

HW 3 Solution

1 Problem 1 (20 points)

We need to apply the backward propagation from the output layer toward the input layer. From the dimensions given in the problem, we note that $\mathbf{W}^1 \in \mathbb{R}^{N \times N_1}$, $\mathbf{W}_0^1 \in \mathbb{R}^{N_1}$, $\mathbf{W}^2 \in \mathbb{R}^{N_1 \times K}$, $\mathbf{W}_0^2 \in \mathbb{R}^K$, and $\hat{\mathbf{y}} \in \mathbb{R}^K$. Since the loss function is not given, I consider a general cost function defined as $\ell(\mathbf{y}, \hat{\mathbf{y}})$. For the output layer we have

$$\nabla \mathbf{W}^2 = \frac{\partial \ell(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{W}^2} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{W}^2} \frac{\partial \ell(\mathbf{y}, \hat{\mathbf{y}})}{\partial \hat{\mathbf{y}}} = \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{W}^2} \nabla \hat{\mathbf{y}} = \sum_{k=1}^K \frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}^2} \nabla \hat{\mathbf{y}}[k] . \quad (1)$$

We know that $\hat{\mathbf{y}}[k] = \sigma_K((\mathbf{W}_k^2)^T + \mathbf{W}_0^2)$ where \mathbf{W}_k^2 is the k -th column of \mathbf{W}_k^2 , and also

$$\frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}^2} = \left[\frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_1^2}, \frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_2^2}, \dots, \frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_K^2} \right] . \quad (2)$$

On the other hand, from the previous homework, we have the partial derivative of softmax function. Therefore,

$$\frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_i^2} = \begin{cases} \hat{\mathbf{y}}[k](1 - \hat{\mathbf{y}}[k])\mathbf{H} & \text{if } i = k \\ -\hat{\mathbf{y}}[k]\hat{\mathbf{y}}[i]\mathbf{H} & \text{if } i \neq k \end{cases} . \quad (3)$$

For the bias vector \mathbf{W}_0^2 the same procedure should be performed where it can be observed the the only difference in that the bias is not multiplied by \mathbf{H} and therefore \mathbf{H} will not appear in the derivatives. Hence,

$$\nabla \mathbf{W}_0^2 = \sum_{k=1}^K \frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_0^2} \nabla \hat{\mathbf{y}}[k] = \sum_{k=1}^K \left[\frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_0^2[1]}, \frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_0^2[2]}, \dots, \frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_0^2[K]} \right] \nabla \hat{\mathbf{y}}[k] , \quad (4)$$

where

$$\frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{W}_0^2[i]} = \begin{cases} \hat{\mathbf{y}}[k](1 - \hat{\mathbf{y}}[k]) & \text{if } i = k \\ -\hat{\mathbf{y}}[k]\hat{\mathbf{y}}[i] & \text{if } i \neq k \end{cases} . \quad (5)$$

In order to calculate the partial derivatives for the weights of the hidden layer, first we need to calculate $\nabla \mathbf{H}$.

$$\nabla \mathbf{H} = \sum_{k=1}^K \frac{\partial \hat{\mathbf{y}}[k]}{\partial \mathbf{H}} \nabla \hat{\mathbf{y}}[k] = \sum_{k=1}^K \left[\mathbf{W}_k^2 \hat{\mathbf{y}}[k] - \sum_{j=1}^K \mathbf{W}_j^2 \hat{\mathbf{y}}[k] \hat{\mathbf{y}}[j] \right] \nabla \hat{\mathbf{y}}[k] . \quad (6)$$

By applying the chain rule, we have

$$\nabla \mathbf{W}^1 = \frac{\partial \mathbf{H}}{\partial \mathbf{W}^1} \nabla \mathbf{H} = \sum_{k=1}^K \frac{\partial \mathbf{H}[k]}{\partial \mathbf{W}^1} \nabla \mathbf{H}[k] \quad (7)$$

$$= \sum_{k=1}^K \left[\frac{\partial \mathbf{H}[k]}{\partial \mathbf{W}_i^1}, \frac{\partial \mathbf{H}[k]}{\partial \mathbf{W}_2^1}, \dots, \frac{\partial \mathbf{H}[k]}{\partial \mathbf{W}_{N_1}^1} \right] \nabla \mathbf{H}[k] , \quad (8)$$

where

$$\frac{\partial \text{ReLU}((\mathbf{W}_k^1)^T \mathbf{X} + \mathbf{W}_0^1)}{\partial \mathbf{W}_i^1} = \begin{cases} \mathbf{X}[k] & \text{if } i = k \text{ and } (\mathbf{W}_i^1)^T \mathbf{X} + \mathbf{W}_0^1[i] > 0 \\ 0 & \text{Otherwise} \end{cases} . \quad (9)$$

