

* Part B

1. • Sigmoid function.

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

we have

$$\begin{aligned} \frac{\partial (\sigma(x))}{\partial x} &= \frac{\partial}{\partial x} \left(\frac{1}{1 + e^{-x}} \right) = -\frac{1}{(1 + e^{-x})^2} \frac{\partial (e^{-x})}{\partial x} \\ &= -\frac{1}{(1 + e^{-x})^2} (-e^{-x}) = \frac{e^{-x}}{(1 + e^{-x})^2} \rightarrow ① \end{aligned}$$

We can reframe ① as

$$\frac{\partial \sigma(x)}{\partial x} = \frac{1 + e^{-x} - 1}{(1 + e^{-x})^2} = \sigma(x)(1 - \sigma(x)) \rightarrow ②$$

Now, applying chain rule to find $\frac{\partial \{\sigma(g(w))\}}{\partial w}$

we get

$$\begin{aligned} \frac{\partial \{\sigma(g(w))\}}{\partial w} &= \underbrace{\frac{\partial \{\sigma(g)\}}{\partial g} \frac{\partial g}{\partial w}}_{\sigma(g)(1 - \sigma(g)) \frac{\partial g}{\partial w}} \quad (\text{by chain rule}) \\ &= \boxed{\sigma(g)(1 - \sigma(g)) \frac{\partial g}{\partial w}} \end{aligned}$$

• Hyperbolic tangent

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

We have

$$\begin{aligned} \frac{\partial(\tanh(x))}{\partial x} &= \frac{\partial}{\partial x} \left\{ \frac{e^x - e^{-x}}{e^x + e^{-x}} \right\} = \frac{(e^x + e^{-x})(e^2 + e^{-2})}{(e^x + e^{-x})^2} \\ &\quad - (e^x - e^{-x})(e^2 - e^{-2}) \\ &= 1 - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} \end{aligned}$$

$\frac{d \tanh(x)}{dx} = 1 - \tanh^2(x)$

→ ①

Applying chain rule on ① with $x = g(w)$, we get

$$\begin{aligned} \frac{\partial \tanh(g(w))}{\partial w} &= \frac{\partial \tanh(g)}{\partial g} \frac{\partial g}{\partial w} \\ &= \left[\left\{ 1 - \tanh^2(g(w)) \right\} \frac{\partial g(w)}{\partial w} \right] \quad (\text{by ①}) \end{aligned}$$

• ReLU

We have $\text{ReLU}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{else.} \end{cases}$

Then, $\frac{\partial \text{ReLU}(x)}{\partial x} = \begin{cases} 1 & x > 0 \\ 0 & \text{else} \end{cases}$

For $x = g(w)$, we have

$\frac{\partial \text{ReLU}(g(w))}{\partial w} = \begin{cases} 1 \cdot \frac{\partial g(w)}{\partial w} & g(w) > 0 \\ 0 & \text{else} \end{cases}$

2. Gradients of common loss functions

• Cross entropy loss

For vector input \bar{a} , we define $S(\bar{a}) : \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_c \end{bmatrix} \rightarrow \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_c \end{bmatrix}$

as $s_i = \frac{e^{a_i}}{\sum_k e^{a_k}}$ (softmax function)

Then, for s_i , we have

$$\frac{\partial s_i}{\partial a_i} = \frac{e^{a_i}}{\sum_k e^{a_k}} + \frac{-1(e^{a_i})}{(\sum_k e^{a_k})^2} e^{a_i} = s_i (1 - s_i) \quad \xrightarrow{\text{---}} \textcircled{1}$$

and (for $i \neq j$)

$$\frac{\partial s_i}{\partial a_j} = \frac{-e^{a_i} e^{a_j}}{(\sum_k e^{a_k})^2} = -s_i s_j \quad \xrightarrow{\text{---}} \textcircled{2}$$

Combining $\textcircled{1}$ and $\textcircled{2}$,

where $\delta_{ij} = \begin{cases} 1 & i=j \\ 0 & \text{else} \end{cases}$

$$\boxed{\frac{\partial s_i}{\partial a_j} = s_i (\delta_{ij} - s_j)} \quad \xrightarrow{\text{---}} \textcircled{3}$$

Now, cross entropy loss is defined as follows

$$J(y, \hat{y}) = - \left(\sum_{k=1}^c y_k \log \hat{y}_k - \vec{y} \cdot \log(S(\vec{\hat{y}})) \right) \quad \xrightarrow{\text{---}} \textcircled{4}$$

where $y \Rightarrow$ one hot encoded vector corresponding to correct class of sample.

