

7/3/23

Module - 3.

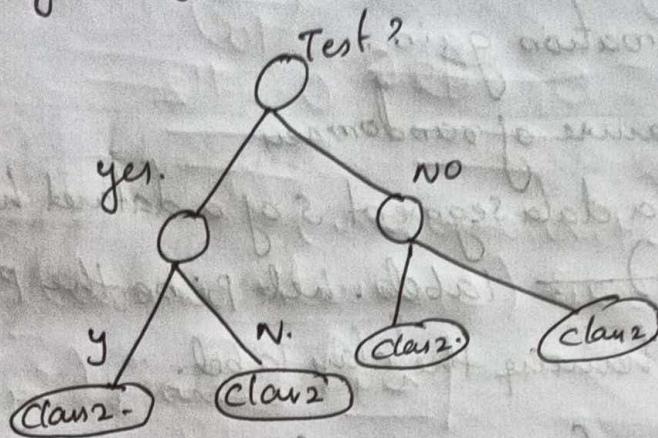
Tuesday Decision tree

clustering: unsupervised

supervised learning method.

It is a hierarchy representation of data.

Simplest type of classification method.



Nodes: Test on attribute.

Branches: outcome for the test.

Leaves - final classes.

Examples of decision tree algorithms.

ID3, C4.5, CART

Two types of decision tree.

1) Classification → o/p are classes.

2) Regression tree (predicting new value) - leaf nodes are numbers.

Consider the following data set.

1)

Name	Features				Class Label
	give birth	airbreath animal	aerial animal	has legs	
Human	yes	no	yes no	yes	mammal
Python	no	no	no	no	reptile
Salmon	no	yes	no	no	fish
frog	yes	semi	yes no	yes	amphibians
bat	yes	no	yes	yes	bird
Pigeon	no	no	yes	yes	bird
cat	yes	no	yes no	yes	mammal
shark	yes	yes	no	yes	fish
turtle	no	semi	yes no	yes	amphibians
Salmander	no	semi	yes no	yes	-

Feature selection measures

Two of the popular selection measures are Information gain & Gini index.

CART \rightarrow Gini index.

ID3 \rightarrow Information gain.

Entropy - measure of randomness.

Consider a data segment S of a dataset having c no. of classes. Let p_i be the proportion of examples in S having the class label i .

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

problem(1) Total no. of samples $|S| = 10$

At. of

$$\text{Probability of mammal} = \frac{2}{10} = \frac{1}{5}$$

$$\text{,, , , } \text{amph} = \frac{3}{10}$$

$$\text{,, , , } \text{bird} = \frac{2}{10}$$

$$\text{,, , , } \text{fish} = \frac{2}{10}$$

$$\text{,, , , } \text{reptile} = \frac{1}{10}$$

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

$$= -P_{\text{mam}} \log_2(P_{\text{mam}}) - P_{\text{amp}} \log_2(P_{\text{amp}}) -$$

$$P_{\text{bird}} \log_2(P_{\text{bird}}) - P_{\text{fish}} \log_2(P_{\text{fish}}) - P_{\text{reptile}} \log_2(P_{\text{reptile}})$$

$$= -\frac{2}{10} \log_2\left(\frac{2}{10}\right) - \frac{3}{10} \log_2\left(\frac{3}{10}\right) - \frac{2}{10} \log_2\left(\frac{2}{10}\right)$$

$$- \frac{2}{10} \log_2\left(\frac{2}{10}\right) - \frac{1}{10} \log_2\left(\frac{1}{10}\right)$$

$$= 0.46 + 0.52 + 0.46 + 0.46 + 0.33$$

$$= 2.232$$

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

A - attribute / feature

v - values of A

S - complete data set which is to be analyzed

S_v - sub-table contains value v.

