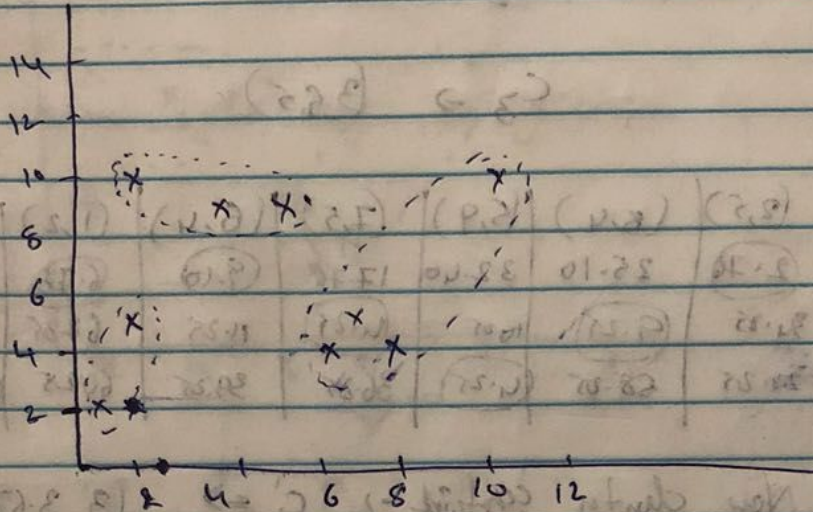Assign-4                                    Sashidhar D.
                                            sxd 173730

Part I:  Clustering.   Points: $(2,10)$ $(2,5)$, $(8,4)$, $(5,9)$, $(7,5)$, $(6,4)$,
                               $(1,2)$, $(4,9)$, $(10,10)$

(i)



Plot data to see appropriate Clusters.

(ii)  Beginning with $(2,5)$, $(5,8)$, $(4,9)$ as initial cluster centroids.

| | $(2,5)$ | $(5,8)$ | $(4,9)$ |
|---|---|---|---|
| $(2,10)$ | $\sqrt{25}$ | $\sqrt{13}$ | $\boxed{\sqrt{5}}$ |
| $(2,5)$ | $\boxed{0}$ | $\sqrt{18}$ | $\sqrt{20}$ |
| $(8,4)$ | $\sqrt{37}$ | $\boxed{25}$ | $\sqrt{41}$ |
| $(5,9)$ | $\sqrt{25}$ | $\boxed{\sqrt{1}}$ | $\sqrt{1}$ |
| $(7,5)$ | $\sqrt{25}$ | $\boxed{\sqrt{13}}$ | $\sqrt{25}$ |
| $(6,4)$ | $\boxed{\sqrt{17}}$ | $\sqrt{17}$ | $\sqrt{29}$ |
| $(1,2)$ | $\boxed{\sqrt{10}}$ | $\sqrt{52}$ | $\sqrt{58}$ |
| $(4,9)$ | $\sqrt{20}$ | $\sqrt{2}$ | $\boxed{\sqrt{0}}$ |
| $(10,10)$ | $\sqrt{89}$ | $\boxed{\sqrt{29}}$ | $\sqrt{37}$ |

$\Rightarrow$   $C_1 \to (2,5), (6,4), (1,2)$
      $C_2 \to (8,4), (5,9), (7,5), (10,10)$
      $C_3 \to (2,10), (4,9)$

⇒ New cluster centroids:-

$$C_1 \to (8/3, 11/3) = (3, 11/3)$$
$$C_2 \to (7.5, 7)$$
$$C_3 \to (3, 9.5)$$

|  | (2,10) | (2,5) | (8,4) | (5,9) | (7,5) | (6,4) | (1,2) | (4,9) | (10,10) |
|---|---|---|---|---|---|---|---|---|---|
| (3,3.67) $C_1$ | 41.66 | 2.76 | 25.10 | 32.40 | 17.76 | 9.10 | 6.78 | 29.40 | 89.60 |
| (7.5,7) $C_2$ | 39.25 | 34.25 | 9.25 | 10.25 | 4.25 | 11.25 | 67.25 | 16.25 | 15.25 |
| (3,9.5) $C_3$ | 1.25 | 21.25 | 58.25 | 4.25 | 36.25 | 39.25 | 60.25 | 1.25 | 49.25 |

