

DMA.

2

PART-A

Problems:

- ✓ Decision Tree
- ✓ KNN [k nearest Neighbor]
- ✓ Linear Regression
- ✓ Bayes

5 out of 4

10m

Part B

1 → 10m comp.

Theory:

- ✓ Multiple logistic regression
- ✓ Logistic regression
- ✓ Generalized linear model (GLM)
- ✓ Simple problems with algo
- ✓ Derivation sigmoidal fun.

DMN

① Decision Tree:

$$\text{Info}(D) = - \sum_{i=1}^m p_i \log_2 (p_i) \quad [\text{Exact info}]$$

$$\text{Gain}(A) = \text{Info}(A) - \text{Info}(D) \quad [\text{Gain}(A)] = \text{Gain}(A)$$

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{info}(D_j) \quad [\text{New info}]$$

Outlook	Temp	Humidity	Windy	PlayTennis
R	H	H	F	No
R	H	H	T	No
R	H	H	F	Yes
S	M	H	F	Yes
S	C	N	F	Yes
S	C	N	T	No
O	C	N	T	Yes
R	M	H	F	No
R	C	N	F	Yes
S	M	N	F	Yes
R	M	N	T	Yes
O	M	H	T	Yes
O	H	N	F	Yes
S	M	H	T	No

$$\text{total} = 14 \quad \text{Yes} = 9$$

$$\text{No} = 5$$

$$\text{calculate Info}(D) = - \sum_{i=1}^m p_i \log_2 (p_i)$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

14 for yes; no;

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

$$\boxed{\text{Info (D)} = 0.940}$$

Total

$$= 0.637$$

1st Attribute outlook

	Y	N
3	2	3 ✓
0	3	2
R	4	0

$$\text{Info outlook} = \frac{5}{14} \left(-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \frac{3}{5} \right) +$$

$$\frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) +$$

$$\frac{4}{14} \left(-\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{0} \right)$$

$$\text{Info outlook (D)} = 0.693$$

2nd Attribute Temperature

	Y	T	N
H	2	2	
M	4	2	
C	3	1	

$$\text{Info outlook (D)} = \frac{4}{14} \left(-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right) +$$

$$\frac{6}{14} \left(-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) +$$

$$\frac{4}{14} \left(-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right)$$

$$\text{Info outlook (D)} = 0.911$$

Attribute	humidity;	Y	N
	H	3	4
	N	6	1

$$\text{Info outlook (D)} = \frac{7}{14} \left(-\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \right) + \frac{7}{14} \left(-\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \right)$$

$$\text{Info outlook (D)} = 0.788$$

Attribute windy;

	Y	N
W	6	2
S	3	3

$$\text{Info outlook (D)} = \frac{8}{14} \left(-\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right) + \frac{6}{14} \left(-\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \right)$$

$$\text{Info outlook (D)} = 0.892$$

$$\text{Gain(A)} = \text{Info (D)} - \text{Info}_A(\text{D})$$

$$\text{Gain(outlook)} = 0.940 - 0.693$$

$$= 0.247$$

$$\text{Gain(Temp)} = 0.940 - 0.911$$

$$= 0.029$$

$$\text{Gain(Hum)} = 0.940 - 0.788$$

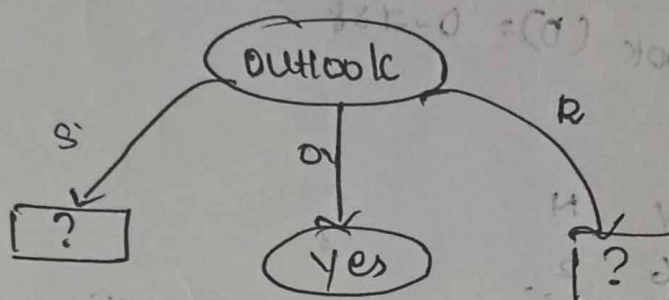
$$= 0.152$$

$$\text{Gain(windy)} = 0.940 - 0.892 = 0.048$$

Attribute	Gain
outlook	0.247 ✓
Temp	0.029
Hum	0.152
windy	0.048

Highest

Root Node: outlook



consider the outlook table.

outlook

Temp

Humid

Windy

Tennis

H

High

W

N

H

High

S

N

M

High

W

N

C

Not

W

Y

M

Not

S

Y

Sunny

outlook

Temp

Humid

Windy

Tennis

M

H

W

Y

C

N

W

N

C

N

S

N

M

N

W

Y

M

H

S

N

Rainy

sunny, yes = 2, no = 3

$$\text{Info (sunny)} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.971 \quad \text{total}$$

