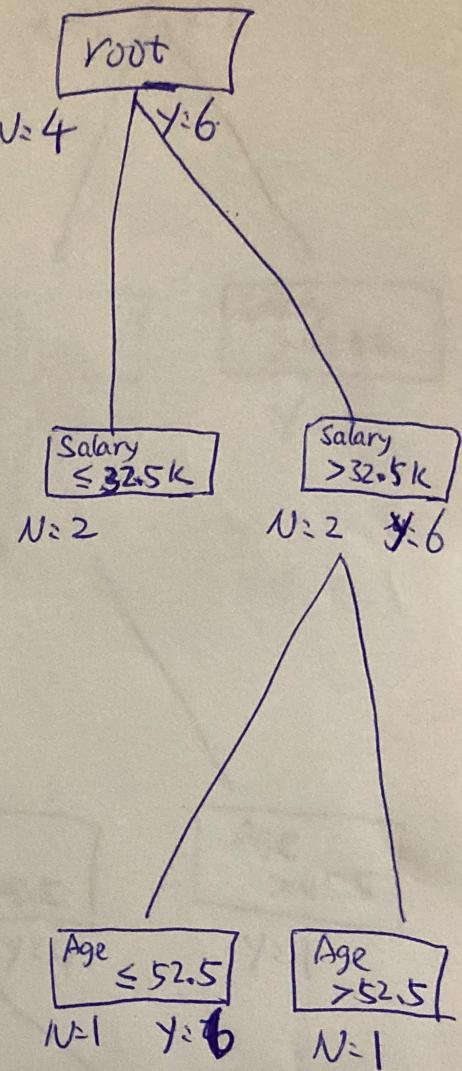


# PS #3.

1. (a) First split

Age:	22	23	24	25	32	43	48	52	52	53	N=4
Degree:	Y	N	Y	Y	Y	N	Y	Y	N	N	
Error:	$\frac{4}{9}$ =0.44	$\frac{1}{2}$ =0.875	$\frac{1}{3}$ =0.76	$\frac{1}{4}$ =0.15	$\frac{1}{5}$ =0.6	$\frac{2}{6}$ =0.83	$\frac{2}{7}$ =0.62	$\frac{1}{8}$ =0.33			

Salary:	25	27	38	40	44	48	52	65	77	110	
Degree:	N	N	Y	Y	N	Y	N	Y	Y	Y	
Error	$\frac{3}{9}$ =0.33	$\frac{2}{8}$ =0.25	$\frac{1}{7}$ <del>=0.14</del>	$\frac{3}{6}$ =0.62	$\frac{2}{5}$ =0.83	$\frac{3}{6}$ =0.6	$\frac{1}{4}$ =0.75	$\frac{3}{7}$ 0.43	$\frac{4}{8}$ 0.5	$\frac{4}{9}$ 0.44	



Second split

Age:	22	24	25	32	43	48	52	53		
Degree:	Y	Y	Y	Y	N	Y	Y	N		
Error:	$\frac{2}{7}$ =0.29	$\frac{2}{6}$ 0.33	$\frac{2}{5}$ 0.4	$\frac{2}{4}$ 0.5	$\frac{1}{3}$ 0.53	$\frac{1}{6}$ 0.67	$\frac{1}{4}$ 0.14			

Salary:	38	40	44	48	52	65	77	110		
Degree:	Y	Y	N	Y	N	Y	Y	Y		
Error	$\frac{2}{7}$ 0.29	$\frac{2}{6}$ 0.33	$\frac{1}{5}$ 0.53	$\frac{1}{4}$ 0.5	$\frac{2}{5}$ 0.4	$\frac{3}{6}$ 0.33	$\frac{2}{7}$ 0.29			

1.(a) cont.

third split.

Age	22	24	25	32	43	48	52
Degree	Y	Y	Y	Y	N	Y	Y
Error	$\frac{1}{6}$	$\frac{1}{5}$	$\frac{1}{4}$	$\frac{1}{3}$	$\frac{1}{5}$	$\frac{1}{6}$	.
=	0.17	0.2	0.25	0.33	0.2	0.17	.

Salary	38	40	44	48	65	77	110
Degree	Y	Y	N	Y	Y	Y	Y
Error	$\frac{1}{6}$	$\frac{1}{5}$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{5}$	$\frac{1}{6}$	.
=	0.17	0.2	0.33	0.25	0.2	0.17	.

Forth split.

Age	22	24	25	32	43	48	.
Degree	Y	Y	Y	Y	N	Y	.
Error	$\frac{1}{5}$	$\frac{1}{4}$	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{5}$	.	.
=	0.2	0.25	0.33	0.5	0.2	.	.

