Name :- Nandish. Bakulkumar. Bhatt.
Student ID :- 200441204

Writing Assignment : 3

## ANS: 1

**1(a)** Given that, Sigmoid function, $\sigma(\alpha) = \dfrac{1}{1+e^{-\alpha}}$

$$\therefore \quad \frac{d\sigma(\alpha)}{d\alpha} = \frac{d}{d\alpha}\left[\left(1+e^{-\alpha}\right)^{-1}\right]$$

$$= (-1)\left(1+e^{-\alpha}\right)^{-2} \cdot \frac{d}{d\alpha}\left(1+e^{-\alpha}\right)$$

$$\left[\because \frac{d[f(x)]^{n}}{dx} = n\left[f(x)\right]^{n-1} \cdot \frac{df(x)}{dx}\right]$$

$$\therefore \quad \frac{d\sigma(\alpha)}{d\alpha} = \frac{(-1)}{\left(1+e^{-\alpha}\right)^{2}} \cdot e^{-\alpha} \cdot (-1)$$

$$\left[\because \frac{d\,e^{f(x)}}{dx} = e^{f(x)} \cdot \frac{df(x)}{dx}\right]$$

$$\therefore \quad \frac{d\sigma(\alpha)}{d\alpha} = \frac{e^{-\alpha}}{\left(1+e^{-\alpha}\right)^{2}}$$

$$\frac{d\sigma(\alpha)}{d\alpha} = \left[\frac{1}{1+e^{-\alpha}}\right]\left[\frac{e^{-\alpha}}{1+e^{-\alpha}}\right]$$

$$\frac{d\,\sigma(\alpha)}{d\alpha} = \left[\frac{1}{1+e^{-\alpha}}\right]\left[1-1+\frac{e^{-\alpha}}{1+e^{-\alpha}}\right]$$

$$(\because \text{Adding and subtracting by } 1)$$

$$= \left[\frac{1}{1+e^{-\alpha}}\right]\left[1 - \left(1 - \frac{e^{-\alpha}}{1+e^{-\alpha}}\right)\right]$$

$$= \left[\frac{1}{1+e^{-\alpha}}\right]\left[1 - \left(\frac{1}{1+e^{-\alpha}}\right)\right]$$

$$\boxed{\therefore \quad \frac{d\,\sigma(\alpha)}{d\alpha} = \sigma(\alpha).\left[1 - \sigma(\alpha)\right]}$$

$\therefore$ Hence, proved.

**1(b)** Now, given that —

Negative Log – Likelihood (NLL) for logistic regression is given by :–

$$NLL = -\sum_{i=1}^{N} \log\left[\mu_i^{1(y_i=1)} \times (1-\mu_i)^{1(y_i=0)}\right]$$

$$= -\sum_{i=1}^{N}\left[y_i \log \mu_i + (1-y_i)\log(1-\mu_i)\right]$$

and $\mu_i = \text{sigm}(w^T x_i) = \dfrac{1}{1+e^{-w^T x_i}}$

$\rightarrow$ Now, $NLL = -\displaystyle\sum_{i=1}^{N}\left[ y_i \log \mu_i + (1-y_i)\log(1-\mu_i) \right]$

$\therefore \dfrac{dNLL(w)}{dw} = (-1)\left\{ \displaystyle\sum_{i=1}^{N}\left[ y_i \cdot \dfrac{d}{dw}(\log\mu_i) + (1-y_i)\dfrac{d}{dw}\log(1-\mu_i) \right] \right.$

$\left( \begin{array}{l} \because \ - \ y_i = \text{const. or independent of } w. \\ \dfrac{d}{dx}\displaystyle\sum_{i=1}^{2} i\,F(x) \ = \ \displaystyle\sum_{i=1}^{2} i\dfrac{d}{dx}[F(x)] \end{array} \right)$

$\therefore \dfrac{dNLL(w)}{dw} = (-1)\displaystyle\sum_{i=1}^{N}\left[ y_i \,\dfrac{1}{\mu_i}\,\dfrac{d\mu_i}{dw} + (1-y_i)\dfrac{1}{1-\mu_i}\,\dfrac{d}{dw}(1-\mu_i) \right]$

$\left( \begin{array}{l} \because \ \dfrac{d}{dx}[\log(F(x))] \ = \ \dfrac{1}{F(x)}\cdot\dfrac{dF(x)}{dx} \\ - \ \mu_i \text{ is a function of } w. \end{array} \right)$

$\therefore \dfrac{dNLL(w)}{dw} = (-1)\displaystyle\sum_{i=1}^{N}\left[ \dfrac{y_i}{\mu_i}\,\dfrac{d\mu_i}{dw} - \dfrac{1-y_i}{1-\mu_i}\cdot\dfrac{d\mu_i}{dw} \right]$

