

Assignment - 2

CS 4783/5783

27th Sept 2021

Nagavathi Bhavani Santhoshi Acharya

A20314248

- ① Derive the Update rule and show how to train a 3-layer (1-input layer, 1-hidden layer, 1-Output layer) Network with Backpropagation for regression using the mean square Error loss. Assume that you are using the Sigmoid activation function for the hidden layer. Explain Briefly how this is different from the Update rule for the network trained for Binary Classification using log loss.

Backpropagation for regression Using the mean Square Error loss

$x \rightarrow \text{input}$ $n \times f$

$y \rightarrow \text{Labels}$ $n \times 1$

First layer weights: w_1 Bias: b_1

Second layer weights: w_2 Bias: b_2

Third layer weights: w_3 Bias: b_3

O/p of first layer: $z_1 = w_1 x + b_1$

O/p of After Activation function: $a_1 = g(z_1)$

O/p of Second layer: $z_2 = w_2 a_1 + b_2$

After Activation function: $a_2 = g(z_2)$

O/p of third layer: $z_3 = w_3 a_2 + b_3$

Final prediction: $\hat{y} = a_3 = g(z_3)$

$$L = - [(1-y) \log(1-\hat{y}) + y \cdot \log \hat{y}]$$

We update the weights using gradient descent

① Start with an initial guess for w_1, w_2, w_3
 b_1, b_2, b_3

② Update weights by

$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$$

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

Where i refers to the i^{th} layer

③ Repeat Until Convergence

\Rightarrow To find $\frac{\partial L}{\partial w_3}$

$$\frac{\partial L}{\partial w_3} = - \frac{\partial}{\partial w_3} ((1-y) \cdot \log(1-\hat{y}) + y \cdot \log \hat{y})$$

$$= -(1-y) \cdot \frac{\partial}{\partial w_3} (1-\hat{y}) - y \cdot \frac{\partial}{\partial w_3} \log \hat{y}$$

$$\because \hat{y} = a_3 = g(z_3)$$

$$\frac{\partial L}{\partial w_3} = -(1-y) \cdot \frac{\partial}{\partial w_3} (1 - \log(g(z_3))) - y \cdot \frac{\partial}{\partial w_3} \log(g(z_3))$$

$$\because \frac{\partial}{\partial z} \log(z) = \frac{1}{z}$$

$$\begin{aligned}\frac{\partial L}{\partial w_3} &= \frac{-(1-y)}{1-g(z_3)} \left(-\frac{\partial}{\partial w_3} (g(z_3)) \right) - \frac{y}{g(z_3)} \frac{\partial}{\partial w_3} g(z_3) \\ &= \frac{(1-y)}{1-g(z_3)} \cdot g'(z_3) - \frac{y}{g(z_3)} \cdot g'(z_3)\end{aligned}$$

\therefore for Sigmoid Activation
 $g'(z) = g(z)(1-g(z))$

$$\begin{aligned}\frac{\partial L}{\partial w_3} &= \frac{(1-y)}{1-g(z_3)} \cdot g(z_3)(1-g(z_3)) \frac{\partial z_3}{\partial w_3} \\ &\quad - \frac{y}{g(z_3)} \cdot g(z_3)(1-g(z_3)) \frac{\partial z_3}{\partial w_3}\end{aligned}$$

$$= ((1-y) \cdot g(z_3) - y(1-g(z_3))) \frac{\partial z_3}{\partial w_3}$$

$$= (g(z_3) - y) \frac{\partial z_3}{\partial w_3}$$

$$= [a_3 - y] \cdot \frac{\partial}{\partial w_3} (w_3 a_2 + b_3)$$

$$\frac{\partial L}{\partial w_3} = (a_3 - y) \cdot a_2^T$$

Similarly, when we diff L w.r.t b_3 we get

$$\boxed{\frac{\partial L}{\partial b_3} = a_3 - y}$$

To derive $\frac{\partial L}{\partial w_2}$, we can use the "Chain rule"

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial w_2}$$

$$= \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_2}$$

$$\because a_3 = g(z_3)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_2} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial w_2}$$

$$\because z_3 = w_3 a_2 + b_3$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial w_2}$$

$$\because a_2 = g(z_2) = g(w_2 a_1 + b_2)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$$

$(a_2 - y)$ $\underbrace{\frac{\partial z_2}{\partial w_2}}_{a_1}$

$$\therefore \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} = a_2 - y$$

$$z_3 = w_3 a_2 + b_3$$

$$\therefore \frac{\partial z_3}{\partial a_2} = w_3$$

$$a_2 = g(z_2)$$

$$\therefore \frac{\partial a_2}{\partial z_2} = g'(z_2)$$

$$\therefore z_2 = w_2 a_1 + b_2$$

$$\frac{\partial z_2}{\partial w_2} = a_1$$

$$\therefore \frac{\partial L}{\partial w_2} = (a_3 - y) \cdot w_3 \cdot g'(z_2) a_1$$

Consider an Example.

But Dimensions of each Matrix don't match

$$a_3 - y \rightarrow 1 \times 1$$

$$w_3 \rightarrow 1 \times h_2$$

$$g'(z_2) \rightarrow h_2 \times 1$$

$$a_1 \rightarrow h_1 \times 1$$

$$\therefore \frac{\partial L}{\partial w_2} = w_3^T \cdot g'(z_2) (a_3 - y) \cdot a_1^T$$

Similarly,

$$\frac{\partial L}{\partial b_2} = w_3^T \cdot g'(z_2) (a_3 - y)$$

Using chain rule,

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial w_1} = (a_3 - y) \cdot w_3 \cdot g'(z_2) \cdot w_2 \cdot g'(z_1) x$$

2.1 Report the Average MSE and the accuracy.

MSE 2858.4015

Accuracy 0.6107

2.2 Plot the loss and accuracy as a function of no. of iterations plotted in .ipynb file

2.3 What is the effect of learning rate on the training process?
Shown in the .ipynb file

2.4 What is the effect of number of neurons in hidden layer?
Shown in the .ipynb file

2.5 What is the effect of Activation function in the network?
Used tanh and ReLU, Shown in .ipynb file