

Neural Network From Scratch

Date: _____

Page: _____

Here, how neural network is implemented from scratch is discussed. Here all mathematical expressions needed to train a simple neural network is explained. This consists of mainly six major parts. They are :-

- 1) Implementation design
- 2) Create the base layer
- 3) Create Dense layer
- 4) Create Activation layer
- 5) Implement ~~the~~ activation functions and loss function
- 6) Solve XOR problem.

Before heading into implementation design, first revising machine learning steps :-

Step 1 :- Feed input: Data flows from layer to layer. Gives output (prediction)

Let $y = \text{network}(x, \omega)$

where y is prediction

x is input / features

ω is parameter to train (weights & bias)
and $\text{network}()$ is a neural network arch.

Step 2: calculate error.

$$\text{eg. } E = \frac{1}{2} (y - \hat{y})^2$$

where y is prediction
 \hat{y} is true.

*Note: These notation might be different on different sources.

Step 3: Adjust the parameter using gradient descent (Back propagation).

$$\omega \leftarrow \omega - \alpha \frac{\partial E}{\partial \omega}$$

where α is learning rate

$\frac{\partial E}{\partial \omega}$ is gradient (first derivative) or partial derivative) of loss function

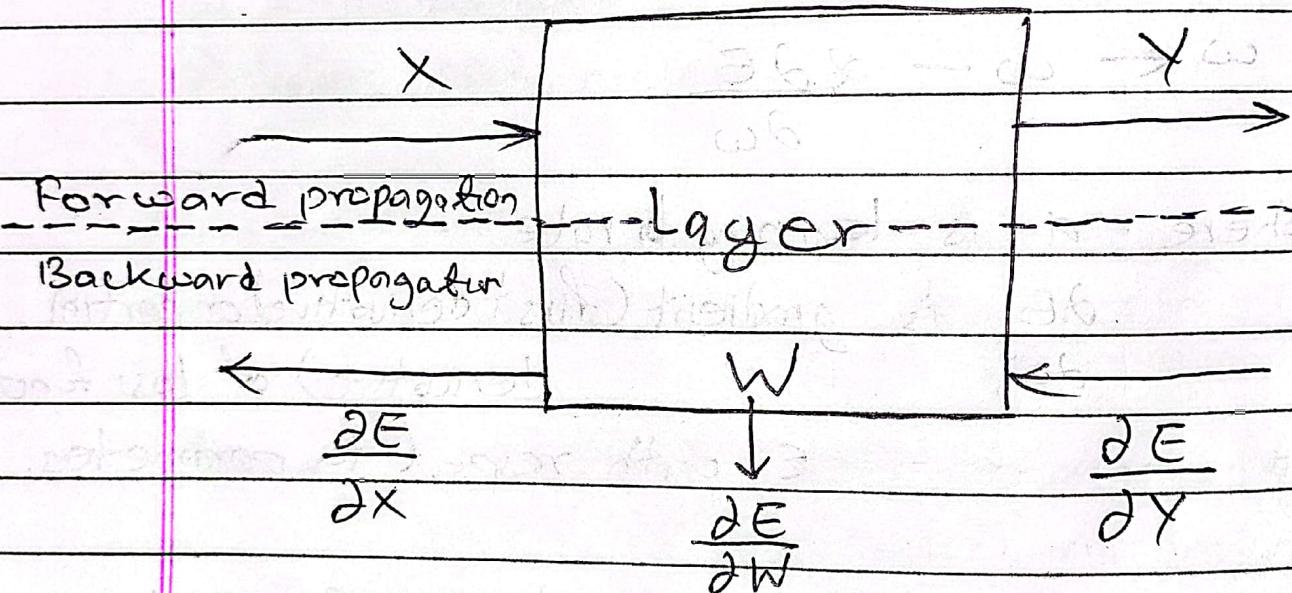
E with respect to parameters.

Step 4: Start again , until error E is minimum.

Now getting into our main purpose.