Salary	38	40	44	48	65	77	.
Degree	Y	Y	N	Y	Y	Y	.
Error	$\frac{1}{5}$	$\frac{1}{4}$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{5}$	.	.
=	0.2	0.25	0.33	0.25	0.2	.	.

Fifth split.

Age	22	24	25	32	43	.
Degree	Y	Y	Y	Y	N	.
Error	$\frac{1}{4}$	$\frac{1}{3}$	$\frac{1}{2}$	$\frac{1}{2}$	0	.

Age  $\leq 52.5$  (from page 1)

N=1, Y=6

Salary  $\leq 93.5k$

N=1, Y=5

Salary  $> 93.5k$

Y=1

Age  $\leq 45.5$

N=1, Y=4

Age  $> 45.5$

Y=1

Age  $\leq 37.5$

Y=4

Age  $> 37.5$

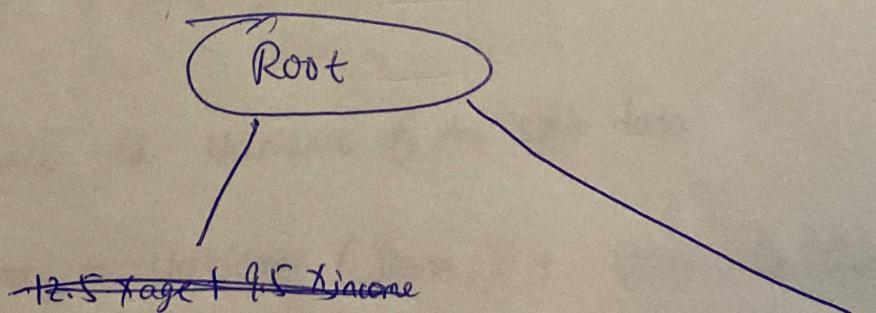
N=1 ~~at~~

1. (b)  
accuracy: for college : 90%  
~~~~~ Iris : ~~93.83%~~

(c). ~~data plot~~  
it's a linear separable data, so there is only 1 split.  
by using hw2 Logistic Regression on the  
~~dataset~~. if learned:

$$\alpha = \cancel{\text{something}} - 1.5$$

$$\beta = \cancel{\text{something}} 2.$$



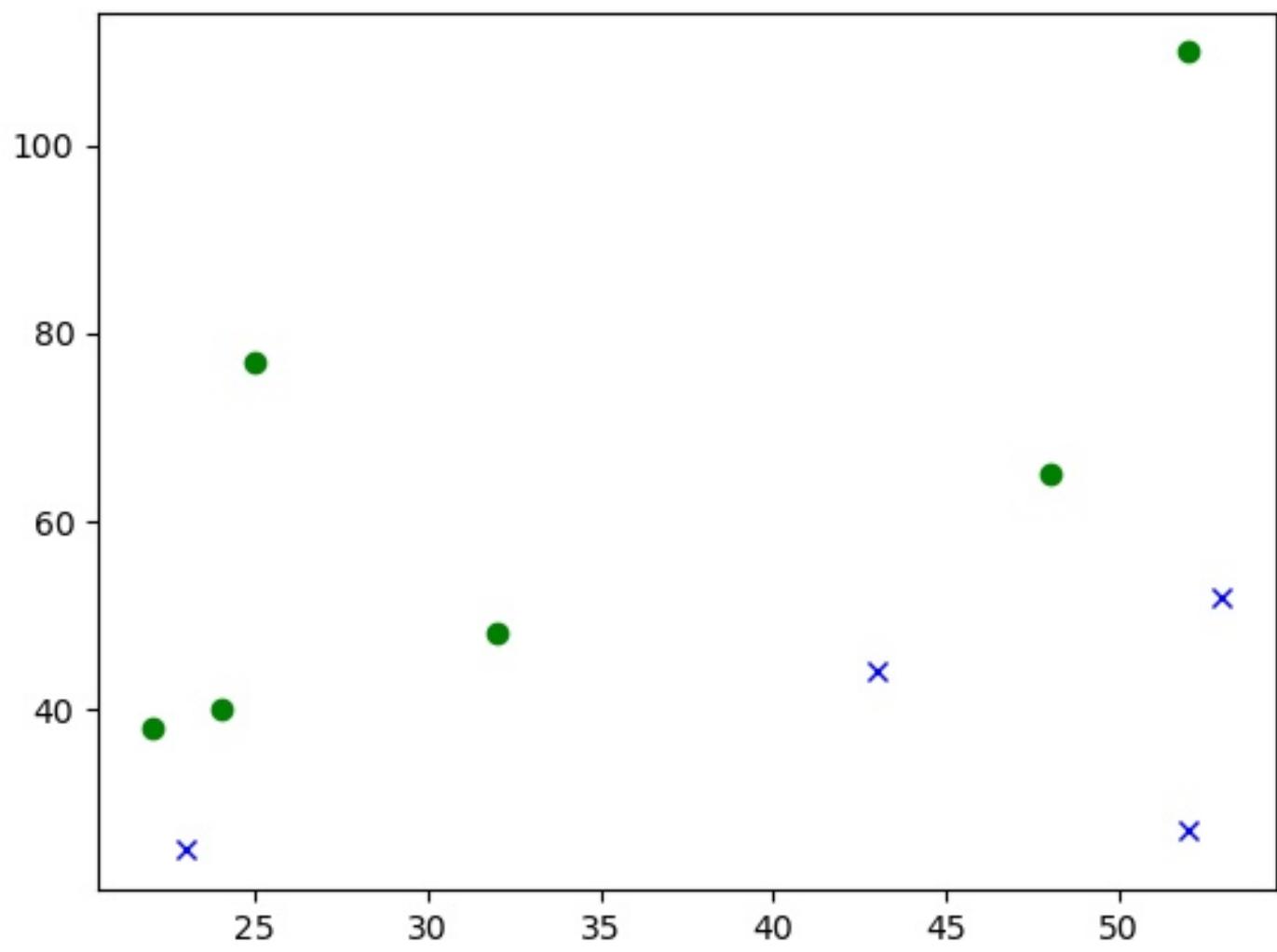
$$\text{Sign}\{-1.5x_{age} + 2x_{income} - 1\} < 0$$

$$N: 4$$

$$\text{Sign}\{-1.5x_{age} + 2x_{income} - 1\} \geq 0$$

$$Y: 6$$

Classification Error: 0.



(1)(d)

Pro: ① tree has less depth. less complicated.

② can cover more decision boundary. eg:

univariate tree can only have decision boundary parallel to the feature axis. multivariate tree can cover other linear boundary.

Con: ① ~~hard~~ hard to reason about. eg.

$-1.5 \times \text{age} + 2 \times \text{income}$  doesn't have any meaning.

② ~~doesn't fit the~~

can much easier to overfit.

(1)(e) i) for regression, the appropriate output would be the average of training data point's  $y$  value.

Line 3: return Leaf( $\frac{1}{n} \sum_{i=1}^n y^{(i)}$ )

ii) minimize the variance of the split data.

Line 8:  $\text{Error}_{\text{L}, \text{U}} = \text{Variance}(\text{Data}_1) + \text{Variance}(\text{Data}_2)$