New clusters:-

⇒ New cluster centroid →  $C_1' \Rightarrow (3, 3.67)$

$C_1' \to (2,5), (6,4), (1,2)$     $C_2' \to (8.33, 6.33)$

$C_2' \to (8,4), (7,5), (10,10)$     $C_3' \to (3.67, 9.33)$

$C_3' \to (2,10), (5,9), (4,9)$

⇒

|  | (2,10) | (2,5) | (8,4) | (5,9) | (7,5) | (6,4) | (1,2) | (4,9) | (10,10) |
|---|---|---|---|---|---|---|---|---|---|
| $C_1'$ | 41.06 | 2.76 | 25.10 | 32.40 | 17.76 | 8.60 | 6.78 | 29.4 | 84.06 |
| $C_2'$ | 53.33 | 41.03 | 5.53 | 18.21 | 3.53 | 10.85 | 72.47 | 25.87 | 16.25 |
| $C_3'$ | 3.24 | 21.5 | 47.13 | 1.87 | 29.8 | 33.54 | 60.85 | 2.17 | 49.32 |

New cluster :-   $C_1'' = (2,5), (6,4), (1,2)$

$C_2'' = (8,4), (7,5), (10,10)$

$C_3'' = (2,10), (5,9), (4,9)$

New Centroid:-

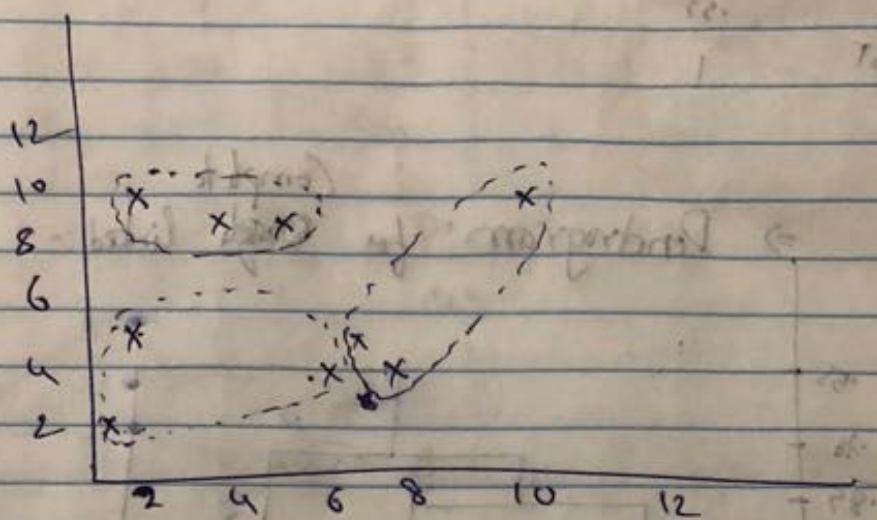⇒ $C_1'' \to (3, 3.67)$

$C_2'' \to (8.33, 6.33)$

$C_3'' \to (3.67, 9.33)$

∴ Converged with $C_1''$, $C_2''$, $C_3''$ as final cluster

⇒ Find cluster:



B)

|    | $P_1$ | $P_2$ | $P_3$ | $P_u$ | $P_5$ |
|----|------|------|------|------|------|
| $P_1$ | 1 | | | | |
| $P_2$ | .10 | 1 | | | |
| $P_3$ | .41 | .64 | 1 | | |
| $P_u$ | .55 | .47 | .44 | 1 | |
| $P_5$ | .35 | (.98) | .85 | .76 | 1 |

