

06/11/22.

Mathematical Intuition of DT :-

• Decision tree Classifier :- DT (classification) \Rightarrow O/p \rightarrow -1 (cat, num)

• Decision tree regression \Rightarrow O/p \rightarrow (cat, num)

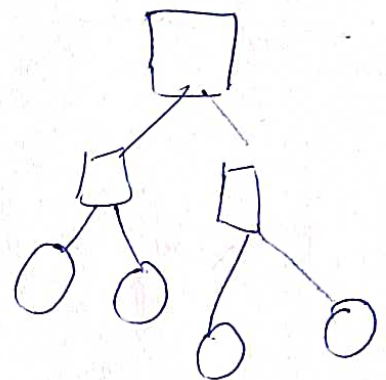
Classification And regression trees

• ID3

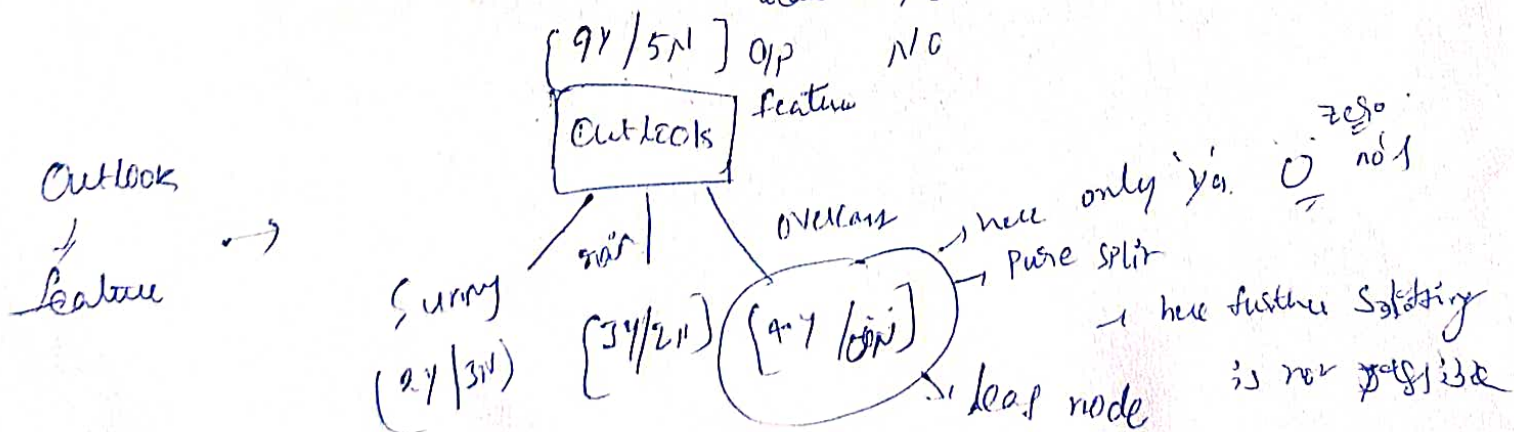
Iterative Decomposes

\rightarrow entropy

• gini impurity



Day	I_1 Outlook	I_2 Temperature	I_3 Humidity	I_4 Wind	\rightarrow O/p Decision
	Sunny	hot	high	weak	No
	Sunny	hot	high	strong	No
	rain	hot	hi	weak	Yes
	overcast	hot		weak	Yes
	Sunny	mid		weak	Yes
				short	Yes
	rain	cool		weak	No
				weak	Yes
				weak	Yes
				weak	No



Pure split :- only one parameter; here further splitting is not possible

Leaf node :- further splitting is not possible here.