$\hat{y} \Rightarrow$ scores ~~for~~ vector for sample y (one score for each class.)

From ④, we have

$$J(y, \hat{y}) = - \sum_{k=1}^c y_k \log(S_k(\hat{y})) \quad \xrightarrow{\text{using definition of softmax}} \quad ⑤$$

Let $\hat{y}_i = g(w_i)$. We want to find the gradient $\frac{\partial J}{\partial w_i}$. This can be found using

⑤ as

$$\begin{aligned} \frac{\partial J}{\partial w_i} &= \frac{\partial J}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial w_i} = \frac{\partial g}{\partial w} \left\{ \frac{\partial}{\partial \hat{y}_i} \left[- \sum_k y_k \log(S_k(\hat{y})) \right] \right\} \\ &= \left\{ - \sum_{k=1}^c \left[\frac{y_k}{S_k(\hat{y})} \frac{\partial S_k(\hat{y})}{\partial \hat{y}_i} \right] \right\} \frac{\partial g(w_i)}{\partial w_i} \\ &= - \sum_{k=1}^c \left\{ \frac{y_k}{S_k(\hat{y})} S_k(\hat{y}) (\delta_{kj} - S_j(\hat{y})) \right\} \frac{\partial g}{\partial w} \\ &\boxed{\frac{\partial J}{\partial w_i} = - \sum_{k=1}^c y_k (\delta_{kj} - S_j(\hat{y})) \frac{\partial g(w_i)}{\partial w_i}} \quad \left[\text{from } ③ \right] \end{aligned}$$

where C is the no. of classes.

and $\hat{y}_i = g(w_i)$

If there are n samples, where for i^{th} sample, $\hat{y}_j = g(w_j, x^{(i)})$, then

$$\boxed{\frac{\partial J}{\partial w_j} = - \sum_{i=1}^n \left\{ \left(\sum_{k=1}^c y_k^{(i)} (\delta_{kj} - S_j(\hat{y}^{(i)})) \right) \frac{\partial g(w, x^{(i)})}{\partial w} \right\}}$$

Basically, loss is summed up over all samples.

* Hinge loss

We assume n samples are given and the classification is binary (classes are -1 or 1).

In this case, we can define the hinge loss as

$$J_{\xi} = \sum_{i=1}^n \max(0, 1 - y_i \hat{y}_i)$$

Assume $\hat{y}_i = g(\omega, x_i)$ [$x_i \rightarrow$ inputs of i^{th} sample]

then

$$\frac{\partial J}{\partial \omega} = \sum_{i=1}^n \frac{\partial \xi}{\partial \omega} \{ \max(0, 1 - y_i \hat{y}_i) \}$$

$$\Rightarrow \frac{\partial J}{\partial \omega} = \sum_{i=1}^n \frac{\partial J}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial \omega} \quad (\text{by chain rule})$$

$$\Rightarrow \boxed{\frac{\partial J}{\partial \omega} = \sum_{i=1}^n \frac{\partial J}{\partial \hat{y}_i} \frac{\partial g(\omega, x_i)}{\partial \omega}} \longrightarrow ①$$

where

$$\frac{\partial J}{\partial \hat{y}_i} = \begin{cases} -y_i & \text{if } y_i \hat{y}_i < 1 \\ 0 & \text{else} \end{cases}$$

Substituting ② in ①, we

can find the gradients wrt the weights.

• L1 loss for ~~a sample~~

We assume y_i^o is the correct output for input x_i^o , and \hat{y}_i is the output of our model.

