

## Homework - 4

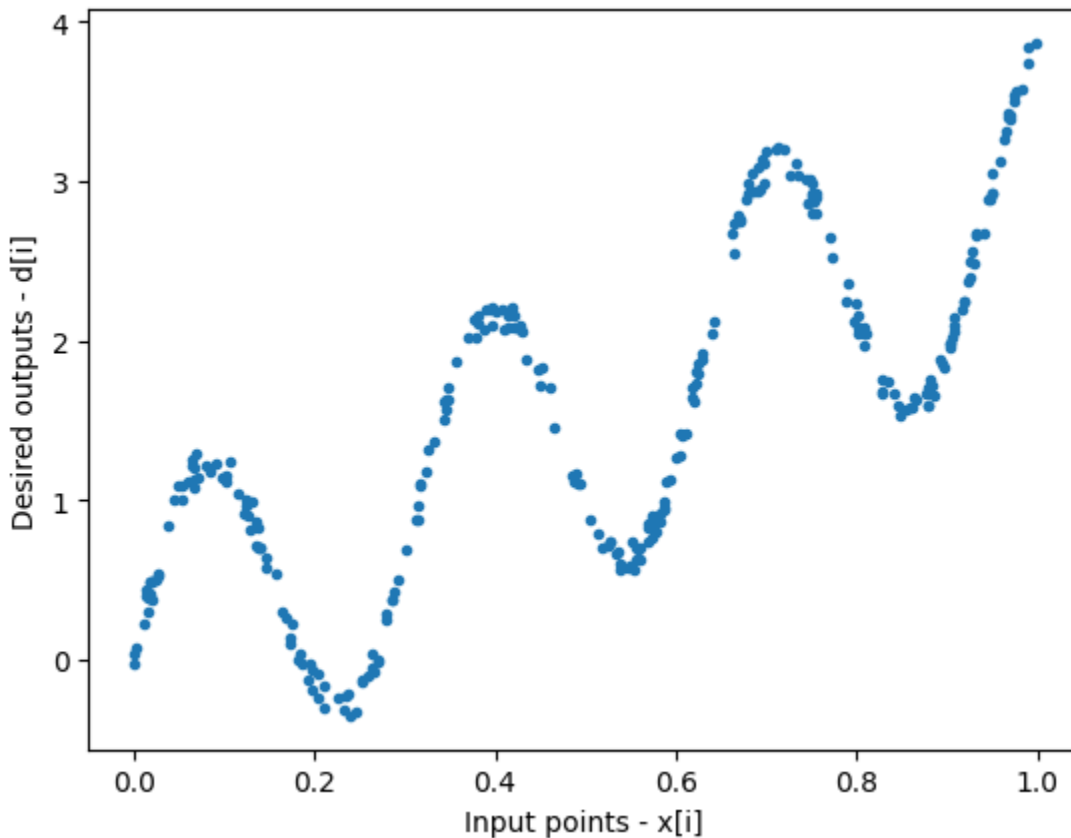
No of training samples are 300, so  $n = 300$

Input of the Neural Network is randomly uniform on  $[0,1]$ , naming as  $x$

There are  $n$  real numbers uniformly at random on  $[-10, 10]$ , naming as  $v$ .

Desired output is calculated using :  $\sin(20x[i]) + 3x[i] + v[i], i = 1, \dots, n$ .

The plot of  $(x_i, d_i), i = 1, \dots, n$  is shown below



Considering  $1 \times N \times 1$  Neural Network. Let's consider  $1 \times N$  as first layer and  $N \times 1$  as a second layer.

Calculating weights for first layer ( $1 \times N$ ):

There are  $N$  weights from  $X$  input to  $N$  neurons and each neuron would be biased. So total weights for first layer is  $2N$ .

Calculating weights for first layer ( $N \times 1$ ):

There are  $N$  weights from  $N$  neurons to 1 output neuron and it would be biased. So total weights for first layer is  $N + 1$ .

Hence, total weights of this  $1 \times N \times 1$  Neural Network is  $3N + 1$ .