1) Implementation design-

First of all, we need a layer model which can fit every types of layers whether it's fully connected, convolution, drop out or activation. Since a layer is a big function, hence, it is represented as a subbox. Below



for every layer, when an input ' x ' is given it will produce output ' y '. This is called forward propagation.

Similarly the other step when training is called backward propagation. when the layer updates its parameters. During back propagation,

Date: _____
Page: _____

If a layer is given the derivative of the error w.r.t its output, it should give back the derivative of the error w.r.t its input.

Ultimately, for every trainable layer there is set of parameters W that need to be updated according to gradient descent or similar other.

In order to update parameters, the layer needs to compute the derivative of error w.r.t. the parameters (parameters are not shown)

Then,

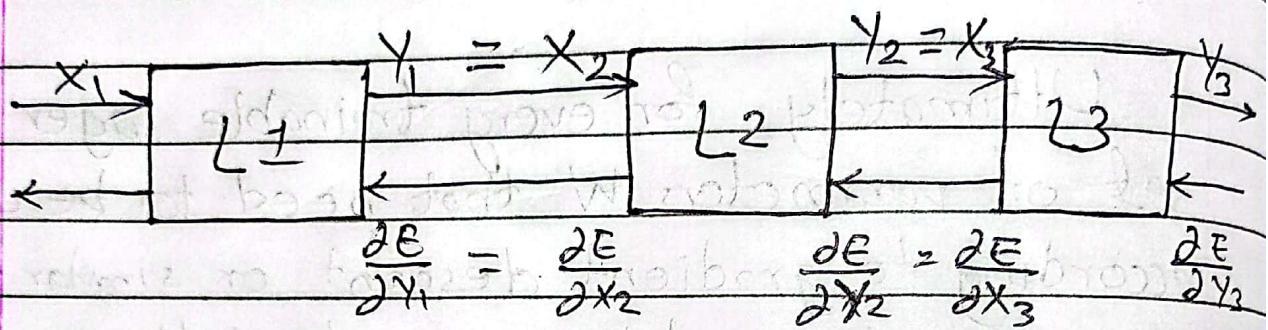
$$\frac{\partial E}{\partial W} = \frac{\partial E}{\partial Y} \cdot \frac{\partial Y}{\partial W}$$

$$\frac{\partial E}{\partial X} = \frac{\partial E}{\partial Y} \cdot \frac{\partial Y}{\partial X}$$

Now suppose

Now suppose there are multiple layers

then,



In these types of sequential model, the output of first layer is input for second layer and goes on.

$$\text{Hence } Y_1 = X_2, \quad Y_2 = X_3$$

Similarly, for layer 1, derivative of error with respect to output of L1 is equal to derivative of error w.r.t. input of L2 and same for other layer.

