

* Problem - 1 :

1) For calculating total entropy, we know that we have 5 datapoints that have a negative class value while we have 4 datapoints that have a positive class value

$$\text{Hence Total Entropy } E = -\frac{4}{9} \log_2 \left(\frac{4}{9}\right) - \frac{5}{9} \log_2 \left(\frac{5}{9}\right)$$

$$\therefore E = 0.9911$$

Let's calculate information gain if we split using attribute a_1 .

We know,

$$\text{gain} = E - \sum_i \frac{n_i}{m} E(i)$$

For a_1 ,

a_1/C	+	-
T	3	1
F	1	4

$$\therefore \text{gain} = E - \frac{4}{9} E(T) - \frac{5}{9} E(F)$$

$$E(T) = -\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \left(\frac{1}{4}\right) = 0.8112$$

$$E(F) = -\frac{1}{5} \log_2 \left(\frac{1}{5}\right) - \frac{4}{5} \log_2 \left(\frac{4}{5}\right) = 0.7219$$

$$\text{gain} = E - \frac{4}{9}(0.8112) - \frac{5}{9}(0.7219)$$

$$= 0.9911 - 0.3605 - 0.4011$$

$$= 0.2295$$

gain for splitting = gain $\underset{a_1}{=} 0.2295$
using attribute a_1

$$\boxed{\text{gain}_{a_1} = 0.2295} \quad \text{--- (1)}$$

Performing binary division at distinct values,

$\rightarrow a_2 (1.0)$

1.0\c	+	-
\geq	4	5
<	0	0

Elapsed Cell Block + (Pseudo)

$$E(\geq 1) = -\frac{4}{9} \log_2(4/9) - \frac{5}{9} \log_2(5/9)$$

$$= 0.9911$$

$$E = 0.9911$$

$$\boxed{\text{gain}_{a_2=1.0} = 0}$$

→ $a_2(3.0)$

2.0/c	+	-
γ=2	3	5
<	1	0

~~log₂(6) = 2.585~~

$$E(\gamma=3.0) = \frac{3}{8} (\log_2(3/8)) + \frac{5}{8} (\log_2(5/8)) \\ \Rightarrow 0.9544$$

$$E(<3.0) = 0$$

$$\text{gain}_{a_2=3.0} \Rightarrow E - \sum_{i=1}^m n_i E(i)$$

$$\Rightarrow E - \frac{8}{9}(0.9544)$$

$$\Rightarrow 0.9911 - 0.8484$$

$$\Rightarrow 0.1427$$

$$\boxed{\text{gain}_{a_2=3.0} = 0.1427}$$

→

Answers

4.0(C)	1	-
$\gamma =$	3	4
<	1	1

Balanced $\Rightarrow E(CC=CC) + EC(CC=CC)$

$$E(\gamma=4) = \frac{3}{2} \log \frac{3}{2} + \frac{4}{2} \log \frac{4}{2}$$

$$\approx 0.9852$$

$$E(C) = \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2}$$

$$= 1$$

$$\text{gain} = E - \sum_i \frac{n_i}{m} E(C_i)$$

$$= 0.9911 - \frac{2}{9}(0.9852) - \frac{2}{9}(1)$$

$$= 0.9911 - 0.2222 - 0.7663$$

$$= \cancel{0.0026}$$

↓

$$\text{gain}_{a_2=4.0} = 0.0026$$

$\rightarrow \alpha_2(5.0)$

$\gamma = 5.0$	C	+	-
$\gamma =$		2	4
$\gamma <$		2	1

$$E(\gamma \geq 5) = -\frac{2}{6} \log_2(2/6) - \frac{4}{6} \log_2(4/6)$$

$$E(\gamma < 5) = -\frac{2}{3} \log_2(2/3) - \frac{1}{3} \log_2(1/3)$$

$$E(\gamma = 5) = 0.9183$$

$$E(\gamma < 5) = 0.9183$$

$$\text{gain}_{\alpha_2=5.0} = E - \frac{6(E(\gamma = 5))}{9} - \frac{3(E(\gamma < 5))}{9}$$

