

## Machine Learning Hw 3 吳偉樂 R09946023

### 1) (b) 30

$$E_D[E_{in}(w_{lin})] = \sigma^2(1 - \frac{d+1}{N}) > 0.006, \quad \sigma = 0.1, d = 11$$

$$0.1^2(1 - \frac{11+1}{N}) > 0.006 \rightarrow \frac{12}{N} < 0.4 \rightarrow N > 30$$

### 2) (a) There exists at least one solution for the normal equation.

For this normal equation,  $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

By Linear Algebra,

If  $\mathbf{X}^T \mathbf{X}$  is invertible, then  $\mathbf{w}$  has unique one solution

If  $\mathbf{X}^T \mathbf{X}$  is not invertible (i.e Singular), then  $\mathbf{w}$  has many optimal solutions

So there exists at least one solution for the normal equation .

### 3) (c)

3. (c) multiplying each of the  $n$ -th row by  $\frac{1}{n}$

Multiplying each of the  $n$ -th row by  $\frac{1}{n}$

$\Rightarrow \exists$  square matrix  $\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \frac{1}{n} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{n} \end{bmatrix}_{n \times n} = \mathbf{X}_i$  s.t.  $\begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \frac{1}{n} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{n} \end{bmatrix}_{n \times n} \begin{bmatrix} -x_1 \\ -x_2 \\ \vdots \\ -x_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} -x_1 \\ -\frac{1}{n}x_2 \\ \vdots \\ -\frac{1}{n}x_n \end{bmatrix}_{n \times 1}$

Note:  $\mathbf{X}_i$  is invertible since it only has values on diagonal.

Given  $\mathbf{H} = \mathbf{X}_A (\mathbf{X}_A^T \mathbf{X}_A)^{-1} \mathbf{X}_A^T$ , where  $\mathbf{X}_A^T \mathbf{X}_A$  invertible

After scaling, does  $(\mathbf{X}_i \mathbf{X}_A) [(\mathbf{X}_i \mathbf{X}_A)^T (\mathbf{X}_i \mathbf{X}_A)]^{-1} (\mathbf{X}_i \mathbf{X}_A)^T = \mathbf{X}_A (\mathbf{X}_A^T \mathbf{X}_A)^{-1} \mathbf{X}_A^T$  ?

$\Rightarrow \mathbf{X}_i \mathbf{X}_A [(\mathbf{X}_i \mathbf{X}_A)^T (\mathbf{X}_i \mathbf{X}_A)]^{-1} (\mathbf{X}_i \mathbf{X}_A)^T = \mathbf{X}_i \mathbf{X}_A [\mathbf{X}_A^T \mathbf{X}_i^T \mathbf{X}_i \mathbf{X}_A]^{-1} \mathbf{X}_A^T \mathbf{X}_i^T \mathbf{X}_i \mathbf{X}_A \leftarrow \mathbf{X}_i^T = \mathbf{X}_i$

$= \mathbf{X}_i \mathbf{X}_A [\mathbf{X}_A^T \mathbf{X}_i^2 \mathbf{X}_A]^{-1} \mathbf{X}_A^T \mathbf{X}_i^T$

Case ①: If  $\mathbf{X}_A$  invertible, then by Linear Algebra

$\mathbf{X}_i \mathbf{X}_A [\mathbf{X}_A^T \mathbf{X}_i \mathbf{X}_i \mathbf{X}_A]^{-1} \mathbf{X}_A^T \mathbf{X}_i^T = \mathbf{X}_i \mathbf{X}_A \mathbf{X}_A^{-1} \mathbf{X}_i^{-1} \mathbf{X}_i^{-1} (\mathbf{X}_A^T)^{-1} \mathbf{X}_A^T \mathbf{X}_i^T$

$= \mathbf{X}_i \mathbf{X}_i^{-1} \mathbf{X}_i^{-1} \mathbf{X}_i = \mathbf{I} = \mathbf{X}_A (\mathbf{X}_A^T \mathbf{X}_A)^{-1} \mathbf{X}_A^T$