Humidity, yes = 1, no = 3

H

N

$$\Rightarrow \frac{3}{5} \left(-\frac{0}{0} - \frac{3}{3} \log_2 \frac{3}{3} \right) +$$

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

$$= 0$$

$$\boxed{\text{Info}_{\text{sun}}^{\text{Hum}} (\text{sunny}) = 0}$$

windy,

W

S

Y

N

1

2

1

$$\text{Info}_{\text{win}} (\text{sunny}) = \frac{3}{5} \left(-\frac{1}{3} \log_2 \frac{1}{3} \right) + \frac{2}{5} \left(-\frac{2}{2} \log_2 \frac{2}{2} \right)$$

$$= \frac{3}{5} \left(-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) +$$

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

$$\boxed{\text{Info}_{\text{win}} (\text{sun}) = 0.951}$$

$$\text{Info}_{\text{tem}} (\text{sun}) =$$

Temp,

	Y	N
H	0	2
C	1	0
M	1	1

$$\Rightarrow \frac{2}{5} \left(-\frac{2}{2} \log_2 \frac{2}{2} \right) + \frac{1}{5} \left(-\frac{1}{1} \log_2 \frac{1}{1} \right) + \frac{2}{5} \left(-\frac{1}{1} \log_2 \frac{1}{1} \right)$$

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

$$\text{Info}_{\text{temp}}(\text{sun}) \Rightarrow 0.4$$

$$\text{Gain}(\text{Temp}) = 0.971 - 0.4$$

$$= 0.571$$

$$\text{Gain}(\text{Hum}) = 0.971 - 0$$

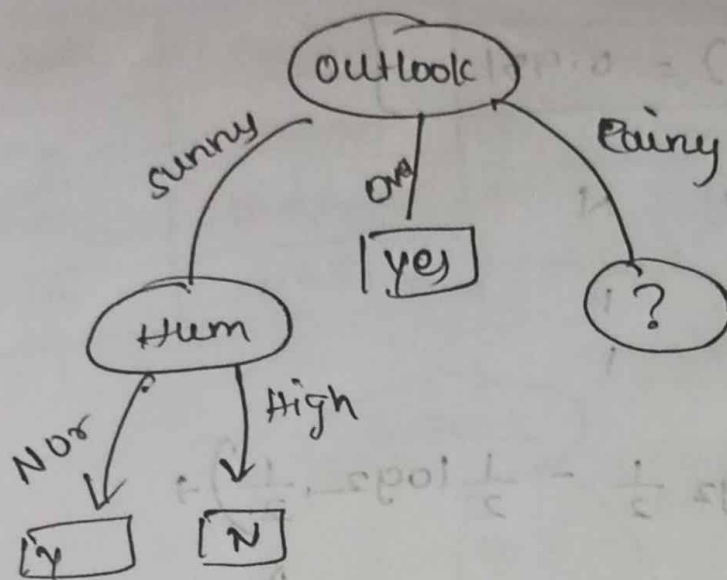
$$= 0.971$$

$$\text{Gain}(\text{Win}) = 0.971 - 0.951$$

$$= 0.02$$

Attribute	Gain
Temp	0.571
Hum	0.971 ✓
Win	0.02

Highest



Rainy;

write the outlook table;

Outlook	temp	Humi	windy	tennis
Rainy	M	H	W 0	Y 4
	C	N	W 2	Y 2
	C	N	S 3	N 1
	M	N	W	Y
	M	H	S	N

Rainy; Y=3, NO=2.

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

	Y	N
M	2	2
C	1	1

$$\Rightarrow 0.971$$

Total

$$\begin{aligned} \text{Info}_{\text{temp}}(\text{Rain}) &= \frac{3}{5} \left(-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) + \\ &\quad \frac{2}{5} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) \end{aligned}$$

$$\text{Info}_{\text{Temp}}(\text{Rain}) = 0.951$$

Hum,

	Y	N
H	1	1
N	2	1

$$\Rightarrow \frac{2}{5} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{3}{5} \left(-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right)$$

$$\text{Hum} \Rightarrow 0.951$$

windy,

	Y	N
W	3	0
S	0	2

$$\text{Info}_{\text{win}}(\text{Rainy}) = \frac{3}{5} \left(-\frac{3}{3} \log_2 \frac{3}{3} - 0 \right) + \frac{2}{5} \left(-0 - \frac{2}{2} \log_2 \frac{2}{2} \right)$$