$$= 0.9911 - \frac{2(0.9183)}{3} - \frac{1(0.9183)}{3}$$

$$= 0.0728$$

$$\boxed{\text{gain}_{\alpha_2=5.0} = 0.0728}$$

$\rightarrow a_2(6 \cdot 0)$

6.0	c	+ -
$\gamma = 2$	2	2
$<$	2	3

$$E(\gamma=6) = -\frac{2}{4} \log(\frac{2}{4}) - \frac{2}{4} \log(\frac{2}{4}) = 1$$

$$E(<6) = -\frac{2}{5} \log(\frac{2}{5}) - \frac{2}{5} \log(\frac{4}{5}) = 0.971$$

$$\text{gain}_{a_2(6.0)} = E - E(\gamma=6) - \frac{4}{9} + \frac{5}{9} E(<6)$$

$$= E - \frac{4}{9}(1) - \frac{5}{9}(0.971)$$

$$= 0.9911 - 0.4444 = 0.5394$$

$$= 0.072$$

$$\boxed{\text{gain}_{a_2=6} = 0.072}$$

→ $a_2 (x=0)$

$a \cdot 0$	C	f	-
γ_2	1	2	
L	3	3	

$$E(\gamma=x) = -\frac{1}{3} \log(\frac{1}{3}) - \frac{2}{3} \log(\frac{2}{3})$$

$$\approx 0.9183$$

$$E(x) = -\frac{1}{2} \log(\frac{1}{2}) - \frac{1}{2} \log(\frac{1}{2})$$

$$\approx 1$$

$$\text{gain } a_2 = x=0 \quad E = \frac{1}{9} (E(\gamma=x)) + \frac{6}{9} E(x)$$

$$= 0.9911 + \frac{1}{3} (0.9183) - \frac{2}{3} (1)$$

$$\approx 0.0183$$

$$\boxed{\begin{array}{l|l} \text{gain} & = 0.0183 \\ a_2 = x=0 & \end{array}}$$

$\rightarrow a_2(8.0)$

$a_2 \setminus c$	+	-
\geq	0	1
$<$	9	4

$$E(\gamma = 8) = 0$$

$$E(C \leq 8) = \frac{1}{2} \log(\frac{1}{2}) - \frac{1}{2} \log(\frac{1}{2})$$

$\Rightarrow 1$

$$\text{gain}_{a_2 = 8.0} = E - \frac{8(1)}{9}$$

$$\Rightarrow 0.1023$$

$\text{gain}_{a_2 = 8.0}$	$= 0.1023$
---------------------------	------------

gain for a_1 is greater than any of the gains for a_2 . Hence we should first select a_1 for splitting for the decision tree.

\rightarrow Hence a_1 is the best attribute to start splitting.

2). Since instances are just used as serial numbers, we should not use it as an attribute for splitting. Suppose we use instance as an attribute. An attribute usually show something about the class distribution. If the attributes are jumbled then there would not be any change to the class distribution since they are linked to the class values. If we jumble the dataset, the values of instances don't change but the class distribution changes instantly. So it is not advisable to use instance as an attribute while splitting.

* Problem - 2 :

		Class Label	
A	B	+	-
T	T	0	20
T	F	20	10
F	T	15	0
F	F	0	35

cost matrix	Attribute value	
Actual class	+	-
+	-1	100
-	0	-10

	A	B	+	-
T	T	0	20	
F	F	20	10	
	T	15	0	
	F	0	35	

For A:

A	+	-	+ 35
T	20	30	- 65
F	15	35	

$$Gini_{\text{orig}} = 1 - \left(\frac{35}{100}\right)^2 - \left(\frac{65}{100}\right)^2 \\ = 0.455$$

$$\boxed{Gini_{\text{orig}} = 0.455}$$

$$Gini_{\text{split}} = \sum_i^n \frac{Gini(i)}{Gini(T)}$$

$$= \frac{50}{100} \left(1 - \left(\frac{20}{50}\right)^2 - \left(\frac{30}{50}\right)^2\right)$$