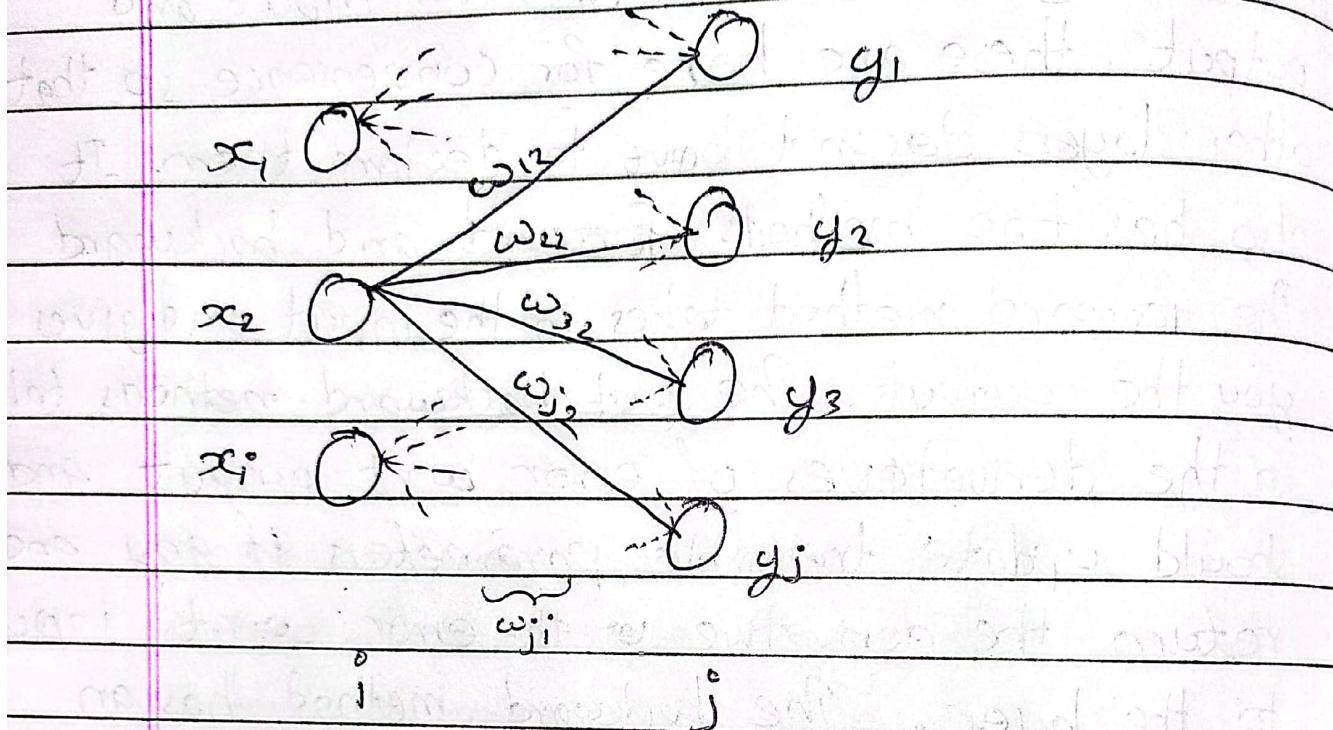
$$\text{i.e. } \frac{\partial E}{\partial Y_1} = \frac{\partial E}{\partial X_2} \quad \text{and} \quad \frac{\partial E}{\partial Y_2} = \frac{\partial E}{\partial X_3}$$

2) Base Layer

This layer has two attributes input and output. These are here for convenience so that other layers don't have to declare them. It also has two methods forward and backward. The forward method takes in the input and gives you the output. The backward method takes in the derivatives of error wrt output and should update trainable parameters if any and return the derivative of the error wrt. input to the layer. The backward method has an additional parameter, the learning-rate.

*Note: Code for these classes are written at ~~very~~ last, ~~forget about them~~.

3) Dense Layer. (Fully-connected layer)



Dense layer, connect a set of ' i ' input neurons to set of ' j ' o/p neurons.

I/P is denoted as x and o/p is denoted as y . In dense layer each input neuron is connected to every o/p neuron. Each connection represent weights. Represented by w_{ij} , which means, weight which connect output neuron j with input neuron i .

Every o/p value is computed as sum of all inputs multiplied by the weights and an additional term called bias denoted as b . Biases are also trainable parameters like weight.

Date: _____
Page: _____

Here,

$$y_1 = x_1 w_{11} + x_2 w_{12} + \dots + x_i w_{1i} + b_1$$

$$y_2 = x_1 w_{21} + x_2 w_{22} + \dots + x_i w_{2i} + b_2$$

$$y_3 = x_1 w_{31} + x_2 w_{32} + \dots + x_i w_{3i} + b_3$$

!

$$y_j = x_1 w_{j1} + x_2 w_{j2} + \dots + x_i w_{ji} + b_j$$

Rewriting in matrix form

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_j \end{bmatrix} = \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1i} \\ w_{21} & w_{22} & \dots & w_{2i} \\ \vdots & \vdots & \ddots & \vdots \\ w_{j1} & w_{j2} & \dots & w_{ji} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_i \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_j \end{bmatrix}$$

$$j \times 1 \quad j \times i \quad j \times 1 \quad j \times 1$$

$$w \quad w \quad w \quad w$$

Weight matrix has a shape of $j \times i$ and bias is a column vector of shape $j \times 1$.