Then for n samples, L1 loss

$$J = \sum_{i=1}^n |y_i^o - \hat{y}_i|$$

We assume $\hat{y}_i = g(w, x_i)$

Then, by chain rule,

$$\frac{\partial J}{\partial w} = \sum_{i=1}^n \frac{\partial J}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial w} \quad \rightarrow ①$$

We have

$$\begin{aligned} \frac{\partial J}{\partial \hat{y}_i} &= \frac{\partial}{\partial \hat{y}_i} \left(\sum_{k=1}^n |y_k^o - \hat{y}_k^o| \right) = -\frac{1}{(y_i^o - \hat{y}_i^o)} \\ &= +\frac{|\hat{y}_i^o - y_i^o|}{\hat{y}_i^o - y_i^o} \quad \rightarrow ② \end{aligned}$$

Substituting ② in ①, we get

$$\boxed{\frac{\partial J}{\partial w} = \sum_{i=1}^n \frac{|\hat{y}_i^o - y_i^o|}{\hat{y}_i^o - y_i^o} \frac{\partial g(w, x_i)}{\partial w}} \quad \leftarrow \text{Ans}$$

where $\hat{y}_i^o = g(w, x_i^o)$

• Huber loss

Assuming same notation as before, for sample i , we assume Huber loss.

$$J_i = \begin{cases} \frac{1}{2} (y_i - \hat{y}_i)^2 & \text{for } |y_i - \hat{y}_i| \leq \delta \\ \delta |y_i - \hat{y}_i| - \frac{\delta^2}{2} & \text{else} \end{cases}$$

and total loss $J = \sum_{i=1}^n J_i$

Assume $y_i = g(w, x_i)$

Then gradient $\frac{\partial J}{\partial w}$ can be found as

$$\begin{aligned} \frac{\partial J}{\partial w} &= \frac{\partial J_1}{\partial w} + \frac{\partial J_2}{\partial w} + \dots + \frac{\partial J_n}{\partial w} \\ &= \sum_{i=1}^n \frac{\partial J_i}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial w} \xrightarrow{\textcircled{1}} (\text{by chain rule}) \end{aligned}$$

We have

$$\frac{\partial J_i}{\partial \hat{y}_i} = \begin{cases} \hat{y}_i - y_i & |y_i - \hat{y}_i| \leq \delta \\ -\frac{\delta(y_i - \hat{y}_i)}{|y_i - \hat{y}_i|} & \text{else} \end{cases} \xrightarrow{\textcircled{2}}$$

Substituting ② ③ in ①, and $\frac{\partial \hat{y}_i}{\partial w} = \frac{\partial g(w, x_i)}{\partial w} \rightarrow \textcircled{3}$
 we get the gradient w.r.t. w.

• L2 loss -

for a batch of n samples, we define
L2 loss as

$$J = \sum_{i=1}^n \frac{1}{2} (y_i - \hat{y}_i)^2$$

we assume $\hat{y}_i = g(w, x_i)$ (w is one of the ~~variable~~
weights)

* Then

$$\frac{\partial J}{\partial w} = - \sum_{i=1}^n \frac{\partial J}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial w} \quad (\text{chain rule})$$

$$= \sum_{i=1}^n \frac{1}{2} (y_i - \hat{y}_i) (-1) \frac{\partial g(w, x_i)}{\partial w}$$

$$\Rightarrow \boxed{\frac{\partial J}{\partial w} = \sum_{i=1}^n (\hat{y}_i - y_i) \frac{\partial g(w, x_i)}{\partial w}}$$

• Cosine similarity &

We assume $y^{(i)}$ is a one-hot encoded vector for C classes, and $\hat{y}^{(i)}$ is a score vector.

Then the loss function $J(y, \hat{y})$ can be defined as:

$$J = \sum_{i=1}^n \left(-\frac{\sum_{k=1}^C y_k^{(i)} \hat{y}_k^{(i)}}{|\hat{y}^{(i)}|} \right) \rightarrow ①$$

where $|\hat{y}^{(i)}|$ denotes absolute value of vector $\hat{y}^{(i)}$. Now assume $y_j^{(i)} = g(w_j, x^{(i)})$.