$Gini(T)$

$$+ \frac{50}{100} \left(1 - \left(\frac{35}{50}\right)^2 - \left(\frac{15}{50}\right)^2\right)$$

$$= \frac{1}{2}(0.48) + \frac{1}{2}(0.42) = 0.24 + 0.21$$

$$\boxed{Gini_{\text{split}} = 0.45}$$

For B:

B_{class}	t	$-$
T	15	20
F	20	45

$\begin{array}{r} 35 \\ - 65 \end{array}$

$$\text{Cini}_{\text{orig}} = 1 - \left(\frac{35}{100} \right)^2 - \left(\frac{65}{100} \right)^2 \\ \Rightarrow 0.455$$

$$\begin{aligned} \text{Cini}_{\text{spur}} &= \sum_i \frac{n_i}{n} \text{Cini}(i) \\ &\quad \downarrow \text{Cini}(T) \\ &= \frac{35}{100} \left(1 - \left(\frac{15}{35} \right)^2 - \left(\frac{20}{35} \right)^2 \right) \\ &\quad \downarrow \text{Cini}(F) \\ &+ \frac{65}{100} \left(1 - \left(\frac{20}{65} \right)^2 - \left(\frac{45}{65} \right)^2 \right) \end{aligned}$$

$$= 0.35(0.4898) + 0.65(0.4260)$$

$$= 0.1814 + 0.2769$$

$$= 0.4483$$

$$\boxed{\text{Cini}_{\text{spur}} = 0.4483}$$

- Since C_{ini}^{split} for B is lesser than C_{ini}^{split} for A, we should choose B for splitting.
- Hence splitting by attribute B, has to be the best way according to C_{ini} indices.

2)

for A:

		A		CM	
		+	-	+	F
T	20	30	-1	100	
	F	15	35	0	-10

From above tables,

$$\text{cost of splitting A} = 20(-1) + 30(0) \\ + 100(15) + 150(-10)$$

$$\text{cost}_{\text{Asplit}} \rightarrow -20 + 1500 + (-350)$$

$$\frac{\text{cost}_{\text{Asplit}} = 1130}{}$$

P.T.O.

For B :

B	+	-
T	15	20
F	20	45

CM	T	F
+	-1	100
-	0	-10

Cost of splitting B ,

$$\text{Coef}_{B \text{ split}} = 15(-1) + 20(0) + 20(100) \\ + (-10)(45)$$

$$= 2000 - 15 - 450$$

$$= 1535$$

Since cost of splitting A is lower,
we will choose attribute A for
splitting .

Problem - 3

ID	1	2	3	4	5	6	7	8	9	10
X	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1
Y	1	1	1	-1	-1	-1	-1	-1	1	1

H1 : if $X \leq 0.35 \rightarrow Y = 1$, else $Y = -1$

H2 : if $X \leq 0.75 \rightarrow Y = -1$, else $Y = 1$

H3 : if $X \leq 0.3$ or $X > 0.95 \rightarrow Y = 1$ else $Y = -1$

a)	X	$0.1 \quad 0.2 \quad 0.3 \quad 0.4 \quad 0.5 \quad 0.6 \quad 0.7 \quad 0.8 \quad 0.9$	$1 \quad 1 \quad 1 \quad -1 \quad -1 \quad -1 \quad -1 \quad -1 \quad 1$
	y		
	$h_1(x)$		

We are going to initialize weights.

X	$0.1 \quad 0.2 \quad 0.3 \quad 0.4 \quad 0.5 \quad 0.6 \quad 0.7 \quad 0.8 \quad 0.9$	$1 \quad 1 \quad 1 \quad -1 \quad -1 \quad -1 \quad -1 \quad -1 \quad 1$
$D_+(i)$	$0.1 \quad 0.1 \quad 0.1$	

always, $D_+(i) = \frac{1}{n}$ (for first step)

incorrect predictions := 0.9 and 1 (x)

hence no. of incorrect predictions = 2

$$\text{error} = \sum D_+(i) \cdot \delta(h_1(x) \neq y)$$

$$= 0.1 + 0.1 = 0.2$$

$$\boxed{\epsilon = 0.2}$$

$$\alpha = \frac{1}{2} \ln \left(\frac{1-\epsilon}{\epsilon} \right) = \frac{1}{2} \ln 4$$