Finally, we can represent it as

$$Y = W \cdot X + B$$

This layer class inherits from the base layer class. The constructor of layer takes two parameters: input size and output size, which are number

of neurons in input and output respectively.
Inside the constructor weights and bias
are initialize randomly.

The forward method simply computes
 $y = w^T x + b$.

For backward propagation

Let us assume we have derivative of
error w.r.t. output of dense layer i.e
 $\frac{dE}{dy}$ is given, then we have to calculate
 $\frac{dE}{dx}$.

~~other~~ two things.

$$\begin{array}{c} \xrightarrow{\quad dE \quad, \quad dE/x_i ?} \\ \frac{dE}{dw}, \frac{dE}{dB} \\ \xrightarrow{\quad dE/dx ? \quad} \end{array}$$

Using $\frac{dE}{dy}$ we need to calculate derivative

of error w.r.t. weights and bias, as they are
trainable parameter, and their value have to
be adjusted according to gradient descent.

On the other, we need derivative of error w.r.t
input, this is because, it will be passed to

the layer before the one computing this.

i) For $\frac{\partial E}{\partial w}$

Starting with derivative of error with respect to w .

we have

$$\frac{\partial E}{\partial w} = \begin{bmatrix} \frac{\partial E}{\partial y_1} \\ \frac{\partial E}{\partial y_2} \\ \vdots \\ \frac{\partial E}{\partial y_j} \end{bmatrix}_{j \times 1}$$

Now

$$\frac{\partial E}{\partial w} = \begin{bmatrix} \frac{\partial E}{\partial w_{11}} & \frac{\partial E}{\partial w_{12}} & \cdots & \frac{\partial E}{\partial w_{1i}} \\ \frac{\partial E}{\partial w_{21}} & \frac{\partial E}{\partial w_{22}} & \cdots & \frac{\partial E}{\partial w_{2i}} \\ \vdots & & & \\ \frac{\partial E}{\partial w_{j1}} & \frac{\partial E}{\partial w_{j2}} & \cdots & \frac{\partial E}{\partial w_{ji}} \end{bmatrix}_{j \times i}$$

Taking first term only

$$\frac{\partial E}{\partial w_{12}} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial w_{12}} + \frac{\partial E}{\partial y_2} \frac{\partial y_2}{\partial w_{12}} + \dots + \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial w_{12}}$$

$\underbrace{}$ $\underbrace{}$ $\underbrace{}$

x_2 0 0

$$\frac{\partial E}{\partial w_{12}} = \frac{\partial E}{\partial y_1} x_2$$

In general

$$\frac{\partial E}{\partial w_{ji}} = \frac{\partial E}{\partial y_j} x_i$$

Now, using this equation

$$\frac{\partial E}{\partial w} = \begin{bmatrix} \frac{\partial E}{\partial y_1} x_1 & \frac{\partial E}{\partial y_1} x_2 & \dots & \frac{\partial E}{\partial y_1} x_i \\ \frac{\partial E}{\partial y_2} x_1 & \frac{\partial E}{\partial y_2} x_2 & \dots & \frac{\partial E}{\partial y_2} x_i \\ \vdots & \vdots & & \vdots \\ \frac{\partial E}{\partial y_j} x_1 & \frac{\partial E}{\partial y_j} x_2 & \dots & \frac{\partial E}{\partial y_j} x_i \end{bmatrix}$$

This can be further more simplified by writing
in matrix multiplication form

Date: _____
Page: _____

$$\therefore \frac{\partial E}{\partial W} = \begin{bmatrix} \frac{\partial E}{\partial y_1} \\ \frac{\partial E}{\partial y_2} \\ \vdots \\ \frac{\partial E}{\partial y_j} \end{bmatrix} \times [x_1 \ x_2 \ \dots \ x_i]$$

$j \times 1 \quad 1 \times i$

$j \times i$

$$\boxed{\therefore \frac{\partial E}{\partial W} = \frac{\partial E}{\partial Y} \cdot X^T}$$