Complete ~~single~~ Link hierarchial clustery :-

⇒

|    | $P_1$ | $P_{25}$ | $P_3$ | $P_u$ |
|----|------|------|------|------|
| $P_1$ | 1 | | | |
| $P_{25}$ | .35 | 1 | | |
| $P_3$ | .41 | (.85) | 1 | |
| $P_u$ | .55 | .76 | .44 | 1 |

$\max (P_{25}, P_1)$

⇒ $\max ((P_1, P_2), (P_1, P_5))$

⇒

|    | $P_1$ | $P_{253}$ | $P_u$ |
|----|------|------|------|
| $P_1$ | 1 | | |
| $P_{235}$ | .41 | 1 | |
| $P_u$ | .55 | (.76) | 1 |

$\max (P_{235}, P_1)$

$\partial \max ((P_{25}, P_3), (P_3, P_1))$

|        | $P_1$ | $P_{2345}$ |
|--------|-------|------------|
| $P_1$  | 1     | .55        |
| $P_{2345}$ | .55 | 1        |

⇒ Dendrogram for ~~Single~~ Complete linked :-



Single ~~complete~~ linked hierarchical clustering :-

|        | $P_1$ | $P_2$ | $P_3$ | $P_4$ | $P_5$ |
|--------|-------|-------|-------|-------|-------|
| $P_1$  | 1     | .10   | .41   | .55   | .35   |
| $P_2$  |       | 1     | .64   | .47   | (.98) |
| $P_3$  |       |       | 1     | .64   | .85   |
| $P_4$  |       |       |       | 1     | .76   |
| $P_5$  |       |       |       |       | 1     |

|        | $P_1$ | $P_{25}$ | $P_3$ | $P_4$ |
|--------|-------|----------|-------|-------|
| $P_1$  | 1     | .10      | .41   | .55   |
| $P_{25}$ |     | 1        | (.64) | .47   |
| $P_3$  |       |          | 1     | .64   |
| $P_4$  |       |          |       | 1     |

$\min(P_{25}, P_1)$
$\Rightarrow \min((P_2, P_1), (P_5, P_1))$
$= \min(.10; .35)$

$\min(P_{25}, P_4)$
$= \min((P_2, P_4),$
$(P_5, P_4))$

$\min(P_{25}, P_3)$
$= \min((P_2, P_3), (P_5, P_3))$

⇒

|  | $P_1$ | $P_{235}$ | $P_4$ |
|---|---|---|---|
| $P_1$ | 1 |  |  |
| $P_{235}$ | .90 | 1 |  |
| $P_4$ | (.55) | .44 | 1 |

⇒

|  | $P_{14}$ | $P_{235}$ |
|---|---|---|
| $P_{14}$ | 1 | .44 |
| $P_{235}$ | .44 | 1 |

⇒ ∴ ~~Complete~~ dendrogram for Single linked:



Ⓠⓒ One-D data points: $\{6, 12, 18, 24, 25, 28, 30, 32, 48\}$.

ⓐ (i) $\{5, 7.5\}$ → assume 2 can of centroid of clusters

⇒ $C_1 \to (5) 6$

$C_2 \to (7.5) 12, 18, 24, 25, 28, 30, 42, 48$

Individual TSE⇒ $C_1 < 1$

$C_2 = (7.5 - 12)^2 + (7.5 - 18)^2 \cdots + (7.5 - 48)^2$

$= 4466$

(Total) TSE ⇒ 4467

② { 15, 25 }

$C_1 \Rightarrow \circled{15} \to 6, 12, 18$

$C_2 \Rightarrow \circled{25} \to 24, 30, 28, 42, 48$

$$TSE_{C_1} = \left( \sum (15 - C_i)^2 \right) + \left. \right\} \; TSE = 952$$

$$TSE_{C_2} = \left( \sum 25 - C_{2i})^2 \right) =$$

③ ⓑ   for set ① ⇒     $C_1 \to \circled{05}$   |   $C_2 \to 7.5$
                         6           |   12, 18, 24, 28, 25, 30, 42, 48

