

Neural Network Maths

Nirbhay Sharma, IISc Bangalore

December 11, 2024

1 Some important activation functions and their derivatives

1. Sigmoid ($\sigma(x) = \frac{1}{1+e^{-x}}$)
2. Tanh ($\tanh(x) = \frac{e^{2x}-1}{e^{2x}+1}$)
3. Relu ($\text{ReLU}(x) = x$ if $x > 0$ else 0)
4. SiLU ($\text{SiLU}(x) = x * \sigma(x)$)
5. Softmax ($\text{softmax}(x) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$)

1.1 Derivative of Sigmoid

$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x)) \quad (1)$$

1.2 Derivative of Tanh

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

1.3 Derivative of Relu

$$\frac{d\text{ReLU}(x)}{dx} = 1 \text{ if } x > 0 \text{ else } 0 \quad (3)$$

1.4 Derivative of SiLU

$$\frac{d\text{SiLU}(x)}{dx} = \frac{dx\sigma(x)}{dx} = x\sigma(x)(1 - \sigma(x)) + \sigma(x) = \sigma(x)(1 + x(1 - \sigma(x))) \quad (4)$$

1.5 Derivative of Softmax

$$x = (x_1, x_2, \dots, x_n) \quad (5)$$

$$\text{Softmax}(x) = y_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}, i = 1 \dots n \quad (6)$$

$$\frac{\partial y_i}{\partial x_i} = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \left(1 - \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}\right) = y_i(1 - y_i) \quad (7)$$

$$\frac{\partial y_i}{\partial x_j} = \frac{-e^{x_i} e^{x_j}}{(\sum_{j=1}^n e^{x_j})^2} = -y_i y_j \quad (8)$$

overall derivative of output of softmax wrt input is as follows

$$\frac{dy}{dx} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \frac{\partial y_n}{\partial x_2} & \dots & \frac{\partial y_n}{\partial x_n} \end{pmatrix} = \begin{pmatrix} y_1(1-y_1) & -y_1y_2 & \dots & -y_1y_n \\ -y_1y_2 & y_2(1-y_2) & \dots & -y_2y_n \\ \vdots & \vdots & \ddots & \vdots \\ -y_1y_n & -y_2y_n & \dots & y_n(1-y_n) \end{pmatrix} \quad (9)$$

2 Forward Propagation

1. The neural network has layers with W, b as weight and biases (learnable parameters)
2. Neural network also uses activation functions to incorporate nonlinearity
3. A typical forward pass looks like **Softmax(Layer(tanh(Layer(sigmoid(Layer(x))))))**, x is the input, sigmoid, tanh, softmax all are activation functions, Layer is a simple hidden layer with W, b as learnable parameters and it uses the equation $Y = WX + b$
4. Forward propagation is relatively straight forward, the major part of neural networks lies in Back propagation which we see in next section

3 BackPropagation

It is the method to backpropagate the loss or the error term E through the network and calculate gradients; the gradients are calculated using the famous chain rule of differentiation.

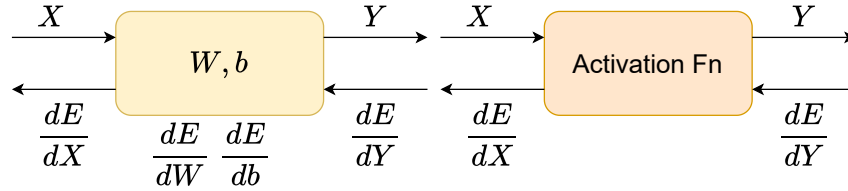


Figure 1: BackPropagation for a layer

1. In this layer, we have gradients coming from next layers, for a particular layer, we calculate the gradient wrt input, learnable parameters, then we flow back the gradients wrt input.
2. The layers having activation function does not have learnable parameters

3.1 Backpropagation for $Y = WX + b$

$$W = \begin{pmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mn} \end{pmatrix} \quad (10)$$

$$Y = (y_1, y_2, \dots, y_m) \quad (11)$$

$$X = (x_1, x_2, \dots, x_n) \quad (12)$$

$$b = (b_1, b_2, \dots, b_m) \quad (13)$$

$$y_i = \sum_{j=1}^n w_{ij}x_j + b_i \quad (14)$$

We now assume $\frac{\partial E}{\partial Y}$ is flowing from the next layer, now we want to calculate $\frac{\partial E}{\partial X}, \frac{\partial E}{\partial W}, \frac{\partial E}{\partial b}$