For the bias term, we have

$$\nabla \mathbf{W}_0^1 = \sum_{k=1}^K \left[\frac{\partial \mathbf{H}[k]}{\partial \mathbf{W}_0^1[1]}, \frac{\partial \mathbf{H}[k]}{\partial \mathbf{W}_0^1[2]}, \dots, \frac{\partial \mathbf{H}[k]}{\partial \mathbf{W}_0^1[N_1]} \right]^T \nabla \mathbf{H}[k] , \quad (10)$$

where

$$\frac{\partial \text{ReLU}((\mathbf{W}_k^1)^T \mathbf{X} + \mathbf{W}_0^1)}{\partial \mathbf{W}_0^1[i]} = \begin{cases} 1 & \text{if } i = k \text{ and } (\mathbf{W}_i^1)^T \mathbf{X} + \mathbf{W}_0^1[i] > 0 \\ 0 & \text{Otherwise} \end{cases} . \quad (11)$$

Correctly calculating each of the gradient of \mathbf{W}^1 and \mathbf{W}^2 has **5 points**, each of the bias terms will earn you **3 points**, and you will get **4 points** for correctly using the chain rule and having the right dimensions for the gradients.

2 Problem 2 (30 points)

2.1

Given the input value $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$, and the initialization of the weight matrix for the first layer W^1 :

$$W^1 = \begin{bmatrix} 3 & 6 \\ 4 & 5 \end{bmatrix}; W_0^1 = \begin{bmatrix} 1 \\ -6 \end{bmatrix}. \quad (12)$$

The value of the hidden nodes can be calculated as:

$$\mathbf{H} = \phi((W^1)^t \mathbf{x} + W_0^1) \quad (13)$$

while $\phi(z) = \frac{1}{1+\exp(-z)}$.

Thus, the value of the hidden nodes are:

$$\begin{aligned} \mathbf{H} &= \phi\left(\begin{bmatrix} 3 & 4 \\ 6 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ -6 \end{bmatrix}\right) \\ &= \begin{bmatrix} 0.982 \\ 0.5 \end{bmatrix} \end{aligned} \quad (14)$$

Using multi-class sigmoid function as the output function, the predicted output value $\hat{\mathbf{y}}$ can be calculated as:

$$\hat{\mathbf{y}} = \sigma_M((W^2)^t \mathbf{H} + W_0^2) \quad (15)$$

while $\sigma_M(z_k) = \frac{\exp((W_k^2)^t \mathbf{H} + W_{0,k}^2)}{\sum_k \exp((W_k^2)^t \mathbf{H} + W_{0,k}^2)}$.

Thus, the value of the predicted output nodes are:

$$\begin{aligned} \hat{\mathbf{y}} &= \sigma_M\left(\begin{bmatrix} 2 & 4 \\ 3 & 3 \end{bmatrix} \begin{bmatrix} 0.982 \\ 0.5 \end{bmatrix} + \begin{bmatrix} -1 \\ -2 \end{bmatrix}\right) \\ &= \begin{bmatrix} 0.6267 \\ 0.3733 \end{bmatrix} \end{aligned} \quad (16)$$

When loss function is chosen as squared loss function, the gradient of the output $\nabla \hat{\mathbf{y}}$ is:

$$\begin{aligned} \nabla \hat{\mathbf{y}} &= \frac{\partial l(\mathbf{y}[m], \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}[m]} \\ &= \frac{1}{2} \frac{\partial (\mathbf{y}[m] - \hat{\mathbf{y}}[m])^t (\mathbf{y}[m] - \hat{\mathbf{y}}[m])}{\partial \hat{\mathbf{y}}[m]} \\ &= -(\mathbf{y}[m] - \hat{\mathbf{y}}[m]) \end{aligned} \quad (17)$$

which is $-\begin{bmatrix} 0 - 0.6267 \\ 1 - 0.3733 \end{bmatrix} = \begin{bmatrix} 0.6267 \\ -0.6267 \end{bmatrix}$ in this case.