$$\text{Info}_{\text{win}}(\text{Rain}) = 0$$

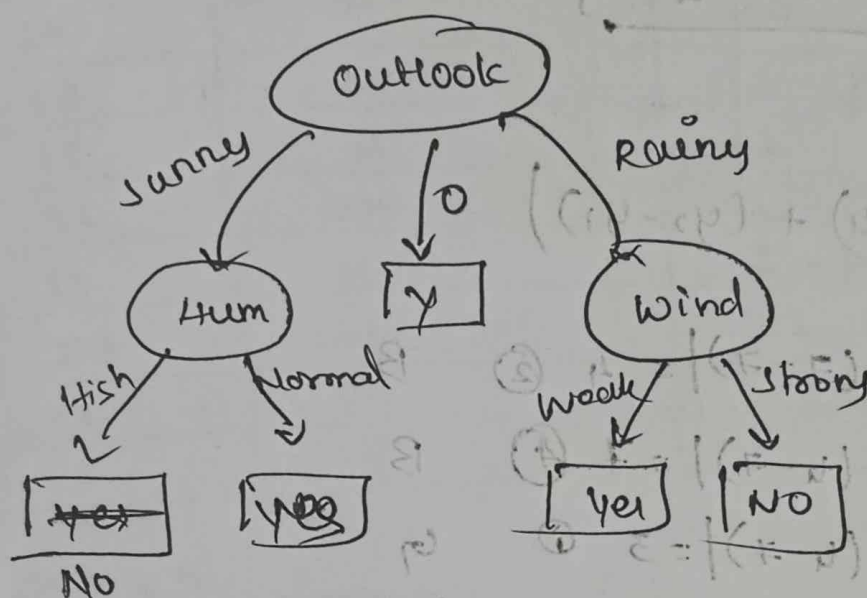
$$\text{Gain}(\text{Temp}) = 0.951 - 0.951 = 0.020$$

$$\text{Gain}(\text{Hum}) = 0.951 - 0.951 = 0.020$$

$$\text{Gain}(\text{wind}) = 0.951 - 0.951 = 0.020$$

Attribute	Gain
Temp	0.020
Hum	0.020
Wind	0.971 ✓

Highest



② K Nearest Neighbors:

↳ Powerful classification algo used in pattern recognition.

Distance formulas:

Eucclidean $\rightarrow \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

Manhattan $\rightarrow |(x_2 - x_1)| + |(y_2 - y_1)|$ (or) $|x_1 - y_1|$

Adv:

- ↳ Simple
- ↳ Applied from any distribution
- ↳ Good classification

Dis:

- ↳ Choosing k may be tough
- ↳ Large no. of samples
- ↳ More time to classify a new example.

x	y	Class
7	7	B
7	4	B
3	4	G
1	4	G

k=3

New instance (3,7)

Predict the class.

Distance

$$|(x_2 - x_1) + (y_2 - y_1)|$$

- B $|(3-7) + (7-7)| = 4$ ② B
 B $|(7-3) + (4-7)| = 4$ ④ B
 G $|(3-3) + (4-7)| = 3$ ① G
 G $|(1-3) + (4-7)| = 3$ ③ G

If, k=1; Good

k=2; Bad

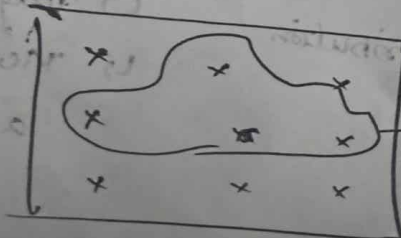
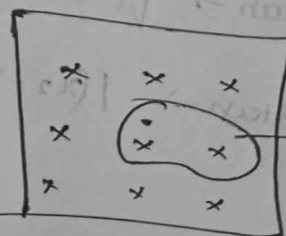
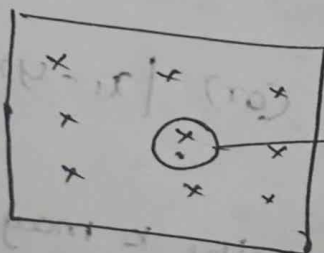
k=3; Good

~~k=4; Bad~~

∴ If k=3 given then

(3,7) is good.

Majority → Good.



$k=5$ gn.

Height	weight	class under/Normal	Distance
167	51	UW	7 - (5)
182	62	N	19
176	64	N	20
173	64	N	12 - (5)
172	65	N	12
169	56	UW	2 - (2)
173	58	N	6 - (4)
170	57	N	2 - (1)
174	56	N	5 - (3)
170	55	?	?