Case ②: If  $\mathbf{X}_A$  not invertible, then we have  $= \mathbf{X}_A \mathbf{X}_A^{-1} (\mathbf{X}_A^T)^{-1} \mathbf{X}_A^T$

$\mathbf{X}_i \mathbf{X}_A [\mathbf{X}_A^T \mathbf{X}_i \mathbf{X}_i \mathbf{X}_A]^{-1} \mathbf{X}_A^T \mathbf{X}_i \neq \mathbf{X}_A [\mathbf{X}_A^T \mathbf{X}_A]^{-1} \mathbf{X}_A^T$

$\Rightarrow$  By ②, the answer is ①.

4) (e) 4

4. (e) 4

① Hoeffding ineq:  $P(|\bar{X} - E(\bar{X})| > \varepsilon) \leq 2e^{-2\varepsilon^2 N}$

Note:  $V = \frac{1}{N} \sum_{n=1}^N Y_n = \bar{Y}$

$\because Y_1, \dots, Y_N \stackrel{iid}{\sim} \text{Bernoulli}(N, \theta)$

$$E(V) = E\left(\frac{1}{N} \sum_{n=1}^N Y_n\right) = \frac{1}{N} (E(Y_1) + \dots + E(Y_N)) = \frac{1}{N} \cdot N E(Y_1) = E(Y_1) = \theta$$

$$\therefore P(|\bar{X} - E(\bar{X})| > \varepsilon) = P(|V - E(V)| > \varepsilon) = \underline{P(|V - \theta| > \varepsilon) \leq 2e^{-2\varepsilon^2 N}}$$

② Let  $L(Y|\theta) = \prod_{i=1}^N P(Y_i|\theta) = \prod_{i=1}^N \theta^{y_i} (1-\theta)^{1-y_i} = \theta^{\sum y_i} (1-\theta)^{N - \sum y_i}$

$$\ell(Y|\theta) = \log L(Y|\theta) = \log [\theta^{\sum y_i} (1-\theta)^{N - \sum y_i}] = \sum y_i \log \theta + (N - \sum y_i) \log (1-\theta)$$

Maximize  $L \Leftrightarrow$  maximize  $\ell$

$$\therefore \frac{d\ell}{d\theta} = \frac{\sum y_i}{\theta} - \frac{N - \sum y_i}{1-\theta} = 0 \Rightarrow \frac{\sum y_i}{\theta} - \frac{N - \sum y_i}{1-\theta} = 0 \Rightarrow \frac{\sum y_i}{\theta} = \frac{N - \sum y_i}{1-\theta} \Rightarrow \underline{\hat{\theta} = \frac{\sum y_i}{N} = V}$$

$$\frac{d^2\ell}{d\theta^2} = -\frac{\sum y_i}{\theta^2} - \frac{N - \sum y_i}{(1-\theta)^2} < 0 \Rightarrow \ell \text{ has max. value at } \frac{d\ell}{d\theta} = 0$$

$\therefore \hat{\theta} = V$  means that  $V$  maximizes likelihood ( $\hat{\theta}$ ) over  $\hat{\theta} \in [0, 1]$

③  $E_{in}(\hat{g}) = \frac{1}{N} \sum_{n=1}^N (\hat{g} - Y_n)^2$

$$\frac{dE}{d\hat{g}} = \frac{2}{N} \sum_{n=1}^N (\hat{g} - Y_n) = 0 \Rightarrow N\hat{g} = N V \Rightarrow \underline{\hat{g} = V}$$

$$\frac{d^2E}{d\hat{g}^2} = \frac{2}{N} \sum_{n=1}^N 1 = 2 > 0 \Rightarrow E \text{ has min. value at } \frac{dE}{d\hat{g}} = 0$$