$$\text{Variance} = \frac{\sum_{i=1}^N (y^{(i)} - \bar{y})^2}{N}, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y^{(i)}$$

2. (a)

$$G(R_m) = P_{m_1}(1-P_{m_1}) + P_{m_2}(1-P_{m_2}).$$

$$P_{m_1} = \frac{m_1}{m_1+m_2} = \begin{array}{l} m_1 \text{ is the \# of class 1} \\ m_2 \text{ is the \# of class 2.} \end{array}$$

$$P_{m_2} = \frac{m_2}{m_1+m_2} = 1 - P_{m_1}. \quad m = m_1 + m_2 \text{ is the \# of all classes}$$

$$G(R_m) = P_{m_1}P_{m_2} + P_{m_2}P_{m_1},$$

$$= 2P_{m_1}P_{m_2}$$

$$= 2 \frac{m_1 \cdot m_2}{m_1 + m_2}$$

$$= 2 \frac{m_1(m-m_1)}{m} = 2 \frac{m_2(m-m_2)}{m}$$

$$\underset{m_1, m_2}{\nabla G(R_m)} = \begin{bmatrix} \frac{2}{m}(m-2m_1) \\ \frac{2}{m}(m-2m_2) \end{bmatrix} = \begin{bmatrix} 2 - \frac{4m_1}{m} \\ 2 - \frac{4m_2}{m} \end{bmatrix}$$

$m \geq 0$ ,  $m_2 \geq 0$  and  $m_1$  and  $m_2$  cannot be 0 at the same time

$\therefore \nabla G(R_m)$  is strictly concave.

$$\nabla^2 G(R_m) = \begin{bmatrix} -\frac{4}{m} & 0 \\ 0 & -\frac{4}{m} \end{bmatrix} \iff -\frac{4}{m} \cdot I$$

negative definite

$\therefore G(R_m)$  is strictly concave.

the let  $P_1 = m_1$ ,  $P_2 = m_2$ ,  $(-\frac{m_1}{m_1+m_2}, \frac{m_2}{m_1+m_2}) \in (0, 1)$

$$G(tP_1 + (1-t)P_2) \geq \frac{m_1}{m_1+m_2} G(m_1) + \frac{m_2}{m_1+m_2} G(m_2)$$

2.(a) cont.

$$\cancel{G(M)} = \cancel{G(0.5M + 0.5M)} > 0.5 G(0.5M) + 0.5 G(0.5M)$$

$$G(M) = 2 \cdot G(0.5M) \quad \text{by definition.}$$

$$\cancel{G(m_1 + m_2)}$$

$$P_1 = m_1, P_2 = m_2, \epsilon = 0.5$$

$$\cancel{G(0.5M)}$$

$$= G(0.5m_1 + 0.5m_2) > 0.5 G(m_1) + 0.5 G(m_2)$$

$$\underline{2} \underline{G(0.5m_1 + 0.5m_2)} > G(m_1) + G(m_2).$$

$$\hookrightarrow = 2G(0.5M)$$

$$= G(M)$$

$$\therefore G(M) > G(m_1) + G(m_2)$$

2.(b) ① when  $G(m) = 0$ : loss is 0, all children are in the same class, so  $G(m_1)$  will be 0  