$$= 0.6931$$

$$\boxed{\alpha = 0.6931}$$

correct classification = $D_2(i) = D_1(i) \exp(-\alpha y_i h_i(x_i))$
(new weights)

$$D_2(i) = 0.1 \times e^{(-0.6931)}$$

correct
classification

$$= 0.05$$

incorrect classification = $D_2(i) = D_1(i) \exp(-\alpha y_i h_i(x_i))$

$$D_2(i) = 0.1 \times e^{(0.6931)}$$

incorrect
classification = 0.2

\therefore new weights	X	0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
	Y	1 1 1 -1 -1 -1 -1 -1 1 1
	$h(x)$	1 1 1 -1 -1 -1 -1 -1 1 -1
	D_t	0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.05 0.2 0.2

P.T.O.

(b) For hypothesis H_2

X	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
y	1	1	1	-1	-1	-1	-1	-1	1	1

Assigning equal weights.

X	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
D_+	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.1

Values of y according to H_2 , ie $h_2(x)$

X	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
$h_2(x)$	-1	-1	-1	-1	-1	-1	-1	-1	1	1
y	1	1	1	-1	-1	-1	-1	-1	-1	1

instances where $h_2(x) \neq y$
predicts incorrectly

$$\therefore \text{error} = \sum_i \delta_i \text{ if } h_2(x) \neq y$$

$$= 0.1 + 0.1 + 0.1 + 0.1 + 0.1$$

$$= 0.5$$

$$\left[\begin{matrix} \varepsilon = 0.4 \\ \end{matrix} \right]$$

$$\alpha = \frac{1}{2} \ln \left(\frac{1-\varepsilon}{\varepsilon} \right) \quad \text{and} \quad \boxed{\alpha = 0.707}$$

$$= \frac{1}{2} \ln \left(\frac{0.6}{0.4} \right) = \frac{1}{2} \ln (3/2)$$

$$= \frac{1}{2} (\ln 0.54) = 0.2027$$

$$\left[\begin{matrix} \alpha = 0.2027 \\ \end{matrix} \right]$$

$$\text{New weights} = D_2(i) = D_i(i) \exp(-\alpha y h(x))$$

$$\text{correct prediction} = 0.1 \times e^{-(0.2027)} \\ \approx 0.0816$$

$$\text{incorrect prediction} = 0.1 \times e^{(0.2027)} \\ \approx 0.1225$$

new weights	X	0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
	y	1 1 1 -1 -1 -1 -1 -1 1 1
	$h(x)$	-1 -1 -1 -1 -1 -1 -1 1 1 1
	D_2	0.1225 0.1225 0.1225 0.08 0.08 0.08 0.08 0.1225 0.08 0.08

(c) For hypothesis h_3

X	0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
y	1 1 1 -1 -1 -1 -1 -1 1 1

initializing weights

X	0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
D_t	0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.1

Predicting y using h_3

X	0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0
h_3	1 1 1 -1 -1 -1 -1 -1 -1 1
y	1 1 1 -1 -1 -1 -1 -1 1 1

number of instances
with incorrect prediction = 1 (0.9)

$$\text{error} = \sum_i \delta(h(x) \neq y)$$

$$E = 0.1$$

$$d = \frac{1}{2} \ln \left(\frac{1-\varepsilon}{\varepsilon} \right) = \frac{1}{2} \ln (0.9/0.1)$$

$$= \frac{1}{2} \ln(9) = (2.1972)/2 = 1.0986$$

$$d = 1.0986$$

$$\text{New weights, } = D_{t+1}^{(i)} \exp(-\alpha y_t h_{t+1}^{(i)})$$

$$\begin{aligned} \text{correct prediction} &= 0.1 \times e^{-1.0986} \\ &\approx 0.0333 \end{aligned}$$

$$\text{incorrect prediction} = 0.1 \times e^{(1-0.98t)}$$
$$= 0.3$$

2) To be precise, for H1 the data instances where $\gamma = 0.9$ and $\gamma = 1$ are reweighted, for H2 the data instances where $\gamma = 0.1, 0.2, 0.3$ and 0.8 are reweighted and for H3 the data instance where $\gamma = 0.9$ is only resampled. Weights of all the data instances are changed but the others are changed for balancing the effect of these particular data instances which were incorrect.