$= (-1)\displaystyle\sum_{i=1}^{N}\left[ \left( \dfrac{y_i}{\mu_i} - \dfrac{1-y_i}{1-\mu_i} \right)\dfrac{d\mu_i}{dw} \right]$

$$. = (-1) \sum_{i=1}^{N} \left[ \left( \frac{y_i - \mu_i y_i - \mu_i + \mu_i y_i}{\mu_i (1-\mu_i)} \right) \cdot \frac{d\mu_i}{dw} \right]$$

$$\frac{dNLL(w)}{dw} = (-1) \sum_{i=1}^{N} \left[ \left( \frac{y_i - \mu_i}{\mu_i (1-\mu_i)} \right) \cdot \frac{d\mu_i}{dw} \right] \underline{\hspace{1cm}} (1)$$

$\longrightarrow$ Now, as proved in $1(a)$, for $\sigma(\alpha) = \frac{1}{1+e^{-\alpha}}$,

$$\frac{d\sigma(\alpha)}{d\alpha} = \sigma(\alpha) \left[ 1 - \sigma(\alpha) \right] \underline{\hspace{1cm}} (2)$$

whereas,

$$\mu_i = sigm(w^T x_i) = \frac{1}{1+e^{-w^T x_i}}$$

$\mu_i$ is a function of $w^T x_i$ or "$wx$" as whole.

→ Let's say, $a = w^T x_i$.

So, $\mu_i$ is fuc$^n$ of $a$.

$\therefore \dfrac{d\mu_i}{dw} = \dfrac{d\mu_i}{da} \times \dfrac{da}{dw}$  (chain Rule)

$= \mu_i[1-\mu_i] \times \dfrac{da}{dw}$

$\left( \begin{array}{c} \because \\ \\ \therefore \end{array} \right.$ From previous $1(a)$,

$\dfrac{d.\sigma(\alpha)}{d\alpha} = \sigma(\alpha)[1-\sigma(\alpha)]$

$\therefore \dfrac{d\mu_i(a)}{da} = \mu_i(a)[1-\mu_i(a)]$ $\left. \vphantom{\begin{array}{c}\\\\\\\end{array}} \right)$

$= \mu_i[1-\mu_i] \cdot x_i$

$\left( \begin{array}{c} \because \\ \\ \\ \end{array} \right.$ $a = w^T x_i$

$\dfrac{da}{dw} = x_i \dfrac{d(w^T)}{dw}$

$= x_i$ $\left. \vphantom{\begin{array}{c}\\\\\\\end{array}} \right)$

$\therefore \dfrac{d\mu_i}{dw} = \mu_i[1-\mu_i] \times x_i$

$\underbrace{\hspace{4cm}}$ (3)

→ So, putting eq$^n$ (3) into eq$^n$ (1), we get —

Scanned by CamScanner

$$\therefore \frac{d\,NLL(w)}{dw} = (-1) \sum_{i=1}^{N} \left[ \left( \frac{y_i - \mu_i}{\mu_i(1-\mu_i)} \right) \mu_i(1-\mu_i)\, x_i \right]$$

$$= (-1) \sum_{i=1}^{N} \left[ (y_i - \mu_i)\, x_i \right]$$

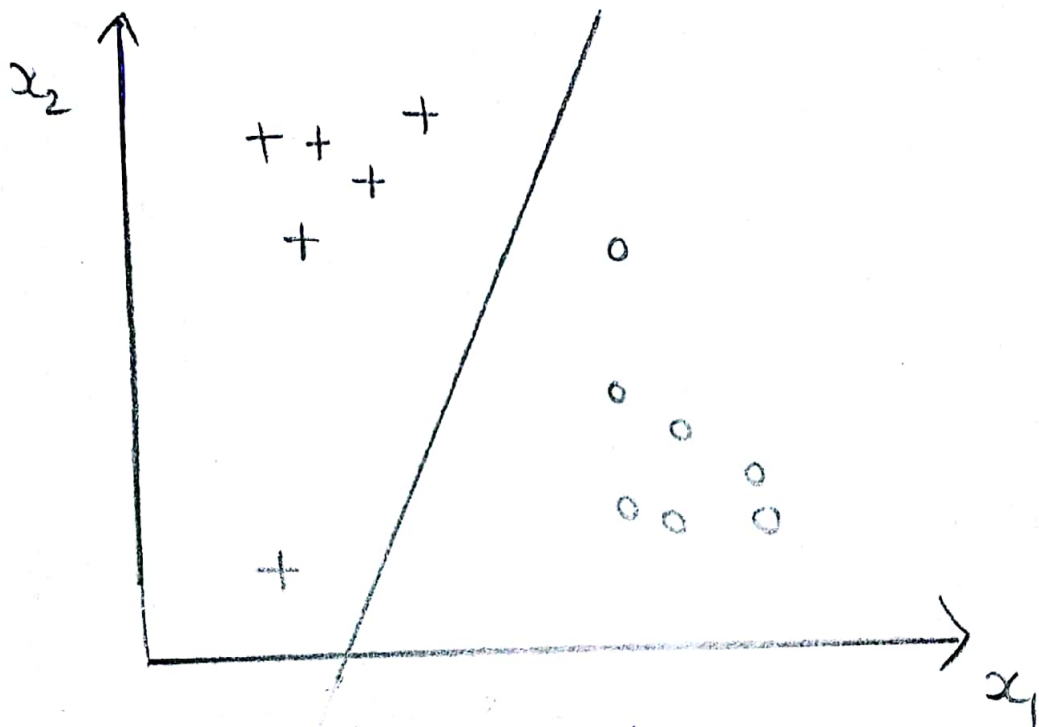$$\therefore \boxed{\frac{d\,NLL(w)}{dw} = \sum_{i=1}^{N} \left[ (\mu_i - y_i)\, x_i \right]}$$

# ANS: 2

## 2 (a)

Given that, $P(y = +1 | x_i, w) = \sigma(w_0 + w_1 x_1 + w_2 x_2)$

— Fitting model by max. likelihood, minimizing the

$$J(w) = -l(w, D_{train})$$

where, $l(w, D_{train})$ is log likehood on training set



→ It is observed from graph, data are separated linearly. Thus, logistic regression will find a line that will fit the data perfect.