$\therefore V$  minimized  $E_{in}(\hat{g}) = \frac{1}{N} \sum_{n=1}^N (\hat{g} - Y_n)^2$  over all  $\hat{g} \in \mathbb{R}$

④  $\nabla E_{in}(\hat{g}) = \frac{dE}{d\hat{g}} = \frac{2}{N} \sum_{n=1}^N (\hat{g} - Y_n)$

$$-\nabla E_{in}(\hat{g})|_{\hat{g}=0} = -\frac{2}{N} \sum_{n=1}^N (-Y_n) = \frac{2}{N} \sum_{n=1}^N Y_n = 2V$$

$\therefore 2V$  is the negative gradient direction  $-\nabla E_{in}(\hat{g})$  at  $\hat{g}=0$

5) (a)  $(\frac{1}{\theta})^N$

5. (a)  $(\frac{1}{\theta})^N$

$$Y_1, \dots, Y_N \stackrel{iid}{\sim} U(0, \theta) \Rightarrow f(y_i) = \begin{cases} \frac{1}{\theta}, & 0 \leq y_i \leq \theta \\ 0, & \text{o.w.} \end{cases} \quad I(0 \leq y_i \leq \theta) = \begin{cases} 1, & 0 \leq y_i \leq \theta \\ 0, & \text{o.w.} \end{cases}$$

$$L(Y|\theta) = \prod_{i=1}^N f(y_i|\theta) = \left(\frac{1}{\theta}\right)^N \prod_{i=1}^N I(0 \leq y_i \leq \theta), \text{ where } 0 \leq y_1, y_2, \dots, y_N \leq \theta$$

$\therefore$  For any  $\hat{\theta} \geq \max(y_1, \dots, y_N)$ , its likelihood is  $(\frac{1}{\hat{\theta}})^N \Leftrightarrow \max(y_1, \dots, y_N) \leq \hat{\theta}$

Additional: MLE.

Since  $(\frac{1}{\theta})$  is decreasing function,  $\therefore \max(f(y_i)) = \max\{y_1, \dots, y_N\}$

$\Rightarrow$  MLE is  $\max\{y_1, \dots, y_N\}$

6) (b)  $\text{err}(w, x, y) = \max(0, -yw^T x)$

6. (b)  $\text{err}(w, x, y) = \max(0, -yw^T x)$   
 $w_{t+1} = w_t + \frac{1}{N} \sum_{n=1}^N \mathbb{I}[\text{sign}(w_t^T x_n) \neq y_n] y_n x_n$   
 - Set  $y_n = 1$ , if  $\text{sign}(w_t^T x_n) = 1$ , then  $w_t^T x_n > 0 \Rightarrow y_n w_t^T x_n > 0 \Rightarrow \text{error} = 0$   
 $y_n = -1$ , if  $\text{sign}(w_t^T x_n) = -1$ , then  $w_t^T x_n < 0 \Rightarrow y_n w_t^T x_n > 0 \Rightarrow \text{error} = 0$   
 $y_n = 1$ , if  $\text{sign}(w_t^T x_n) = -1$ , then  $w_t^T x_n < 0 \Rightarrow y_n w_t^T x_n < 0 \Rightarrow \text{error} = y_n w_t^T x_n$   
 $y_n = -1$ , if  $\text{sign}(w_t^T x_n) = 1$ , then  $w_t^T x_n > 0 \Rightarrow y_n w_t^T x_n < 0 \Rightarrow \text{error} = y_n w_t^T x_n$   
 combine the statements above, we have  $\text{error}(w, x, y) = \max(0, -y_n w_n^T x_n)$