* Problem - 4

Part I :

data points : { $(1, 1, -)$, $(2, 2, -)$, $(1, 3, -)$,
 $(4, 1, -)$, $(4, 4, -)$, $(5, 1, -)$,
 $(3, 0, -)$, $(8, 1, -)$, $(9, 1, +)$,
 $(0, 5, +)$, $(2, 8, +)$, $(7, 3, +)$,
 $(9, 4, +)$, $(9, 6, +)$, $(8, 7, +)$,
 $(7, 7, +)$, $(6, 8, +)$, $(5, 9, +)$,
 $(8, 9, +)$ }

Let's suppose one test point belongs to class p which can be + or -.

data point

to test $p = (5, 4, p)$

Calculating distance of pt from different datapoints.

By Pythagoras

$$\begin{aligned} d((1, 1, -)) &= \sqrt{(5-1)^2 + (4-1)^2} \\ &= \sqrt{25} \\ &= 5 \end{aligned} \quad \therefore f((1, 1, -)) = (5, -)$$

$$\begin{aligned} d((2, 2, -)) &= \sqrt{13} \\ &= 3.605 \end{aligned} \quad \therefore f((2, 2, -)) = (3.6, -)$$

$$\begin{aligned} d((1, 3, -)) &= \sqrt{18} \\ &= 4.123 \end{aligned} \quad \therefore f((1, 3, -)) = (4.123, -)$$

$$\begin{aligned} d((4, 1, -)) &= \sqrt{10} \\ &= 3.16 \end{aligned} \quad \therefore f((4, 1, -)) = (3.16, -)$$

$$\begin{aligned} d((4, 4, -)) &= \sqrt{18} \\ &= 1 \end{aligned} \quad \therefore f((4, 4, -)) = (1, -)$$

$$\begin{aligned} d((5, 1, -)) &= \sqrt{9} \\ &= 3 \end{aligned} \quad \therefore f((5, 1, -)) = (3, -)$$

$$\begin{aligned} d((2, 6, -)) &= \sqrt{13} \\ &= 3.605 \end{aligned} \quad \therefore f((2, 6, -)) = (3.6, -)$$

$$\begin{aligned} d((8, 1, -)) &= \sqrt{18} \\ &= 4.24 \end{aligned} \quad \therefore f((8, 1, -)) = (4.24, -)$$

$$\begin{aligned} d((9, 1, -)) &= \sqrt{25} \\ &= 5 \end{aligned} \quad \therefore f((9, 1, -)) = (5, -)$$

$$d((6, 5, +)) \approx \sqrt{2}$$
$$\approx 1.41$$
$$\therefore f((6, 5, +)) = (1.41, +)$$

$$d((2, 8, +)) \approx \sqrt{25}$$
$$\approx 5.0$$
$$\therefore f((2, 8, +)) = (5.0, +)$$

$$d((2, 3, +)) \approx \sqrt{5}$$
$$\approx 2.23$$
$$\therefore f((2, 3, +)) = (2.23, +)$$

$$d((9, 4, +)) \approx \sqrt{16}$$
$$\approx 4$$
$$\therefore f((9, 4, +)) = (4, +)$$

$$d((9, 6, +)) \approx \sqrt{20}$$
$$\approx 4.47$$
$$\therefore f((9, 6, +)) = (4.47, +)$$

$$d((8, 2, +)) \approx \sqrt{18}$$
$$\approx 4.24$$
$$\therefore f((8, 2, +)) = (4.24, +)$$

$$d((7, 7, +)) \approx \sqrt{13}$$
$$\approx 3.60$$
$$\therefore f((7, 7, +)) = (3.6, +)$$

