

CSCE 636: Homework 2

Aryan Sharma
UIN: 326006767

Question 1

1a. We need to flatten the data set loaded as the data that is loaded from MNIST are 2D images. We convert this to a vector because it makes the computation easier and also conveniently fits into the vectorized computation as required in tensorflow.

1b. Implemented in code

1c. Implemented in code

1d. Implemented in code

1e. Epoch was trained as the hyperparameter. The best setting came at

Question 2

$$\begin{aligned}\frac{\partial E_{in}}{\partial w} &= \frac{\partial}{\partial w} \left(\frac{1}{N} \sum_{n=1}^N (\tanh(\mathbf{w}^T \mathbf{x}_n) - y_n)^2 \right) \\ &= \frac{2}{N} \sum_{n=1}^N (\tanh(\mathbf{w}^T \mathbf{x}_n) - y_n) \frac{\partial}{\partial w} (\tanh(\mathbf{w}^T \mathbf{x}_n) - y_n) \\ &= \frac{2}{N} \sum_{n=1}^N (\tanh(\mathbf{w}^T \mathbf{x}_n) - y_n)(1 - \tanh^2(\mathbf{w}^T \mathbf{x}_n)) \frac{\partial}{\partial w} (\mathbf{w}^T \mathbf{x}_n) \\ &= \frac{2}{N} \sum_{n=1}^N (\tanh(\mathbf{w}^T \mathbf{x}_n) - y_n)(1 - \tanh^2(\mathbf{w}^T \mathbf{x}_n)) \mathbf{x}_n\end{aligned}$$

When the weights go to infinity, the gradient goes to 0. This exploding gradient problem would then cause the perceptron to always have large positive activation and thus would fail to classify negative data points. Also since the gradient is 0, no extra training iteration would update its weight and bring down. This is a common problem in perceptron which classifies all the data

points as positive when the gradient explodes. To solve this the model should have a regularization parameter.

Question 3

Given, the values we have,

$x^{(0)}$	$s^{(1)}$	$x^{(1)}$	$s^{(2)}$	$x^{(2)}$	$s^{(3)}$	$x^{(3)}$
$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$	$\begin{bmatrix} 0.7 \\ 1 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 0.6 \\ 0.76 \end{bmatrix}$	$\begin{bmatrix} -4.48 \end{bmatrix}$	$\begin{bmatrix} 1 \\ -0.9 \end{bmatrix}$	$\begin{bmatrix} -0.8 \end{bmatrix}$	-0.8

$$g^{(3)} = [2(-0.8 - 1)] = [-3.6]$$

$$g^{(2)} = [(1 - 0.9^2) * 2 * -3.6] = [-1.368]$$

$$\begin{aligned} g^{(1)} &= \begin{bmatrix} 1 - 0.6^2 \\ 1 - 0.76^2 \end{bmatrix} \otimes \left(\begin{bmatrix} 1 \\ -3 \end{bmatrix} \begin{bmatrix} -1.368 \end{bmatrix} \right) \\ &= \begin{bmatrix} -0.875 \\ 1.733 \end{bmatrix} \end{aligned}$$

$$\frac{\partial e}{\partial w^{(3)}} = x^{(2)}(g^{(3)})^T = \begin{bmatrix} -3.6 \\ 3.24 \end{bmatrix}$$

$$\frac{\partial e}{\partial w^{(2)}} = x^{(1)}(g^{(2)})^T = \begin{bmatrix} -1.368 \\ -0.821 \\ -1.039 \end{bmatrix}$$

$$\frac{\partial e}{\partial w^{(1)}} = x^{(0)}(g^{(1)})^T = \begin{bmatrix} -0.875 & 1.733 \\ -1.75 & 3.466 \end{bmatrix}$$

Figure 1: Question 3

Question 4

Given,

$$E(w) = (w - w^*)^T Q(w - w^*)$$

$$\frac{\partial}{\partial w} E(w) = (Q + Q^T)(w - w^*)$$

Since Q is positive definite matrix,

$$\frac{\partial}{\partial w} E(w) = (2Q)(w - w^*) \quad (1)$$

Putting $w = 0$, we have

$$\frac{\partial}{\partial w} E(w) = -2Qw^*$$

Setting Equation 1 to 0, we get the $w = w^*$. Hence the optimal value is w^* and these are the weights that minimize w .

No, the gradient does not move in the direction of the optimal weights as when w is 0, we can see the gradient to negative when w^* is positive and positive otherwise. Hence, owing to the gradient, the optimization does not move in the correct direction.

The more big the step the more oscillation it causes while reaching the optimal value. In cases with very high learning rate, it may diverge.