$$\frac{\partial E}{\partial X} = \left(\frac{\partial E}{\partial x_1}, \frac{\partial E}{\partial x_2}, \dots, \frac{\partial E}{\partial x_n} \right) \quad (15)$$

$$\frac{\partial E}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} w_{ij} \quad (16)$$

$$\frac{\partial E}{\partial X} = \left(\sum_{i=1}^n \frac{\partial E}{\partial y_i} w_{i1}, \sum_{i=1}^n \frac{\partial E}{\partial y_i} w_{i2}, \dots, \sum_{i=1}^n \frac{\partial E}{\partial y_i} w_{in} \right) = W^T \frac{\partial E}{\partial Y} \quad (17)$$

$$\frac{\partial E}{\partial b} = \left(\frac{\partial E}{\partial b_1}, \frac{\partial E}{\partial b_2}, \dots, \frac{\partial E}{\partial b_m} \right) \quad (18)$$

$$\frac{\partial E}{\partial b_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial b_j} = \frac{\partial E}{\partial y_j} \quad (19)$$

$$\frac{\partial E}{\partial b} = \frac{\partial E}{\partial Y} \quad (20)$$

$$\frac{\partial E}{\partial W} = \begin{pmatrix} \frac{\partial E}{\partial w_{11}} & \frac{\partial E}{\partial w_{12}} & \dots & \frac{\partial E}{\partial w_{1n}} \\ \frac{\partial E}{\partial w_{m1}} & \frac{\partial E}{\partial w_{m2}} & \dots & \frac{\partial E}{\partial w_{mn}} \end{pmatrix} \quad (21)$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial w_{ij}} = \frac{\partial E}{\partial y_i} x_j \quad (22)$$

$$\frac{\partial E}{\partial W} = \begin{pmatrix} \frac{\partial E}{\partial y_1} x_1 & \frac{\partial E}{\partial y_1} x_2 & \dots & \frac{\partial E}{\partial y_1} x_n \\ \frac{\partial E}{\partial y_m} x_1 & \frac{\partial E}{\partial y_m} x_2 & \dots & \frac{\partial E}{\partial y_m} x_n \end{pmatrix} = \frac{\partial E}{\partial Y} X^T \quad (23)$$

In summary,

$$\frac{\partial E}{\partial X} = W^T \frac{\partial E}{\partial Y}, \quad \frac{\partial E}{\partial b} = \frac{\partial E}{\partial Y}, \quad \frac{\partial E}{\partial W} = \frac{\partial E}{\partial Y} X^T \quad (24)$$

3.2 Backpropagation for Sigmoid $Y = \sigma(X)$

for activation function we only calculate $\frac{\partial E}{\partial X}$ as there are no learnable parameters

$$\frac{\partial E}{\partial X} = \left(\frac{\partial E}{\partial x_1}, \frac{\partial E}{\partial x_2}, \dots, \frac{\partial E}{\partial x_n} \right) \quad (25)$$

$$\frac{\partial E}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_j} = \frac{\partial E}{\partial y_j} y_j(1 - y_j) \quad (26)$$

$$\frac{\partial E}{\partial X} = \left(\frac{\partial E}{\partial y_1} y_1(1 - y_1), \frac{\partial E}{\partial y_2} y_2(1 - y_2), \dots, \frac{\partial E}{\partial y_n} y_n(1 - y_n) \right) \quad (27)$$

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} \odot Y(1 - Y) \quad (28)$$

\odot is element wise multiplication

3.3 Backpropagation for Tanh $Y = \tanh(X)$

for activation function we only calculate $\frac{\partial E}{\partial X}$ as there are no learnable parameters

$$\frac{\partial E}{\partial X} = (\frac{\partial E}{\partial x_1}, \frac{\partial E}{\partial x_2}, \dots, \frac{\partial E}{\partial x_n}) \quad (29)$$

$$\frac{\partial E}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_j} = \frac{\partial E}{\partial y_j} (1 - y_j^2) \quad (30)$$

$$\frac{\partial E}{\partial X} = (\frac{\partial E}{\partial y_1} (1 - y_1^2), \frac{\partial E}{\partial y_2} (1 - y_2^2), \dots, \frac{\partial E}{\partial y_n} (1 - y_n^2)) \quad (31)$$

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} \odot (1 - Y^2) \quad (32)$$