2.2

Output Layer For the output layer, the gradient of the weight matrix can be calculated as:

$$\begin{aligned} \nabla W^2 &= \frac{\partial \hat{\mathbf{y}}}{\partial W^2} \nabla \mathbf{y} \\ &= \sum_{i=1}^2 \frac{\partial \hat{y}[i]}{\partial W^2} \nabla y[i] \\ &= \sum_{i=1}^2 \begin{bmatrix} \frac{\partial \hat{y}[i]}{\partial W_1^2} & \frac{\partial \hat{y}[i]}{\partial W_2^2} \end{bmatrix} \nabla y[i] \end{aligned} \quad (18)$$

Thus, the gradient of W^2 is:

$$\begin{aligned} \nabla W^2 &= [\hat{y}[1](1 - \hat{y}[1])\mathbf{H} \quad -\hat{y}[1]\hat{y}[2]\mathbf{H}] \nabla \hat{y}[1] + [-\hat{y}[1]\hat{y}[2]\mathbf{H} \quad \hat{y}[2](1 - \hat{y}[2])\mathbf{H}] \nabla y[2] \\ &= \begin{bmatrix} 0.6267 * 0.3733 \begin{bmatrix} 0.982 \\ 0.5 \end{bmatrix} & -0.6267 * 0.3733 \begin{bmatrix} 0.982 \\ 0.5 \end{bmatrix} \end{bmatrix} * 0.6267 \\ &+ \begin{bmatrix} -0.6267 * 0.3733 \begin{bmatrix} 0.982 \\ 0.5 \end{bmatrix} & 0.6267 * 0.3733 \begin{bmatrix} 0.982 \\ 0.5 \end{bmatrix} \end{bmatrix} * (-0.6267) \\ &= \begin{bmatrix} 0.2880 & -0.2880 \\ 0.1466 & -0.1466 \end{bmatrix} \end{aligned} \quad (19)$$

The gradient of the weight bias vector:

$$\nabla W_0^2 = \begin{bmatrix} 0.2932 \\ -0.2932 \end{bmatrix} \quad (20)$$

So, the gradient of the Hidden nodes $\nabla \mathbf{H}$ is:

$$\begin{aligned}
\nabla \mathbf{H} &= \frac{\partial \sigma_M(\mathbf{z})}{\partial \mathbf{H}} \nabla \mathbf{y} \\
&= \sum_{i=1}^2 \frac{\partial \hat{y}[i]}{\partial \mathbf{H}} \nabla y[i] \\
&= \hat{y}[1][W_1^2 - \sum_{j=1}^2 \hat{y}[j]W_j^2] \nabla y[1] + \hat{y}[2][W_2^2 - \sum_{j=1}^2 \hat{y}[j]W_j^2] \nabla y[2] \\
&= 0.6267 * \left[\begin{bmatrix} 2 \\ 4 \end{bmatrix} - (0.6267 * \begin{bmatrix} 2 \\ 4 \end{bmatrix} + 0.3733 * \begin{bmatrix} 3 \\ 3 \end{bmatrix}) \right] \\
&\quad + 0.3733 * \left[\begin{bmatrix} 3 \\ 3 \end{bmatrix} - (0.6267 * \begin{bmatrix} 2 \\ 4 \end{bmatrix} + 0.3733 * \begin{bmatrix} 3 \\ 3 \end{bmatrix}) \right] \\
&= \begin{bmatrix} -0.2932 \\ 0.2932 \end{bmatrix}
\end{aligned} \tag{21}$$

Hidden Layer Given $\nabla \mathbf{H}$, the gradient of the first weight matrix is:

$$\begin{aligned}
\nabla W^1 &= \frac{\partial \mathbf{H}}{\partial W^1} \nabla \mathbf{H} \\
&= \sum_{i=1}^2 \frac{\partial H[i]}{\partial W^1} \nabla H[i] \\
&= \sum_{i=1}^2 \left[\frac{\partial H[i]}{\partial W_1^1} \quad \frac{\partial H[i]}{\partial W_2^1} \right] \nabla H[i] \\
&= \begin{bmatrix} H[1](1 - H[1]) \begin{bmatrix} 1 \\ 0 \end{bmatrix} & 0 \end{bmatrix} * (-0.2932) + \begin{bmatrix} 0 & H[2](1 - H[2]) \begin{bmatrix} 1 \\ 0 \end{bmatrix} \end{bmatrix} * 0.2932 \\
&= \begin{bmatrix} -0.0052 & 0.0733 \\ 0 & 0 \end{bmatrix}
\end{aligned} \tag{22}$$