$(170, 55) \rightarrow ?$ Instance Normal for $k=1 \rightarrow N$
 $k=2 \rightarrow UW$

$$|(170 - 167) + (55 - 51)| = 3 + 4 = 7$$

$k=3 \rightarrow N$
 $k=4 \rightarrow N$
 $k=5 \rightarrow UW$

- if $k=1 \rightarrow N$
 $k=2 \rightarrow UW$ ✓ Majority Normal
 $k=3 \rightarrow N$
 $k=4 \rightarrow ~~N~~ UW$
 $k=5 \rightarrow UW$ ✓
- $(170, 55) \rightarrow \text{Normal}$

③ Baye's theorem:

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

B \rightarrow evidence
A \rightarrow hypothesis.

$$P(\text{playtennis} = \text{"yes"}) = \frac{9}{14} = 0.643$$

$$P(\text{playtennis} = \text{"No"}) = \frac{5}{14} = 0.357$$

conditional probability;

Outlook	Y	N
S	2/9	3/5
R	3/9	2/5
O	4/9	0

Temp	Y	N
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

Hum	Y	N
High	3/9	4/5
Not	6/9	1/5

Windy	Y	N
S	3/9	3/5
W	6/9	2/5

compute $P(x|c_i)$ for each class;

$$P(\text{outlook} = \text{"sunny"} | \text{playtennis} = \text{"yes"}) = \frac{2}{9} = 0.222$$

④ Linear Regression 1.

least square method.

$$y = b_0 + b_1x$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\text{Slope (m)} = \frac{S_{xy}}{S_{xx}}$$

Simple linear regression

$$\hat{y} = b_0 + b_1x + e$$

No. of ty	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$
1	14	-1	-6	1
3	24	1	4	1
2	18	0	-2	0
1	17	-1	-3	1
3	27	1	7	1
10/2=5	20			4

$\bar{x} = 2$ $\bar{y} = 20$

slope

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\Rightarrow (-1 \times -6) + (1 \times 4) + (0 \times -2) + (-1 \times -3) + (1 \times 7)$$

4

$$b_1 \Rightarrow \frac{6 + 4 + 0 + 3 + 7}{4} = \frac{20}{4} = 5$$

y-intercept

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$= 20 - 5(2)$$

$$= 20 - 10$$

$$= 10$$

$$\hat{y} = 10 + 5x$$

Least square method:

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
1	2	-2	-2	4	4
2	4	-1	0	0	1
3	5	0	1	0	0
4	4	1	0	0	1
5	5	2	1	2	4
Σ	4			6	10

$$m, \text{slope} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} = \frac{6}{10} = \frac{3}{5}$$

$y = mx + c \rightarrow \text{intercept}$

$$\bar{y} = 4 \quad \bar{x} = 3$$

$$4 = \frac{3}{5} \times 3 + c$$

$$4 = \frac{9}{5} + c$$

$$4 = \frac{9}{5} + c$$

$$\frac{9}{5} + c = 4$$

$$\frac{9}{5} - \frac{9}{5} = c$$

$$\frac{9 - 20}{5} = c$$

$$\boxed{c = -2.2}$$

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
2	50	-3.2	-25	80	10.24
3	60	-2.2	-15	33	4.84
5	80	-0.2	5	-1	0.04
7	90	2.2	15	33	4.84
9	95	3.8	20	76	14.44
				221	34.4

$$\boxed{5.2} \quad \boxed{75}$$

$$m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$m = \frac{221}{34.4} = 6.424$$

$$y = mx + c$$

$$75 = 6.424 \times 5.2 + c$$

$$\boxed{c = 41.6}$$

Predict for $x = 4$

$$y = 6.424 \times 4 + 41.6 = 67.3$$

$$\boxed{y = 67.3}$$

$$y = mx + c$$

$$y = 6.424 \times 4 + 41.6$$

R^2 coefficient of determination:

$R^2 = \frac{SSR}{SST} \rightarrow$ sum of square reg.
 $SST \rightarrow$ Total sum of reg.