→ From fig, with the given data and estimation of model, Classification error on training set = 0.

→ Line is actually not unique. (we can even change it normally, i.e., twitch it little bit)

It is indicating decision boundary.

## 2(b)

Now, by regularizing only $w_0$ and all others parameters are un-regularized. Then, the training error increases. So, Then the boundary will eventually go through origin.
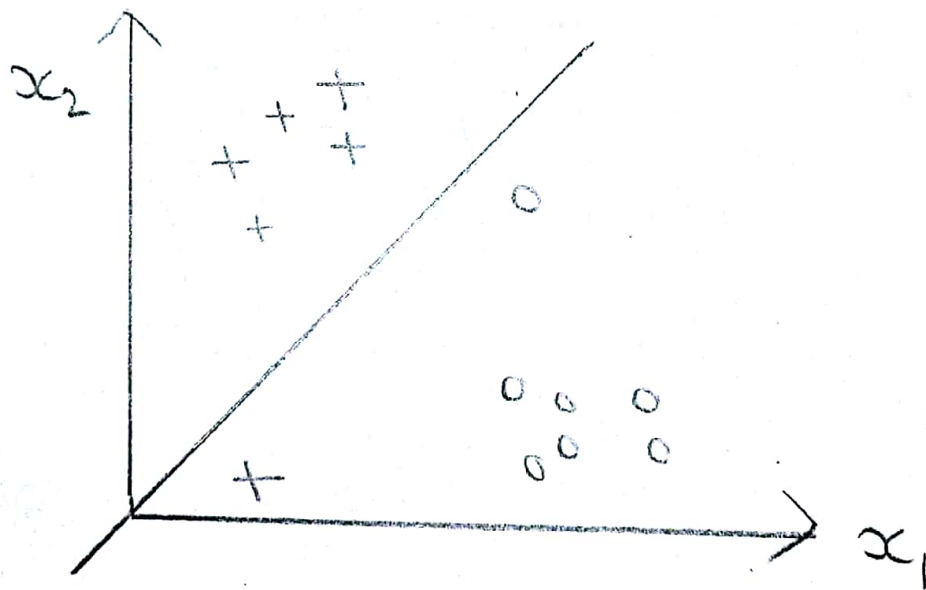
→ As, $w_0 = 0$, point $(0,0)$ origin will must be on the decision boundary.

Since, at that point,

$$\sigma(w_0 + w_1 x_1 + w_2 x_2) = \sigma(0) = 0.5$$

→ So, regularized logistic regression will find best decision boundary as plotted which passes through $(0,0)$.

→ It will make only one mistake on training data.



→ Moreover, due to this reason in regularized logistic or linear regression, we generally don't penalize the bias term.
(i.e, the weight which corresponding to feature) which is always $1$

**2(c)**
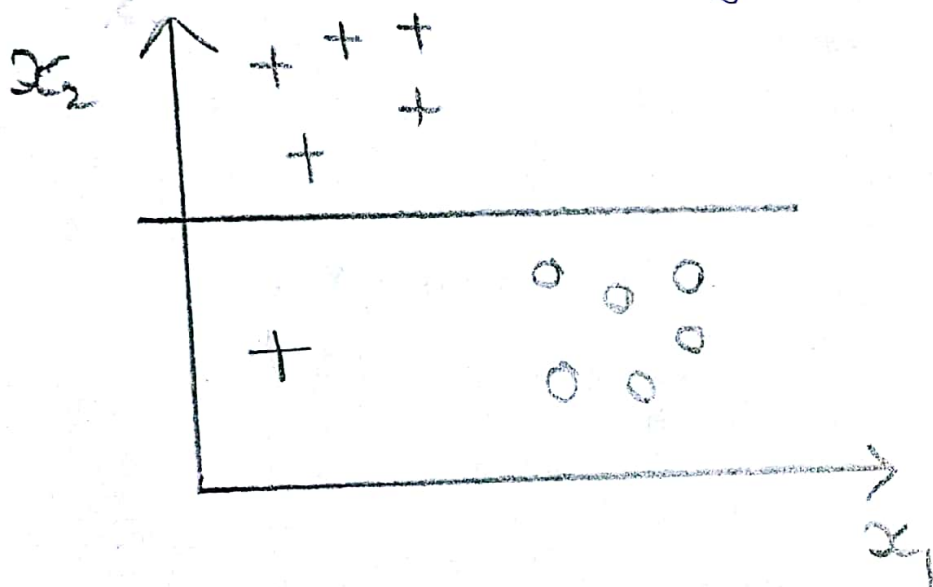
Minimizing, $J_1(w) = -l(w, D_{train}) + \lambda w_1^2$

So, by heavily regularizing $w_1$, the resulting boundary can rely less and less on values of $x_i$ and thus, it becomes more horizontal.