Let A = "gives birth"

Values of A = {yes, no}.

$$\text{Information Gain}(S, \text{"gives birth"}) = \text{Entropy}(S) - \frac{|S_y|}{|S|} \text{Entropy}(S_{\text{yes}}) - \frac{|S_n|}{|S|} \text{Entropy}(S_{\text{no}})$$

$$(S_{\text{yes}}) = 0$$

$$\text{Entropy}(S, \text{"gives birth = yes"}) = -P_{\text{man}} \log_2(P_{\text{man}}) - P_{\text{bird}} \log_2(P_{\text{bird}}) - P_{\text{jinh}} \log_2(P_{\text{jinh}})$$

$$= -\frac{1}{4} \log_2\left(\frac{2}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right)$$

$$= \underline{\underline{1.5}}$$

$$\text{Entropy}(S, \text{"gives birth = no"}) = -P_{\text{man}} \log_2(P_{\text{man}}) - P_{\text{bird}} \log_2(P_{\text{bird}}) - P_{\text{jinh}} \log_2(P_{\text{jinh}}) - P_{\text{rep}} \log_2(P_{\text{rep}}) - P_{\text{amp}} \log_2(P_{\text{amp}})$$

$$= -\frac{1}{6} \log_2 \left(\frac{1}{6}\right) - \frac{1}{6} \log_2 \left(\frac{1}{6}\right) - \frac{1}{6} \log_2 \left(\frac{1}{3}\right)$$

$$\frac{3}{6} \log_2 \left(\frac{3}{6}\right) = \underline{\underline{1.792}}$$

$$\text{Information gain} = 2.24 - \left(\frac{4}{10}\right) \times 1.5 - \left(\frac{6}{10}\right) \times 1.792$$

$$= \underline{\underline{0.5768}} \quad \underline{\underline{0.5712}}$$

$A = \text{aerobic}$

values of $A = \{y, n, \text{aerosomi}\}$

Information Gain (S , aerobic animal) = Entropy (S)

$$= \frac{|S_y|}{S} \text{Entropy}(S_y) + \frac{|S_n|}{S} \text{Entropy}(S_n)$$

$$\text{Entropy}(S_y) = \frac{|S_y|}{S} \text{Entropy}(S_y)$$

$$\text{Entropy}(S, \text{"aerobic animal = yes"}) = -P_{fish} \log_2(P_{fish})$$

$$= 1.1092(1)$$

$\therefore \underline{\underline{0}}$ (pure class)

$$\text{Entropy}(S, \text{"aerobic animal = no"}) = -P_{mamm} \log_2(P_{mamm}) - P_{rep} \log_2$$

$$P(rep) - P_{bird} \log_2(P_{bird})$$

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

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

$$= \underline{\underline{0.0643}} \quad \underline{\underline{1.5}}$$

$$\text{Entropy}(S, \text{"aerobic animal = semi"}) = \underline{\underline{0.1092(2)}} \quad \underline{\underline{0}}$$

$$\text{Information gain} = 2.24 - 0 - \cancel{0.5} \times \frac{5}{10} \times 1.5 - 0.$$

$$= \underline{\underline{1.49}}$$

$A = \text{aerial}$

values of $A = \{Y, N\}$.

$$\text{Entropy}(S, "aerial \text{ animal} = \text{Yes}") = \underline{\underline{0}}$$

$$\begin{aligned} \text{Entropy}(S, "aerial \text{ animal} = \text{No}") &= -\frac{2}{8} \log_2 \left(\frac{2}{8}\right) - \frac{3}{8} \log_2 \left(\frac{3}{8}\right) \\ &\quad - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) = \frac{3}{8} \log_2 \left(\frac{3}{8}\right) \\ &= \underline{\underline{1.905}} \end{aligned}$$

$$\text{Information gain}(S, \text{aerial animal}) = \text{Entropy}(S) - [S_Y] \text{Entropy}(S_Y) -$$

$$\begin{matrix} | S_Y \\ | S_N \end{matrix}$$

$$= 2.24 - \cancel{0} \times 0 - \frac{8}{10} \times 1.905.$$

$$= \underline{\underline{0.716}}$$

$A = \text{has legs}$.