$$\boxed{SST = SSR + SSE}$$

x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})(y - \bar{y})$	$(x - \bar{x})^2$
1	14	-1	-6	6	1
3	24	1	4	4	1
2	18	0	-2	0	0
1	17	-1	-3	3	1
3	27	1	7	7	1
<u>2</u>	<u>20</u>			<u>20</u>	<u>4</u>

$$m = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$= \frac{20}{4} = 5$$

$$y = mx + c$$

$$20 = 5 \times 2 + c$$

$$c = 10$$

$$\hat{y} = mx + c \Rightarrow \hat{y} = 5x + 10$$

	$y - \bar{y}$	$\hat{y} - \bar{y}$	$(\hat{y} - \bar{y})^2$
1	36	5	25
3	16	5	25
2	4	0	0
1	9	5	25
3	49	5	25
	<u>114</u>	<u>X</u>	<u>100</u>

$$R^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2}$$

$$R^2 = \frac{100}{114}$$

$$R^2$$

④ Derivation for sigmoid fun

$$B_0 = 0.07 \quad B_1 = 0.28 \quad B_2 = 0.36$$

	B ₀	B ₁	B ₂
0	0	0	0
1	1	1	1
2	0	0	0
3	1	1	1
4	0	0	0

Derivation of
B₀ = 150
B₁ = 50
B₂ = 50

$$B(x_1 + B_0 x_2 + \dots + B_n x_n)$$

$$+ 0.28 \times 2 + 0.28 \times 50 + 0.36 \times 150$$

$$FF - FF =$$

⑤ Multiple logistic reg:

Given:

$$\beta_0 = 0.67$$

$$\beta_2 = 0.47$$

$$\beta_1 = 0.58$$

$$\beta_3 = 0.36$$

Age x_1	BMI x_2	BP x_3	Dia y
45	28	140	1
50	25	135	0
40	30	150	1
35	24	125	0

predict age 25

BMI 20

BP = 120

Dial = ?

$$x_1 = 25$$

$$x_2 = 20$$

$$x_3 = 120$$

$$y = ?$$

$$Z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

$$Z = 0.67 + 0.58 \times 25 + 0.58 \times 20 + 0.36 \times 120$$

$$Z = 67.77$$

$$\hat{p} = \frac{1}{1 + e^{-z}} = \frac{1}{1 + 0} = 1$$

$$5(1) \rightarrow 10 \times 4 = 40$$

$$1(10) \rightarrow 1 \times 10 = \frac{10}{50}$$

Theory:

* Multiple logistic Regression

* logistic regression

* Generalized linear model

* Derivation sigmoidal fun.

same

- DMP:
- Problem:
- ✓ * Decision Tree
 - ✓ * KNN
 - ✓ * Linear Regression
 - * Bayes

① Decision Tree:

$$\text{Info}(D) = - \sum P_i \log_2(P_i)$$

$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

SNo	outlook	temperature	humidity	windy	Play tennis
1	S	H	High	weak	No
2	S	H	H	Strong	No
3	O	H	H	W	Yes
4	R	M	H	W	Yes
5	R	C	Normal	W	Yes
6	R	C	N	S	No
7	O	C	N	S	Yes
8	S	M	H	W	No
9	S	C	N	W	Yes
10	R	M	N	W	Yes
11	S	M	N	S	Yes
12	O	M	H	S	Yes
13	O	H	N	W	Yes
14	R	M	H	S	No

Total = 14.

Yes = 9, No = 5

0.41

$$Info(D) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$Info(D) = 0.940$$

Attributes

outlook

	Y	N
S	2	3
O	4	0
R	3	2

0.44

0.53

Info outlook (D) =

$$\begin{aligned} & \frac{5}{14} \left(-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right) + \\ & \frac{4}{14} \left(-\frac{4}{4} \log_2 \frac{4}{4} - 0 \right) + \\ & \frac{5}{14} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) \end{aligned}$$

$$\begin{aligned} &= \frac{5}{14} (0.97) + \frac{4}{14} (1) + \frac{5}{14} (0.97) \\ &= 0.35 + 0.29 + 0.35 \\ &= 0.70 \end{aligned}$$

Info temp (D) =

temp

	Y	N
H	2	2
M	4	2
C	3	1

Info temp (D) =

$$\begin{aligned} & \frac{4}{14} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right] + \\ & \frac{6}{14} \left[-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right] + \end{aligned}$$

$$\frac{4}{14} \left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right]$$

$$= \frac{4}{14} (1) + \frac{6}{14} (0.92) + \frac{4}{14} (0.81)$$

$$= 0.29 + 0.39 + 0.23$$

$$= 0.91$$

Humidity:

	Y	N
H	3	4
N	6	1

$$0.52 +$$

$$0.19$$

$$\text{Info}_{\text{Hum}}(D) = \frac{7}{14} \left[-\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} \right] +$$

$$\frac{7}{14} \left[-\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} \right]$$

$$= \frac{7}{14} (0.98) + \frac{7}{14} (0.59)$$

$$= 0.49 + 0.30$$

$$= 0.788$$

$$0.1029$$

$$0.4284$$

Info windy (D) =

	Y	N
windy	6	2
	3	3

$$\text{Info}_{\text{windy}}(D) = \frac{8}{14} \left[-\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right] +$$

$$\frac{6}{14} \left[-\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6} \right]$$

$$= \frac{8}{14} (0.8113) + \frac{6}{14} (1)$$

$$0.8113$$

$$\approx 0.8922$$

$$0.4636$$

$$0.3113$$

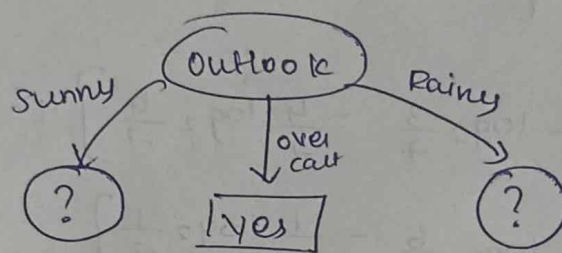
$$\text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\text{Gain}(\text{outlook}) = 0.940 - 0.70 \\ = 0.24 \checkmark$$

$$\text{Gain}(\text{Temp}) = 0.940 - 0.911 \\ = 0.029$$

$$\text{Gain}(\text{Hum}) = 0.940 - 0.788 \\ = 0.152$$

$$\text{Gain}(\text{wind}) = 0.940 - 0.892 \\ = 0.048.$$



Sunny, Rainy table.

Outlook	Temp	Hum	windy	PT
Sunny	Hot	High	Weak	No
	Hot	High	Strong	No
	Mild	High	Weak	No
	Cold	Normal	Weak	Yes
	Mild	Normal	Strong	Yes

Yes = 2, No = 3

$$\text{Info}(\text{sunny}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \\ = 0.5288 + 0.4422$$

$$\text{sunny} = 0.941$$

Attribute, Temp

	Y	N
H	0	2
M	1	1
C	1	0

$$\begin{aligned}
 \text{Info}_{\text{Temp}}(D) &= \frac{2}{5} \left(-\frac{2}{2} \log_2 \frac{2}{2} \right) + \\
 &\quad \frac{2}{5} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \\
 &\quad \frac{1}{5} \left(-\frac{1}{1} \log_2 \frac{1}{1} \right) \\
 &= \frac{2}{5} (0) + \frac{2}{5} (1) + \frac{1}{5} (0)
 \end{aligned}$$

$$\text{Info}_{\text{Temp}}(D) = 0.4$$

Hum,

	Y	N
H	0	3
N	2	0

$$\begin{aligned}
 \text{Info}_{\text{Hum}}(D) &= \frac{3}{5} \left(-\frac{3}{3} \log_2 \frac{3}{3} \right) + \frac{2}{5} \left(-\frac{2}{2} \log_2 \frac{2}{2} \right) \\
 &= 0
 \end{aligned}$$

windy

	Y	N
W	0	2
S	1	1

0.5283

$$\begin{aligned}
 \text{Info}_{\text{wind}}(D) &= \frac{3}{5} \left(-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) + \\
 &\quad \frac{2}{5} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) \\
 &= \frac{3}{5} (0.9183) + \frac{2}{5} (1) \\
 &= 0.5510 + 0.4 \\
 &= 0.9510
 \end{aligned}$$

$$\text{Gain}(\text{Temp}) = 0.971 - 0.4$$

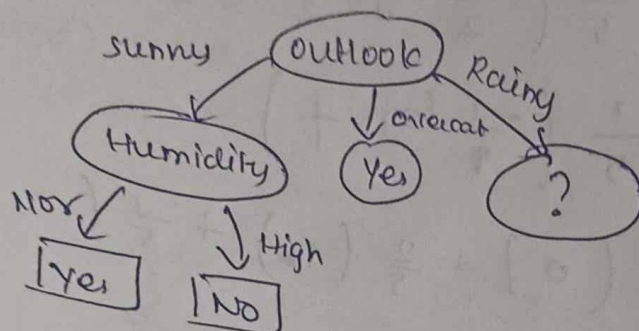
$$= 0.5710$$

$$\text{Gain}(\text{Hum}) = 0.971 - 0$$

$$= 0.971 \checkmark$$

$$\text{Gain}(\text{wind}) = 0.971 - 0.9510$$

$$= 0.02$$



Rainy table:-

Outlook	Temp	Hum	Windy	PT
Rainy	Mild	High	w	y
	Cool	Nor	w	y
	Cool	Nor	S	N
	Mild	Nor	w	y
	Mild	High	S	N

$$0.3900$$

$$0.4422$$

$$\text{Rainy} = \text{yes} = 3, \text{NO} = 2$$

$$\text{Info}(\text{Rainy}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.9710$$