7) (a)  $+y_n x_n \exp(-y_n w^T x_n)$

7. (a)  $+y_n x_n e^{-y_n w^T x_n}$   
 $\text{err}_{\text{exp}}(w, x, y) = e^{-y w^T x} = e^{-y(w_0 x_0 + \dots + w_n x_n)}$   
 $\nabla \text{err}_{\text{exp}}(w, x, y) = \begin{pmatrix} \frac{d \text{err}}{d w_0} \\ \frac{d \text{err}}{d w_1} \\ \vdots \\ \frac{d \text{err}}{d w_n} \end{pmatrix} = \begin{pmatrix} -y e^{-y w^T x} x_0 \\ -y e^{-y w^T x} x_1 \\ \vdots \\ -y e^{-y w^T x} x_n \end{pmatrix} = -y e^{-y w^T x} \begin{pmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{pmatrix} = -y e^{-y w^T x} x$   
 $\Rightarrow -\nabla \text{err}_{\text{exp}}(w, x, y) = y x e^{-y w^T x}$

8) (b)  $-(A_E(u))^{-1} b_E(u)$

8. (b)  $-(A_E(u))^{-1} b_E(u)$   
 $w \leftarrow u + v \Rightarrow v = w - u$   
 $\Rightarrow E(w) \approx E(u) + b_E(u)^T \cdot v + \frac{1}{2} v^T A_E(u) v$   
 $\nabla_v E(w) = b_E(u)^T + \frac{1}{2} \cdot 2 A_E(u) v = b_E(u)^T + v A_E(u) = 0$   
 $\Rightarrow v A_E(u) = -b_E(u)^T \Rightarrow v = -b_E(u)^T \cdot (A_E(u))^{-1}$   
 $A_E(u)$  is symmetric,  $\therefore A_E(u)^T$  is also symmetric (By Linear Algebra)  
 $\Rightarrow v = \underbrace{-b_E(u)^T}_{1 \times d} \underbrace{(A_E(u))^{-1}}_{d \times d} = -[A_E(u)^{-1}]^T b_E(u)^T = -\underbrace{(A_E(u))^{-1}}_{d \times d} b_E(u)$   
 same elements but different interpretation method: (0,0) or  $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$

9) (b)

9. (b)  $\frac{2}{N} X^T X$

$$E = E_{in} = \frac{1}{N} \sum_{i=1}^N (w^T x_i - y_i)^2, \text{ where } x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{pmatrix}, X = \begin{bmatrix} -x_1^T & - \\ & \vdots \\ -x_N^T & - \end{bmatrix}_{N \times d}$$

$$\Rightarrow E_{in} = \frac{1}{N} \left[ (w_1 x_{11} + w_2 x_{12} + \dots + w_d x_{1d} - y_1)^2 + (w_1 x_{21} + w_2 x_{22} + \dots + w_d x_{2d} - y_2)^2 + \dots \right. \\ \left. + (w_1 x_{N1} + w_2 x_{N2} + \dots + w_d x_{Nd} - y_N)^2 \right]$$

$$\Rightarrow \frac{dE_{in}}{dw_1} = \frac{2}{N} \left[ (w_1 x_{11} + w_2 x_{12} + \dots + w_d x_{1d} - y_1) x_{11} + \dots + (w_1 x_{N1} + w_2 x_{N2} + \dots + w_d x_{Nd} - y_N) x_{N1} \right]$$

$$\Rightarrow \frac{dE_{in}}{dw_k} = \frac{2}{N} \left[ (w_1 x_{11} + w_2 x_{12} + \dots + w_d x_{1d} - y_1) x_{1k} + \dots + (w_1 x_{N1} + w_2 x_{N2} + \dots + w_d x_{Nd} - y_N) x_{Nk} \right] \text{ for } k=1, 2, \dots, d$$

$$\frac{d^2 E_{in}}{dw_k^2} = \frac{2}{N} \left[ x_{1k}^2 + x_{2k}^2 + \dots + x_{Nk}^2 \right] \text{ for } k=1, \dots, d$$

$$\frac{d^2 E_{in}}{dw_k dw_j} = \frac{2}{N} \left[ x_{1j} x_{1k} + x_{2j} x_{2k} + \dots + x_{Nj} x_{Nk} \right] \text{ for } j=1, \dots, d, j \neq k$$