⇒ Center$_1$ = 6          |          Center$_2$ = 28.375

| | 6 | 28.375 |
|---|---|---|
| ⇒ 6 | ⓪ | 22.375 |
| 12 | ⑥ | 16.375 |
| 18 | 12 | ⑩.375 |
| 24 | 18 | ④.375 |
| 25 | 19 | ③.375 |
| 28 | 22 | ·375 |
| 30 | 24 | ①.625 |
| 42 | 36 | ⑬.625 |
| 48 | 42 | ⑲.625 |

New clusters :-

ⓑ $C_1' \to 6, 1\circled{0}$

Center$_1'$ = 9

$\circled{28.375}\, C_2' \to 18, 24, 25, 30, 42, 48$

Center$_2'$ = 30.714

| | 9 | 30.714 | New cluster : |
|---|---|---|---|
| ⇒ 6 | ③ | 24.714 | $9 \to C_1'' \to$ ~~29~~ 6, 18, 12 |
| 12 | ⑥ | 18.714 | $30.714 \to C_2'' \to 24, 25, 28, 30,$ |
| 18 | ⑨ | 12.714 | 42, 48 |
| 24 | 15 | ⑥.714 | |
| 25 | 16 | ⑤.714 | Center$_1''$ = 12 |
| 28 | 18 | ②.714 | Center$_2''$ = 32.8 |
| 30 | 21 | ·714 | |
| 42 | 33 | ⑪.28 | |
| 48 | 39 | ⑰.28 | |

|      | 12   | 32.8   |
|------|------|--------|
| 6    | (6)  | 26.8   |
| 12   | (0)  | 20.8   |
| 18   | (6)  | 26.8   |
| 24   | 12   | (8.8)  |
| 25   | 13   | (7.8)  |
| 30   | 18   | (2.8)  |
| 28   | 16   | (4.8)  |
| 42   | 20   | (9.2)  |
| 48   | 36   | (3.2)  |

New cluster :-

$$C_1''' \rightarrow (6, 12, 18)$$

$$C_2''' \rightarrow \; (24, 25, 28, 42, 30, 48)$$

Here,

$$C_1'' \equiv C_1'''$$

$$C_2'' \equiv C_2'''$$

∴ there converged.

Now consider Set ② :- and repeat above procedure with initial centroids at {15, 25}

Then, it converges at ⟹ $C_1''' = (6, 12, 18) \rightarrow 12$

$$C_2^{VI} = (24, 25, 28, 30, 42, 48)$$

$$\downarrow$$

$$32.8$$

∴ Both are stable solutions since they converged and also have same clustering.

ⓒ Two clusters produced by MIN (single clustering) link.

⟹



24   25   28   30  18  6   12

MIN
single
Linked ⟹



24    25    28    30    18    6    12    42    48

(d) MIN gives most natural clustering than K-means in this situation.

(e) K-means clustering depends on selection of initial centroids.

Part II :- Classification:-

(i) ① P(A=0|+) = 2/5          P(B=0|+) = 4/5
      P(A=1|+) = 3/5          P(B=1|+) = 1/5
      P(A=0|-) = 3/5          P(B=0|-) = 3/5
      P(A=1|-) = 2/5          P(B=1|-) = 2/5

