SaiParthish Mandumula
(100202 2847)

1. Let Anthony be represented as 'a'
   Brutus as 'b'
   Caesar as 'C₁'
   Calpurnia as 'C₂'
   cleopatra as 'C₃'

∴ log weighted $tf \rightarrow \begin{cases} 1 + \log_{10} tf & , tf > 0 \\ 0 & , \text{otherwise} \end{cases}$

\* log weighted

|  | Anthony & cleopatra | Julius Caesar | The tempest | Hamlet |
|---|---|---|---|---|
| a → | 3.195 | 2.86 | 2.69 | 3 |
| b → | 1.60 | 3.19 | 2 | 1.30 |
| C₁ → | 3.36 | 1.30 | 3.35 | 3 |
| C₂ → | 0 | 2 | 0 | 0 |
| C₃ → | 2.75 | 0 | 0 | 2.77 |

\* log. weighted IDF
$$\frac{}{IDF} \rightarrow \log_{10}\left(\frac{N}{df}\right)$$

| a | 0.30 |
|---|---|
| b | 1.30 |
| C₁ | 0.22 |
| C₂ | 1.30 |
| C₃ | 1.69 |

Tf-IDF vector representation

Anthony & Cleopatra (A)
$$\rightarrow 6.9577a + 2.02b + 0.73c_1 + 0c_2 + 4.64c_3$$

Julius Caesar (J) $= 0.86a + 4.15b + 0.29c_1 + 26c_2 + 0c_3$

The tempest (T) $\rightarrow 0.81a + 2.6b + 0.34c_1 + 0.c_2 + 0c_3$

Hamlet (H) $\rightarrow 0.9a + 1.6b + 0.66c_1 + 0c_2 + $ ~~$c_3$~~
$$4.84c_3$$

$$||A|| = 5.236$$
$$||J|| = 4.921$$
$$||T|| = 2.82$$
$$||H|| = 4.341$$

∴ Normalized Tf-IDF vector representation.

A $\rightarrow$ $0.12a + 0.46 + 0.14c + 0c_2 + 0.89c_3$

J $\rightarrow$ $0.17a + 0.83b + 0.06c_1 + 0.52c_2 + 0c_3$

T $\rightarrow$ $0.29a + 0.92b + 0.26c_1 + 0c_2 + 0c_3$

H $\rightarrow$ $0.21a + 0.39b + 0.15c_1 + 0c_2 + 0.80c_3$
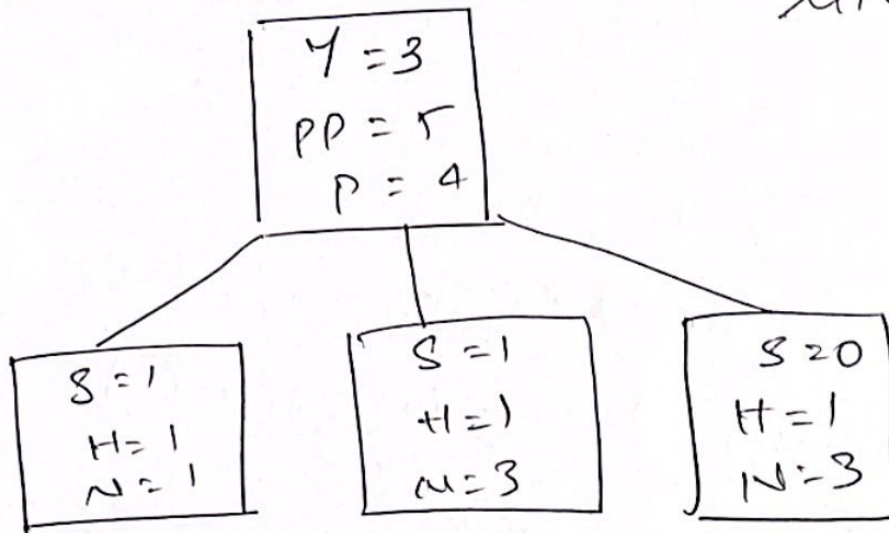
✗ Cosine similarities b/w 'Anthony & Cleopatra'

and other  A-J = 0.371
         A-T = 0.4566
         A-H = 0.998

⇒ Anthony & cleopetra' is most similar
to 'Hamlet'.