Value of $A = \{Y, N\}$

$$\begin{aligned} \text{Entropy}(S, "has \text{ legs} = \text{Yes}") &= -\frac{2}{7} \log_2 \left(\frac{2}{7}\right) - \frac{3}{7} \log_2 \left(\frac{3}{7}\right) \\ &\quad - \frac{2}{7} \log_2 \left(\frac{2}{7}\right) \\ &= \underline{\underline{-1.5}} \end{aligned}$$

$$\text{Entropy}(S, \text{"has legs = } \text{Y}") = P_{\text{no leg}} \log_2(P_{\text{no leg}}) - P_{\text{has leg}} \cdot$$

$(P_{\text{has leg}}).$

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

$$= \underline{\underline{0.918}}$$

$$\text{Information gain} = \text{Entropy}(S) - \frac{1}{|S|} \text{Entropy}$$

$$= 2.24 - \frac{7}{10} \times 1.5 - \frac{3}{10} \times 0.918$$

$$= \underline{\underline{0.8856}}$$

Aquatic animal is the root node because information gain of aquatic animal is high.

Find the root node of Decision tree using ID3 algorithm for the following data

Age	Completion	Type	Class (Profit.)
old	yes	Software	Down.
old	no	H/W	Down.
old	no	Hardware	Down.
mid	yes	Software	Down.
mid	yes	Hardware	Down.
mid	no	Hardware	Up.
mid	no	Software	U.P.
new	yes	S/W S	U.P.
new	no	H/W	Up
new	no	S/W L	Up

$$\text{Total Entropy } H = \sum -P_i \log_2 P_i$$

$$= -P_{\text{down}} \log_2(P_{\text{down}}) - P_{\text{up}} \log_2(P_{\text{up}})$$

$$= -\frac{5}{10} \log_2 \left(\frac{5}{10}\right) - \frac{5}{10} \log_2 \left(\frac{5}{10}\right)$$

$$= \underline{\underline{1.5}}$$

A = competition age age.

values of A = {old, mid, new}.

$$\text{Entropy}(s, "age = \cancel{\text{old}}") = -0$$

$$\text{Entropy}(s, "age = \cancel{\text{mid}}") = \cancel{\frac{2}{4} \log_2 \left(\frac{2}{4}\right)} - \cancel{\frac{2}{4} \log_2 \left(\frac{2}{4}\right)}$$

$$\text{Entropy}(s, "age = \cancel{\text{new}}") = \underline{\underline{1}}$$

$$\text{Entropy}(s, "age = \cancel{\text{mid}}") = -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right)$$

$$= (1.5)$$

$$\text{Entropy}(s, "age = \cancel{\text{new}}") = \underline{\underline{0}}$$

$$\text{Information gain } I = 1 - 0 - \frac{2}{10} * 1 - \frac{3}{10} * 0$$

$$= 0.6$$

A = competition.

values of A = {yes, no}.

$$\text{Entropy}(s, "competition = yes") = \frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \left(\frac{1}{4}\right)$$

$$= \underline{\underline{0.8112}}$$

$$\text{Entropy } (S, \text{"Comp = no"}) = \frac{2}{6} \log_2 \left(\frac{2}{6}\right) - \frac{4}{6} \log_2 \left(\frac{4}{6}\right)$$

$$= 0.918$$

$$\text{Information gain}_{(S, \text{comp})} = 1 - \frac{4}{10} \times 0.8112 - \frac{6}{10} \times 0.918$$

$$= \underline{\underline{0.124}}$$

~~Ex~~ $A = \text{software type}$

Values of $A = \{S/w, H/w\}$

$$\text{Entropy } (S, \text{"Type = "S/w"}) = \frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right)$$

$$= \underline{\underline{0}}$$

$$\text{Entropy } (S, \text{"Type = "H/w"}) = \frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right)$$

$$(S, \text{Type}) = \underline{\underline{1}}$$

$$\text{Info gain} = 1 - \frac{6}{10} \times 1 - \frac{4}{10} \times 1$$

$$= \underline{\underline{0}}$$

Age is the root node. $\text{outgoing} = A$

$\text{pos} = \frac{5}{10} = 0.5$

Instance no.	Class label	x_1	x_2
1	T	T	T
2	T	T	F
3	F	T	F
4	F	F	F
5	T	F	T
6	F	F	T

$$\text{Entropy}(S) = \sum p_i \cdot \log_2(p_i)$$

$$= p_1 \log_2(p_1) - p_0 \log_2(p_0).$$

$$= \frac{-3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right)$$

$$= \underline{\underline{1}}$$

$$A = x_1.$$

Values of $A = \{T, F\}$.