The gradient of the bias weight vector of the first layer is:

$$\nabla W_0^1 = \begin{bmatrix} -0.0052 \\ 0.0733 \end{bmatrix} \tag{23}$$

2.3

Weight Update The weight update equation is:

$$W = W - r * \nabla W \tag{24}$$

Update the weight matrix accordingly:

$$\begin{aligned}
W_{new}^2 &= W^2 - r * \nabla W^2 \\
&= \begin{bmatrix} 2 & 3 \\ 4 & 3 \end{bmatrix} - 0.5 * \begin{bmatrix} 0.2880 & -0.2880 \\ 0.1466 & -0.1466 \end{bmatrix} \\
&= \begin{bmatrix} 1.8560 & 3.1440 \\ 3.9267 & 3.0733 \end{bmatrix}
\end{aligned} \tag{25}$$

The updated weight bias vector is:

$$\begin{aligned}
W_{0,new}^2 &= W_0^2 - r * \nabla W_0^2 \\
&= \begin{bmatrix} -1 \\ -2 \end{bmatrix} - 0.5 * \begin{bmatrix} 0.2932 \\ -0.2932 \end{bmatrix} \\
&= \begin{bmatrix} -1.1466 \\ -1.8534 \end{bmatrix}
\end{aligned} \tag{26}$$

The updated weight matrix for the first layer is:

$$\begin{aligned}
W_{new}^1 &= W^1 - r * \nabla W^1 \\
&= \begin{bmatrix} 3 & 6 \\ 4 & 5 \end{bmatrix} - 0.5 * \begin{bmatrix} -0.0052 & 0.0733 \\ 0 & 0 \end{bmatrix} \\
&= \begin{bmatrix} 3.0026 & 5.9633 \\ 4 & 5 \end{bmatrix}
\end{aligned} \tag{27}$$

Weight bias vector of the first layer:

$$\begin{aligned}
W_{0,new}^1 &= W_0^1 - r * \nabla W_0^1 \\
&= \begin{bmatrix} 1 \\ -6 \end{bmatrix} - 0.5 * \begin{bmatrix} -0.0052 \\ 0.0733 \end{bmatrix} \\
&= \begin{bmatrix} 1.0026 \\ -6.0367 \end{bmatrix}
\end{aligned} \tag{28}$$

Update Output Value

Case 1 After the updating of the weight matrix of the first layer, the value of hidden nodes are:

$$\begin{aligned}
\mathbf{H} &= \phi((W_{new}^1)^t \mathbf{x} + W_{0,new}^1) \\
&= \begin{bmatrix} 0.9821 \\ 0.4817 \end{bmatrix}
\end{aligned} \tag{29}$$

The updated predicted output value $\hat{\mathbf{y}}$ is:

$$\begin{aligned}\hat{\mathbf{y}} &= \sigma_M((W_{new}^2)^t \mathbf{H} + W_{0,new}^2) \\ &= \begin{bmatrix} 0.4633 \\ 0.5367 \end{bmatrix}\end{aligned}\quad (30)$$

The loss function value before updating is:

$$\begin{aligned}& \frac{1}{2}(\mathbf{y}[m] - \hat{\mathbf{y}}[m])^t(\mathbf{y}[m] - \hat{\mathbf{y}}[m]) \\ &= \frac{1}{2} \sum_{i=1}^2 \nabla y[i]^2 \\ &= 0.3927\end{aligned}\quad (31)$$

After updating, the loss function value is:

$$\begin{aligned}& \frac{1}{2} \sum_{i=1}^2 \nabla y[i]^2 \\ &= 0.2146\end{aligned}\quad (32)$$

It is clear that the loss function value has been reduced.

Case 2 If the weight matrix of the first layer is fixed, the value of the hidden nodes are also fixed, the updated predicted output value should be:

$$\begin{aligned}\hat{\mathbf{y}} &= \sigma_M((W_{new}^2)^t \mathbf{H} + W_{0,new}^2) \\ &= \begin{bmatrix} 0.4672 \\ 0.5328 \end{bmatrix}\end{aligned}\quad (33)$$

From Case 1, we know the loss function value before updating is 0.3927, After updating weight matrix in the second layer, the loss function value has been changed to:

$$\begin{aligned}& \frac{1}{2} \sum_{i=1}^2 \nabla y[i]^2 \\ &= 0.2183\end{aligned}\quad (34)$$

Also, the loss function value has been reduced.