$$d((6, 8, +)) \approx \sqrt{17}$$
$$\approx 4.123$$
$$\therefore f((6, 8, +)) = (4.123, +)$$

$$d((5, 9, +)) \approx \sqrt{25}$$
$$\approx 5$$
$$\therefore f((5, 9, +)) = (5, +)$$

$$d((8, 9, +)) \approx \sqrt{34}$$
$$\approx 5.83$$
$$\therefore f((8, 9, +)) = (5.83, +)$$

point	class	distance
(1, 1)	-	5
(2, 2)	-	3.6
(1, 3)	-	4.12
(4, 1)	-	3.16 ←
(4, 4)	-	1.0 ←—
(5, 1)	-	3.0 ←—
(2, 6)	-	3.6
(8, 1)	-	4.24
(9, 1)	-	5
(6, 5)	+	5 1.41 ←
(2, 8)	+	5 5
(4, 3)	+	5 2.23 ←
(9, 4)	+	5 4.47
(9, 6)	+	5 4.24
(8, 7)	+	3.6
(7, 8)	+	4.12
(6, 8)	+	5
(5, 9)	+	5.83
(8, 9)	+	

5 nearest :

(4, 4, -), (6, 5, +), (7, 3, +), (4, 4, -), (4, 1, -)

majority of the classes are negative
from the 5-nearest neighbour,
so the prediction is class '-'.

2) points

Manhattan

distance (d)

$(1, 1, ' - ')$	7
$(2, 2, ' - ')$	5
$(1, 3, ' - ')$	5
$(4, 1, ' + ')$	4
$(4, 4, ' - ')$	7 → min
$(5, 1, ' - ')$	3 → min
$(2, 6, ' - ')$	5
$(8, 1, ' - ')$	6
$(9, 1, ' - ')$	7
$(6, 5, ' + ')$	2 → min
$(2, 8, ' + ')$	7
$(7, 3, ' + ')$	3 → min
$(9, 4, ' + ')$	4
$(9, 6, ' + ')$	6
$(8, 7, ' + ')$	6
$(7, 7, ' + ')$	5
$(6, 8, ' + ')$	5
$(5, 9, ' + ')$	5
$(8, 9, ' + ')$	8

$(4, 4, ' - ')$

+

-

↓

$$w = \frac{1}{d^2} = \frac{1}{1} = 1$$

$(5, 1, ' - ')$

↓

-

↓

$$w = \frac{1}{d^2} = \frac{1}{9}$$

$(4, 5, ' + ')$

↓

+

↓

$$w = \frac{1}{d^2} = \frac{1}{4}$$

$(7, 3, ' + ')$

+

+

↓

$$w = \frac{1}{d^2} = \frac{1}{9}$$

since the weights are $\{1, 2, 3, 3\}$ we consider both possibilities for 3.

Suppose we pick the 3 which has a '+' class, then the weights.

$$\begin{array}{cc} - & + \\ | & | \\ & 1/4 \\ & 1/9 \end{array}$$

T: 1 13/36

$$w(-) > w(+)$$

Hence prediction would be -.

Suppose we pick the 3 which has a '-' class, then the weights

$$\begin{array}{cc} - & + \\ | & | \\ & 1/4 \\ 1/9 & \end{array}$$

T: 10/9 1/4

$$w(-) > w(+)$$

Hence prediction would be -.

\rightarrow Hence in both cases prediction is -.

▼ Problem 4. K-nearest Neighbor Classification

Part II.

▼ Preparing the dataset

▼ Importing Libraries

```
import pandas as pd
```

▼ Importing the training sample

```
train_sample = pd.read_csv('train.csv')
```

```
train_values = train_sample.values  
tn_X, tn_y = train_values[:, :-1], train_values[:, -1]
```

▼ Importing the testing sample

```
test_sample = pd.read_csv('test.csv')
```

▼ Removing unwanted columns from the testing sample

```
final_cols = ['x', 'y', 'z', 'actual-class',]  
test_sample = test_sample[final_cols]
```

```
test_values = test_sample.values  
tt_X, tt_y = test_values[:, :-1], test_values[:, -1]
```