→ Also, training data can be seperated with zero training error with horizontal linear separator.

→ From fig., it is clear that,

~~those~~ Classification errors = 2
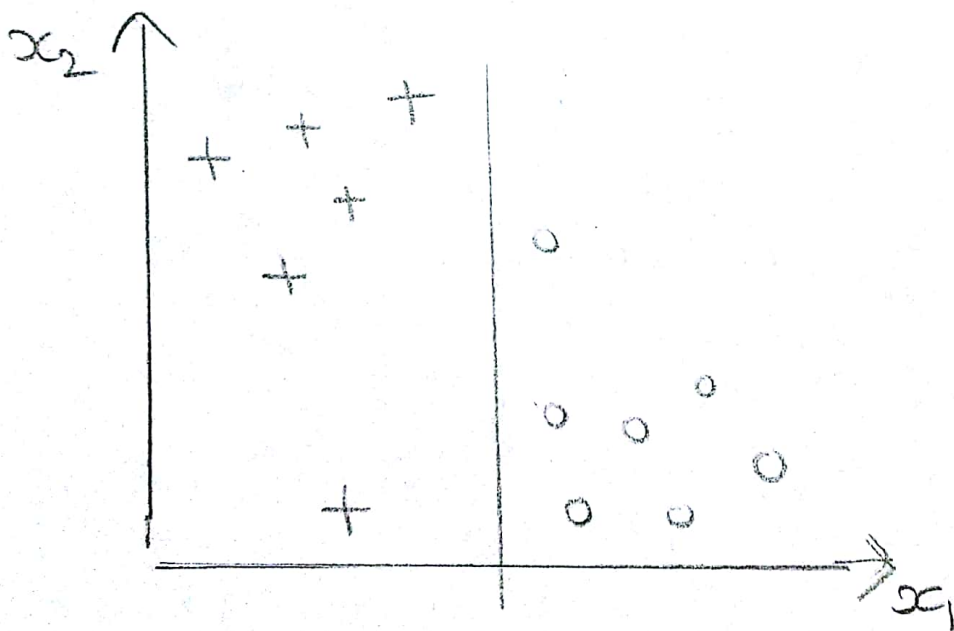on training set



$x_2$

**2(d)**

→ Now, if we regularize only $w_2$ parameter, then the resulting boundary will rely less and less on value of $x_2$ and therefore, it becomes more vertical.

→ So, the decision boundary will become vertical line.

→ Classification error on training set = 0

# ANS: 3

Now, T = 3.

| Credit | Term | Income | $\alpha$ | | | $y$ | $\hat{y}_1$ | $\hat{y}_2$ | $\hat{y}_3$ |
|--------|------|--------|----------|------|--------|------|-------|-------|-------|
| Good | 3 yrs. | High | $\frac{1}{9}$ | 0.5485 | 0.05 | 0.0245 | Safe | Safe | Safe | Safe |
| Good | 5 yrs. | Low | $\frac{1}{9}$ | 0.5485 | 0.05 | 0.1011 | Risky | Risky | Safe | Risky |
| Good | 5 yrs. | High | $\frac{1}{9}$ | 0.2078 | 0.207 | 0.1014 | Safe | Risky | Safe | Risky |
| Bad | 5 yrs. | High | $\frac{1}{9}$ | 0.5485 | 0.05 | 0.0245 | Risky | Risky | Risky | Risky |
| Bad | 3 yrs | Low | $\frac{1}{9}$ | 0.5485 | 0.05 | 0.1011 | Safe | Safe | Risky | Safe |
| Good | 5 yrs. | Low | $\frac{1}{9}$ | 0.5485 | 0.05 | 0.1011 | Risky | Risky | Safe | Risky |
| Bad | 3 yrs | High | $\frac{1}{9}$ | 0.2078 | 0.207 | 0.1014 | Risky | Safe | Risky | Safe |
| Bad | 5 yrs. | Low | $\frac{1}{9}$ | 0.5485 | 0.05 | 0.0245 | Risky | Risky | Risky | Risky |
| Good | 3 yrs. | High | $\frac{1}{9}$ | 0.5485 | 0.05 | 0.0245 | Safe | Safe | Safe | Safe |