Sai Parthish·M

② ① <u>Age as the root</u>

$$\text{Gini(Parent)} : 1 - \left[\frac{2^2 + 3^2 + 7^2}{12^2}\right]$$

$$= 0.569$$

```
              ┌─────────┐
              │  Y = 3  │
              │  PP = 5 │
              │  P = 4  │
              └─────────┘
         ┌────────┼────────┐
   ┌─────────┐ ┌───────┐ ┌────────┐
   │  S = 1  │ │ S = 1 │ │ S = 0  │
   │  H = 1  │ │ H = 1 │ │ H = 1  │
   │  N = 1  │ │ M = 3 │ │ N = 3  │
   └─────────┘ └───────┘ └────────┘
```

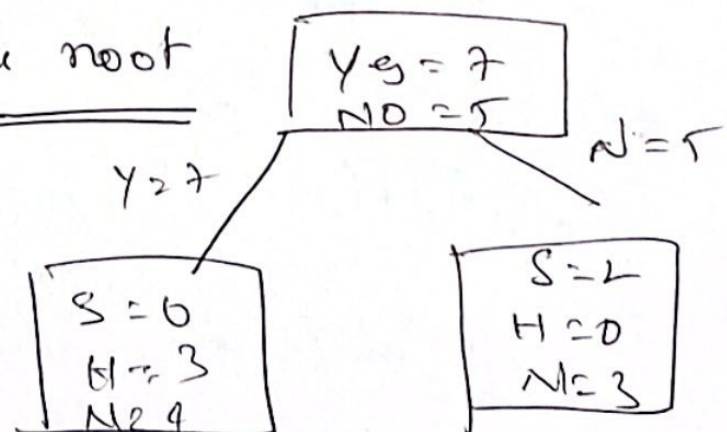$$\text{GINI}(Y) = 1 - \left[\frac{1^2 + 1^2 + 1}{3^2}\right] = 0.67$$

$$\text{GINI}(PP) = 1 - \left[\frac{1^2 + 1^2 + 3^2}{5^2}\right] = 0.56$$

$$\text{GINI}(P) = 1 - \left[\frac{1^2 + 3^2}{4^2}\right] = 0.375$$

child split $= \frac{3}{12}(0.67) + \frac{5}{12}(0.56) + \frac{4}{12}(0.375)$

$$= 0.525$$

Gain $= \text{GINI(Parent)} - \text{Given(child)}$

$$= 0.569 - 0.525$$

$$= 0.044$$

② <u>Affirmation as the root</u>

```
                        ┌─────────┐
                        │ Yg = 7  │
                        │ NO = 5  │
                        └─────────┘
                   Y=7 /           \ N = 5
            ┌─────────┐          ┌─────────┐
            │  S = 6  │          │  S = 2  │
            │  H = 3  │          │  H = 0  │
            │  N = 4  │          │  N = 3  │
            └─────────┘          └─────────┘
```

Scanned with CamScanner

$$GINI(Y) = 1 - \left[\frac{3^2 + 4^2}{7^2}\right] = 0.489$$

$$GINI(N) = 1 - \left[\frac{2^2 + 3^2}{5^2}\right] = 0.48$$

Child split $= \frac{7}{12}(0.489) + \frac{5}{12}(0.48) = 0.48525$

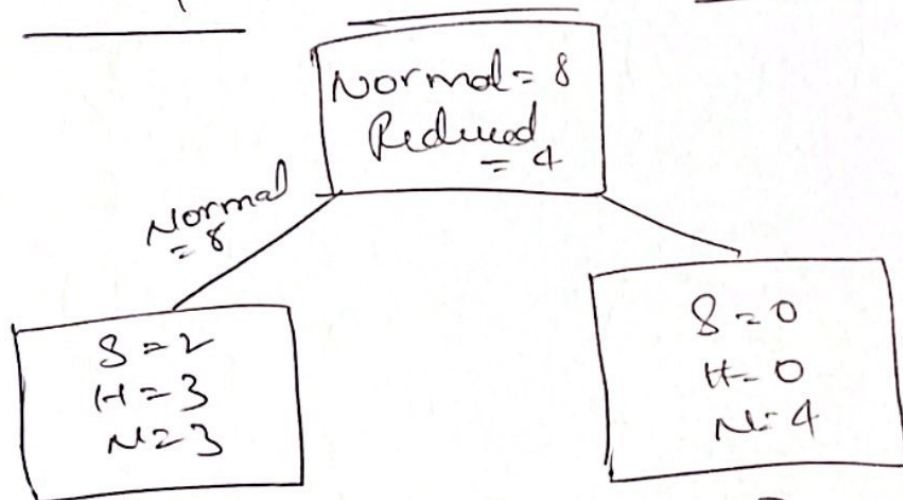Gain $=$ GINI(Parent) $-$ GINI(child)