to check the purity of the feature, 2 ways are there

Purity \rightarrow 1. Entropy

2. Gini Coefficient, Gini impurity

Entropy :- $H(S) = - \sum_{i=1}^n P_i \times \log(P_i)$

2 class $\rightarrow Y/N$

Gini Coefficient $\rightarrow 1 - \sum_{i=1}^n P_i^2$

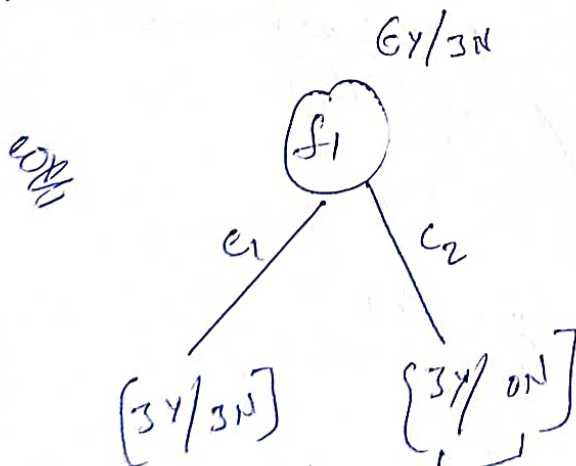
for binary $1 - [P_Y^2 + P_N^2]$

Binary ^{classification} $\rightarrow -P_Y \log(P_Y) - P_N \log(P_N)$

for 3 class - C_1, C_2, C_3

$-P_{C_1} \log(P_{C_1}) - P_{C_2} \log(P_{C_2}) - P_{C_3} \log(P_{C_3})$

Multiclass classification



f1	C/P
C1	Y
C2	Y
C1	Y
C2	Y
C1	Y
C1	N
C2	Y
C1	N
C1	N

$$\text{Entropy } H(S) = - \sum_{i=1}^N P_i \log(P_i)$$

$$= - P_Y \log(P_Y) - P_N \log(P_N)$$

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

← from $(3Y/3N)$ →

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

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

$$= - \frac{1}{2} (0-1) - \frac{1}{2} (0-1)$$

$$= \frac{1}{2} + \frac{1}{2} = 1$$

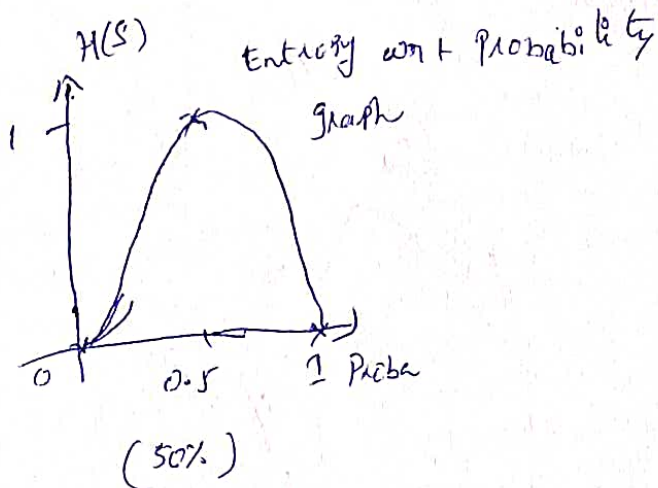
Entropy $H(S) = 1$. \Rightarrow This is for impurity data = 1 .

↑ whenever there is high impurity

Ex 3Y/0N \Rightarrow full pure split entropy = 0 . , entropy value is 1 :

$$- \frac{3}{3} \log_2\left(\frac{3}{3}\right) - \frac{0}{3} \log_2\left(\frac{0}{3}\right)$$

$$\text{Entropy} = 0$$



$H(S) = 1 \Rightarrow$ very impure split

$H(S) = 0 \Rightarrow$ pure split

Ex. 2 Yes / 3 No.

$$H(S) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$= 0.97$$

② Cini coefficient / Cini impurity

$$1 - \sum_{i=1}^n p_i^2$$

1. $3Y/3N \rightarrow$ Entropy $\rightarrow H(S) = 1$ (very impure split)

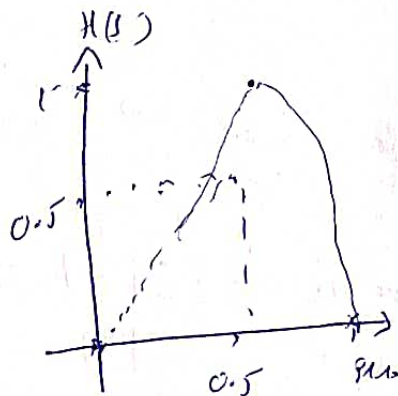
2. $3Y/0N =$

3. $2Y/3N =$

$$\begin{aligned} C_{ii} &= 1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2 \right] \\ &= 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right] \\ &= 1 - \left[\frac{1}{4} + \frac{1}{4} \right] \end{aligned}$$