**For t = 1**    N = 9
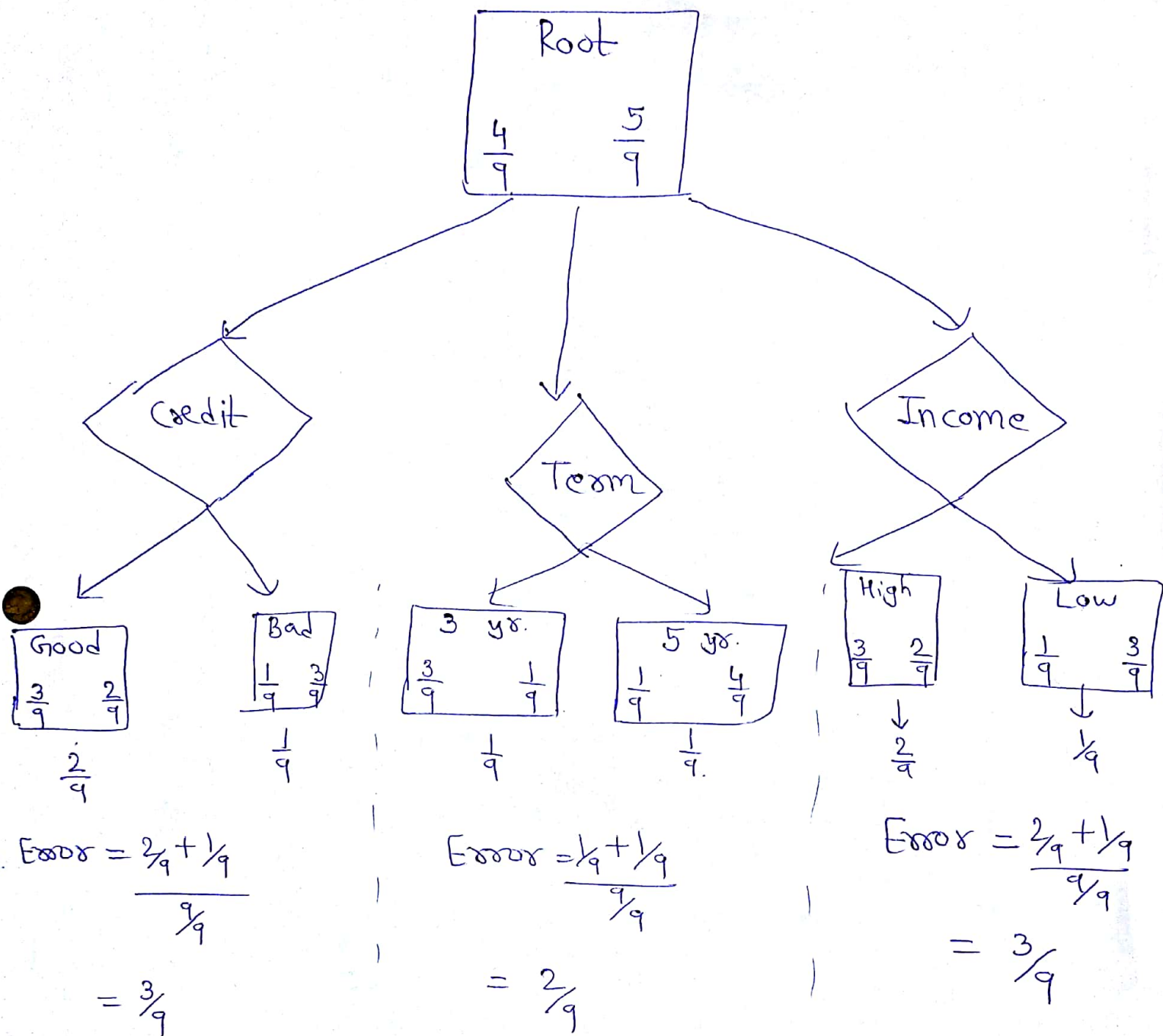
$$\therefore \alpha_i = \frac{1}{N} = \frac{1}{9}.$$

Firstly, same weights for all points.

So, making trees as follows:—

Root

$\frac{4}{9}$ $\frac{5}{9}$

Credit

Term

Income

Good
$\frac{3}{9}$ $\frac{2}{9}$
$\frac{2}{9}$

Bad
$\frac{1}{9}$ $\frac{3}{9}$
$\frac{1}{9}$

3 yr.
$\frac{3}{9}$ $\frac{1}{9}$
$\frac{1}{9}$

5 yr.
$\frac{1}{9}$ $\frac{4}{9}$
$\frac{1}{9}$

High
$\frac{3}{9}$ $\frac{2}{9}$
$\frac{2}{9}$

Low
$\frac{1}{9}$ $\frac{3}{9}$
$\frac{1}{9}$

$\therefore$ Error $= \dfrac{\frac{2}{9} + \frac{1}{9}}{\frac{9}{9}}$

$= \frac{3}{9}$

Error $= \dfrac{\frac{1}{9} + \frac{1}{9}}{\frac{9}{9}}$

$= \frac{2}{9}$

Error $= \dfrac{\frac{2}{9} + \frac{1}{9}}{\frac{9}{9}}$

$= \frac{3}{9}$

So, "Term" is winner for $t = 1$.

→ So, for Term → majorities are:

$\boxed{\begin{array}{cc} 3 \text{ yr.} \\ \frac{3}{9} & \frac{1}{9} \end{array}}$ → Safe

$\boxed{\begin{array}{cc} 5 \text{ yr.} \\ \frac{1}{9} & \frac{4}{9} \end{array}}$ → Risky.

Scanned by CamScanner

→ So, all with term -3 yr are predicted to be safe and all with term of 5 yr. are predicted to be risky.