 $G(m_2)$

② When  $P_1 = P_2$ , so the children region has the exact same proportion of classes examples as the parent.

because we aim to reduce the loss, we always looks for a split that has smaller loss than the parent. therefore, it will converge

2(c) let  $M$  be the root,  $N_{M_1}$  be # of class 1 in  $M$   
 $N_{M_2} \swarrow \searrow \dots \swarrow \searrow$

assume  $M$  split at some feature and result in  $X, Y$ , child node

let  $N_{X_1}$  be # of class 1 in  $X$

$$\begin{array}{ccccccc} N_{X_2} & - & - & - & 2 & - & - \\ N_{Y_1} & - & - & - & 1 & - & Y \\ N_{Y_2} & - & - & - & 2 & - & Y \end{array}$$

further, assume  $N_{M_1} \leq N_{M_2}$

by definition:

$$L(M) = \frac{N_{M_1}}{N_M}, \quad N_{M_1} + N_{M_2} = N_M$$

$$N_{X_1} + N_{Y_1} = N_{M_1}, \quad N_{X_2} + N_{Y_2} = N_{M_2}, \quad N_X + N_Y = N_M$$

$$L(X) + L(Y) = \frac{\min(N_{X_1}, N_{X_2})}{N_X} + \frac{\min(N_{Y_1}, N_{Y_2})}{N_Y}$$

$$= \frac{N_Y \min(N_{X_1}, N_{X_2}) + N_X \cdot \min(N_{Y_1}, N_{Y_2})}{N_X \cdot N_Y}$$

$$= \frac{N_X \min(N_{Y_1}, N_{Y_2}) + (N_M - N_X) \min(N_{X_1}, N_{X_2})}{N_X \cdot (N_M - N_X)}$$

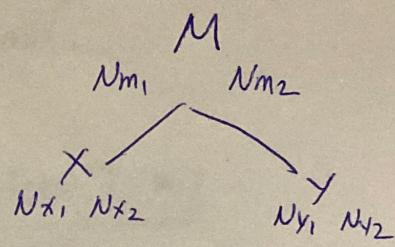
$$L(M) = \frac{N_{M_1}}{N_M} = \frac{N_{X_1} + N_{Y_1}}{N_M}$$

$$\text{let } L(M) = L(X) + L(Y)$$

$$\frac{N_{X_1} + N_{Y_1}}{N_M} = \frac{N_M \cdot \min(N_{X_1}, N_{X_2}) - [\min(N_{X_1}, N_{X_2}) - \min(N_{Y_1}, N_{Y_2})] \cdot N_X}{N_X \cdot (N_M - N_X)}$$

$$N_{M_1} \cdot N_X (N_M - N_X) = N_M^2 \cdot \min(N_{X_1}, N_{X_2}) - N_M \cdot N_X \cdot [\min(N_{X_1}, N_{X_2}) - \min(N_{Y_1}, N_{Y_2})]$$

$$N_{M_1} \cdot N_X N_M - N_{M_1} N_X^2 = N_M^2 \cdot \min(N_{X_1}, N_{X_2}) - N_M \cdot N_X [\min(N_{X_1}, N_{X_2}) - \min(N_{Y_1}, N_{Y_2})]$$



2.(c) cont.

set:  $N_m \cdot N_x = N_m^2 \cdot \min(N_{x_1}, N_{x_2})$