                P(C=0|+) = 1/5
                P(C=1|+) = 4/5
                P(C=0|-) = 0
                P(C=1|-) = 1

②     Using Naive Bayes Calculate: $\quad P(A=1, B=1, C=0)$

$$P(+|A=1, B=1, C=0) = \frac{P(A=1, B=1, C=0|+) \; P(+)}{P(A=1, B=1, C=0)}$$

$$= \frac{P(A=1|+) \; P(B=1|+) \; P(C=0|+) \; P(+)}{K}$$

$$= (3/5)(1/5)(1/5)(1/2) / K = \frac{3}{250K} \quad Ⓐ$$

$$P(-|A=1, B=1, C=0) = \frac{P(A=1, B=1, C=0|-) \; P(-)}{K}$$

$$= \frac{P(A=1|-) \; P(B=1|-) \; P(C=0|-) \; P(-)}{K}$$

$$= (2/5)(4/5)(0)(1/2) / K = 0 \quad Ⓑ$$

$$\therefore \quad Ⓐ > Ⓑ \implies \text{The class of sample } P(A=1, B=1, C=0)$$
$$\text{is } ⊕$$

③   Using m-estimate approach $\implies n = 1/2, \quad m = 4$

$$P(A=1|+) = \left\{ \text{apply} \; \frac{\text{B2}}{\left(\frac{n_c + mp}{n+m}\right)} \right\} = \frac{3 + (4)(1/2)}{5 + 4} \quad \begin{cases} 5/9 \\ 4/9 \\ 3/9 \\ 4/9 \\ 3/9 \\ 2/9 \end{cases}$$

$$P(A=1|-)$$
$$P(B=1|+)$$
$$P(B=1|-)$$
$$P(C=0|+)$$
$$P(C=0|-)$$

④ repeat part ③ :-

$$P(A=1, B=1, C=0) \implies \text{which class?}$$

$$\implies P(+ \mid A=1, B=1, C=0) = P(+) \frac{\left( P(A=1 \mid +) \, P(B=1 \mid +) \, P(C=0 \mid +) \right)}{P(A=1, B=1, C=0)}$$

$$= (5/9)(3/9)(3/9)(1/2) / \kappa = \frac{0.30}{\kappa} \, Ⓐ$$

$$\implies P(- \mid A=1, B=1, C=0) = \frac{P(-) \left( P(A=1 \mid -) \, P(B=1 \mid -) \, P(C \mid -) \right)}{\kappa}$$

$$= (4/9)(4/9)(2/4)(1/2) / \kappa$$

$$= .219 / \kappa \, Ⓑ$$

∴ ⟹ Ⓐ > Ⓑ   The classification class of sample is +

⑤ Considering above two cases, when the conditional probability of one of the class is zero, without smoothing the contribution of all other probabilities is not taken into consideration.

    Say, then the contribution of other probabilities contribute too much to the given class, then the zero prob. value.

    Hence, it is better to consider all probabilities by smoothing because it given minimal probability weight to all the classes.

## (ii) Adaboosting:-

**Iteration① $x_2 \uparrow$**



$w_1 \to w_8$

$\Rightarrow w_i = 1/8$

$= .125$

Let $x_1 z .25 \Rightarrow 2$ points $(x_5, x_6)$ are misclassified.

$\Rightarrow \varepsilon_1 = 2 \times 1/8 = 0.25$

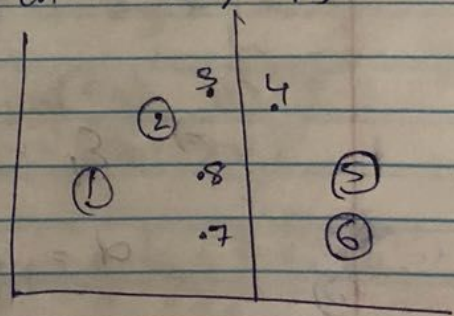$$\therefore \alpha = \frac{1}{2} \ln\left(\frac{1-0.25}{0.25}\right) = \ln\sqrt{3} = .55$$

$$Z_1 = 2\sqrt{(.25)(1-0.25)} = .866$$

**Iteration②:=** let boundary be at $x \geq .75$

2 points are misclassified $(x_5, x_6)$

$$D_2 = \frac{(1/8) e^{0.55}}{.866} = .25$$



weights of $x_1, x_2, x_3, x_4, x_7, x_8$ are correctly classified

$$\Rightarrow D_2 = \frac{(1/8) e^{-0.55}}{.866} = .083$$

$$\varepsilon_2 = 0.083 \times 2 = .167$$

$$\alpha = \frac{1}{2} \ln\left(\frac{1-\varepsilon_2}{\varepsilon_2}\right) = .804$$

$$Z_2 = .746$$

Iteration ③:

consider $y < .75$

$\Rightarrow$ $D_1$ (wrongly classified)