$$\text{Entropy}(S, "x_1 = T") = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right)$$

$$= \underline{\underline{0.918}}$$

$$\text{Entropy}(S, "x_1 = F") = \frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right)$$

$$= \underline{\underline{0.918}}$$

$$\text{Information gain}(S, x_1) = 1 - \frac{3}{6} \times 0.918 - \frac{3}{6} \times 0.918$$

$$= \underline{\underline{0.082}}$$

$$A = x_2$$

Values of $A = \{T, F\}$.

$$= -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right)$$

$$= 1$$

$$= -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \times \log_2 \left(\frac{1}{2}\right).$$

$$= 1$$

$$\text{Information}(s, \cancel{s}) = 1 - \cancel{\frac{1}{2} \times 1} - \cancel{\frac{1}{2} \times 1}$$

$$= \cancel{0}$$

$$= 1 - \cancel{\frac{4}{6} \times 1} - \cancel{\frac{2}{6} \times 1}$$

$$= 1 - 1 = 0$$

~~2~~

$$I_{\text{pol}} \frac{1}{\varepsilon} - \left(\frac{1}{\varepsilon}\right) I_{\text{pol}} \frac{1}{\varepsilon} = \left\{ T = 1 \text{ for } \varepsilon \right\} \text{ parts}$$

$$\frac{8}{10} P \cdot 0.5$$

$$I_{\text{pol}} \frac{1}{\varepsilon} - \left(\frac{1}{\varepsilon}\right) I_{\text{pol}} \frac{1}{\varepsilon} = \left\{ T = 1 \text{ for } \varepsilon \right\} \text{ parts}$$

$$0.12 \cdot 0.5 = 0.06$$

$$8 \cdot P \cdot 0.5 \cdot \varepsilon - 1 = (x \cdot 0.5) \text{ a good control}$$

13/3/23 Crini Index
Monday

It's a probability of random a particular class being
incorrectly classified.

Crini index is used in CART algorithm.

Consider a dataset S having class labels c_1, c_2, \dots, c_r
Let p_i is the probability of examples having class label c_i
Then Crini index is calculated by

$$\text{Gini}(S) = 1 - \sum_{i=1}^r (p_i)^2$$

$$P_{amp} = \frac{3}{10}$$

$$P_{aep} = \frac{1}{10}$$

$$P_{bird} = \frac{2}{10}$$

$$P_{fish} = \frac{2}{10}$$

$$P_{mama} = \frac{2}{10}$$

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

$$= 0.78$$

Crini split index

Let S be a set of examples, A - attribute / feature.

S_A - Subtable with $A = v$

Values(A) - set of all possible values of A .

Crini split index of an attribute A is defined by,

$$\text{Crini split}(S, A) =$$

$$C_{\text{minisplit}}(S, A) = \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \times C_{\text{ini}}(S_v)$$

C_{ini} split (S , $A = \text{"gives birth"}$).