→ Thus, weighted error $[f_1] = \dfrac{\sum\limits_{i=1}^{N} \alpha_i \,(y_i \neq \hat{y}_i)}{\sum\limits_{i=1}^{N} \alpha_i}$

$$= \frac{\tfrac{1}{9} + \tfrac{1}{9}}{1} \qquad = \frac{2}{9}$$

∴ $\hat{W}_1 = \dfrac{1}{2} \ln\left[\dfrac{1 - \text{weighted error}[f_1]}{\text{weighted error}[f_1]}\right] = \dfrac{1}{2}\ln\left[\dfrac{\tfrac{7}{9}}{\tfrac{2}{9}}\right]$

$$\boxed{\hat{W}_1 = 0.6263} \quad \leftarrow \text{coefficient}$$

→ Now, recomputing weights such that,

$$\alpha_i \leftarrow \begin{cases} \alpha_i \, e^{-\hat{w}_t} & , \text{ if } y_i = \hat{y}_i \\ \alpha_i \, e^{+\hat{w}_t} & , \text{ if } y_i \neq \hat{y}_i \end{cases}$$

→ So,
$$e^{-\hat{w_1}} = 0.5345 \quad \text{and} \quad e^{+\hat{w_1}} = 1.8706$$
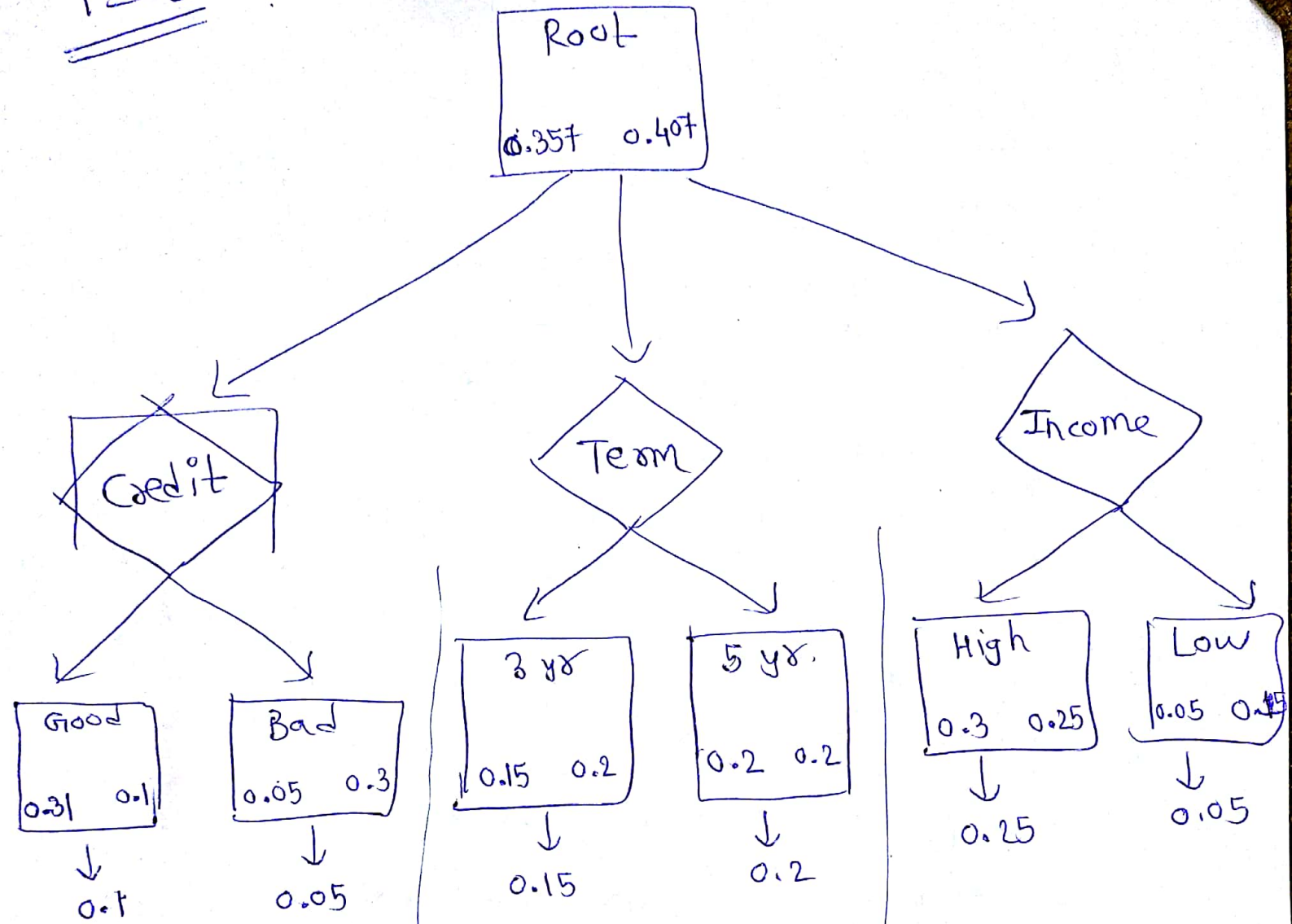
→ Since, $\alpha_i$ is same for all $t=1$.