$$\therefore A_E(w) = \begin{bmatrix} \frac{\partial^2 E}{\partial w_1^2}(w) & \frac{\partial^2 E}{\partial w_1 \partial w_2}(w) & \dots & \frac{\partial^2 E}{\partial w_1 \partial w_d}(w) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 E}{\partial w_d \partial w_1}(w) & \frac{\partial^2 E}{\partial w_d \partial w_2}(w) & \dots & \frac{\partial^2 E}{\partial w_d^2}(w) \end{bmatrix} dy$$

$$= \frac{2}{N} \begin{bmatrix} x_{11}^2 + x_{21}^2 + \dots + x_{N1}^2 & x_{12}x_{11} + x_{22}x_{21} + \dots + x_{N2}x_{N1} & \dots & x_{1d}x_{11} + x_{2d}x_{21} + \dots + x_{Nd}x_{N1} \\ \vdots & \vdots & \ddots & \vdots \\ x_{1d}x_{11} + \dots + x_{Nd}x_{N1} & x_{1d}x_{12} + x_{2d}x_{22} + \dots + x_{Nd}x_{N2} & \dots & x_{1d}^2 + x_{2d}^2 + \dots + x_{Nd}^2 \end{bmatrix}$$

$$= \frac{2}{N} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Nd} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Nd} \end{bmatrix}_{N \times d} = \underline{\underline{\frac{2}{N} X^T X}}$$

10) (b)

10. (b)  $(h_k(x) - \mathbb{I}[y=k]) x_i$

$$w = \begin{bmatrix} w_1 & w_2 & \dots & w_k & \dots & w_K \end{bmatrix}_{(d+1) \times K} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1k} & \dots & w_{1K} \\ w_{21} & w_{22} & \dots & w_{2k} & \dots & w_{2K} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{d+1,1} & w_{d+1,2} & \dots & w_{d+1,k} & \dots & w_{d+1,K} \end{bmatrix}, h_y(y) = \frac{e^{w_y^T x}}{\sum_{i=1}^K e^{w_i^T x}}$$

$$\text{err}(w, x, y) = -\ln h_y(x) = -\sum_{j=1}^K \mathbb{I}[y=j] \ln h_j(x)$$

$$\begin{aligned} &= -\sum_{j=1}^K \mathbb{I}[y=j] \ln \frac{e^{w_j^T x}}{\sum_{i=1}^K e^{w_i^T x}} \\ &= -\mathbb{I}[y=1] \ln \frac{e^{w_1^T x}}{\sum_{i=1}^K e^{w_i^T x}} - \dots - \mathbb{I}[y=k] \ln \frac{e^{w_k^T x}}{\sum_{i=1}^K e^{w_i^T x}} - \dots - \mathbb{I}[y=K] \ln \frac{e^{w_K^T x}}{\sum_{i=1}^K e^{w_i^T x}} \end{aligned}$$

$$\frac{\partial \text{err}(w, x, y)}{\partial w_{ik}} = -\mathbb{I}[y=1] \cdot \frac{1}{h_1(x)} \cdot \frac{-e^{w_1^T x} \cdot e^{w_k^T x} x_i}{\left(\sum_{i=1}^K e^{w_i^T x}\right)^2} - \dots - \mathbb{I}[y=k] \cdot \frac{1}{h_k(x)} \cdot \frac{e^{w_k^T x} \cdot e^{w_k^T x} x_i}{\left(\sum_{i=1}^K e^{w_i^T x}\right)^2} - \dots - \mathbb{I}[y=K] \cdot \frac{1}{h_K(x)} \cdot \frac{-e^{w_k^T x} \cdot e^{w_K^T x} x_i}{\left(\sum_{i=1}^K e^{w_i^T x}\right)^2}$$

$$= \mathbb{I}[y=1] \cdot \frac{1}{h_1(x)} \cdot h_1(x) h_k(x) x_i + \dots + \mathbb{I}[y=k] \cdot \frac{1}{h_k(x)} [-h_k(x) + h_k(x) h_k(x)] x_i + \dots$$