$$= 1 - \left[\frac{2}{4} \right] \rightarrow 1 - \frac{1}{2} = \frac{1}{2}$$

$$= 0.5$$



$$\begin{aligned} C_{12} &= 8Y/2N \\ &= 1 - \left[\left(\frac{8}{10}\right)^2 + \left(\frac{2}{10}\right)^2 \right] \end{aligned}$$

$$= 1 - \left[\left(\frac{4}{5}\right)^2 + \left(\frac{1}{5}\right)^2 \right]$$

$$= 1 - \left[\frac{16}{25} + \frac{1}{25} \right]$$

$$= 1 - \frac{17}{25} = \frac{25-17}{25}$$

$$= 0.32$$

Ex. 3 $4Y/8N$

$$\begin{aligned} &= 1 - \left[\left(\frac{4}{12}\right)^2 + \left(\frac{8}{12}\right)^2 \right] \\ &= 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] \end{aligned}$$

$$= 1 - \left[\frac{1}{9} + \frac{4}{9} \right]$$

$$= 1 - \left[\frac{5}{9} \right] = \frac{9-5}{9} = \frac{4}{9}$$

$$= 0.44$$

$3Y/0N$

$$1 - \left[\left(\frac{3}{5}\right)^2 + \left(\frac{0}{5}\right)^2 \right]$$

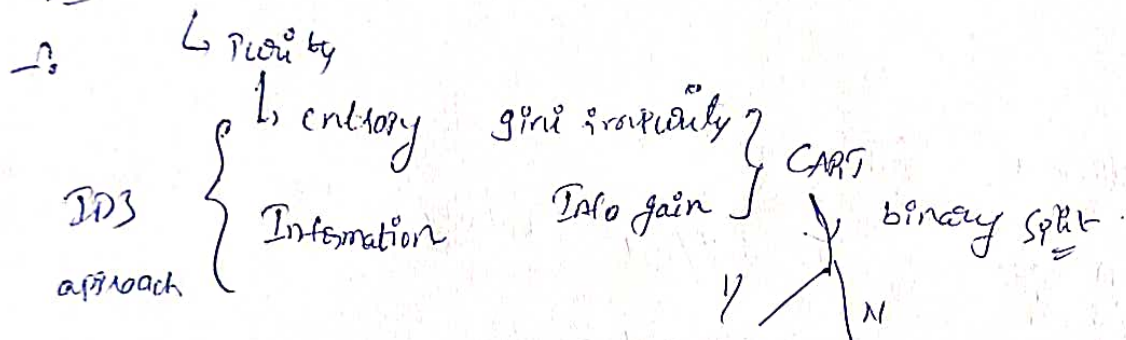
$$= 1 - 1$$

$$= 0$$

Information Gain :- ID3 approach
 ↓
 entropy

CART
 ↓
 Gini impurity
 ↓
 Binary Classification

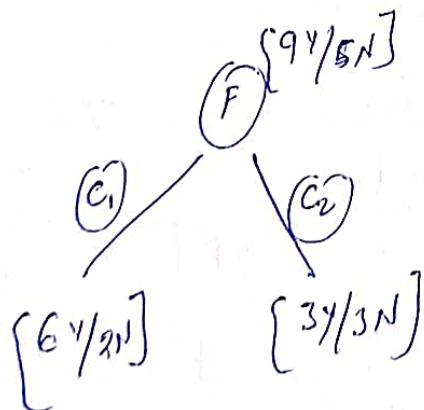
feature 1, feature 2, feature 3



Gini impurity is very fast, when compare with entropy mathematically.

IG :- Information Gain → we can cal using entropy & Gini impurity

$$\boxed{\text{Gain}(S, f_1) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)} \rightarrow \text{w.r.to entropy}$$