∴ So, updating weights —

$$\alpha e^{-\hat{w_1}} = 0.059. \quad \text{and} \quad \alpha e^{+\hat{w_1}} = 0.207$$

$$\simeq 0.05 \qquad\qquad\qquad \simeq 0.2$$

→ Now, we have updated weights. Thus, this will form another decision stump for these new values of $\alpha$.

# T = 2

Root
| 0.357 | 0.407 |



**Credit** (crossed out)
- Good: | 0.31 | 0.1 | → 0.1
- Bad: | 0.05 | 0.3 | → 0.05

Error = $\dfrac{0.15}{0.764}$ = 0.1963

**Term**
- 3 yr: | 0.15 | 0.2 | → 0.15
- 5 yr: | 0.2 | 0.2 | → 0.2

Error = $\dfrac{0.35}{0.764}$ = 0.4581

**Income**
- High: | 0.3 | 0.25 | → 0.25
- Low: | 0.05 | 0.05 | → 0.05

Error = $\dfrac{0.30}{0.764}$ = 0.3926

→ So, Credit has lowest error, so it is winner for T = 2

majority Credit → Good | 0.31 | 0.1 | → Safe
are → Bad | 0.05 | 0.3 | → Risky

So,

∴ weighted error for $f_2$ $= \dfrac{0.05 + 0.05 + 0.05}{0.764}$

$$= 0.1963$$

∴ $\hat{w_2} = \dfrac{1}{2} \ln \left[ 1 - \dfrac{\text{weighted} - \text{error}}{\text{weighted} - \text{error}} \right]$

$$= \dfrac{1}{2} \ln \left[ \dfrac{1 - 0.1963}{0.1963} \right]$$
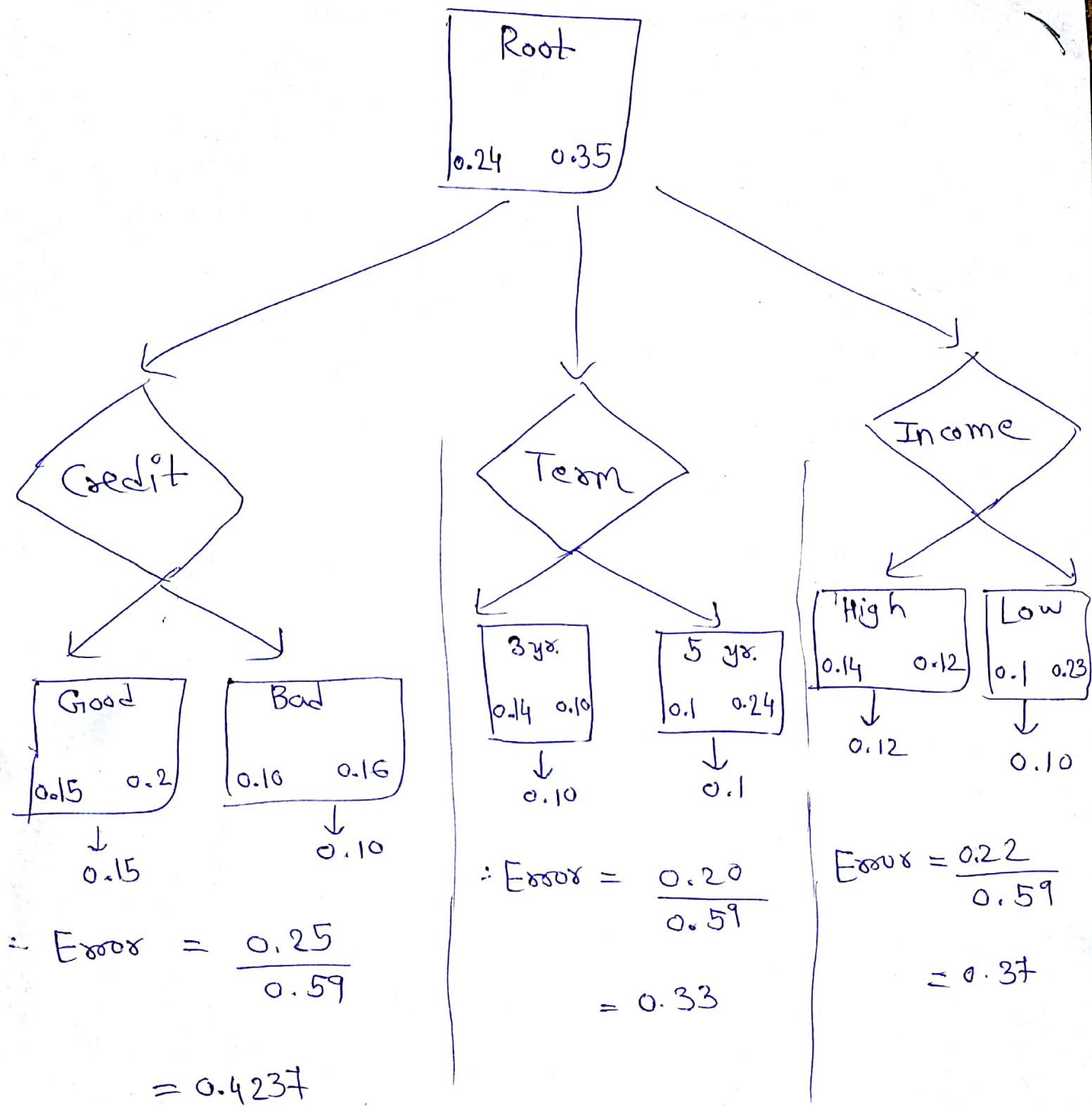
$$= 0.7047.$$

↳ Recomputing weights, so,

$e^{-\hat{w_2}} = 0.4942$      $e^{+\hat{w_2}} = 2.0232.$

$(y_i = \hat{y_i})$          $(y_i \neq \hat{y_i})$

→ So, we have updated the values of $\alpha$. So, we need to form one more decision stump of T=3.

