

# 1 Neural Network

Assume we have the following data point  $\mathbf{x}$ :

$$\mathbf{x} = \begin{bmatrix} 2 \\ 5 \end{bmatrix}$$

with corresponding binary label:

$$y = [1]$$

Which is part of a larger dataset  $\mathbf{X}$  with binary labels  $\mathbf{y}$ .

We want to create a neural network to solve this classification problem. We want to use squared error as the loss function:  $\ell = \frac{1}{2}(y - \hat{y})^2$ . (Please take a moment and think about why do we use this loss function.) We will use stochastic gradient descent to train our network by minimizing the loss. Suppose we want only one hidden layer with two neurons, sigmoid activation functions and we want to include all bias terms.

1. Draw what our neural network will look like.

2. What will be the shape of our weight matrices (Hint: we view  $\mathbf{x}$  as a column vector and compute the product as  $\mathbf{W}^T \mathbf{x}$ .)

(a)  $\mathbf{W}^1 =$

(b)  $\mathbf{W}^2 =$

$$\mathbf{W}^1 = 3 \times 2$$

$$\mathbf{W}^2 = 3 \times 1$$

Assume our weights are as follows:

$$\mathbf{W}^1 = \begin{bmatrix} -1 & 2 \\ -2 & 1 \\ 1 & 4 \end{bmatrix} \quad \mathbf{W}^2 = \begin{bmatrix} 1 \\ 2 \\ -4 \end{bmatrix}$$

3. What are the values being passed into our hidden neurons? (Hint: the bias term is considered the first feature.)

Recall that we need to add a one to our example to represent the bias term.

$$\boldsymbol{\alpha}^1 = \mathbf{W}^{1T} \mathbf{x} = \begin{bmatrix} -1 & -2 & 1 \\ 2 & 1 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} -1 - 4 + 5 \\ 2 + 2 + 20 \end{bmatrix} = \begin{bmatrix} 0 \\ 24 \end{bmatrix}$$

4. What are the outputs of our hidden neurons?

We need to apply the Sigmoid function to this vector to get the output, recall the Sigmoid function is:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\text{So our output will be } \boldsymbol{\beta}^1 = \begin{bmatrix} 0.5 \\ 0.999 \end{bmatrix}$$

Recall we need to add a +1 to the start of this vector to represent the Bias term in this layer!

5. What is the value being passed into our output layer?

$$\alpha^2 = \mathbf{W}^{2T} \boldsymbol{\beta}^1 = \begin{bmatrix} 1 & 2 & -4 \end{bmatrix} \begin{bmatrix} 1 \\ 0.5 \\ 0.999 \end{bmatrix} = 1 + 1 - 3.996 = -1.996$$

6. What is our final output  $\hat{y}$ ?

Recall we need to apply the Sigmoid function to this value  $\alpha^2$ .

$$\hat{y} = \frac{1}{1+e^{1.996}} = 0.119623$$

7. What does our model classify  $\mathbf{x}_1$  as?

0

8. What is the loss  $\ell$  on this example?

$$\ell = \frac{1}{2}(y - \hat{y})^2 = \frac{1}{2}(1 - 0.119623)^2 = 0.38753$$

Now we do backpropagation to update the weights. I hope you stored your values because you will be reusing them!

First you need to calculate the derivative of  $\ell$  with respect to all of the weights in  $\mathbf{W}^1$  and  $\mathbf{W}^2$ , then you use gradient descent to update them. Be sure to do it in this order because if you update the weights in  $\mathbf{W}^2$  before calculating the derivative with respect to  $\mathbf{W}^1$  then your calculations will be incorrect.

In the interest of time let's just calculate the derivatives:  $\frac{d\ell}{d\mathbf{W}_{11}^2}$  and  $\frac{d\ell}{d\mathbf{W}_{21}^1}$ .

Let's start with  $\frac{d\ell}{d\mathbf{W}_{11}^2}$ . It helps to examine the path that this weight effects  $\mathbf{W}_{11}^2$ :

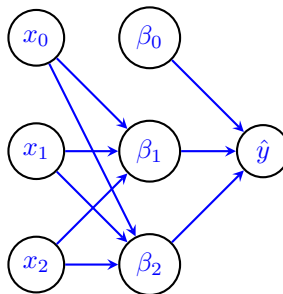


Figure 1: Neural Network

$\mathbf{W}_{11}^2$  follows a simple propagation path from  $\beta_1^1$  to  $\hat{y}$ .