ISB 36 - 36
 X^T is transpose of X

(ii) For $\frac{\partial E}{\partial B}$

$$\frac{\partial E}{\partial B} = \begin{bmatrix} \frac{\partial E}{\partial b_1} \\ \frac{\partial E}{\partial b_2} \\ \vdots \\ \frac{\partial E}{\partial b_3} \end{bmatrix}$$

ISB

ISB

ISB = ISB

ISB = ISB

ISB

ISB

For

$$\frac{\partial E}{\partial b_1} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial b_1} + \frac{\partial E}{\partial y_2} \frac{\partial y_2}{\partial b_1} + \dots + \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial b_1}$$

↓ ↓ ↓
 1 0 0

$$\therefore \frac{\partial E}{\partial b_1} = \frac{\partial E}{\partial y_1}$$

In general,

$$\frac{\partial E}{\partial b_j} = \frac{\partial E}{\partial y_j}$$

$$\therefore \boxed{\frac{\partial E}{\partial b} = \frac{\partial E}{\partial Y}}$$

(iii) For $\frac{\partial E}{\partial x}$.

$$\frac{\partial E}{\partial x} = \begin{vmatrix} \frac{\partial E}{\partial x_1} \\ \frac{\partial E}{\partial x_2} \\ \vdots \\ \frac{\partial E}{\partial x_i} \end{vmatrix}$$

Date: _____
 Page: _____

$$\frac{\partial E}{\partial x_i} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial x_i} + \frac{\partial E}{\partial y_2} \frac{\partial y_2}{\partial x_i} + \dots + \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

$\underbrace{\hspace{1cm}}$ $\underbrace{\hspace{1cm}}$ $\underbrace{\hspace{1cm}}$
 ω_{1i} ω_{2i} ω_{ji}

In general,

$$\frac{\partial E}{\partial x_i} = \frac{\partial E}{\partial y_1} \omega_{1i} + \frac{\partial E}{\partial y_2} \omega_{2i} + \dots + \frac{\partial E}{\partial y_j} \omega_{ji}$$

	$\frac{\partial E}{\partial y_1} \omega_{1i} + \frac{\partial E}{\partial y_2} \omega_{2i} + \dots + \frac{\partial E}{\partial y_j} \omega_{ji}$
	$\frac{\partial E}{\partial y_1} \omega_{12} + \frac{\partial E}{\partial y_2} \omega_{22} + \dots + \frac{\partial E}{\partial y_j} \omega_{j2}$
$\frac{\partial E}{\partial x}$	⋮
	$\frac{\partial E}{\partial y_1} \omega_{1i} + \frac{\partial E}{\partial y_2} \omega_{2i} + \dots + \frac{\partial E}{\partial y_j} \omega_{ji}$

Expressing it into matrix multiplication

$$\frac{\partial E}{\partial x} = \begin{bmatrix} \omega_{11} & \omega_{21} & \dots & \omega_{j1} \\ \omega_{12} & \omega_{22} & \dots & \omega_{j2} \\ \vdots & \ddots & & \\ \omega_{ij} & \omega_{2j} & \dots & \omega_{jj} \end{bmatrix} \begin{bmatrix} \frac{\partial E}{\partial y_1} \\ \frac{\partial E}{\partial y_2} \\ \vdots \\ \frac{\partial E}{\partial y_j} \end{bmatrix}_{j \times 1}$$

$$\frac{\partial E}{\partial x} = W^T \cdot \frac{\partial E}{\partial y}$$

Now we can implement these three equations for backward propagation

4) Activation Layer.

$$x_1 \text{---} O \text{---} y_1 = f(x_1)$$

$$x_2 \text{---} O \text{---} y_2 = f(x_2)$$

$$x_3 \text{---} O \text{---} y_3 = f(x_3)$$

$$x_i \text{---} O \text{---} y_i = f(x_i)$$

The activation layer takes some input neurons and passes them through activation function. Thus, for that layer the output has the same shape as the input.

The forward propagation is expressed as

$$\bullet Y = f(E)$$

where f is activation function

Date: _____
Page: _____

The activation layer inherits from the base layer class. The constructor takes two parameters the activation and its derivatives.