$v = \text{value}$

feature $f_1 = (94/51)$

$H(S) = \text{root entropy}$

$$= -P_v \log_2 P_v - P_n \log_2 P_n$$

$$= -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$

$$= -0.64 \log_2(0.64) - 0.35 \log_2(0.35)$$

$$H(S) = 0.94$$

root feature entropy

$$0.4 (0.15) - 0.35 \cdot 0.96$$

$$= 0.1216 - 0.3376$$

$$C_1 = 64/2N \Rightarrow$$

$$C_1 = -\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right)$$

$$= -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right)$$

$$= 0.75(0.12) + 0.25(0.60)$$

$$= 0.81$$

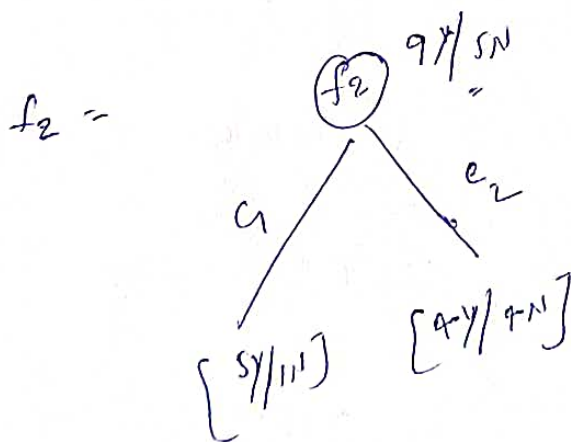
$$C_2 = 3/3N = 1$$

$$H(S) = 1$$

$$Gain(S, f_1) = 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$= 0.94 [0.462 + 0.42]$$

$$f_1 = 0.049$$



$$39 = ?$$

$$Gain = 0.94 - \left(\frac{6}{14}\right)(0.65) - \frac{8}{14} \times 1$$

$$f_2 H(S) = 0.94$$

$$Gain(f_2) = 0.04$$

$$C_1 = 54/1N$$

$$= -\frac{5}{6} \log_2\left(\frac{5}{6}\right) - \frac{1}{6} \log_2\left(\frac{1}{6}\right)$$

$$= 0.65$$

$$I(f_2) > I(f_1)$$

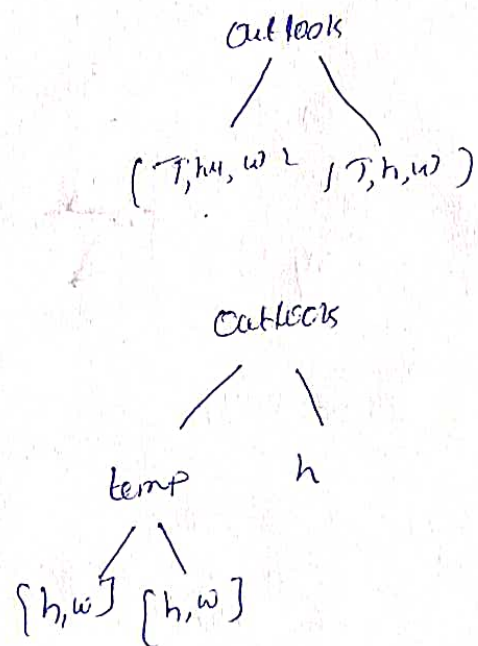
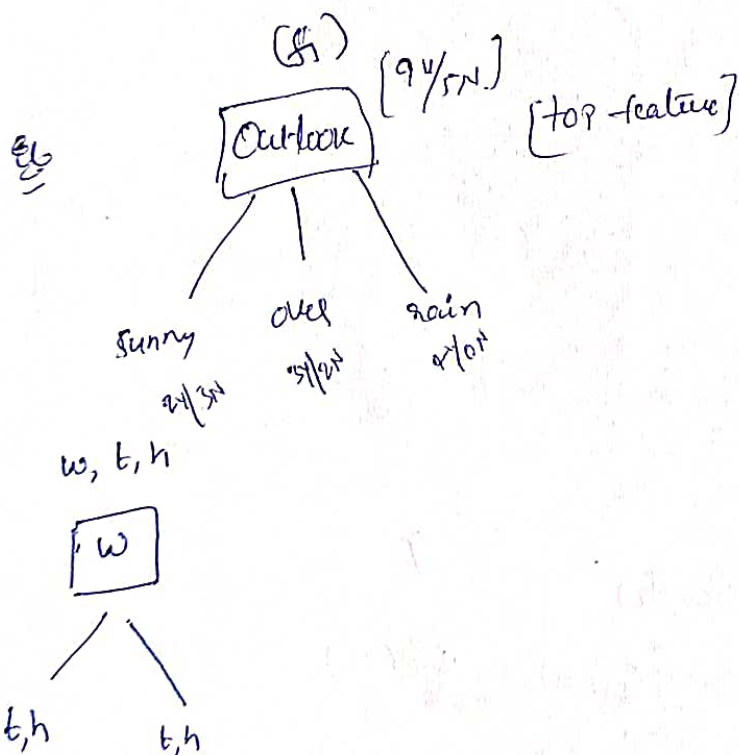
↓
This is greater and it is providing more info.