~~From ①, we have~~

Then

$$\frac{\partial J}{\partial w_j} = \sum_{i=1}^n \frac{\partial}{\partial w_j} \left\{ -\frac{\sum_{k=1}^C y_k^{(i)} \hat{y}_k^{(i)}}{|\hat{y}^{(i)}|} \right\}$$

$$\Rightarrow \frac{\partial J}{\partial w_k} = \sum_{i=1}^n \frac{\partial}{\partial \hat{y}_j^{(i)}} \left\{ -\frac{\sum_{k=1}^C y_k^{(i)} \hat{y}_k^{(i)}}{|\hat{y}^{(i)}|} \right\} \frac{\partial \hat{y}_j^{(i)}}{\partial w_j}$$

$$\frac{\partial J}{\partial w_k} = -\sum_{i=1}^n \left[\left\{ \frac{y_j^{(i)}}{|\hat{y}^{(i)}|} - \left(\frac{\sum_{k=1}^C y_k^{(i)} \hat{y}_k^{(i)}}{|\hat{y}^{(i)}|^3} \right) \hat{y}_j^{(i)} \right\} \frac{\partial g(w_j, x^{(i)})}{\partial w_j} \right] \rightarrow ②$$

② is the final gradient after chain rule application.

3. Hand calculation of gradients.

for my set:

$$\text{inp}1 = 0.18$$

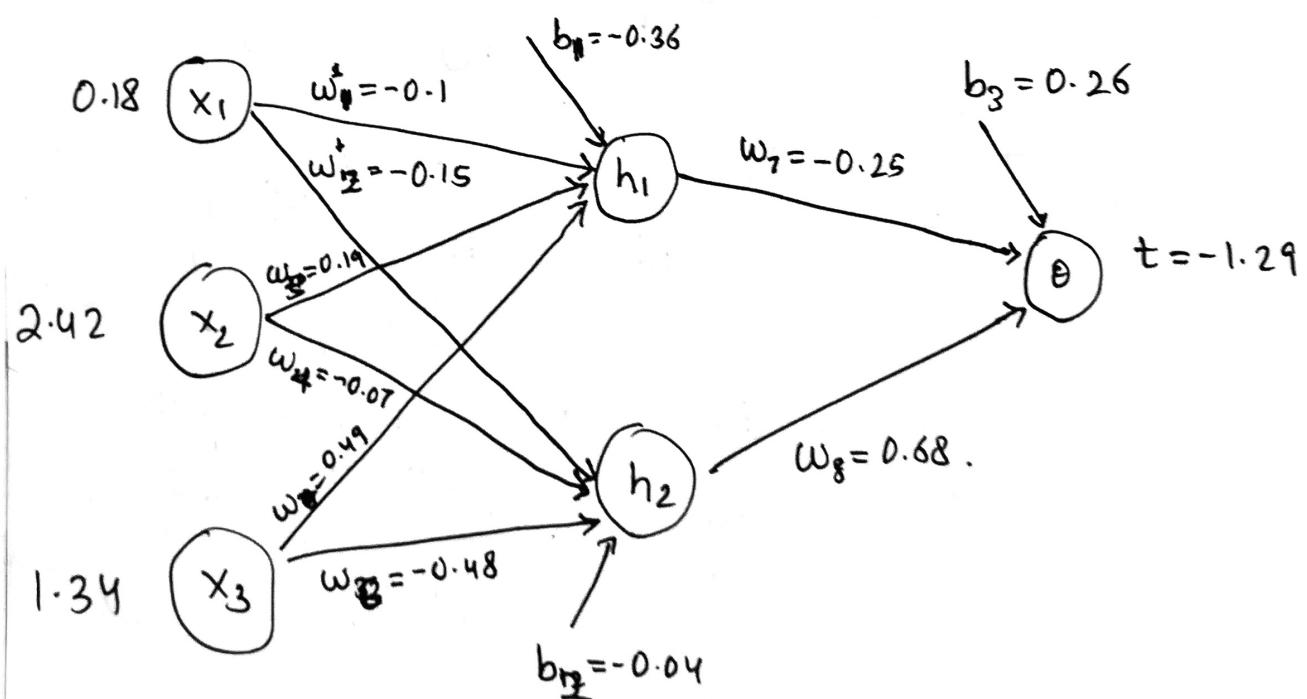
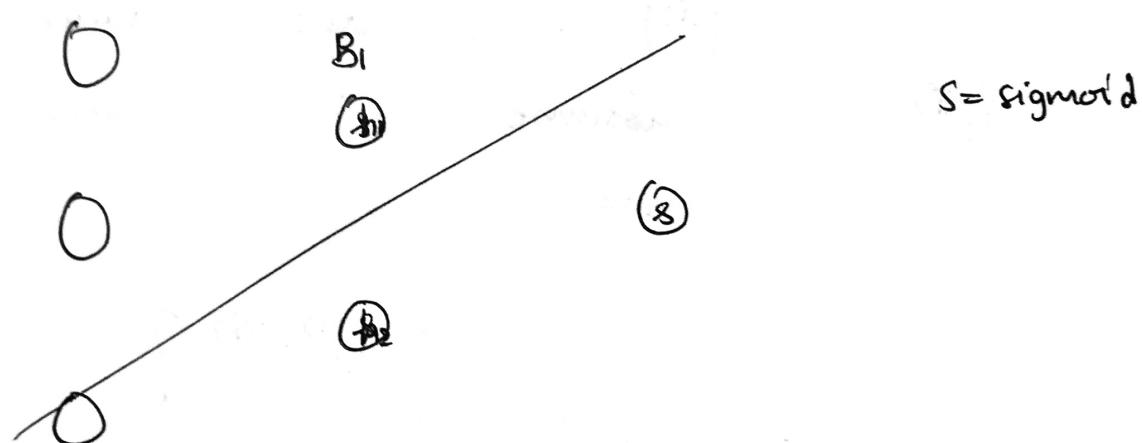
$$\text{inp}2 = 2.42$$