~~$N_m \cdot N_x^2 = N_m \cdot N_x \cdot [\min(N_{x_1}, N_{x_2}) - \min(N_{y_1}, N_{y_2})]$~~

let  $L(x,y) = L(x) + L(y)$

$$\frac{N_m}{N_m} = \frac{\min(N_{x_1}, N_{x_2})}{N_x} + \frac{\min(N_{y_1}, N_{y_2})}{N_y}$$

so the proportion doesn't need to be exactly the same.  
as long the sum adds up to be equal.

2.(d)

$$\rho = \frac{\text{cov}(T_i(x), T_j(x))}{\sqrt{\text{Var}(T_i(x))} \cdot \sqrt{\text{Var}(T_j(x))}} = \frac{\text{cov}(T_i(x), T_j(x))}{\sigma^2}$$

$$\text{cov}(T_i(x), T_j(x)) = \rho \cdot \sigma^2$$

$$\text{Var}\left(\frac{1}{B} \sum_{i=1}^B T_i(x)\right) = \frac{1}{B^2} \text{Var}\left(\sum_{i=1}^B T_i(x)\right) = \frac{1}{B^2} \left[ \sum_{i=1}^B \text{Var}(T_i(x)) + 2 \sum_{i < j} \text{cov}(x_i, x_j) \right]$$

$$= \frac{1}{B^2} \left[ B \cdot \sigma^2 + 2 \sum_{i < j} \rho \sigma^2 \right]$$

$$\hookrightarrow = \underbrace{\text{cov}(x_1, x_2) + \dots + \text{cov}(x_1, x_B)}_{B} + \text{cov}(x_2, x_3) + \dots + \text{cov}(x_2, x_B) + \text{cov}(x_{B-1}, x_B) \quad \} B-1$$

∴

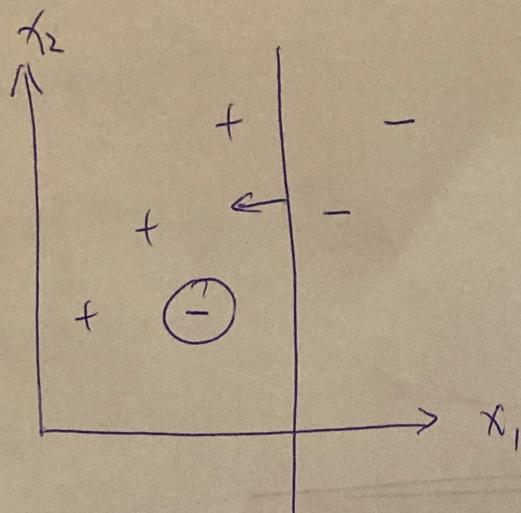
$$= \frac{1}{B^2} \left[ B \sigma^2 + 2 \cdot \frac{B(B-1)}{2} \cdot \rho \sigma^2 \right]$$

$$= \frac{\sigma^2}{B} + \frac{B-1}{B} \rho \sigma^2$$

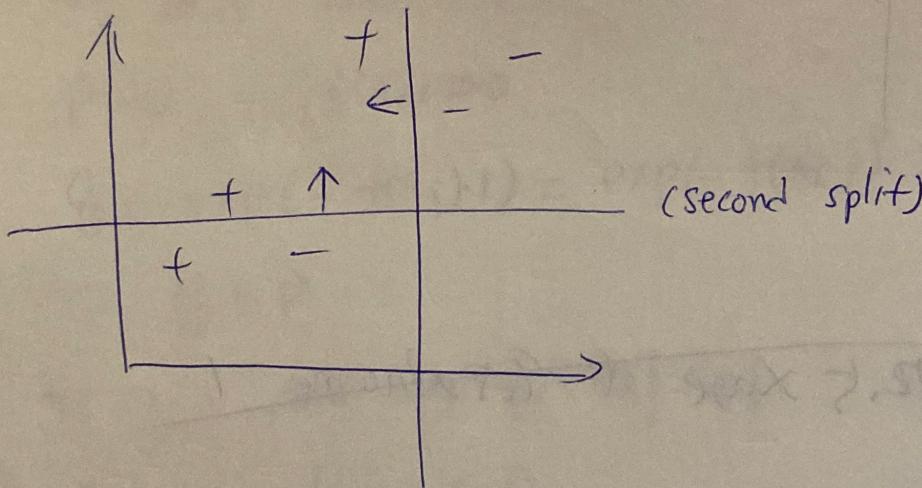
$$= \rho \sigma^2 + \frac{1-\rho}{B} \sigma^2$$

3.

(a)



(b)



$$(c) \quad \hat{w}_2 < \hat{w}_1 < 0.5$$

Weighted error of  $f_2$  is smaller

than  $f_1$  ~~an~~

$$\text{so } \hat{w}_2 < \hat{w}_1$$

4. a)

$$1\{f(x_i) \neq y_i\}$$

$$= 1(f(x_i) \cdot y_i > 0)$$

for a data point  $i$

$$\exp(-f(x_i) \cdot y_i) \stackrel{\Delta}{=} q_i$$

$$1(f(x_i) \cdot y_i > 0)$$

$$= 1(\text{sign}\{f(x_i)\} \cdot y_i > 0) \stackrel{\Delta}{=} P$$

if  $y_i = 1$  and  $\text{sign}\{f(x_i)\} = 1$

$$\Rightarrow P=1, f(x_i) \geq 0$$