$A = \text{gives birth}$

Values of $A = \{\text{"yes"}, \text{"no"}\}$

$$\begin{aligned} C_{\text{ini}} \text{ split } (S, \text{Agives birth}) &= \frac{|S_{\text{yes}}|}{|S|} \times C_{\text{ini}}(S_{\text{gives birth}}) \\ &= \frac{|S_{\text{yes}}|}{|S|} \times C_{\text{ini}}(S_{\text{gives birth}=\text{yes}}) \end{aligned}$$

$$C_{\text{ini}}(S_{\text{gives birth}=\text{yes}}) = 1 - \sum p_i^2 \quad \text{Total no.}$$

$$P(\text{mam}) = \frac{2}{4}$$

$$P(\text{bird}) = \frac{1}{4}$$

$$P(\text{fish}) = \frac{1}{4}$$

$$\begin{aligned} &= 1 - \left[\left(\frac{2}{4}\right)^2 + \left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right] \\ &\approx 0.625 \end{aligned}$$

$$\begin{aligned} C_{\text{ini}}(S_{\text{gives birth}=\text{no}}) &= 1 - \left[\left(\frac{1}{6}\right)^2 + \left(\frac{1}{6}\right)^2 + \left(\frac{3}{6}\right)^2 + \left(\frac{1}{6}\right)^2 \right] \\ &= 0.667 \end{aligned}$$

$$C_{\text{ini}} \text{ split } (S, \text{gives birth}) = \frac{4}{10} \times 0.625 + \frac{6}{10} \times 0.667$$

$$\approx 0.6502$$

$\text{Gini } A = \text{aauatic animal}$

values of $A = \{\text{"yes"}, \text{"no"}, \text{"semi"}\}$

Gini :

$$\text{Crini}(S, \text{aauatic animal} = \text{yes}) = 1 - \frac{0}{10} = 0$$

$$\begin{aligned} \text{Crini}(S, \text{aauatic} = \text{no}) &= 1 - \left[\left(\frac{2}{5}\right)^2 + \left(\frac{1}{5}\right)^2 + \left(\frac{2}{5}\right)^2 \right] \\ &= 1 - \left[\frac{2}{5} + \frac{1}{5} + \frac{2}{5} \right] = 1 - \frac{5}{5} = 0.64 \end{aligned}$$

$$\text{Crini}(S, \text{aauatic} = \text{semi}) = 0$$

$$\begin{aligned} \text{Crini}^{\text{split}}(S, \text{aauatic animal}) &= \frac{|S_{\text{yes}}|}{|S|} \times \text{Crini}(S, \text{aa} = \text{yes}) \\ &\quad + \frac{|S_{\text{no}}|}{|S|} \times \text{Crini}(S, \text{aa} = \text{no}) \\ &\quad + \frac{|S_{\text{semi}}|}{|S|} \times \text{Crini}(S, \text{aa} = \text{semi}) \\ &= \frac{0.64 \times 5}{10} = 0.32 \end{aligned}$$

$A = \text{aerial animal}$

values of $A = \{\text{"yes"}, \text{"no"}\}$

$$\text{Crini}(S, \text{aerial} = \text{yes}) = 1 - \frac{0}{10} = 1$$

$$\begin{aligned} \text{Crini}(S, \text{aerial} = \text{no}) &= 1 - \left[\left(\frac{2}{8}\right)^2 + \left(\frac{1}{8}\right)^2 + \left(\frac{2}{8}\right)^2 + \left(\frac{3}{8}\right)^2 \right] \\ &= 1 - \left[\frac{2}{8} + \frac{1}{8} + \frac{2}{8} + \frac{3}{8} \right] = 1 - \frac{8}{8} = 0.71875 \end{aligned}$$

$$\begin{aligned} \text{Crini}^{\text{split}}(S, \text{aerial animal}) &= \frac{|S_{\text{yes}}|}{|S|} \times \text{Crini}(S, \text{aa} = \text{yes}) \\ &\quad + \frac{|S_{\text{no}}|}{|S|} \times \text{Crini}(S, \text{aa} = \text{no}) \\ &= \frac{1}{10} \times 1 + \frac{9}{10} \times 0.71875 = 0.71875 \end{aligned}$$

$$0 + \frac{8}{10} \times 0.71875 \\ = 0.575$$

$A = \text{has legs}$

Values of $A = \{\text{"yes", "no"}\}$.

