Problem 1.1 :-

1)

$$X < 60$$

$$Y < 40 \qquad X < 80 \quad R4$$

$$R3$$

$$X < 20 \qquad Y < 30$$

$$R1 \qquad R2 \qquad R5 \qquad R6$$



3)



→ Entropy

→ Gini

→ Misclassification error

Attached pdf with Sample code

4) a)

$X_1 < 1$

```
        X₁<1
         /\
        /  \
      X₂<1  (5)
       /\
      /  \
    X₁<0 (15)
     /\
    /  \
  (3)  X₂<0
        /\
       /  \
     (10)  (0)
```

b)



5)

Majority approach:-

we have 6 out of 10 samples which has probability greater than 0.5

so predicted class = Red

Average approach:-

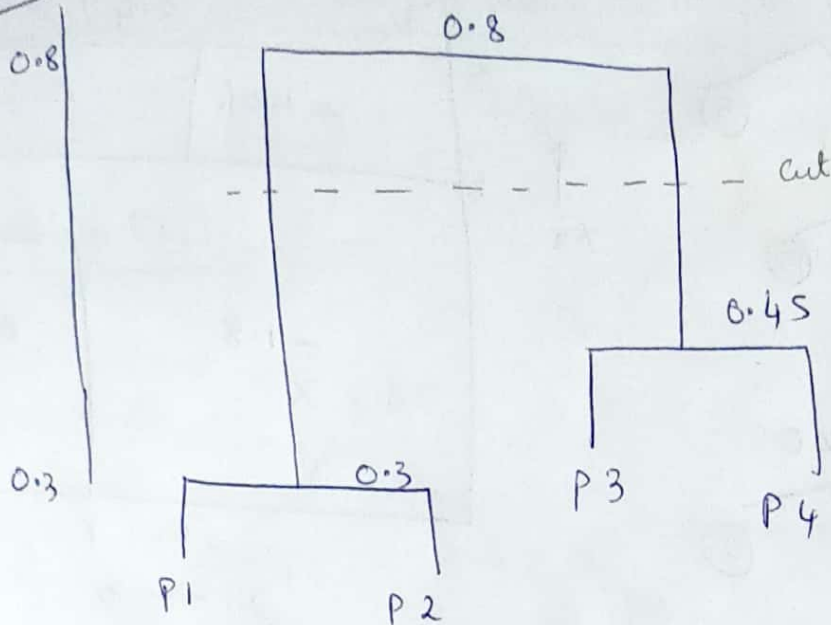mean (observed probabilities) = 0.45

so the predicted class = Green    since the average probability is less than 0.5

Problem 1.2

2)
$$\begin{bmatrix} & 0.3 & 0.4 & 0.7 \\ 0.3 & & 0.5 & 0.8 \\ 0.4 & 0.5 & & 0.45 \\ 0.7 & 0.8 & 0.45 & \end{bmatrix}$$

complete

a)



|     | $P_1$ | $P_2$ | $P_3$ | $P_4$ |
|-----|-------|-------|-------|-------|
| $P_1$ | 0 |  |  |  |
| $P_2$ | 0.3 Min | 0 |  |  |
| $P_3$ | 0.4 | 0.5 | 0 |  |
| $P_4$ | 0.7 | 0.8 | 0.45 | 0 |

Max operation

|           | $P_1 - P_2$ | $P_3$ | $P_4$ |
|-----------|-------------|-------|-------|
| $P_1 - P_2$ | 0 |  |  |
| $P_3$ | 0.5 | 0 Min |  |
| $P_4$ | 0.8 | 0.45 | 0 |

$\Rightarrow$

Max operation

|           | $P_1 - P_2$ | $P_3 - P_4$ |
|-----------|-------------|-------------|
| $P_1 - P_2$ | 0 |  |
| $P_3 - P_4$ | 0.8 Min | 0 |

b) single

|     | $P_1$ | $P_2$ | $P_3$ | $P_4$ |
|-----|-------|-------|-------|-------|
| $P_1$ | 0 |  |  |  |
| $P_2$ | 0.3 Min | 0 |  |  |
| $P_3$ | 0.4 | 0.5 | 0 |  |
| $P_4$ | 0.7 | 0.8 | 0.45 | 0 |

Min operation

$P_1 - P_2 \quad P_3 \quad P_4$

$P_1 - P_2$    0    min

$P_3$    0.4    0

$P_4$    0.7    0.45   0

Min operation

$P_1 - P_2 - P_3 \quad P_4$

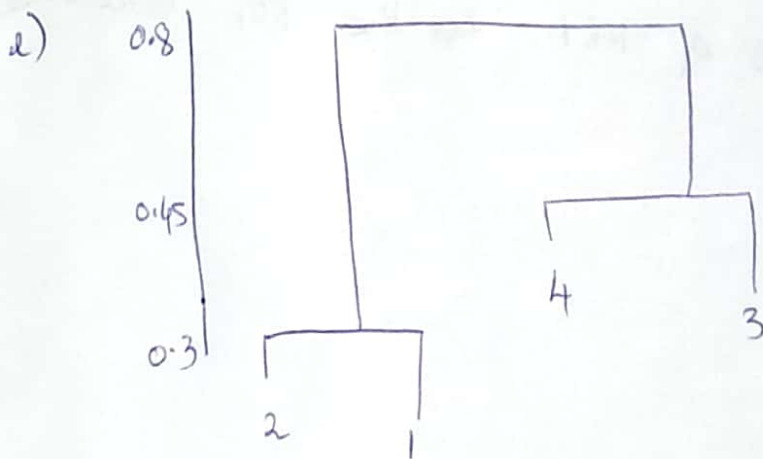$\Rightarrow P_1 - P_2 - P_3$    0

$P_4$    min    0.45    0

$\Rightarrow P_1 - P_2 - P_3 - P_4$



c) (1, 2) (3, 4)

d) (4) (1, 2, 3)

e)



3) Attached Assignment - Theory . Rmd file.

6) a) "explains 10% of the variation"

me ans the first principle component

has captured only 10% of the total variance.
and 90% of information is lost by
Projecting tissue sample and gene data set.

b) Since the first principle component
explained only 10% of the variance. I would
suggest the researcher to use (A vsB) as a
feature of the data set. This inturn increases
the variance before applying the two-
Sample t-test.

c) Attached R-code in Assigment4_Theory.Rmd
Analysis:- After performing A vsB there is
an Improvement of 1.6% in the PC₁ variance.