$(x_1, x_2)$

$$= \frac{(0.83)\, e^{-0.804}}{.746} = .249$$

$D_2$ (correctly classified) $\Rightarrow (x_3, x_4, x_8, x_7)$

$$= \frac{(0.083)\, e^{.804}}{.746} = 0.50$$
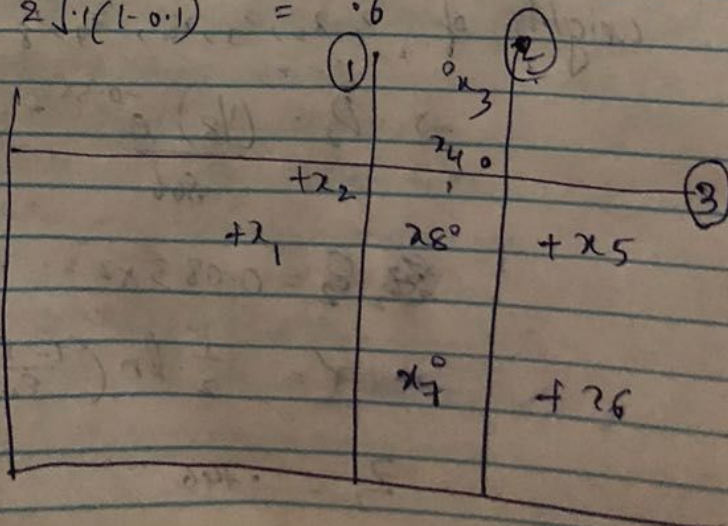
$D_3$ (correctly classified) $= (x_5, x_6)$

$$= \frac{(.25)\, e^{.804}}{.746} = .150$$

$\Rightarrow$ $\varepsilon_4 = .050 + .050 = .10$

$$\alpha = \frac{1}{2} \ln\left(\frac{1-0.1}{0.1}\right) = 1.099$$

$$Z_3 = 2\sqrt{.1(1-0.1)} = .6$$

@ Decision Stumps:

| t | $\varepsilon_t$ | $\alpha_t$ | $z_t$ | $w_t(1)$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ | $w_7$ | $w_8$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 025 | .55 | .886 | .125 | .125 | .125 | .125 | .125 | .125 | .125 | .125 |
| 2 | .167 | 804 | .746 | .083 | .083 | .083 | .083 | .25 | .25 | .083 | .083 |
| 3 | .10 | 1099 | .6 | .249 | .249 | .05 | .068 | .150 | .05 | .05 | .050 |

## ② Training Error of adaboost :-

Adaboost outperform a single decision stump because when single decision stump is used, the training error is ~~the~~ mostly greater than zero and also has high variance which is not the same using adaboost.

## Q⑤

$$X_1 \in \{a, b\} \qquad X_2 \in \{c, g, u, w\} \qquad X_3 \in \{k, sv\}$$

$$\text{parent entropy} \Rightarrow - \left( \sum \left( \frac{n_{yi}}{n_{y1i} + n_{y2i}} \right) \log \left( \frac{n_{yi}}{n_{y2i} + n_{y2i}} \right) \right)$$

$$= - \left( \frac{6}{11} \log \frac{6}{11} + \frac{5}{11} \log \frac{5}{11} \right) = .994$$