$= 0.569 - 0.48525$

$= 0.08325$

0.569
0.48525
.0.08325

③ Tear Production Rate as the root



$$GINI(N) = 1 - \left[\frac{2^2 + 3^2 + 3^2}{8^2}\right] = 0.656$$

$$GINI(R) = 1 - \left[\frac{4^2}{4^2}\right] = 0$$

Child split $= \frac{8}{12}[0.656] + \frac{4}{12}[0] = 0.437$

Sri·Parthish·Mandumula

∴ Gain = GINI ( Parent) - GINI ( child )

$$= 0.569 - 0.437$$

$$= 0.132$$

Since we get the highest gain with the Tear Production as the root, we will proceed with 'Tear production' as the root node.

③ Test Instance attribute

| outlook | Temperature | Humidity | windy |
|---------|-------------|----------|-------|
| Rainy | Cool | High | True. |

$P(play\ Golf = Y/A) = P(A/play\ Golf = Yes) \cdot P(Play\ Golf = Yes)$

$P(play\ Golf = No/A) = P(A/play\ Golf = NO) \cdot P(play\ Golf = NO)$

we calculate the following

$P(play\ Golf = Yes) = 9/14$

$P(play\ Golf = NO) = 5/14$ → $P(outlook = Rainly/play\ golf = y) = 2/9$

$P(Temperature = cool/play\ Golf = y) = \frac{3}{9}$

$P(Temperature = cool/play\ Golf = N) = 1/5$

$P(windy = True/play\ Golf = N) = 3/9$

$P(\text{windy} = True / Play\ Golf = N) = 3/5$

$P(outlook) = Rainly / play\ Golf = N) = 3/5$

$P(Humidity = High / play\ Golf = y) = 3/9$

$P(Humidity = High / play\ Golf = N) = 4/5$

$P(play\ Golf = Yes / A) =$
$$\left(\frac{2}{9}\right)\left(\frac{3}{9}\right)\left(\frac{3}{9}\right)\left(\frac{3}{9}\right)\left(\frac{9}{14}\right) = 5.29 \times 10^{-3}$$

$P(play\ Golf = No / A) =$
$$\left(\frac{3}{5}\right)\left(\frac{1}{5}\right)\left(\frac{4}{5}\right)\left(\frac{3}{5}\right)\left(\frac{5}{14}\right) = 20 \times 10^{-3}$$

Rest . attributes . class $\longrightarrow$ play golf = No

4) a Keeping the Samples in mind,

Support vector ( +ve class) $\rightarrow (0, 2), (2, 0)$

Support Vector (-ve class) $\rightarrow (-2, -2)$

For the . +ve hyperplane; $\vec{\omega} \cdot \vec{x} + b = 1$

For the -ve hyperplane ; $\vec{\omega} \cdot \vec{x} + b = -1$

$\rightarrow$ $0\omega_1 + 2\omega_2 + b = 1$ ————①

$2\omega_1 + 0\omega_2 + b = 1$ ————②

Sai Parthish. M
Co020230847

$- 2\omega_2 - 2\omega_3 + b = -1$

$\Rightarrow 2\omega_1 + 2\omega_2 - b = 1$ ——— ③

By ① in ③, and ② in ③

$(1 - b) + (1 - b) - b = 1$

$2 - 3b = 1$

$$\boxed{b = 1/3}$$ ——— ④

with ④ in ① and ②, we get

$$\boxed{\begin{array}{l} \omega_1 = 1/3 \\ \omega_2 = 1/3 \end{array}}$$

ⓑ margin $= \dfrac{2}{||\omega||} = \dfrac{2}{0.471} = \underline{\underline{4.242}}$

$||\omega|| = \sqrt{\left(\dfrac{1}{3}\right)^2 + \left(\dfrac{1}{3}\right)^2}$

$= \underline{\underline{0.471}}$

ⓒ $\omega_1 x_1 + \omega_2 x_2 + b$

$(1/3 \, x - 1) + (1/3 \, x - 1) + 1/3$

$= -1/3 - 1/3 + 1/3 = \left(\dfrac{-1/3}{0}\right)$

So, the above then 0, it will lie ~~above~~ below
the Hyperplane, Hence, it will be of
the ~~positive class~~ negative class.