Attribute;	Temp	y	N
M	2	1	
C	1	1	

$$\begin{aligned} \text{Info temp}(D) &= \frac{3}{5} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right] + \\ &\quad \frac{2}{5} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] \\ &= \frac{3}{5} (0.9183) + \frac{2}{5} (1) \\ &= 0.9510 \end{aligned}$$

Hum,		Y	N
H	1	1	
N	2		1

$$\begin{aligned} \text{Info Hum (Red)} &= \frac{2}{5} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] + \\ &\quad \frac{3}{5} \left[-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right] \\ &= \frac{2}{5} (1) + \frac{3}{5} (0.9183) \\ &= 0.9510 \end{aligned}$$

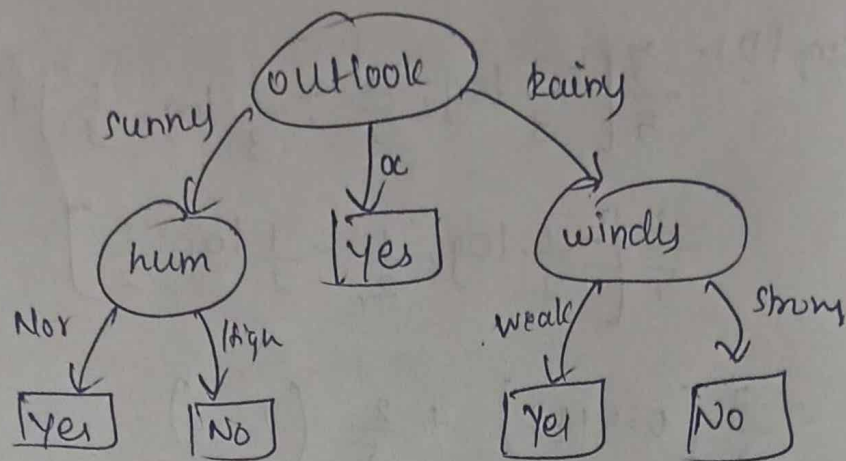
windy,		Y	N
W	3	0	
S	2		2

$$\begin{aligned} \text{Info windy (Rainy)} &= \frac{3}{5} \left(-\frac{3}{3} \log_2 \frac{3}{3} - 0 \right) + \frac{2}{5} \left(-\frac{2}{2} \log_2 \frac{2}{2} \right) \\ &= 0 \end{aligned}$$

$$\text{Gain (temp)} = 0.9710 - 0.9510 = 0.020$$

$$\text{Gain (Hum)} = 0.9710 - 0.9510 = 0.020$$

$$\text{Gain (wind)} = 0.9710 - 0 = 0.9710 \checkmark$$



2. KNN:

- Non-parametric
- Classification & Regression
- lazy learner algo.

x	y	class
7	7	Bad
7	4	Bad
3	4	Good
1	4	Good

New instance $(3, 7)$
 with $K=3$.

$$= \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$= \sqrt{(x_2 - x_1) + (y_2 - y_1)}$$

↓
Euclidean

$$= |(x_2 - x_1) + (y_2 - y_1)|$$

↓
Manhattan

B $| (7 - 3) + (7 - 7) | = | 4 + 0 | = 4$ ②

B $| (7 - 3) + (4 - 7) | = | 4 + (-3) | = 1$ ④

G $| (3 - 3) + (4 - 7) | = | 0 + (-3) | = 3$ ①

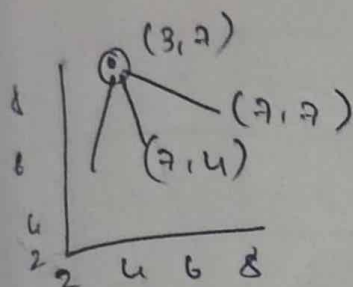
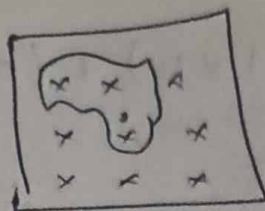
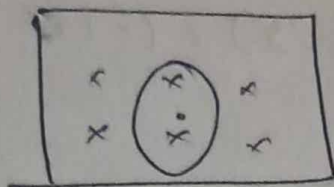
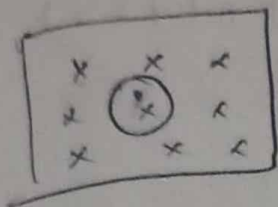
G $| (1 - 3) + (4 - 7) | = | -2 + (-3) | = 5$ ③

$K=1$: Good

$K=2$: Bad

$K=3$: Good

$K=3 \therefore (3, 7)$ the good.