$$\text{criterion}(S, \text{has legs} = \text{yes}) = 1 - \left[\left(\frac{2}{7} \right)^2 + \left(\frac{3}{7} \right)^2 + \left(\frac{2}{7} \right)^2 \right] \\ = 0.653$$

$$\text{criterion}(S, \text{has legs} = \text{no}) = 1 - \left[\left(\frac{1}{3} \right)^2 + \left(\frac{2}{3} \right)^2 \right] \\ = 0.44$$

$$\text{criterion}_{\text{split}}(S, \text{has legs}) = \frac{|S_{\text{yes}}|}{|S|} \times \text{criterion}(S, \text{has legs} = \text{yes}) +$$

$$\frac{|S_{\text{no}}|}{|S|} \times \text{criterion}(S, \text{has legs} = \text{no})$$

$$= \frac{7}{10} \times 0.653 + \frac{3}{10} \times 0.44 \\ = 0.589$$

split attribute

Root node (split node) is selected as the attribute with lowest criterion split index.

Here, aquatic animal is the lowest criterion split index.

Criterion split index is used in the CART algorithm.

$$P_1 = \frac{3}{6}$$

$$P_Q = \frac{3}{6}$$

$$\text{Crini}(S) = 1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2 \right].$$

$$= \underline{\underline{0.5}}.$$

$$A = x_1$$

values of A = {T, F}.

$$\text{Crini}(S, x_1 = T) = 1 - \left[\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right]$$

$$= \underline{\underline{0.44}}$$

$$\text{Crini}(S, x_1 = F) = 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right]$$

$$= \underline{\underline{0.44}}$$

$$\text{Crini}(S, x_1) = \frac{|S_T|}{|S|} \times \text{Crini}(S, x_1 = T) + \frac{|S_F|}{|S|} \times \text{Crini}(S, x_1 = F)$$

$$\text{Crini}(S, x_2 = F)$$

$$= \frac{3}{6} \times 0.44 + \frac{3}{6} \times 0.44$$

$$= \underline{\underline{0.44}}$$

$$A = x_2$$

values of A = {T, F}.

$$\text{Crini}(S, x_2 = T) = 1 - \left[\left(\frac{2}{4}\right)^2 + \left(\frac{2}{4}\right)^2 \right]$$

$$= \underline{\underline{0.5}}$$

$$\text{Crini}(S, x_2 = F) = 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]$$

$$= \underline{\underline{0.5}}$$

$$gini(S, x_2) = \frac{4}{6} \times 0.5 + \frac{2}{6} \times 0.5 \\ = \underline{\underline{0.5}}$$

x_1 is the root node

14/3/23
Tuesday

Measuring Performance of a classifier Algorithm

- o confusion matrix

		Actual class	
		Yes	No
predicted class	Yes	TP	FP
	No	FN	TN

TP → True positive

FP → False positive

FN → False negative

TN → True negative

2 class classification problem.

Confusion matrix is a model that is used to describe the performance of a classification model. It is a table that categorises predictions according to whether they match the actual values.

Precision

Recall

accuracy.

$$\text{Precision} = \frac{TP}{TP+FP}$$

$$\text{Recall} = \frac{TP}{TP+FN}$$

$$\text{Accuracy} = \frac{TP+TN}{TP+FP+TN+FN}$$

1) Suppose a computer program for recognizing dogs in a photograph identifies 8 dogs in a picture containing 12 dogs and some cats. Of the 8 dogs identified, 5 actually are dogs, while the rest are cats. Compute the precision and recall for the computer program.

		Actual	
		dog	cat
predicted	dog	5	3
	cat	7	

Sum = 12

$$\begin{aligned} TP &= 5 \\ FP &= 3 \\ FN &= 7 \end{aligned}$$

$$\text{Precision} = \frac{TP}{TP+FP} = \frac{5}{5+3} = \frac{5}{8} = 0.625$$

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{5}{5+7} = \frac{5}{12} = 0.416$$

- 2) Let there be 10 balls (6 white & 4 red balls) in a box and it is required to pick up red balls from them. Suppose we pick up balls as red of which 2 are actually red balls. What are the values of precision & recall in picking red ball.