$$\text{inp}3 = 1.34$$

$$W_1 = \begin{bmatrix} -0.1 & -0.15 \\ 0.19 & -0.07 \\ 0.49 & -0.48 \end{bmatrix}$$

$$B_1 = \begin{bmatrix} -0.36 \\ -0.04 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} -0.25 \\ 0.68 \end{bmatrix} \quad B_2 = 0.26 \cdot \text{target} = -1.29.$$



$$b_1 + \omega_1 x_1 + \omega_3 x_2 + \omega_5 x_3 = z_{h_1}$$

$$b_2 + \omega_2 x_1 + \omega_4 x_2 + \omega_6 x_3 = z_{h_2}$$

$$h_1 = \sigma(z_{h_1})$$

$$h_2 = \sigma(z_{h_2})$$

$$b_3 + \omega_7 h_1 + \omega_8 h_2 = z_\theta$$

$$\theta = \sigma(z_\theta)$$

$$t = -1.29$$

Substituting in the above formulas,
we get.

$$z_{h_1} = -0.36 + (-0.1)(0.18) + (0.19)(2.42) + (1.34)(0.49)$$

$$= \underline{\underline{0.7384}}$$

$$z_{h_2} = -0.04 + (-0.15)(0.18) + (-0.07)(2.42) + (-0.48)(1.34)$$

$$= \underline{\underline{-0.8796}}$$

$$h_1 = \sigma(0.7384) = \underline{\underline{0.6767}} = \underline{\underline{0.68}}$$

$$h_2 = \sigma(-0.8796) = \underline{\underline{0.2933}} = \underline{\underline{0.29}}.$$

$$z_\theta = 0.26 + (-0.25)(0.68) + (0.68)(0.29)$$

$$= \underline{\underline{0.2872}}.$$

$$\theta = \sigma(z_\theta) = \underline{\underline{0.57}}.$$

$$E = \frac{1}{2} (\theta - t)^2 = \underline{\underline{1.73}}$$

$$\frac{\partial E}{\partial \theta} = \theta - t = \underline{\underline{1.86}}.$$

$$\frac{\partial z_0}{\partial w_1} = h_1 \quad ; \quad \frac{\partial z_0}{\partial w_8} = h_2 \quad ; \quad \frac{\partial z_0}{\partial b_3} = 1$$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial \theta} \frac{\partial \theta}{\partial z_0} \frac{\partial z_0}{\partial w_1} = (\theta - t) \theta(1-\theta) h_1$$

$$\frac{\partial E}{\partial w_1} = (1.86)(0.57(1-0.57)) (0.68) = \boxed{0.31}$$

$$\frac{\partial E}{\partial w_8} = \frac{\partial E}{\partial \theta} \frac{\partial \theta}{\partial z_0} \frac{\partial z_0}{\partial w_8} = (\theta - t) \theta(1-\theta) h_2$$

