

Assignment 2 Songhwan Kim

Q1

Input \rightarrow Hidden 1 \rightarrow output 2 layer.

Let x be the input

y be the output

First layer $\rightarrow w_1$, bias $\rightarrow b_1$

Second layer $\rightarrow w_2$, bias $\rightarrow b_2$

the output layer 1 is:

$$z_1 = w_1 x + b_1$$

$$a_1 = \sigma(z_1) \text{ where } \sigma \text{ is activation function. (Sigmoid)}$$

The output layer 2 is:

$$\bar{z}_2 = w_2 a_1 + b_2$$

$$a_2 = \bar{z}_2$$

The output of the neural network is

$$\hat{y} = a_2$$

Task: Back propagation for regression where loss function is the Mean Square Error loss.

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

To learn the parameters.

① Provide random values for weights w_1, w_2 and biases b_1, b_2

② Update the weights using gradient descent by

$$w_i := w_i - \alpha \cdot \frac{\partial L}{\partial w_i} \text{ where } i \text{ refers to the } i\text{th layer}$$

$$b_i := b_i - \alpha \cdot \frac{\partial L}{\partial b_i}$$

③ Repeat until convergence.

To find $\frac{\partial L}{\partial w_2}$ to update the second layer.

$$\frac{\partial L}{\partial w_2} = \frac{\partial}{\partial w_2} \left[\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right] \quad \bullet \text{ L is MSE Loss function.}$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial w_2} (y_i - \hat{y}_i)^2$$

$$= \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \frac{\partial}{\partial w_2} (y_i - \hat{y}_i) = \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \left(\frac{\partial}{\partial w_2} y_i - \frac{\partial}{\partial w_2} \hat{y}_i \right)$$

$$= \frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \left(0 - \frac{\partial}{\partial w_2} \hat{y}_i \right) = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \left(\frac{\partial}{\partial w_2} \hat{y}_i \right)$$

$$= -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \cdot a_1$$

$$\rightarrow \frac{\partial}{\partial w_2} \hat{y} = \frac{\partial \hat{y}}{\partial w_2} \quad \hat{y} = a_2 = z_2$$

$$\text{Similarly } \frac{\partial L}{\partial b_2} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$

$$= \frac{\partial a_2}{\partial w_2} = \frac{\partial z_2}{\partial w_2}$$

↳ identity function for output layer

$$= \frac{\partial}{\partial w_2} (w_2 a_1 + b_2)$$

- $z_2 = w_2 a_1 + b_2$ from previous page

$$= a_1$$

To find $\frac{\partial L}{\partial w_1}$, we use chain rule.

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$a_2 = z_2$$

a_2 is dependent on z_2

$$z_2 = w_2 a_1 + b_2$$

z_2 is dependent on a_1

$$a_1 = \sigma(z_1)$$

a_1 is dependent on z_1

$$z_1 = w_1 x + b_1$$

Rewrite it as

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

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

$$w_2$$

$$\sigma'(z_1)$$

$$x$$

Therefore,

$$\frac{\partial L}{\partial w_1} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \cdot w_2 \cdot \sigma'(z_1) \cdot x$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$$

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

$$a_1$$

Similarly

$$\frac{\partial L}{\partial b_1} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \cdot w_2 \cdot \sigma'(z_1)$$

The difference from the update rule for the network trained for binary classification using log loss.

For classification problem

$$\frac{\partial L}{\partial w_1} = (a_2 - y) \cdot \underline{w_2 \cdot \sigma'(z_1)} \cdot x$$

For Regression problem //

$$\frac{\partial L}{\partial w_1} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \underline{w_2 \cdot \sigma'(z_1)} \cdot x$$

The difference of update rule is that the classification problem uses log loss to output a probability value between 0 and 1. In the other hands, the regression problem uses MSE loss to find an optimal value.