$$q_i = \exp(-f(x_i) \cdot 1) = \exp(-f(x_i)) > 0$$

$$\therefore q_i > P.$$

if  $y_i = 1$  and  $\text{sign}\{f(x_i)\} = -1$

$$\Rightarrow P=0, f(x_i) < 0$$

$$q_i = \exp(-f(x_i) \cdot 1) \geq \exp(0) = 1$$

$$\therefore q_i \geq P$$

if  $y_i = -1$  and  $\text{sign}\{f(x_i)\} = -1$

$$\Rightarrow P=0, f(x_i) < 0$$

$$q_i = \exp(-f(x_i) \cdot -1) > 0 \quad \text{by definition of exp}$$

$$\therefore q_i > P$$

if  $y_i = -1$  and  $\text{sign}\{f(x_i)\} = 1$

$$\Rightarrow P=1, f(x_i) > 0$$

$$q_i = \exp(-f(x_i) \cdot -1) \geq \exp(0) = 1 = P$$

4.(a) Continue.

$\therefore$  for any  $i$ .  $P \leq q$

$$\therefore \frac{1}{n} \sum_{i=1}^n P_i \leq \frac{1}{n} \sum_{i=1}^n q.$$

4.(b)

$$Z_t = \sum_{i=1}^n a_{i,t} \exp(-\hat{w}_t^\top f_t(x_i) y_i)$$

$$(let \quad \exp(-\hat{w}_t^\top f_t(x_i) y_i) \triangleq \exp_{i,t}(c))$$

$$Z_t = \sum_{i=1}^n a_{i,t} \cdot \exp_{i,t}(c)$$

$$\begin{aligned} \prod_{t=1}^T Z_t &= Z_1 \cdot Z_2 \cdot \dots \cdot Z_{t-1} \cdot Z_t \\ &= Z_1 \cdot Z_2 \cdot \dots \cdot Z_{t-1} \cdot \left( \sum_{i=1}^n a_{i,t} \exp_{i,t}(c) \right) \\ &= Z_1 \cdot Z_2 \cdot \dots \cdot \cancel{Z_{t-1}} \cdot \left( \sum_{i=1}^n \frac{a_{i,t-1} \cdot \exp_{i,t-1}(c)}{\cancel{Z_{t-1}}} \cdot \exp_{i,t}(c) \right) \end{aligned}$$

$$\Rightarrow \prod_{t=1}^T Z_t = \sum_{i=1}^n a_{i,0} \cdot \exp_{i,1}(c) \cdot \exp_{i,2}(c) \cdot \dots \exp_{i,T}(c)$$

$$a_{i,0} = \frac{1}{n} \# \text{as the unweighted dataset}$$

$$\therefore \prod_{t=1}^T Z_t = \sum_{i=1}^n \frac{1}{n} \exp_{i,1}(c) \cdot \exp_{i,2}(c) \dots \exp_{i,T}(c).$$

$$= \frac{1}{n} \sum_{i=1}^n \exp(-\hat{w}_1^\top f_1(x_i) y_i - \hat{w}_2^\top f_2(x_i) y_i - \dots - \hat{w}_T^\top f_T(x_i) y_i)$$

$$= \frac{1}{n} \sum_{i=1}^n \exp(-f(x_i) \cdot y_i)$$

4(c) (i)