$$\frac{\partial E}{\partial w_8} = (1.86)(0.57)(0.43)(0.29) = \boxed{0.13}$$

$$\frac{\partial E}{\partial b_3} = \frac{\partial E}{\partial \theta} \frac{\partial \theta}{\partial z_0} \frac{\partial z_0}{\partial b_3} = (1.86)(0.57)(0.43) \\ = \boxed{0.46} .$$

$$\frac{\partial E}{\partial w_1} = \frac{\partial E}{\partial h_1} \frac{\partial h_1}{\partial z_{h_1}} \frac{\partial z_{h_1}}{\partial w_1}$$

$$\frac{\partial E}{\partial h_1} = \frac{\partial E}{\partial \theta} \frac{\partial \theta}{\partial z_0} \frac{\partial z_0}{\partial h_1} = (1.86)(0.57)(0.43)(-0.25) \\ = - \underline{\underline{0.11}}$$

$$\frac{\partial E}{\partial w_1} = (-0.11)(0.68)(0.32)(0.18) = \boxed{0.0043}$$

$$\frac{\partial E}{\partial w_3} = \frac{\partial E}{\partial h_1} \frac{\partial h_1}{\partial z_{h_1}} \frac{\partial z_{h_1}}{\partial w_3} = (-0.11)(0.68)(0.32)(2.42) \\ = \boxed{-0.0579}$$

$$\frac{\partial E}{\partial w_5} = \frac{\partial E}{\partial h_1} \frac{\partial h_1}{\partial z_{h_1}} \frac{\partial z_{h_1}}{\partial w_5} = (-0.11)(0.68)(0.32)(1.34) \\ = \boxed{-0.0321}$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial h_2} \frac{\partial h_2}{\partial z_{h_2}} \frac{\partial z_{h_2}}{\partial w_2}$$

$$\frac{\partial E}{\partial h_2} = \frac{\partial E}{\partial \theta} \frac{\partial \theta}{\partial z_0} \frac{\partial z_0}{\partial h_2} = (1.86)(0.57)(0.43)(0.68) \\ = \boxed{0.31//}$$

$$\frac{\partial E}{\partial w_2} = \frac{\partial E}{\partial h_2} \frac{\partial h_2}{\partial z_{h_2}} \frac{\partial z_{h_2}}{\partial w_2} = (0.31)(0.29)(0.71)(0.18) = \boxed{0.01148}$$

$$\frac{\partial E}{\partial w_4} = \frac{\partial E}{\partial h_2} \frac{\partial h_2}{\partial z_{h_2}} \frac{\partial z_{h_2}}{\partial w_4} = (0.31)(0.29)(0.71)(2.42) = \boxed{0.1544}$$

$$\frac{\partial E}{\partial w_6} = \frac{\partial E}{\partial h_2} \frac{\partial h_2}{\partial z_{h_2}} \frac{\partial z_{h_2}}{\partial w_6} = (0.31)(0.29)(0.71)(1.34) = \boxed{0.0855}$$

$$\frac{\partial E}{\partial b_1} = \frac{\partial E}{\partial h_1} \frac{\partial h_1}{\partial z_{h_1}} \frac{\partial z_{h_1}}{\partial b_1} = (-0.11)(0.68)(0.32)(1) = \boxed{-0.0239}$$

$$\frac{\partial E}{\partial b_2} = \frac{\partial E}{\partial h_2} \frac{\partial h_2}{\partial z_{h_2}} \frac{\partial z_{h_2}}{\partial b_2} = (0.31)(0.29)(0.71)(1) = \boxed{0.0638}$$

Thus summarizing all results :

$$\nabla W_1 = \begin{pmatrix} 0.0043 & 0.0115 \\ -0.0579 & 0.1544 \\ -0.0321 & 0.0855 \end{pmatrix}$$

$$X = \begin{pmatrix} 0.18 \\ 2.42 \\ 1.31 \end{pmatrix}$$

$$\nabla B_1 = \begin{pmatrix} -0.0239 \\ 0.0638 \end{pmatrix}$$

$$\nabla B_2 = (0.46)$$

$$\text{Error} = 1.73$$

(mean square error)

$$\nabla W_2 = \begin{pmatrix} 0.31 \\ 0.13 \end{pmatrix}$$