$$+ \mathbb{I}[y=K] \cdot \frac{1}{h_K(x)} \cdot h_K(x) h_k(x) x_i$$

$$= \sum_{j=1}^K \mathbb{I}[y=j] h_k(x) x_i - \mathbb{I}[y=k] x_i$$

$$= \left[ \sum_{j=1}^K \mathbb{I}[y=j] h_k(x) - \mathbb{I}[y=k] \right] x_i$$

$$\stackrel{y=k}{=} \left[ h_k(x) - \mathbb{I}[y=k] \right] x_i \quad \times$$

11) (e)

11. (e)  $w_2^* - w_1^*$

When  $K=2$ , error in MLR is equivalent to error in logistic reg.

$\therefore \text{err}(w, x, y)_{\text{MLR}} = -\ln h_y(x) = \text{err}(w, x, y)_{\text{LR}} = -\ln \theta(yw^T x)$

$\Rightarrow h_y(x) = \theta(yw^T x) = \frac{1}{1+e^{-yw^T x}} \Leftrightarrow h_{y_n}(x_n) = \theta(y_n w^T x_n)$

relabel:

$K=y_n=1 \Rightarrow y_n'=-1$

$\therefore h_1(x_n) = \frac{e^{w_1^T x_n}}{e^{w_1^T x_n} + e^{w_2^T x_n}} = \frac{1}{1+e^{(w_2^T - w_1^T)x_n}} = \theta(y_n' w^T x_n) = \frac{1}{1+e^{-w^T x_n}}$

$\Rightarrow w^T = w_2^{*T} - w_1^{*T} \Rightarrow w = \underline{w_2^* - w_1^*}$

$K=2=y_n \Rightarrow y_n'=1$

$\therefore h_2(x_n) = \frac{e^{w_2^T x_n}}{e^{w_2^T x_n} + e^{w_1^T x_n}} = \frac{1}{1+e^{(w_1^T - w_2^T)x_n}} = \theta(y_n' w^T x_n) = \frac{1}{1+e^{-w^T x_n}}$

$\Rightarrow -w^T = (w_1^* - w_2^*)^T \Rightarrow \underline{w = w_2^* - w_1^*}$

12) (e) [-7,0,0,2,-2,3]

Draw each 5 curve on the graph, and mark those points, obviously is (e)

13) (b)

13.  $(x_1, \dots, x_d) \in \mathbb{R}^d \xrightarrow{\mathcal{I}_{(k)}} (1, x_k) \in \mathbb{R}^2$  Each  $\mathcal{I}_{(k)}$  is a decision stump.

$\mathcal{I}_{(1)} : (1, x_1) \Rightarrow w_0 + w_1 x_1 \Rightarrow m_{H_1}(N) = 2N$

$\mathcal{I}_{(2)} : (1, x_2) \Rightarrow w_0 + w_2 x_2 \Rightarrow m_{H_2}(N) = 2N$

$\vdots$

$\mathcal{I}_{(d)} : (1, x_d) \Rightarrow w_0 + w_d x_d \Rightarrow m_{H_d}(N) = 2N$

$\Rightarrow m_{H_k}(N) = 2Nd$

$\text{dec}(U_{k+1}^d, H_k) = \{\text{largest } N \text{ for which } m_{H_k}(N) = 2^N\} = N$

$2^N \leq m_{H_k}(N) = 2Nd \quad \leftarrow \text{we cannot shatter } X \text{ if } 2^N > m_{H_k}(N)$

$2^{N-1} \leq Nd$

$N-1 \leq \log_2 d + \log_2 N \leq \log_2 d + \frac{N}{2} \Rightarrow N \leq 2(\log_2 d + 1)$

