunny}}|}{|S|} \text{Entropy}(S_{\text{sunny}}) \right] + \left[\frac{|S_{\text{overcast}}|}{|S|} \text{Entropy}(S_{\text{overcast}}) \right] + \left[\frac{|S_{\text{rain}}|}{|S|} \text{Entropy}(S_{\text{rain}}) \right]$$

$$S = 14 \quad (9 \text{ yes } 5 \text{ NO})$$

$$S_{\text{sunny}} = 5 \quad (2 \text{ yes } 3 \text{ NO})$$

$$\Rightarrow P_{\text{yes}} = \frac{2}{5} \quad P_{\text{no}} = \frac{3}{5}$$

$$S_{\text{overcast}} = 4 \quad (4 \text{ yes } 0 \text{ NO}) \quad P_{\text{yes}} = \frac{4}{4} \quad P_{\text{no}} = \frac{0}{4}$$

$$S_{\text{rain}} = 5 \quad (3 \text{ yes } 2 \text{ NO})$$

$$P_{\text{yes}} = \frac{3}{5} \quad P_{\text{no}} = \frac{2}{5}$$

$$\text{Entropy}(S_{\text{sunny}}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}}$$

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

$$= 0.971$$

$$\text{Entropy}(S_{\text{overcast}}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}}$$

$$= -1 \log_2 1 - 0 \log_2 0$$

$$= 0$$

$$\text{Entropy}(S_{\text{rain}}) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}}$$

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

$$= 0.971$$

$$\text{Gain}(S, \text{outlook}) = 0.940 - \left[\frac{5}{14} \times 0.971 \right] +$$

$$\left[\frac{4}{14} \times 0 \right] + \left[\frac{5}{14} \times 0.971 \right]$$

$$\text{Gain}(S, \text{outlook}) = 0.247$$

~~Gross~~ Temperature (no calculation for numerical data values)

Gain (S, Humidity) =

$$= \text{Entropy (S)} - \left(\frac{7}{14}\right) \text{Entropy (S}_{\text{High}}) - \frac{7}{14} \text{Entropy (S}_{\text{Normal}})$$

[∵ values of Humidity is High or Normal]
 High = 7 Normal = 7
 (6 Yes 1 No)

Entropy S_{High} = 7 (3 Yes 4 No)

$$\begin{aligned} \text{Entropy (S}_{\text{High}}) &= -P_{\text{Yes}} \log_2 P_{\text{Yes}} - P_{\text{No}} \log_2 P_{\text{No}} \\ &= -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \\ &= 0.985 \end{aligned}$$

$$\begin{aligned} \text{Entropy (S}_{\text{Normal}}) &= -P_{\text{Yes}} \log_2 P_{\text{Yes}} - P_{\text{No}} \log_2 P_{\text{No}} \\ &= -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \\ &= 0.592 \end{aligned}$$

Gain (S, Humidity) = 0.94

$$\begin{aligned} 0.940 - \frac{7}{14} \times 0.985 - \frac{7}{14} \times 0.592 \\ = 0.152 \end{aligned}$$

Wind

Page No. 5
Date: 11

values of wind's weak, strong

$$S_{\text{weak}} = 8 \text{ (6 Yes 2 No)}$$

$$S_{\text{strong}} = 6 \text{ (3 Yes 3 No)}$$

$$\begin{aligned} \text{Entropy}(S_{\text{weak}}) &= -P_{\text{Yes}} \log_2 P_{\text{Yes}} - P_{\text{No}} \log_2 P_{\text{No}} \\ &= -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \\ &= 0.811 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(S_{\text{strong}}) &= -P_{\text{Yes}} \log_2 P_{\text{Yes}} - P_{\text{No}} \log_2 P_{\text{No}} \\ &= -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Gain}(S, \text{wind}) &= \text{Entropy}(S) - \left(\frac{8}{14} \times 0.811 \right) \\ &\quad - \left(\frac{6}{14} \times 1 \right) \\ &= 0.940 - \left(\frac{8}{14} \right) \cdot 0.811 - \frac{6}{14} \\ &= 0.48 \end{aligned}$$

Among all the information Gain values the $\text{Gain}(S, \text{outlook})$ is high.

→ So outlook Attribute will be chosen as root.

→ The decision tree is expanded to cover outlook's possible values.