So our derivative chain will be:

$$\frac{d\ell}{d\mathbf{W}_{11}^2} = \frac{d\ell}{d\hat{y}} \frac{d\hat{y}}{d\alpha^2} \frac{d\alpha^2}{d\mathbf{W}_{11}^2}$$

1. What is  $\frac{d\ell}{d\hat{y}}$ ?

$$-(y - \hat{y})$$

2. What is  $\frac{d\hat{y}}{d\alpha^2}$ ?

Recall that the derivative of the Sigmoid function has a very simple form!  $\sigma(\alpha^2)(1 - \sigma(\alpha^2))$   
 Notice that we can simplify this further as  $\hat{y}(1 - \hat{y})$

3. What is  $\frac{d\alpha^2}{d\mathbf{W}_{11}^2}$ ?

$$\text{Recall: } \alpha^2 = \mathbf{W}_{01}^2 \beta_0^1 + \mathbf{W}_{11}^2 \beta_1^1 + \mathbf{W}_{21}^2 \beta_2^1$$

$$\beta_1^1 = 0.5$$

4. Finally, what is  $\frac{d\ell}{d\mathbf{W}_{11}^2}$ ?

$$\frac{d\ell}{d\mathbf{W}_{11}^2} = -(y - \hat{y})\hat{y}(1 - \hat{y})0.5 = -(1 - 0.119623)0.119623(1 - 0.119623)0.5 = -0.04635772$$

Now we will calculate the derivative with respect to  $\mathbf{W}_{21}^1$ :  $\frac{d\ell}{d\mathbf{W}_{21}^1}$

$$\frac{d\ell}{d\mathbf{W}_{21}^1} = \frac{d\ell}{d\hat{y}} \frac{d\hat{y}}{d\alpha^2} \frac{d\alpha^2}{d\beta_1^1} \frac{d\beta_1^1}{d\alpha_1^1} \frac{d\alpha_1^1}{d\mathbf{W}_{21}^1}$$

We already had the first two derivatives calculated from the previous question.

5. What is  $\frac{d\alpha^2}{d\beta_1^1}$ ?

Recall:  $\alpha^2 = \mathbf{W}_{01}^2 \beta_0^1 + \mathbf{W}_{11}^2 \beta_1^1 + \mathbf{W}_{21}^2 \beta_2^1$   
 The answer is  $\mathbf{W}_{11}^2 = 2$

6. What is  $\frac{d\beta_1^1}{d\alpha_1^1}$ ?

This is again the derivative of a sigmoid function, so our answer is:  $\sigma(\alpha_1^1)(1 - \sigma(\alpha_1^1)) = \beta_1^1(1 - \beta_1^1)$

7. What is  $\frac{d\alpha_1^1}{d\mathbf{W}_{21}^1}$ ?

Recall  $\alpha_1^1 = \mathbf{W}_{01}^1 x_0 + \mathbf{W}_{11}^1 x_1 + \mathbf{W}_{21}^1 x_2$   
 Hence:  $x_2 = 5$

8. Finally what is  $\frac{d\ell}{d\mathbf{W}_{21}^1}$ ?

$$\frac{d\ell}{d\mathbf{W}_{21}^1} = -(y - \hat{y})\hat{y}(1 - \hat{y})2\beta_1^1(1 - \beta_1^1)5 = -(1 - 0.119623)0.119623(1 - 0.119623)2 \cdot 0.5 \cdot 0.5 \cdot 5 = -0.23178$$

9. What is our updated  $\mathbf{W}_{11}^2$  and  $\mathbf{W}_{21}^1$  if we use learning rate  $\lambda = 2$ ?

Recall gradient descent rule:  $w_{new} = w_{old} - \lambda \frac{d\ell}{dw_{old}}$   
 $\mathbf{W}_{11}^2 = 2 - 2 \cdot -0.04635772 = 2.0927$   
 $\mathbf{W}_{21}^1 = 1 - 2 \cdot -0.23178 = 1.4636$