(b)  $2 \lceil \log_2 d + 1 \rceil$

14) (d) 0.60

```
import numpy as np
import random
import math

dat_file = r'C://Users//USER//Desktop//hw3_train.dat.txt'
with open(dat_file, 'r') as f:
    text = f.read()

data = text.split() #split string into a list

D = []
for i in range(1000):
    K = list(data[i*11:(i+1)*11])
    D.append(K)

for i in range(len(D)):
    for j in range(len(D[0])):
        D[i][j] = float(D[i][j])

# 14 (d)
X = []
Y = []
for i in range(len(D)):
    X.append([1.0] + D[i][0:10]) # xi = [x0, x1, ..., x10]
    Y.append(D[i][-1])

XTX = np.transpose(X).dot(X)
W_LIN = np.linalg.inv(XTX).dot(np.transpose(X)).dot(Y)
E_in = np.square(np.subtract(np.array(X).dot(W_LIN), Y)).mean()
print(E_in)

0.6053223804672918
```

15) (c) 1800

```
# 15 (c)
iteration = 0
random.seed(40)
SUM = 0
while iteration != 1000:
    E_wt = 10 # initialize
    W = [0]*11 # initialize
    i=0

    while E_wt>1.01*E_in:
        r = random.randint(0, 999)
        neg_gra_err = 2*(Y[r] - np.array(X[r]).dot(W))*np.array(X[r])
        W = W + 0.001* neg_gra_err
        E_wt = np.square(np.subtract(np.array(X).dot(W),Y)).mean()
        i+=1
    SUM = SUM+i
    iteration +=1

SUM = SUM/1000
print(SUM)
## 20 mins runtime

1772.577
```

16) (c) 0.56

```
# 16 (c) 0.56
iteration = 0
SUM = 0
random.seed(50)
while iteration != 1000:
    E_in_c = 0
    W = [0]*11
    i = 0

    while i != 500:
        r = random.randint(0, 999)
        Exp = math.exp(-Y[r]*np.array(W).dot(X[r]))
        neg_gra_err = Y[r]* np.array(X[r]) *Exp/(1+Exp)
        W = W + 0.001* neg_gra_err
        for j in range(1000):
            Ep = math.exp(-Y[j]*np.array(W).dot(X[j]))
            E_in_c += math.log(1+Ep)
        E_in_c = E_in_c/1000
        i+=1
    SUM = SUM + E_in_c
    iteration += 1
SUM = SUM/1000
print(SUM)
# Almost 30 mins runtime

0.5694796431407935
```

17) (b) 0.50

```
# 17 (b) 0.50
iteration = 0
SUM = 0
random.seed(50)
while iteration != 1000:
    E_in_c = 0
    W = W_LIN
    i = 0

    while i != 500:
        r = random.randint(0, 999)
        Exp = math.exp(-Y[r]*np.array(W).dot(X[r]))
        neg_gra_err = Y[r]* np.array(X[r]) *Exp/(1+Exp)
        W = W + 0.001* neg_gra_err
        for j in range(1000):
            Ep = math.exp(-Y[j]*np.array(W).dot(X[j]))
            E_in_c += math.log(1+Ep)
        E_in_c = E_in_c/1000
        i+=1
    SUM = SUM + E_in_c
    iteration += 1
SUM = SUM/1000
print(SUM)
# Almost 30 mins runtime
```