∴ T=3

**Root**

$0.24 \quad 0.35$

**Credit**

- **Good**: $0.15 \quad 0.2 \rightarrow 0.15$
- **Bad**: $0.10 \quad 0.16 \rightarrow 0.10$

$\therefore$ Error $= \dfrac{0.25}{0.59} = 0.4237$

**Term**

- **3 yr.**: $0.14 \quad 0.10 \rightarrow 0.10$
- **5 yr.**: $0.1 \quad 0.24 \rightarrow 0.1$

$\therefore$ Error $= \dfrac{0.20}{0.59} = 0.33$

**Income**

- **High**: $0.14 \quad 0.12 \rightarrow 0.12$
- **Low**: $0.1 \quad 0.23 \rightarrow 0.10$

Error $= \dfrac{0.22}{0.59} = 0.37$

$\therefore$ So, lowest classification error $= 0.33$
  is of Term.

$\rightarrow$ So, Term is winner for $T = 3$.

So, for Term
majorities
are:

$\rightarrow$ 3 yr. $0.14 \quad 0.10 \rightarrow$ Safe

$\rightarrow$ 5 yr. $0.1 \quad 0.24 \rightarrow$ Risky

→ So, weighted error of $f_3 = \dfrac{0.1 + 0.1.}{0.59.}$

$$= 0.33$$

∴ $\hat{W}_3 = \dfrac{1}{2} \ln \left[ \dfrac{1 - \text{weighted\_error}}{\text{weighted\_error}} \right]$

$$= \dfrac{1}{2} \ln \left[ \dfrac{1 - 0.33}{0.33} \right]$$

∴ $\boxed{\hat{W}_3 = 0.3540}$

→ So, recomputing weights, $e^{-\hat{W}_3} = 0.7018$     $\boxed{e^{+\hat{W}_3} = 1.4247}$

| old weights | update weights |
|---|---|
| 0.0245 | $0.0245 \times 0.7018 = 0.017$ |
| 0.1011 | $0.1011 \times 0.7018 = 0.071$ |
| 0.1014 | $0.1014 \times 1.4247 = 0.14247$ |
| 0.0245 | $0.0245 \times 0.7018 = 0.017$ |
| 0.1011 | $0.1011 \times 0.7018 = 0.071$ |
| 0.1011 | $0.1011 \times 0.7018 = 0.071$ |
| 0.1014 | $0.1014 \times 1.4247 = 0.144$ |
| 0.0245 | $0.0245 \times 0.7018 = 0.017$ |
| 0.0245 | $0.0245 \times 0.7018 = 0.017$ |

$\rightarrow$ So, finally,

$$\hat{y} = \text{sign}\left(\sum_{t=1}^{T} \hat{w}_t \hat{f}_t(x)\right).$$

$$= \text{sign}\left(\sum_{t=1}^{3} \hat{w}_t \hat{f}_t(x)\right)$$

$$= \text{sign}\left[\hat{w}_1 \hat{f}_1(x) + \hat{w}_2 f_2(x) + \hat{w}_3 f_3(x)\right]$$

So, ~~$= \text{sign} [(0.6263)]$~~ for each datapoint, $\hat{y}$

$0.6263(+1) + 0.70(+1) + 0.35(+1) = 1.6763 \Rightarrow$ Safe

$0.6263(-1) + 0.70(+1) + 0.35(-1) = (-0.2763) \Rightarrow$ Risky

$0.6263(-1) + 0.70(+1) + 0.35(-1) = (-0.2763) \Rightarrow$ Risky

$0.6263(-1) + 0.70(-1) + 0.35(-1) = (-1.6763) \Rightarrow$ Risky

$0.6263(+1) + 0.70(-1) + 0.35(+1) = 0.2763 \Rightarrow$ Safe

$0.6263(-1) + 0.70(+1) + 0.35(-1) = (-0.2763) \Rightarrow$ Risky

$0.6263(+1) + 0.70(-1) + 0.35(+1) = 0.2763 \Rightarrow$ Safe

$0.6263(-1) + 0.70(-1) + 0.35(-1) = (-1.6763) \Rightarrow$ ~~Safe~~ Risky

$0.6263(+1) + 0.70(+1) + 0.35(+1) = +1.6763 \Rightarrow$ Safe

$\rightarrow$ Thus, we can say from true obs$^n$ and predicted obs$^n$ of target variable are same for all, except $3^{rd}$ and $7^{th}$, which has incorrect classification.