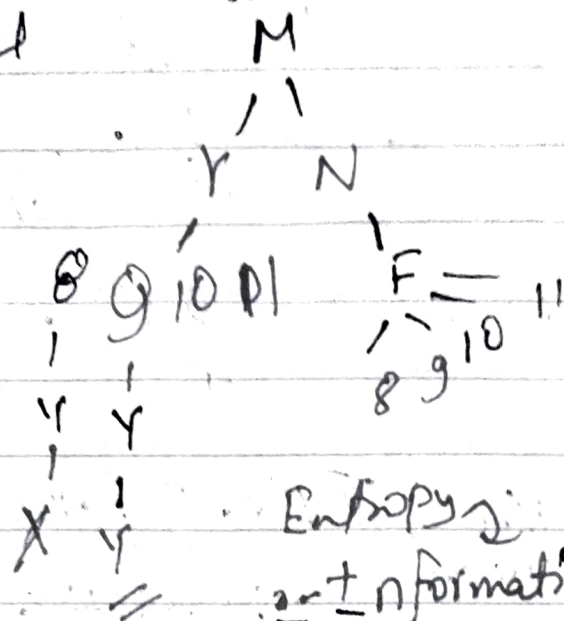


Decision Tree

classifier. as well as Regression task.
Powerful than both Linear Regr & Logic. Regr.

Class	Gender	Stay in Hostel
9	M	Y
10	F	N
8	F	Y
8	F	N
8	F	Y
9	M	N
10	M	N
11	M	Y
11	M	Y
8	F	Y
9	M	Y
11	M	N
10	F	N
10	M	Y
8	F	?



Entropy
+ Information Gain
Genny.

PS How to build a tree?

$$Ginni^0 = 1 - \sum_{i=1}^C (P_i)^2$$

Class	SIH	Total Value	$P(V)$	$P(N)$
8	$V=2, N=1$	3	$2/3$	$1/3$
9	$V=2, N=1$	3	$2/3$	$1/3$
10	$V=1, N=3$	4	$1/4$	$3/4$
11	$V=3, N=1$	4	$3/4$	$1/4$
		$T = 14$		

$$Ginni^0 \text{ for 8th class} = 1 - [P(V)^2] - [P(N)^2]$$

$$= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 4/9$$

$$G(9) = 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 = 4/9$$

$$G(10) = 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = \frac{6}{16} = \frac{3}{8}$$

$$G(11) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = \frac{6}{16} = \frac{3}{8}$$

$$G(\text{entire column}) = \frac{\sum \text{No. of instance of class}}{\text{Total No. of instance}} \times G(\text{respected class})$$

$$G(C) = \frac{n_8}{T} \cdot G(8) + \frac{n_9}{T} \cdot G(9) + \frac{n_{10}}{T} \cdot G(10) + \frac{n_{11}}{T} \cdot G(11)$$

$$= \frac{3}{14} \cdot \frac{4}{9} + \frac{3}{14} \cdot \frac{4}{9} + \frac{4}{14} \cdot \frac{6}{16} + \frac{4}{14} \cdot \frac{6}{16}$$

$$G(C) = 0.404$$

Ginni \equiv Ginni Impurity.

Gender	SIH	TV	P(Y)	P(N)
M	$Y=5, N=3$	8	$5/8$	$3/8$
F	$Y=3, N=3$	6	$1/2$	$1/2$
		14		

$$G(M) = 1 - \left(\frac{5}{8}\right)^2 - \left(\frac{3}{8}\right)^2 = 15/32 = 0.468$$

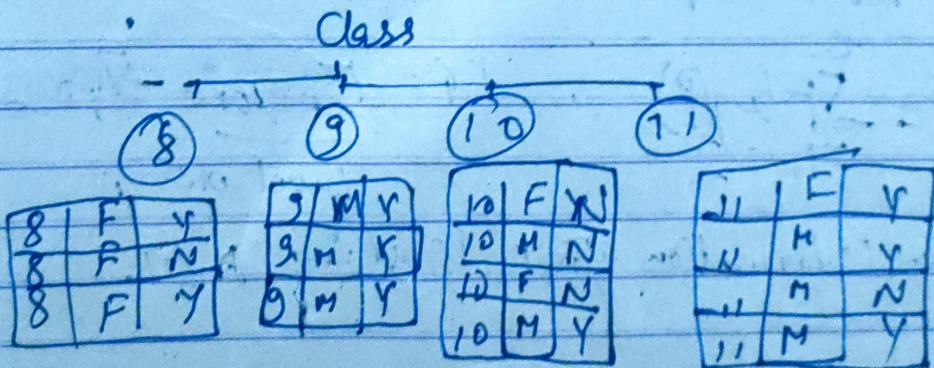
$$G(F) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

~~G(c) =~~

$$G(\text{Gender}) = \frac{8}{14} \cdot 0.468 + \frac{6}{14} \cdot \frac{1}{2} = 0.482$$

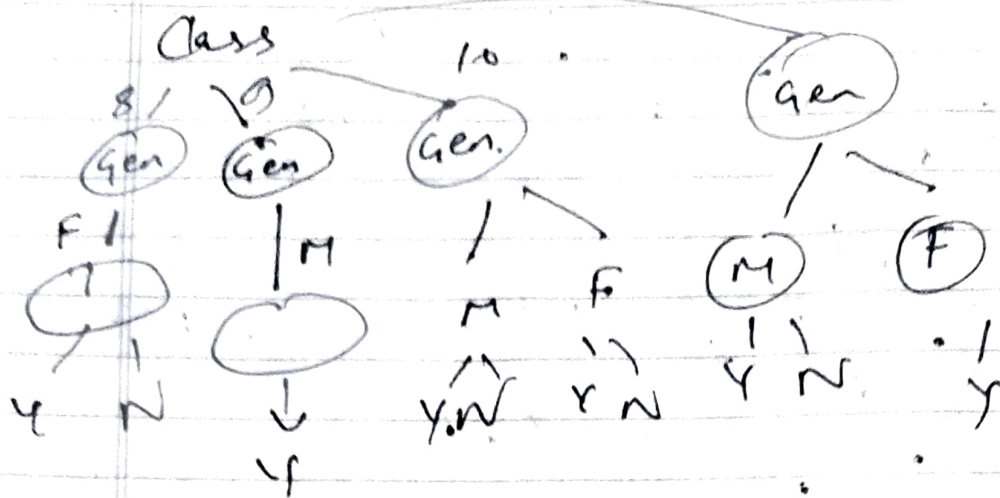
= Select Ginni Impurity which is low.
(select class column).