① Now consider sorting on basis of $X_1$

$$\Rightarrow \text{entropy} (S_a) = - \left( \frac{3}{5} \log 3/5 + 2/5 \log 2/5 \right)$$

$$= 0.97$$

$$S_b = - \left( 3/6 \log 3/6 + 3/6 \log 3/6 \right)$$

$$= 1$$

$$IG = .994 - \left( \left( \frac{5}{11} \right) (.97) + \left( \frac{6}{11} \right) 1 \right) = .0076$$

② Sort on basis of $X_2$:

$$S_C = - \left( \cancel{2} \ \frac{3}{5} \log 3/5 + \frac{2}{5} \log 2/5 \right)$$

$$= 0.97$$

$$S_{CO} = - \left( \frac{2}{4} \log 2/4 + \frac{2}{4} \log 2/4 \right)$$

$$= - \left( \log 1/4 \right) = \log 4.$$

$$S_g = - \left( 1 \cdot \log 1 \right)^{3} = 0 \quad ; \quad S_4 = 0.$$

$$IG = .994 - \left( (.97)\left(\frac{5}{11}\right) + (1.386)\left(\frac{4}{11}\right) \right)$$

$$= .994 - \left( \frac{(.4.85) +}{11} \boxed{5.545} \right)$$

$$= .601$$

③ Sort on basis of $X_3$:

$$S_K = - \left( \frac{3}{3} \log 3 \right) = 0$$

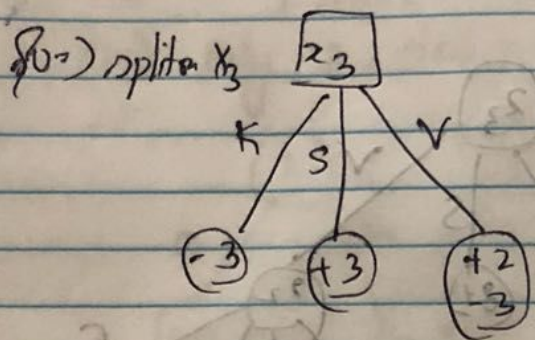$$S_V = - \left( \frac{2}{5} \log 2/5 + 3/5 \log 3/5 \right) = \frac{-1}{5} \left( \log 36 \right)$$

$$= 3.583$$

$$S_3 = 0$$

$$IG = .994 - \left( \left(\frac{5}{11}\right)(3.583) \right) = 1.6288$$

$$IG = .634$$

IG of ③ is highest

$S(.)$ splits $x_3$ | $z_3$ |



At Node ③ ⟹ entropy $= \dfrac{-2}{5}\log\dfrac{2}{5} + - \dfrac{3}{5}\log 3/5$
(parent now)

$= .97$

based on $z_1$

$S_a = -\dfrac{1}{2}\log 1/2 = (\log 1/2)\dfrac{1}{2} = 1$

$S_b = -\dfrac{2}{3}\log 2/3 - \dfrac{1}{3}\log 1/3 = .918$

$IG = .97 - \left(\dfrac{2}{5}(1) + \dfrac{3}{5}(.918)\right)$

$= 0.192$

based on $x_2$

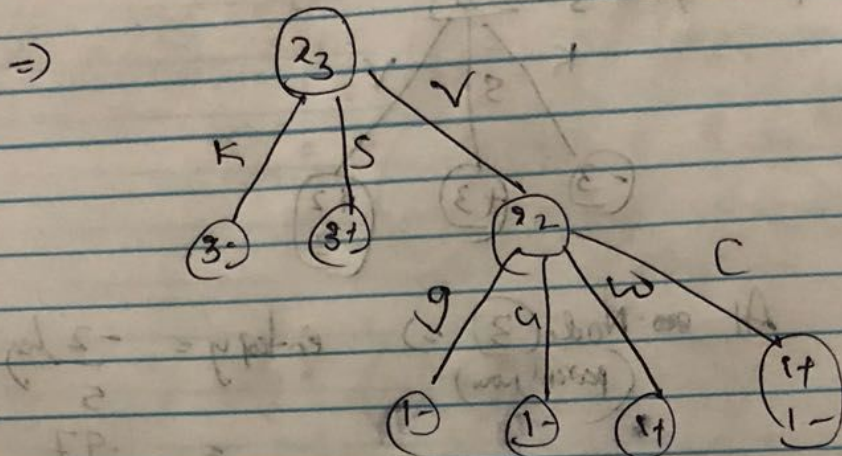$S_c = -\dfrac{1}{2}\log\dfrac{1}{2} - \dfrac{1}{2}\log 1/2 = 1$