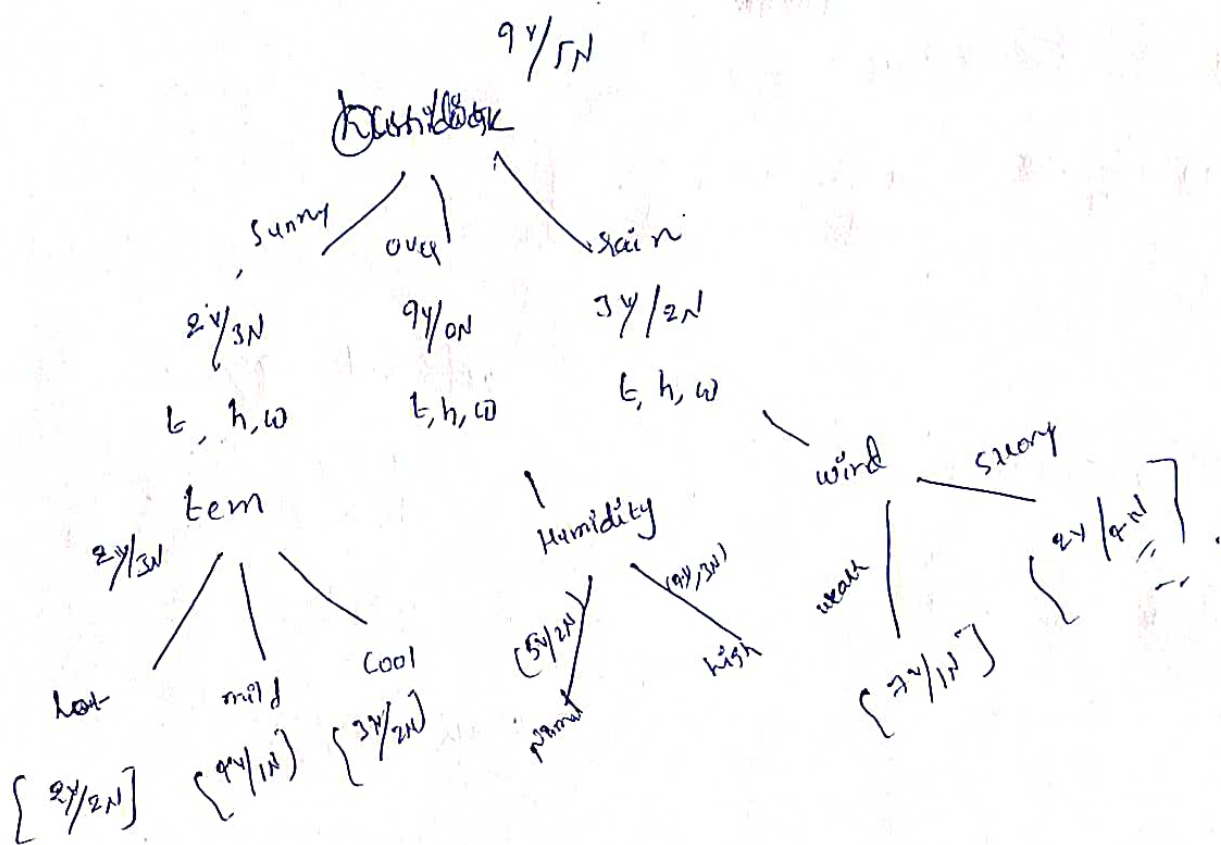
$$c_2 = 44/4N = -\frac{1}{4} \log_2\left(\frac{1}{4}\right) + \frac{3}{4} \log_2\left(\frac{3}{4}\right)$$

$$= -\left\{ \frac{1}{4} \log_2\left(\frac{1}{4}\right) + \frac{3}{4} \log_2\left(\frac{3}{4}\right) \right\} = -\frac{1}{2} \log 1$$