Root node \equiv Less Ginni value \equiv class column.



↓
Calculate Ginni
for all

9



So.

$$8 \quad F \quad ? \quad \equiv \quad 8 \quad F$$

$$11 \quad F \quad ? \quad \equiv \quad 11 \quad F \quad \cdot \quad Y$$

$$11 \quad M \quad ? \quad \equiv \quad 11 \quad F \quad \cdot \quad Y$$

(According to Weightage)

Entropy & Information Gain

$$E = - \sum_{i=1}^L P \log_2(P)$$

$$IG = \frac{E_{\text{before}} - E_{\text{After}}}{\text{Label Column } n.}$$

Gaining information after making a DT.

14 record $n(Y) = 8$ $n(N) = 6$

$$E(L) = -P(Y) \log_2(Y) - P(N) \log_2(N)$$

$$= -\frac{8}{14} \log_2 \frac{8}{14} - \frac{6}{14} \log_2 \frac{6}{14} = 0.98522$$

$$E(8) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$$

$$E(9) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$$

$$E(10) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.811$$

$$E(11) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.811$$

$$I(\text{class}) = \frac{3}{14} \times 0.918 + \frac{3}{14} \times 0.918 + \frac{4}{14} \times 0.811 + \frac{4}{14} \times 0.811 = 0.8574$$

$$IG = 0.98522 - 0.8574 = 0.1278$$

(class
column)

$$IG_{\text{Gender}} = 0.98522 - 0.974 = 0.01$$

$$E(M) = -\frac{5}{8} \log_2 \frac{5}{8} - \frac{3}{8} \log_2 \frac{3}{8} = 0.954$$

$$E(N) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} = 1$$

$IG(\text{gender})$

$$IG_{\text{Before Gender}} = \frac{8}{14} \times 0.954 + \frac{6}{14} \times 1 = 0.974 \text{ Final.}$$

ID3 C4.5 CART

$IG_{\text{Gender}} = 0.01$

$IG_{\text{class}} = 0.1278$

* More $IG \equiv$ that is root node \equiv Class Column

\Rightarrow Which Class has high IG less randomness

3 possibility.

Feature \rightarrow categorical
Feature \rightarrow conti.
Feature \rightarrow conti.

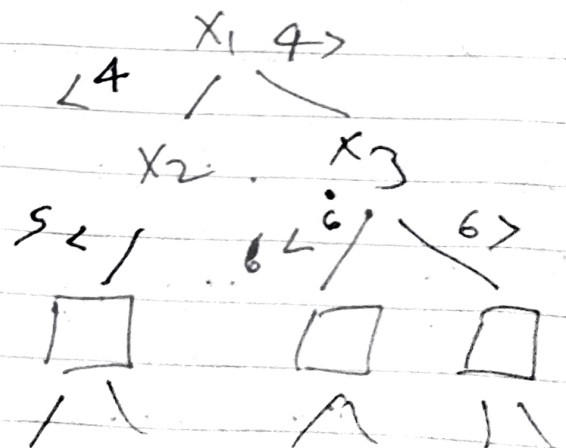
Outcome \rightarrow cat \rightarrow class
Outcome \rightarrow cat \rightarrow class
Outcome \rightarrow conti. \rightarrow Regression

ID3 \leftrightarrow Iterative Dichotomizer (Classifier)
C4.5 \leftrightarrow

CART \leftrightarrow Classi. & Regression Tree.

avg. $\leq 4 < \leq 5 < \leq 6 <$

x_1	x_2	x_3	y
1.1	7.5	2.5	A
2.2	8.8	5.5	B
3	9.2	6	A
3.6	5.1	6.7	A
5	5.4	7	B
5.8	2	8.9	B
8	1	9.1	A



\hat{y} is avg. of Record(y)

x_1	x_2	x_3	y
1	2	1	3
2.4	5.8	6	10
1.9	6.1	7.2	12.8
5.8	2.3	9.6	13.9
2.2	7.8	8.9	14.8

x_1	x_2	x_3
5	1	6
1	5	10

$$\sum (y - \hat{y})^2 \neq \sum (y - \bar{y})^2$$

CART According to errors both side rearrangement of Tree occurs

In all these above algos. we have to determine threshold first (different mechanism for that).

To control under/over fitting we using pruning.
(backward)

Prepruning
Prepruning
Post pruning
Post pruning.

Post pruning :- Build a complete DT then identify which branch is contributing to over/under fit issues. using Cross validation.

Cross validation :- Tune a parameter & removal of branches.

Backward pruning (Pre) :- Stop DT to create insignificant branches.

*

Ensemble - Techniques

* Combination of multiple algos. { Multiple.
(Multidecision
Maker)

Q/ Bagging Boosting Stacking

Q/ Bootstrap Aggregation (Bagging).

Q/ Pasting?

Q/ Random Forest?

Bootstrapping

Bagging

- SVM

- KNN

- NB

- DT

- RF