The forward method simply applies the activation to the input.

The backward method which should return the derivative of the error wrt the input

As always let we have derivative of E wrt. Y then we have to calculate derivative of E wrt X

$$\frac{\partial E}{\partial x} = \begin{vmatrix} \frac{\partial E}{\partial x_1} \\ \frac{\partial E}{\partial x_2} \\ \vdots \\ \frac{\partial E}{\partial x_i} \end{vmatrix}$$
$$y_1 = f(x_1)$$
$$y_2 = f(x_2)$$
$$\vdots$$
$$y_i = f(x_i)$$

$$\frac{\partial E}{\partial x_1} = \frac{\partial E}{\partial y_1} \frac{\partial y_1}{\partial x_1} + \frac{\partial E}{\partial y_2} \frac{\partial y_2}{\partial x_1} + \dots + \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial x_1}$$

$\underbrace{f'(x_1)}$ $\underbrace{\circ}$

$$\therefore \frac{\partial E}{\partial x_i} = \frac{\partial E}{\partial y_i} f'(x_i)$$

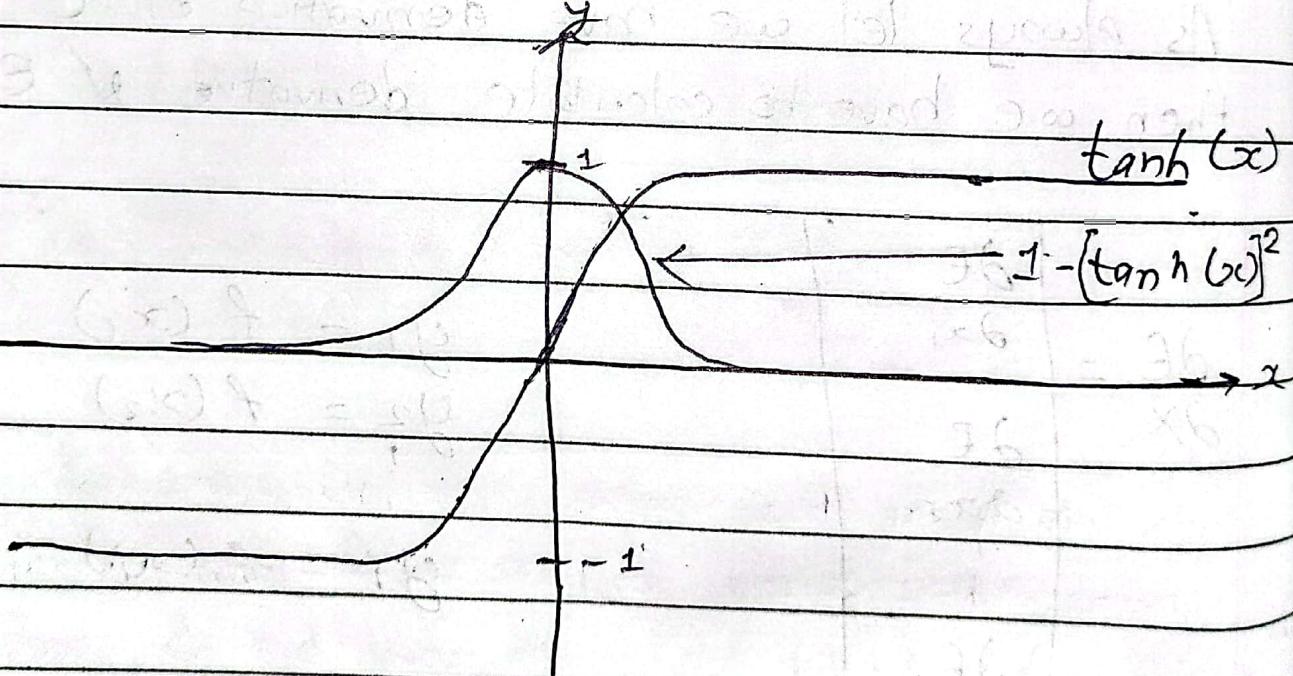
Hence,

$$\frac{\partial E}{\partial x} = \frac{\partial E}{\partial y} \odot f'(x)$$

① element wise
multiplication

5) Implementation of activation function and loss function

Hyperbolic Tangent as activation function



$$f(x) = \tanh(x)$$

$$f'(x) = 1 - [\tanh(x)]^2$$

For implementation of tanh activation, we inherit from activation class and pass hyperbolic tangent function and its derivatives to super constructor.