$$\frac{\partial Z_t}{\partial \varepsilon_t} = (1 - \varepsilon_t) \cdot (-\exp(-\hat{w}_t)) + \varepsilon_t \cdot \exp(\hat{w}_t)$$

$$\text{Set } \frac{\partial Z_t}{\partial \varepsilon_t} = 0$$

$$(1 - \varepsilon_t) \cdot \exp(-\hat{w}_t) = \varepsilon_t \exp(\hat{w}_t)$$

$$\exp(\hat{w}_t)^2 = \frac{1 - \varepsilon_t}{\varepsilon_t}$$

$$\exp(\hat{w}_t) = \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$$

$$Z_t^{\text{opt}} = (1 - \varepsilon_t) \exp(-\hat{w}_t) + \varepsilon_t \exp(\hat{w}_t)$$

$$= \frac{1 - \varepsilon_t}{\sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}} + \varepsilon_t \cdot \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$$

$$= \sqrt{(1 - \varepsilon_t)(\varepsilon_t)} + \sqrt{\varepsilon_t(1 - \varepsilon_t)}$$

$$= 2\sqrt{\varepsilon_t(1 - \varepsilon_t)}$$

(c) (ii) by definition  $0 < \varepsilon_t \leq 1$

$$\Rightarrow 0 < \frac{1}{2} - r_t \leq 1$$

$$\Rightarrow -\frac{1}{2} \leq r_t < \frac{1}{2}$$

$$\Rightarrow \frac{1}{4} > r_t^2 \geq 0 \Rightarrow 1 - 4r_t^2 \geq 0$$

$$\Rightarrow 0 \leq 1 - 4r_t^2 \leq 1$$

$$Z_t = 2\sqrt{\varepsilon_t(1 - \varepsilon_t)} = 2\sqrt{\frac{1}{4} - r_t^2} = \sqrt{1 - 4r_t^2}$$

$$\log(Z_t) = \frac{1}{2} \log(1 - 4r_t^2) \leq \frac{1}{2} \cdot (-4r_t^2) = -2r_t^2$$

$$\therefore Z_t \leq \exp(-2r_t^2)$$

4.(c) (iii)

$$\begin{aligned} \mathcal{E}_{\text{training}} &\leq \frac{1}{n} \sum_{i=1}^n \exp(-f(x_i)g_i) \\ &= \prod_{t=1}^T Z_t \\ &= Z_1 \cdot Z_2 \cdot \dots \cdot Z_T \\ &\leq \exp(-2r_1^2) \cdot \exp(-2r_2^2) \cdot \dots \cdot \exp(-2r_T^2) \\ &= \exp(-2r_1^2 - 2r_2^2 - \dots - 2r_T^2) \\ \because r_t > r \text{ for all } t \text{ and } \delta > 0 \\ \Rightarrow r_t^2 > r^2 \\ \Rightarrow -r_t^2 < -r^2 \\ \Rightarrow \underbrace{\exp(-2r^2 - 2r^2 - \dots - 2r^2)}_T \\ &= \exp(-2Tr^2) \end{aligned}$$