3.4 Backpropagation for Relu $Y = ReLU(X)$

for activation function we only calculate $\frac{\partial E}{\partial X}$ as there are no learnable parameters

$$\frac{\partial E}{\partial X} = (\frac{\partial E}{\partial x_1}, \frac{\partial E}{\partial x_2}, \dots, \frac{\partial E}{\partial x_n}) \quad (33)$$

$$\frac{\partial E}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_j} = \frac{\partial E}{\partial y_j} \quad y_j \neq 0 \quad (34)$$

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} \quad Y < 0 = 0 \quad (35)$$

3.5 Backpropagation for SiLU $Y = SiLU(X)$

for activation function we only calculate $\frac{\partial E}{\partial X}$ as there are no learnable parameters

$$\frac{\partial E}{\partial X} = (\frac{\partial E}{\partial x_1}, \frac{\partial E}{\partial x_2}, \dots, \frac{\partial E}{\partial x_n}) \quad (36)$$

$$\frac{\partial E}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_j} = \frac{\partial E}{\partial y_j} \sigma(x_j) (1 + x_j (1 - \sigma(x_j))) \quad (37)$$

$$\frac{\partial E}{\partial X} = (\frac{\partial E}{\partial y_1} \sigma(x_1) (1 + x_1 (1 - \sigma(x_1))), \dots, \frac{\partial E}{\partial y_n} \sigma(x_n) (1 + x_n (1 - \sigma(x_n)))) \quad (38)$$

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} \odot \sigma(X) (1 + X \odot (1 - \sigma(X))) \quad (39)$$

3.6 Backpropagation for Softmax $Y = Softmax(X)$

for activation function we only calculate $\frac{\partial E}{\partial X}$ as there are no learnable parameters

$$\frac{\partial E}{\partial X} = (\frac{\partial E}{\partial x_1}, \frac{\partial E}{\partial x_2}, \dots, \frac{\partial E}{\partial x_n}) \quad (40)$$

$$\frac{\partial E}{\partial x_j} = \sum_{i=1}^n \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_j} \quad (41)$$

Recall that

$$\frac{\partial y_i}{\partial x_j} = \begin{cases} y_i(1 - y_i) & i = j \\ -y_i y_j & i \neq j \end{cases} \quad (42)$$

$$\frac{\partial E}{\partial X} = \begin{pmatrix} \frac{\partial E}{\partial y_1} y_1(1-y_1) - \frac{\partial E}{\partial y_2} y_1 y_2 + \dots - \frac{\partial E}{\partial y_n} y_1 y_n \\ -\frac{\partial E}{\partial y_1} y_1 y_2 + \frac{\partial E}{\partial y_2} y_2(1-y_2) + \dots - \frac{\partial E}{\partial y_n} y_2 y_n \\ \vdots \\ -\frac{\partial E}{\partial y_1} y_1 y_n - \frac{\partial E}{\partial y_2} y_2 y_n + \dots + \frac{\partial E}{\partial y_n} y_n(1-y_n) \end{pmatrix} \quad (43)$$

$$= \begin{pmatrix} y_1(1-y_1) & -y_1 y_2 & \dots & -y_1 y_n \\ -y_1 y_2 & y_2(1-y_2) & \dots & -y_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ -y_1 y_n & -y_2 y_n & \dots & y_n(1-y_n) \end{pmatrix} \frac{\partial E}{\partial Y} \quad (44)$$

$$\text{define } M = \begin{pmatrix} y_1 & y_2 & \dots & y_n \\ y_1 & y_2 & \dots & y_n \\ \vdots & \vdots & \ddots & \vdots \\ y_1 & y_2 & \dots & y_n \end{pmatrix} \quad (45)$$

$$= \begin{pmatrix} y_1(1-y_1) & -y_1 y_2 & \dots & -y_1 y_n \\ -y_1 y_2 & y_2(1-y_2) & \dots & -y_2 y_n \\ \vdots & \vdots & \ddots & \vdots \\ -y_1 y_n & -y_2 y_n & \dots & y_n(1-y_n) \end{pmatrix} \frac{\partial E}{\partial Y} = M \odot (I - M^T) \frac{\partial E}{\partial Y} \quad (46)$$

3.7 Backpropagation for Cross Entropy Loss $E = -\sum y_i \log(\hat{y}_i)$

we only calculate derivative wrt \hat{y}_i which is straight forward

$$\frac{\partial E}{\partial \hat{y}_i} = \begin{cases} \frac{-1}{\hat{y}_i} & y_i = 1 \\ 0 & \text{otherwise} \end{cases} \quad (47)$$

The gradients from these loss will flow backwards