loss Function

Here mean squared error is implemented. Given a vector y that represents the desired output denoted as y and another vector that represents the actual output represented by \hat{y} .

Error is defined as

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

Now,

From the beginning we have supposed derivative of E w.r.t. y is given during back propagation. For L1 layer, it is given by L2 layer and for L2 it is given by L3 and so on. At very last layer the output of last layer is the output of whole neuron network. Hence it is the output used to calculate error.

Hence derivative of E is calculated using above function.

Date: _____
Page: _____

$$\begin{aligned}\therefore \frac{\partial E}{\partial y_1} &= \frac{\partial}{\partial y_1} \left[(\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 + \dots + (\hat{y}_n - y_n)^2 \right] \\ &= \frac{\partial}{\partial y_1} \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \\ &= \frac{2}{n} (\hat{y}_1 - y_1)\end{aligned}$$

In general

$$\frac{\partial E}{\partial y_i} = \frac{2}{n} (\hat{y}_i - y_i)$$

In matrix form

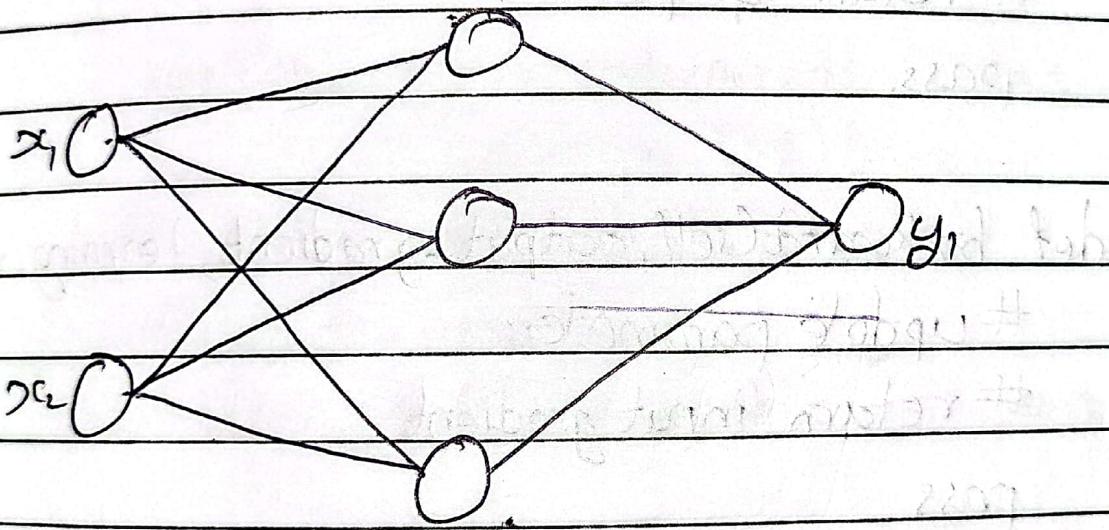
$$\boxed{\frac{\partial E}{\partial Y} = \frac{2}{n} (Y - \hat{Y})}$$

Now two different function will be made for these two task i.e. to calculate error and to calculate derivative of E w.r.t. Y.