0.5033221372400456



18) (a) 0.32

```
# 18 (a) 0.32
test_file = r'C://Users//USER//Desktop//hw3_test.dat.txt'
with open(test_file, 'r') as f:
    text = f.read()

tdata = text.split() #split string into a list

T = []
for i in range(int(len(tdata)/11)):
    K = list(tdata[i*11:(i+1)*11])
    T.append(K)

for i in range(len(T)):
    for j in range(len(T[0])):
        T[i][j] = float(T[i][j])

X_test = []
Y_test = []
for i in range(len(T)):
    X_test.append([1.0] + T[i][0:10]) # xi = [x0, x1, ..., x10]
    Y_test.append(T[i][-1])

def sign(y):
    h = y
    if h>0: h = 1
    else: h = -1
    return h

E_in_binary = 0
Y_temporary = np.array(X).dot(W_LIN)
Y_hat = [sign(Y_temporary[j]) for j in range(len(Y_temporary))]
for j in range(len(Y)):
    if Y_hat[j] != Y[j]:
        E_in_binary+=1
E_in_binary = E_in_binary/len(Y)

E_out_binary = 0
Y_test_temp = np.array(X_test).dot(W_LIN)
Y_test_hat = [sign(Y_test_temp[j]) for j in range(len(Y_test_temp))]
for j in range(len(Y_test)):
    if Y_test_hat[j] != Y_test[j]:
        E_out_binary+=1
E_out_binary = E_out_binary/len(Y_test)

Diff_E = abs( E_out_binary - E_in_binary )
print(Diff_E)

0.32266666666666666
```

### 19) (b) 0.36

```
# 19 (b) 0.36
X_new = X.copy()
for i in range(2):
    for j in range(len(X_new)):
        X_new[j] = X_new[j] + [np.power(X[j][k+1], i+2) for k in range(len(X[0])-1)]

XTX_t = np.transpose(X_new).dot(X_new)
W_LIN_t = np.linalg.inv(XTX_t).dot(np.transpose(X_new)).dot(Y)

E_in_bi = 0
Y_train_hat_temp = np.array(X_new).dot(W_LIN_t)
Y_hatrain = [sign(Y_train_hat_temp[j]) for j in range(len(Y_train_hat_temp))]
for j in range(len(Y)):
    if Y_hatrain[j] != Y[j]:
        E_in_bi += 1
E_in_bi = E_in_bi / len(Y)

X_test_new = X_test.copy()
for i in range(2):
    for j in range(len(X_test_new)):
        X_test_new[j] = X_test_new[j] + [np.power(X_test[j][k+1], i+2) for k in range(len(X_test[0])-1)]

E_out_bi = 0
Y_test_hat_temp = np.array(X_test_new).dot(W_LIN_t)
Y_hatest = [sign(Y_test_hat_temp[j]) for j in range(len(Y_test_hat_temp))]
for j in range(len(Y_test)):
    if Y_hatest[j] != Y_test[j]:
        E_out_bi += 1
E_out_bi = E_out_bi / len(Y_test)

Diff_E_new = abs(E_out_bi - E_in_bi)
print(Diff_E_new)

0.37366666666666665
```

## 20) (d) 0.44

```
# 20 (d) 0.44
X_new = X.copy()
for i in range(9):
    for j in range(len(X_new)):
        X_new[j] = X_new[j] + [np.power(X[j][k+1], i+2) for k in range(len(X[0])-1)]
XTX_t = np.transpose(X_new).dot(X_new)
W_LIN_t = np.linalg.inv(XTX_t).dot(np.transpose(X_new)).dot(Y)

E_in_bi = 0
Y_train_hat_temp = np.array(X_new).dot(W_LIN_t)
Y_hatrain = [sign(Y_train_hat_temp[j]) for j in range(len(Y_train_hat_temp))]
for j in range(len(Y)):
    if Y_hatrain[j] != Y[j]:
        E_in_bi += 1
E_in_bi = E_in_bi / len(Y)

X_test_new = X_test.copy()
for i in range(9):
    for j in range(len(X_test_new)):
        X_test_new[j] = X_test_new[j] + [np.power(X_test[j][k+1], i+2) for k in range(len(X_test[0])-1)]

E_out_bi = 0
Y_test_hat_temp = np.array(X_test_new).dot(W_LIN_t)
Y_hatest = [sign(Y_test_hat_temp[j]) for j in range(len(Y_test_hat_temp))]
for j in range(len(Y_test)):
    if Y_hatest[j] != Y_test[j]:
        E_out_bi += 1
E_out_bi = E_out_bi / len(Y_test)

Diff_E_new = abs(E_out_bi - E_in_bi)
print(Diff_E_new)

0.4466666666666666
```