

# Neural Network from Scratch

Your Name

## 1 Introduction

In this document I am going to go over my implementation of a very simple Neural Network and the math behind it.

## 2 Network Architecture

My simple neural network consists of:

1. **Input Layer:** 3-dimensional input
2. **Hidden Layer:** One neuron
3. **Output Layer:** 1-dimensional output

## 3 Initialization

We generate/initialize the weights using the Xavier Initialization:

$$W \sim N(0, \frac{2}{n_{in} + n_{out}}) \quad (1)$$

Where:

- $n_{in}$  is the number of input features
- $n_{out}$  is the number of neurons in the layer
- $N$  is the normal distribution function

In our Python implementation:

```
1 n_in = 3
2 n_out = 1
3 w = np.random.normal(0, np.sqrt(2 / (n_in + n_out)), size=(n_in,
4   n_out))
5 b = 0
```

## 4 Forward Propagation

In the Forward Propagation we pass the input data through the Neural Network to generate an output.

1. **Linear Combination:** Calculate the weighted sum of inputs plus a bias:

$$Z = W \cdot X + b \quad (2)$$

2. **Activation Function:** Apply a non-linear function, in this case Sigmoid, to introduce some non-linearity:

$$A = \frac{1}{1 + e^{-Z}} \quad (3)$$

In our code:

```
1 z = np.dot(w.T, x) + b
2 a = sigmoid(z)
```

Where `sigmoid` is defined as:

```
1 def sigmoid(z):
2     return 1 / (1 + np.exp(-z))
```

## 5 Loss Function

We use the Mean Squared Error (MSE) to calculate the loss of this epoch and evaluate how the model is doing:

$$L = \frac{1}{2}(A - Y)^2 \quad (4)$$

Where:

- $A$  is the predicted output
- $Y$  is the actual target value

## 6 Backpropagation

In the Backpropagation process we calculate the gradients of the loss function with respect to the weights and the biases to then be able to update them accordingly.

### 6.1 Derivatives

We compute the following partial derivatives:

### 6.1.1 Derivative of Loss with respect to Activation

We start with the Mean Squared Error loss function:

$$L = \frac{1}{2}(A - Y)^2 \quad (5)$$

To find  $\frac{\partial L}{\partial A}$ , we apply the chain rule:

$$\frac{\partial L}{\partial A} = \frac{\partial}{\partial A} \left[ \frac{1}{2}(A - Y)^2 \right] \quad (6)$$

$$= \frac{1}{2} \cdot 2(A - Y) \cdot \frac{\partial}{\partial A}(A - Y) \quad (7)$$

$$= (A - Y) \cdot 1 \quad (8)$$

$$= A - Y \quad (9)$$

Thus we derive:  $\frac{\partial L}{\partial A} = (A - Y)$

### 6.1.2 Derivative of Activation with respect to Linear Combination

We use the sigmoid function as our activation function:

$$A = \frac{1}{1 + e^{-Z}} = \sigma(Z) \quad (10)$$

To find  $\frac{\partial A}{\partial Z}$ , we use the derivative of the sigmoid function:

$$\frac{\partial A}{\partial Z} = \frac{\partial}{\partial Z} \sigma(Z) \quad (11)$$

$$= \sigma(Z)(1 - \sigma(Z)) \quad (12)$$

$$= A(1 - A) \quad (13)$$

$A = \sigma(Z)$  Thus, we derive:  $\frac{\partial A}{\partial Z} = A(1 - A)$

### 6.1.3 Derivative of Linear Combination with respect to Weights

Our linear combination is:

$$Z = W \cdot X + b \quad (14)$$

To find  $\frac{\partial Z}{\partial W}$ , we differentiate with respect to W:

$$\frac{\partial Z}{\partial W} = \frac{\partial}{\partial W}(W \cdot X + b) \quad (15)$$

$$= \frac{\partial}{\partial W}(W \cdot X) + \frac{\partial}{\partial W}b \quad (16)$$

$$= X + 0 \quad (17)$$

$$= X \quad (18)$$

Thus, we derive:  $\frac{\partial Z}{\partial W} = X$

#### 6.1.4 Derivative of Linear Combination with respect to Bias

$$Z = W \cdot X + b \quad (19)$$

To find  $\frac{\partial Z}{\partial b}$ , we differentiate with respect to b:

$$\frac{\partial Z}{\partial b} = \frac{\partial}{\partial b}(W \cdot X + b) \quad (20)$$

$$= \frac{\partial}{\partial b}(W \cdot X) + \frac{\partial}{\partial b}b \quad (21)$$

$$= 0 + 1 \quad (22)$$

$$= 1 \quad (23)$$

Thus, we derive:  $\frac{\partial Z}{\partial b} = 1$

## 6.2 Chain Rule

We use the chain rule to calculate the gradients of the loss with respect to the weights and bias:

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial A} \cdot \frac{\partial A}{\partial Z} \cdot \frac{\partial Z}{\partial W} = (A - Y) \cdot A(1 - A) \cdot X \quad (24)$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial A} \cdot \frac{\partial A}{\partial Z} \cdot \frac{\partial Z}{\partial b} = (A - Y) \cdot A(1 - A) \quad (25)$$

In our code:

```
1 dl_dw = (a - y) * a * (1 - a) * x
2 dl_db = (a - y) * a * (1 - a)
```

## 7 Gradient Descent

In the last step we apply an algorithm called Gradient Descent to find the direction for which the weights and biases will be minimal. Afterwards we update them to minimize the loss.

$$W = W - \eta \cdot \frac{\partial L}{\partial W} \quad (26)$$

$$b = b - \eta \cdot \frac{\partial L}{\partial b} \quad (27)$$

Where  $\eta$  is the learning rate, which controls the step size of each update.

In our implementation:

```
1 learning_rate = 0.005
2 w = w - learning_rate * dl_dw
3 b = b - learning_rate * dl_db
```

## 8 Training Loop

In the full code this process gets repeated multiple times to hopefully minimize the loss and get as accurate as possible. In the full code it looks like this:

```
1 def feedforward_nn(epochs, x, y):
2     # ... (initialization code)
3
4     for epoch in range(epochs):
5         # Forward propagation
6         z = np.dot(w.T, x) + b
7         a = sigmoid(z)
8
9         # Backpropagation
10        dl_dw = (a - y) * a * (1 - a) * x
11        dl_db = (a - y) * a * (1 - a)
12
13        # Gradient descent
14        w = w - learning_rate * dl_dw
15        b = b - learning_rate * dl_db
16
17        # Print progress (optional)
18        if epoch % 100 == 0:
19            print(f"Epoch {epoch}: w = {list(map('{:.3f}%'.format,
20            w.flatten()))}, b = {list(map('{:.3f}%'.format, b))}, loss = {
21            list(map('{:.3f}%'.format, (a - y) ** 2))}")
22
23    return w, b
```

Hope you enjoyed my implementation:)