Date: _____
Page: _____

6) Solving for X-OR Problem.

x_1	x_2	\hat{y}
0	0	0
0	1	1
1	0	1
1	1	0



There two input neuron is kept as there are two inputs for each row and one output neuron for one output value. Also three neurons are used for one hidden layer.

layer.py

class Layer:

def __init__(self):

self.input = None

self.output = None

def forward(self, input):

return output

pass.

def backward(self, output-gradient, learning-rate):

update parameters

return input gradient

pass

$$\text{output-gradient} = \frac{\partial E}{\partial Y} \quad \text{input} = X$$

$$\quad \quad \quad \text{output} = Y$$

$$\text{input-gradient} = \frac{\partial E}{\partial X} \quad \text{input.T} = X^T$$

$$\text{weight-gradient} = \frac{\partial E}{\partial W} \quad \text{weight.T} = W^T$$

$$\text{bias-gradient} = \frac{\partial E}{\partial B}$$

`np.dot()` → matrix multiplication

dense.py

```
import numpy as np  
from layer import Layer
```

```
class Dense(Layer):
```

```
    def __init__(self, input_size, output_size):
```

```
        self.weights = np.random.randn(output_size,
```

```
                                input_size)
```

```
        self.bias = np.random.randn(output_size, 1)
```

```
    def forward(self, input):
```

```
        self.input = input
```

```
        return np.dot(self.weights, self.input) + self.bias
```

```
    def backward(self, output_gradient, learning_rate):
```

```
        weights_gradient = np.dot(output_gradient, self.input.T)
```

```
        input_gradient = np.dot(self.weights.T, output_gradient)
```

```
        self.bias =
```

```
            self.weight -= learning_rate * weights_gradient
```

```
            self.bias -= learning_rate * output_gradient
```

```
        return input_gradient
```

activation.py

```
import numpy as np
```

```
from layer import Layer
```

```
class Activation(Layer):
```

```
    def __init__(self, activation, activation_prime):
```

```
        self.activation = activation
```

```
        self.activation_prime = activation_prime
```

```
    def forward(self, input):
```

```
        self.input = input
```

```
        return self.activation(self.input)
```

```
    def backward(self, output_gradient, learning_rate):
```

```
        return np.multiply(output_gradient, self.activation_
```

```
prime(self.input))
```

$$\text{activation} = f(x)$$

$$\text{activation-prime} = f'(x)$$

`np.multiply()` → element wise multiplication.

```
class Tanh(Activation):
```

```
    def __init__(self):
```

```
        def tanh(x):
```

```
            return np.tanh(x)
```

```
        def tanh_prime(x):
```

```
            return 1 - np.tanh(x) ** 2
```

```
    super().__init__(tanh, tanh_prime)
```

losses.py

①

```
import numpy as np
```

```
def mse(y_true, y_pred):
```

```
    return np.mean(np.power(y_true - y_pred, 2))
```

```
def mse_prime(y_true, y_pred):
```

```
    return 2 * (y_pred - y_true) / np.size(y_true)
```

$y_{\text{true}} = \hat{y}$

$y_{\text{pred}} = y$

$\text{np.size}() = n$

train.py

```
def predict(network, input):  
    output = input  
    for layer in layers network:  
        output = layer.forward(output)  
    return output.
```

```
def fit(network, loss, loss-prime, x-train, y-train,  
       epochs=1000, learning-rate=0.01):
```

```
for e in range(epochs):  
    error = 0  
    for x, y in zip(x-train, y-train):  
        # forward propagation  
        output = predict(network, x)  
  
        # error calculation  
        error += loss(y, output)  
  
        # backward propagation.  
        grad = loss-prime(y, output)  
        for layer in reverse(network):  
            grad = layer.backward(grad, learning-rate)  
        error /= len(x-train)  
    print(f"Epoch {e+1}/{epochs}, Error = {error}")
```

progr.py

```
import numpy as np
from dense import Dense
from activation import Tanh
from losses import mse, mse_prime
from network import train, predict, fit.
```

```
x = np.reshape([ [0,0], [0,1], [1,0], [1,1] ], (4,2,1))
y = np.reshape([ [0], [1], [1], [0] ], (4,1,1))
```

```
network = [ Dense(2, 3),
            Tanh(),
            Dense(3, 1),
            Tanh() ]
```

```
fit(network, mse, mse_prime, x, y, epochs=1000,
     learning_rate=0.1)
```