		Actual	
		red	white
predicted	red	2	5
	white	2	1

Sum = 10

$$\begin{aligned} TP &= 2 \\ FP &= 5 \\ FN &= 8 \end{aligned}$$

$$\text{Precision} = \frac{TP}{TP+FP} = \frac{2}{2+5} = \frac{2}{7} = \frac{2}{7} = 0.2857$$

$$\text{Recall} = \frac{TP}{TP+FN} = \frac{2}{2+8} = \frac{2}{10} = \frac{2}{10} = \frac{2}{10} = 0.2$$

3) A database contains 80 records on a particular topic of which 55 are relevant to certain investigations. A search was conducted on that topic and 50 records are retrieved, 40 were relevant. construct the confusion matrix for the search and calculate the precision & recall and finally the accuracy whole classification problem

		Actual		
		relevant	not relevant	
predicted	relevant	40	10	50
	not relevant	15	15	

$$\text{TP} = 40$$

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{FP} = 10$$

$$\text{TN} = 15$$

$$\text{FN} = 15$$

$$\text{Precision} = \frac{40}{40 + 10} = \frac{40}{50}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{40}{40 + 15} = \frac{40}{55}$$

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{TN} + \text{FN}} = \frac{40 + 10}{40 + 15 + 10 + 15} = \frac{50}{80}$$

The ID3 Algorithm

- Developed by Ross Quinlan
- Iterative Dichotomiser 3

Assumption:

- The algorithm uses information gain to select the

most useful attribute for classification.

Assume that there are only two class labels namely, "+" and "-". The examples with class labels "+" are called positive examples and others negative examples.

Notations

The following notations are used in the algorithms-

S The set of examples.

C The set of class labels.

F The set of features.

A An arbitrary feature (attribute).

$\text{values}(A)$ The set of values of the feature A .

v An arbitrary value of A .

S_v The set of examples with $A = v$.

Algorithm ID3 (S, F, C)

1. Create a root node for the tree.
2. if (all examples in S are positive) then
3. returns single node tree Root with label "+"
4. end if
5. if (all examples are negative) then
6. returns single node tree Root with label "-".
7. end if.
8. if (number of feature is 0) then
9. returns single nodes & root tree.
10. else
11. let A be the feature in F with the highest information gain.
12. Assign A to the root node in decision tree.

13. for all (values v of A) do.
 14. Add a new tree branch below root node.
 corresponding to v .
 15. if (S_v is empty) then
 16. Below this branch add a leaf node with label
 equal to the most common class label in the set
 17. else
 18. Below this branch add the subtree formed by applying
 the same algorithm ID3 with the values
 $ID3(S_v, C, T - \{A\})$.
 19. end if
 20. end for
 21. end if.
 4) Consider a two class classification problem of predicting
 whether a photograph contains a man or a woman. Suppose
 we have a test dataset of 10 records with expected
 outcomes and a set of predictions from our classification
 algorithm.

	Expected	Predicted
1	man	woman.
2	man	man
3	woman	woman
4	man	man
5	woman	man
6	woman	woman
7	woman.	woman
8	man.	man
9	man	woman
10	woman.	women.

- a) Compute the confusion matrix for the data.
 b) Compute the accuracy, precision recall.

		Actual		
		Man	woman	
predicted	Man	3	1	4
	woman	2	4	6

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN} = \frac{3}{3+1} = \frac{3}{4}$$

$$\text{Accuracy} = \frac{3+4}{3+2+1+4} = \frac{7}{10}$$

- 5) Suppose 10000 patients tested for flu out of them, 9000 are actually healthy and 1000 are actually sick. For the sick people a test was positive for 620 and negative for 380. For the healthy people, the same test was positive for 180 and negative for 8820. Construct a confusion matrix for the data and compute the accuracy, precision and recall for the data.