K-Nearest Neighbors Algorithm

Problem 4_Part II

(1)

- ▼ For k = 3:

```
from sklearn.neighbors import KNeighborsClassifier #Importing KNN and CM

#preparing a model
model = KNeighborsClassifier(n_neighbors = 3)

#training the model on the train sample
model.fit(tn_X, tn_y)

#predicting using the model
prediction = model.predict(tt_X)
print(prediction)
print(tt_y)

[1. 0. 0. 1. 0. 1. 0. 0. 0. 1. 0. 0. 0. 1. 1. 1. 0. 0. 0.]
[1. 0. 0. 1. 0. 1. 0. 0. 0. 1. 0. 0. 0. 0. 1. 1. 0. 0. 0.]
```

- ▼ Probability of the final prediction

```
probab = model.predict_proba(tt_X)
print(probab)

[[0. 1.
 [0.66666667 0.33333333]
 [1. 0.]
 [0. 1.]
 [1. 0.]
 [0. 1.]
 [1. 0.]
 [1. 0.]
 [1. 0.]
 [1. 0.]
 [0. 1.]
 [1. 0.]
 [1. 0.]
 [1. 0.]
 [0.33333333 0.66666667]
 [0. 1.]
 [0. 1.]
 [1. 0.]
 [1. 0.]
 [1. 0.]
 [1. 0.]]]
```

▼ Confusion Matrix, accuracy, precision and F-measure

```

tp = 0
tn = 0
fp = 0
fn = 0

for x in range(len(prediction)):
    if(prediction[x] == 0 and tt_y[x] == 0):
        tn += 1
    if(prediction[x] == 1 and tt_y[x] == 1):
        tp += 1
    if(prediction[x] == 0 and tt_y[x] == 1):
        fn += 1
    if(prediction[x] == 1 and tt_y[x] == 0):
        fp += 1

total = tn + fn + tp + fp
cm = [[tp,fp],[fn,tn]]
accuracy = (tp + tn)/total
precision = tp/(tp + fp)
recall = tp/(tp + fn)
f_score = 2*precision*recall/(precision + recall)

print("Accuracy:",accuracy)
print("Precision:",precision)
print("Recall:",recall)
print("F-Measure:",f_score)

Accuracy: 0.95
Precision: 0.8571428571428571
Recall: 1.0
F-Measure: 0.923076923076923

```

▼ Problem 4_Part II

(2)

```

#Implementing weighted KNN
w_knn = KNeighborsClassifier(n_neighbors=3, weights='distance', p=2, metric='euclidean')

# Training the KNN model
w_knn.fit(tn_X, tn_y)

KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='euclidean',
                     metric_params=None, n_jobs=None, n_neighbors=3, p=2,

```

```
weights='distance')

#Predicting the Output
w_prediction = w_knn.predict(tt_X)
print(w_prediction)

[1. 0. 0. 1. 0. 0. 0. 0. 1. 0. 0. 0. 1. 1. 1. 0. 0. 0.]

tp = 0
tn = 0
fp = 0
fn = 0

for x in range(len(w_prediction)):
    if(w_prediction[x] == 0 and tt_y[x] == 0):
        tn += 1
    if(w_prediction[x] == 1 and tt_y[x] == 1):
        tp += 1
    if(w_prediction[x] == 0 and tt_y[x] == 1):
        fn += 1
    if(w_prediction[x] == 1 and tt_y[x] == 0):
        fp += 1

total = tn + fn + tp + fp
cm = [[tp,fp],[fn,tn]]
accuracy = (tp + tn)/total
precision = tp/(tp + fp)
recall = tp/(tp + fn)
f_score = 2*precision*recall/(precision + recall)

print("Accuracy:",accuracy)
print("Precision:",precision)
print("Recall:",recall)
print("F-Measure:",f_score)

Accuracy: 0.95
Precision: 0.8571428571428571
Recall: 1.0
F-Measure: 0.923076923076923
```

Since the dataset extremely small, both of the models perform same. If dataset was huge maybe one of them would have been better, and my guess would have been the weighted KNN would've been more accurate



✓ 0s completed at 5:11 PM