$S_g = \dfrac{1}{1}\log 1/1 = 0$
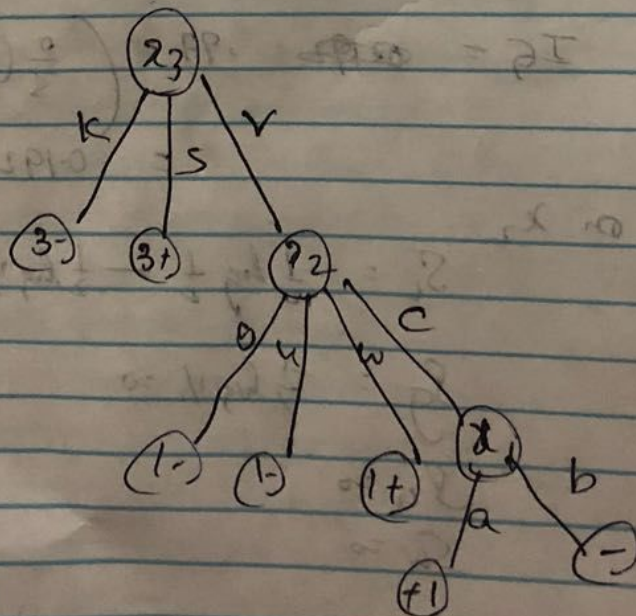
$S_u = 0$

$S_w = 0$

$IG = .97 - \left(\dfrac{2}{5}\times 1 + \dfrac{1}{5}\times 0 + 0 + 0\right) = .57$

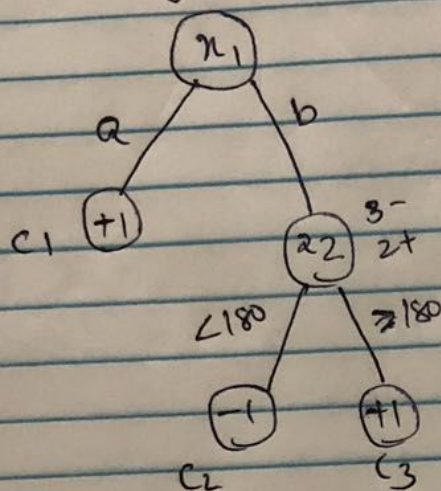Build stump based on $z_2$ ( $IG_{z_1} < IG_{z_2}$ )

$\Rightarrow$



$z_2 \to$ new parent

but no need to calculate IG because
we are left with $z_1$ only

$\Rightarrow$ Find Decision Tree :-

(b) For given data,
Decision Tree :-

(i) only 2 attribute node and 3 leaf nodes with
100% accuracy.



$c_i \rightarrow$ class nodes
$z_i \rightarrow$ attribute nodes.

(ii) Calculating accuracy of test sample using above model :-

From model

$\Rightarrow$ $x_1 = b$, $x_2 = 170$    expected value $\rightarrow -1$
Obtained value $\Rightarrow$ $Y = -1$    ✓

$x_1 = a$, $x_2 = 150$    expected value $\rightarrow +1$
Obtained $\Rightarrow Y = +1$    ✓

$x_1 = b$, $x_2 = 60$    expected value $\Rightarrow +1$
Obtained $\Rightarrow Y = -1$    ✗

Accuracy $= 2/3 = \cdot 667$

$\Rightarrow 66.7 \%$