for  $T \geq 1$

when  $T=0$ ,  $\mathcal{E}_{\text{training}} \leq 1 = \exp(0)$

$\therefore \mathcal{E}_{\text{training}} \leq \exp(-2Tr^2)$

(5) (a) for a single data point.  
 let  $a_1 = w^{[1]T} \cdot x + b_1$

$$(b_1 \triangleq \begin{bmatrix} w_{01}^{[1]} \\ w_{02}^{[1]} \\ w_{03}^{[1]} \end{bmatrix})$$

$$h = \sigma(a_1)$$

~~$$\text{let } a_2 = w^{[2]T} \cdot h + b_2 \quad (b_2 \triangleq w_0^{[2]})$$~~

$$o = \sigma(a_2)$$

$$\frac{\partial o}{\partial a_2} = \sigma(a_2) \cdot (1 - \sigma(a_2))$$

$$\frac{\partial a_2}{\partial h_2} = w_2^{[2]} \quad \frac{\partial a_2}{\partial h_1} = w_1^{[2]}, \quad \frac{\partial a_2}{\partial b_2} = w_0^{[2]}$$

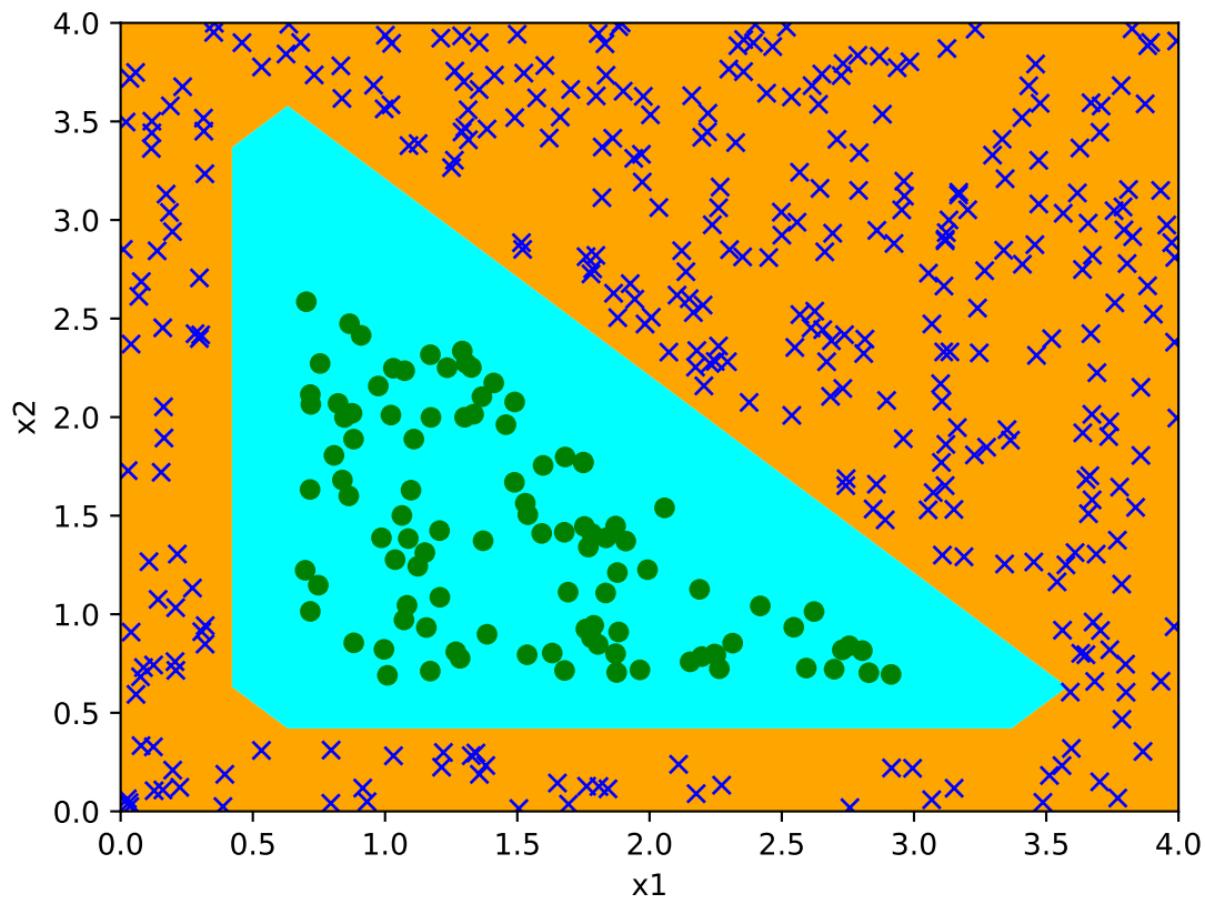
$$\frac{\partial h_2}{\partial a_1} = \sigma(a_1)(1 - \sigma(a_1))$$

$$\frac{\partial h}{\partial a_1} = \begin{bmatrix} \sigma(a_{11})(1 - \sigma(a_{11})) \\ \sigma(a_{12})(1 - \sigma(a_{12})) \\ \sigma(a_{13})(1 - \sigma(a_{13})) \end{bmatrix}, \quad \frac{\partial h_2}{\partial a_{12}} = \sigma(a_{12})(1 - \sigma(a_{12}))$$

$$\frac{\partial a_{12}}{\partial w_{1,2}^{[2]}} = x_1$$

$$\frac{\partial l}{\partial o} = 2(o - y)$$

$$\begin{aligned} \frac{\partial l}{\partial w_{1,2}^{[2]}} &= \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial a_2} \cdot \frac{\partial a_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial a_{12}} \cdot \frac{\partial a_{12}}{\partial w_{1,2}^{[2]}} \\ &= 2(o - y) \cdot \sigma(a_2)(1 - \sigma(a_2)) \cdot w_2^{[2]} \cdot \sigma(a_{12})(1 - \sigma(a_{12})) \cdot x_1 \end{aligned}$$



5.(c)

No.

linear function can only result in ~~linear~~  
~~regression~~. representing linear relationship among  
the data.

the training data is non-linear  
thus cannot be represented by linear  
function.