$$= 1$$

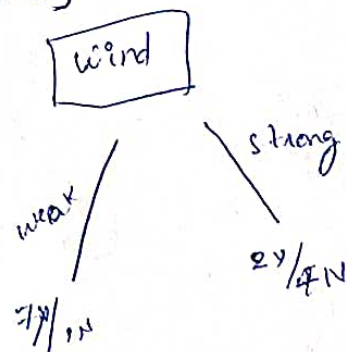
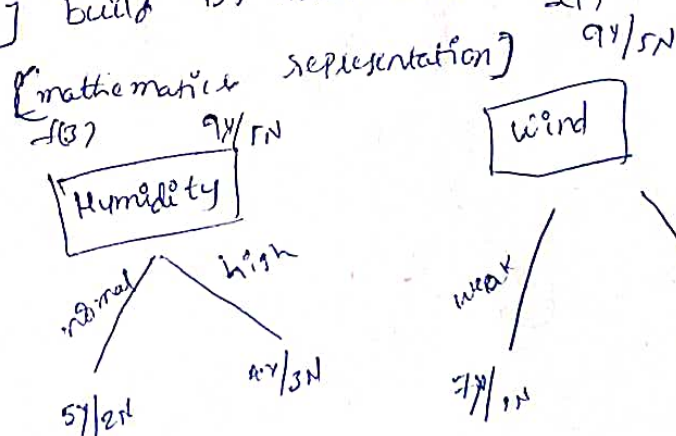
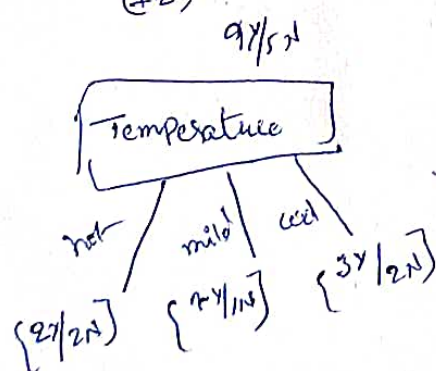
Day	Outlook	temperature	humidity	wind	Decision
Sp	Sunny	hot	high	weak	No
	Sunny	hot	high	strong	No
	Sunny	hot	high	weak	YES
	Overcast	hot	high	weak	YES
	Rainfall	mild	normal	weak	No
	Rainfall	cool	normal	strong	YES
	Rainfall	cool	normal	strong	No
	Rainfall	cool	normal	weak	YES
	Overcast	cool	high	weak	YES
	Sunny	mild	normal	weak	YES
	Rainfall	cool	normal	strong	YES
	Sunny	mild	normal	strong	YES
	Overcast	mild	high	weak	No
	Overcast	hot	normal	strong	
	Rainfall	mild	high		





• Entropy vs Gini-impurity [5 dif]

• Take this dataset [tennis play] build DT from scratch (f1)



f2 → Temperature

$$9y/5N = 0.1233 \quad 0.1576$$

$$H(S) = -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$H(S) = 0.94$$

Cool (3y/2N)

$$= -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right)$$

$$\text{Hot: } 2y/2N$$

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

$$H(S) = 1$$

$$\text{Mild: } 4y/1N$$

$$H(S) = -\frac{4}{5} \log_2 \left(\frac{4}{5} \right) - \frac{1}{5} \log_2 \left(\frac{1}{5} \right)$$

=

$I_9 \rightarrow f_d \rightarrow f_2$ (Temperature)

$$= 0.94 \left[\frac{4}{14} \cdot + \frac{5}{14} \cdot + \frac{5}{14} \cdot \right]$$

f_3 (Humidity) = $9\% / 5N$

$$-\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

=

Normal = $5\% / 2N$

$$= -\frac{5}{7} \log_2\left(\frac{1}{7}\right) - \frac{2}{7} \log_2\left(\frac{2}{7}\right)$$

=

High = $4\% / 3N$

$$= -\frac{4}{7} \log_2\left(\frac{4}{7}\right) - \frac{3}{7} \log_2\left(\frac{3}{7}\right)$$

=

^(f3)
Gain = $0.94 \left[\frac{7}{14} \cdot + \frac{7}{14} \cdot \right]$

f_4 (wind) = $9\% / 5N$

$$-\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

$$= 0.94$$

Weak = $7\% / 1N$

$$= -\frac{7}{8} \log_2\left(\frac{7}{8}\right) - \frac{1}{8} \log_2\left(\frac{1}{8}\right)$$

=

Strong = $2\% / 4N$

$$= -\frac{2}{6} \log_2\left(\frac{2}{6}\right) - \frac{4}{6} \log_2\left(\frac{4}{6}\right)$$

=

I_9 = $0.94 \left[\frac{7}{14} \cdot + \frac{1}{14} \cdot \right]$
(f4)