③ linear regression:

a) least square method:

$$y = b_0 + b_1 x$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \left. \begin{array}{l} \text{slopes} \\ m \end{array} \right\}$$

x	y	$x - \bar{x}$	$(x_i - \bar{x})^2$	$y_i - \bar{y}$
1	14	-1	1	-6
3	24	1	1	4
2	18	0	0	-2
1	17	-1	1	-3
3	27	1	1	7
<u>10/5</u>	<u>100/5</u>	<u>-2</u>	<u>4</u>	<u>4</u>
$\bar{x} = 2$	$\bar{y} = 20$			

$\hat{y} =$

11-2
= 11+2

$$= (-1 \times -6) + (1 \times 4) + (0 \times -2) + (-1 \times -3) + (1 \times 7)$$

$$= 6 + 4 + 3 + 7$$

$$= \frac{20}{4}$$

$$= 5$$

$$\boxed{b_1 = 5} \quad \boxed{\text{slope}}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

$$b_0 = 20 - 5 \times 2$$

$$= 20 - 10$$

$$\boxed{b_0 = 10}$$

$$\boxed{\hat{y} = b_0 + b_1 x + \epsilon}$$

$$\hat{y} = 10 + 5x$$

$$\boxed{y = b_0 + b_1 x}$$

\hat{y}

$(x_i - \bar{x})(y_i - \bar{y})$	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})^2$	$(y_i - \bar{y})^2$
-6	-1	-6	1	36
4	1	4	1	16
-2	0	-2	0	4
3	-1	3	1	9
7	1	7	1	49

b) R^2

$$R^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2}$$

$$\hat{y} = mx + c$$

	y	x - \bar{x}	y - \bar{y}	(x - \bar{x})(y - \bar{y})	(x - \bar{x}) ²	(y - \bar{y}) ²	\hat{y}	$\hat{y} - \bar{y}$	($\hat{y} - \bar{y}$) ²
1	14	-1	-6	6	1	36	15	5	25
3	24	-1	4	-4	1	16	25	5	25
2	18	0	-2	2	0	4	20	0	0
1	17	-1	-3	3	1	9	15	5	25
3	27	1	7	-7	1	49	25	5	25
Σ	20			20		114			100

$\frac{20}{4}$

$$\hat{y} = 5x + 10$$

$$5(1) + 10$$

$$\hat{y} = mx + c$$

$$y = 5x + 10$$

$$\bar{y} = m\bar{x} + c$$

$$20 = 5 \times 2 + c$$

$$20 = 10 + c$$

$$10 + c = 20$$

$$c = 20 - 10$$

$$c = 10$$

$$R^2 = \frac{\sum (\hat{y} - \bar{y})^2}{\sum (y - \bar{y})^2}$$

$$R^2 = \frac{100}{114} = 0.877$$

$$R^2 = 88\%$$

Sigmoid Function

$$\text{odds}(p) = \frac{\text{Prob of event happening}}{\text{Prob of event Not happening}}$$

$$\text{odds}(p) = \frac{p}{1-p}$$

Value of Odds range from 0 to ∞

$$y = \beta_0 + \beta_1 x$$

Take log on both sides

$$\log \frac{p(x)}{1-p(x)} = \beta_0 + \beta_1 x$$

Exponentiating both sides

$$e^{\ln \frac{p(x)}{1-p(x)}} = e^{\beta_0 + \beta_1 x}$$

$$\frac{p(x)}{1-p(x)} = e^{\beta_0 + \beta_1 x}$$

$$[; e^{\ln} = 1]$$

$$\text{Let } y = e^{\beta_0 + \beta_1 x}$$

$$\text{Then, } \frac{p(x)}{1-p(x)} = y$$

$$p(x) = y [1-p(x)]$$

$$= y - y p(x)$$

$$p(x) + y p(x) = y$$

$$p(x) (1+y) = y$$

$$p(x) = y / (1+y)$$

$$p(x) = \frac{e^{B_0 + B_1 x}}{1 + e^{B_0 + B_1 x}}$$

divide $e^{B_0 + B_1 x}$

$$= \frac{e^{B_0 + B_1 x} / e^{B_0 + B_1 x}}{1 / e^{B_0 + B_1 x} + \frac{e^{B_0 + B_1 x}}{e^{B_0 + B_1 x}}}$$

$$= \frac{1}{1 + e^{-(B_0 + B_1 x)}}$$