		Actual		
		Sick	Healthy	
predicted	Sick	620	180	
	Healthy	380	8820	

$$\text{Predicted} = \frac{620}{620+180} = \frac{620}{800} = \frac{31}{40}$$

1620
 5820
 790

$$\text{Recall} = \frac{620}{620+380} = \frac{620}{1000} = \frac{31}{50}$$

$$\text{Accuracy} = \frac{620+8820}{620+180+380+8820} = \frac{9440}{10000} = \underline{\underline{94\%}}$$

Q3/23 Suppose that the data for analysis includes the attribute age. The age values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 43, 46, 52, 70.

- what is the mean of the data? what is its median?
- use smoothing by bin means to smooth the data, using a bin depth of 3. illustrate your steps comment on the effect of this technique for the given data.
- how might you determine outliers in the data?
- what other methods are there for data smoothing?
- use min-max algorithm normalization to transform value 35 for age on to the range [0:1]
- comment on which method of normalization you would prefer to use for the given data, giving reasons as to why?
- use Z-score normalization to transform the value 35 for age, where the standard deviation of age is 12.94 years.
- use normalization by decimal scaling to transform the value 35

a) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25,
30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70

$$\text{mean} = \underline{\underline{29.96}}$$

b) median = 25.

b) smoothing by bin means of b =

Bin 1: 13, 15, 16 {

Bin 2: 16, 19, 20 }

Bin 3: 20, 21, 22

" 4: 22, 25, 25.

" 5: 25, 25, 30.

" 6: 33, 33, 35

" 7: 35, 35, 35

" 8: 36, 40, 45

" 9: 46, 52, 70.

smoothing by bin means

Bin 1: 14.6, 14.6, 14.6. Replaced by mean value

Bin 2: 18.3, 18.3, 18.3.

Bin 3: 21, 21, 21

" 4: 24, 24, 24.

" 5: 26.6, 26.6, 26.6

" 6: 33.6, 33.6, 33.6

" 7: 35, 35, 35

" 8: 40.3, 40.3, 40.3

" 9: 56, 56, 56

c) eliminating outliers (values fall outside the dots)

d) smoothing by bin medians

" " boundaries

$$(e) \text{new_max} = 1.0 \\ \text{new_min} = 0.0$$

$$V = 35^{\circ}$$

$$\text{min}_A = 13$$

$$\text{max}_A = 70$$

$$V' = \frac{V - \text{min}_A}{\text{max}_A - \text{min}_A} (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$$

$$= \frac{35 - 13}{70 - 13} (1.0 - 0) + 0 \\ = 0.385$$

$$(f) V' = \frac{V - \bar{A}}{\sigma_A}$$

$$\bar{A} = 29.96$$

$$V' = \frac{35 - 29.96}{12.94}$$

$$= 0.389$$

$$gg) V' = \frac{V}{10^3} = \frac{35}{100} = \underline{\underline{0.35}}$$

h) Z-score normalization because it is in units of standard deviation.

complete data 2 10^3
divided by 10^3
mean dev & me
several dev
value
complete
dataset

2) Suppose a hospital tested the ages and body fat data for 18 randomly selected adults with the following result.

age	23	23	27	27	38	41	47	49	50	50	54	54	57	57	57	58
%fat	9.5	26.5	7.8	17.8	31.4	25.9	27.4	27.2	31.2	31.6	42.5	28.8	33.1	33.1	33.1	33.1

58	58	60	61	56	57	58	58	60	61						
				33.4	30.2	34.1	32.9	41.2	35.9						
				32.9	41.2	35.7									

- a) calculate the mean, median and SD of age and fat
- b) normalize the two variables based on z-score normalization
- c) calculate the correlation coefficient (Pearson's product moment coefficient) Are there two variables positively or negatively correlated

a) Mean of age = $\frac{46.4}{18}$

Median of age = $\frac{50 + 52}{2} = 51$

~~Mean~~

7.8, 9.5, 17.8, 25.9, 26.5, 27.4, 27.2, 31.2, 31.4, 31.6, 42.5, 28.8, 33.1, 33.1, 33.1, 33.1, 33.1

Mean of fat = $\frac{26.9}{18}$

~~decimals
tree~~