Some initializations of the Neural Network :

$N = 24$  (hidden neurons)

$\eta = 0.01$

$w_{\text{first\_layer}} = \text{np.array}(\text{np.random.uniform}(-0.5, 0.5, \text{size} = n))$

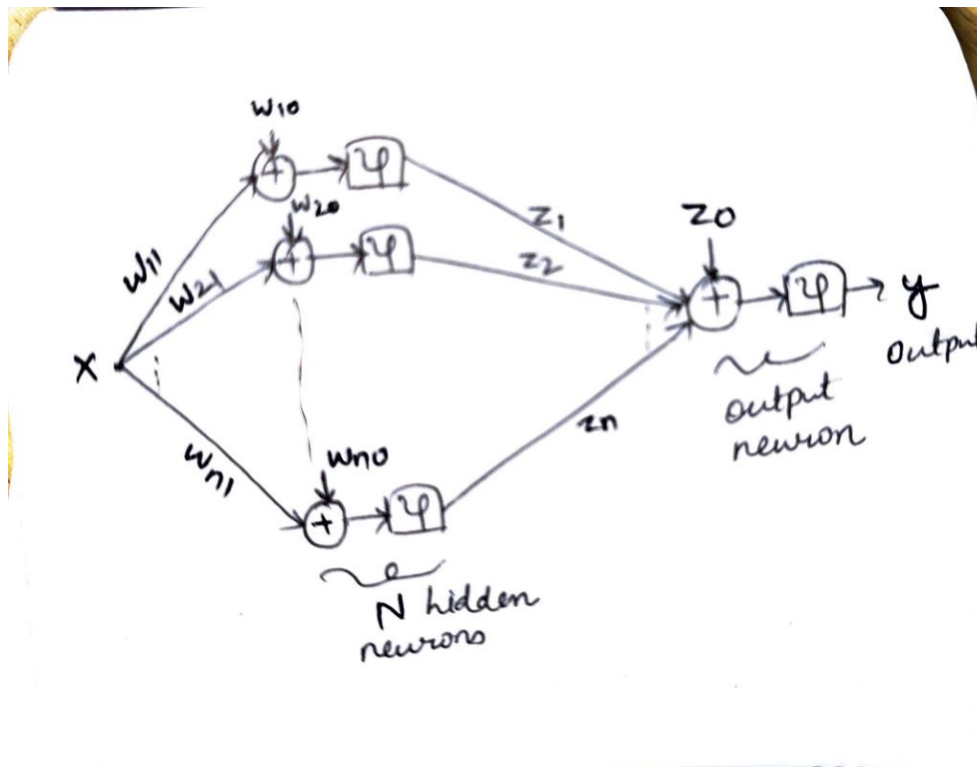
$w_{\text{first\_layer\_bias}} = \text{np.array}(\text{np.random.uniform}(-0.5, 5, \text{size} = n))$

$w_{\text{second\_layer}} = \text{np.array}(\text{np.random.uniform}(-0.5, 5, \text{size} = n))$

$w_{\text{second\_layer\_bias}} = \text{np.array}(\text{np.random.uniform}(-0.5, 5, \text{size} = 1))$

Here,  $w$  represents weights, hence  $w_{\text{first\_layer}}$  is weights from  $X$  input to  $N$  neurons in the first layer.

Please find the below diagram of  $1 \times N \times 1$  Neural Network



Here,  $w_{i1}$  is  $w_{\text{first\_layer}}$ ,  $i = 1, \dots, n$

$w_{i0}$  is  $w_{\text{first\_layer\_bias}}$ ,  $i = 1, \dots, n$

$z_i$  is  $w_{\text{second\_layer}}$ ,  $i = 1, \dots, n$

$z_0$  is  $w_{\text{second\_layer\_bias}}$

The output neuron has activation function  $\phi(v) = v$ ; all other neurons have activation function  $\phi(v) = \tanh v$ .

## Pseudo Code for Backpropagation algorithm :

Initialization is already done for below variables in the above pages.

$w\_first\_layer, w\_first\_layer\_bias, w\_second\_layer, w\_second\_layer\_bias, d, x, N, n, \eta$

epoch = 0

$m\_s\_e\_array = []$  // It will store MSE for all epochs

while(1):

    Initializing  $sum\_first\_layer[n, N] = 0$ ,  $activation\_first\_layer[n, N] = 0$ ,  $sum\_second\_layer[n] = 0$  and  $y[n] = 0$

    // Doing forward propagation

    // Calculating local field for first layer

    for i = 1 to n:

        for j = 1 to N:

$sum\_first\_layer[i][j] = (w\_first\_layer[j] * x[i][j]) + w\_first\_layer\_bias[j]$

    // Calculating the output of first layer from the activation function, here  $\phi(v) = \tanh v$

    for i = 1 to n:

        for j = 1 to N:

$activation\_first\_layer[i][j] = \tanh(sum\_first\_layer[i][j])$

    // Calculating local field for second layer

    for i = 1 to n:

$sum\_second\_layer[i] = 0$

        for j = 1 to N:

$sum\_second\_layer[i] += (activation\_first\_layer[i][j] * w\_second\_layer[j]) + w\_second\_layer\_bias[j]$

    // Calculating the output from the activation function, here  $\phi(v) = v$

    for i = 1 to n:

$y[i] = sum\_second\_layer[i]$

    // Calculating the MSE for that epoch

$m\_s\_e = 0$

    for i = 1 to n:

$m\_s\_e += (d[i] - y[i])^2$

$m\_s\_e = m\_s\_e / n$

    // Modifying  $\eta \leftarrow 0.9\eta$  whenever the MSE has increased

    if(  $m\_s\_e > m\_s\_e\_array[epoch - 1]$ ):

$\eta = \eta * 0.9$

```

// Storing MSE for all epochs
m_s_e_array.append(m_s_e)

// breaking the loop after MSE is less than or equal to threshold. Here threshold is 0.01
if(m_s_e <= 0.01):
    break

// Updating the epoch
epoch = epoch + 1

// Doing backward propagation and updating the weights

for i = 1 to n:
    // Calculating y with the given weights

    // Calculating local field for first layer
    for j = 1 to N:
        sum_first_layer[i][j] = (w_first_layer[j] * x[i][j]) + w_first_layer_bias[j]

    // Calculating the output of first layer from the activation function, here  $\phi(v) = \tanh v$ 
    for j = 1 to N:
        activation_first_layer[i][j] = tanh(sum_first_layer[i][j])

    // Calculating local field for second layer
    sum_second_layer = 0
    for j = 1 to N:
        sum_second_layer += (activation_first_layer[i][j] * w_second_layer[j]) +
            w_second_layer_bias[j]

    // Calculating the output from the activation function, here  $\phi(v) = v$ 
    y = sum_second_layer

    // Updating w_first_layer – gradient descent equation
    gradient_first_layer_weights[N] = 0
    for j = 1 to N:
        gradient_first_layer_weights[j] = x[i] * (d[i]-y) * (1 - (tanh((w_first_layer[j] * x[i][j]) +
            w_first_layer_bias[j] )2 ) * w_second_layer[j]

    for j = 1 to N:
        w_first_layer[j] = w_first_layer[j] +  $\eta$ *gradient_first_layer_weights[j]

```

```

// Updating w_first_layer_bias – gradient descent equation
gradient_first_layer_weights_bias[N] = 0
for j = 1 to N:
    gradient_first_layer_weights_bias[j] = 1 * (d[i]-y) * (1 - (tanh((w_first_layer[j] *
        x[i][j]) + w_first_layer_bias[j] )2 ) * w_second_layer[j]

for j = 1 to N:
    w_first_layer_bias[j] = w_first_layer_bias[j] + η*gradient_first_layer_weights_bias[j]

// Updating w_second_layer – gradient descent equation
gradient_second_layer_weights[N] = 0
for j = 1 to N:
    gradient_second_layer_weights[j] = (d[i]-y) * (tanh((w_first_layer[j] * x[i][j]) +
        w_first_layer_bias[j] )

for j = 1 to N:
    w_second_layer[j] = w_second_layer[j] + η*gradient_second_layer_weights[j]

// Updating w_second_layer_bias – gradient descent equation
gradient_second_layer_weights_bias = (d[i]-y) * 1
w_second_layer_bias[j] = w_second_layer_bias[j] +
    η*gradient_second_layer_weights[j]

```

After the algorithm function is completed

`W_first_layer`, `w_first_layer_bias`, `w_second_layer`, `w_second_layer_bias` are the final optimal weights

`m_s_e_array` would have the MSEs for each epoch, and size of `m_s_e_array` is the number of epochs required.

Gradient Descent Equations are :

```
for i = 1 to n:
    // Updating w_first_layer – gradient descent equation
    gradient_first_layer_weights[N] = 0
    for j = 1 to N:
        gradient_first_layer_weights[j] = x[i] * (d[i]-y) * (1 - (tanh((w_first_layer[j] * x[i][j]) +
            w_first_layer_bias[j] )2 ) * w_second_layer[j]

    for j = 1 to N:
        w_first_layer[j] = w_first_layer[j] + η*gradient_first_layer_weights[j]

// Updating w_first_layer_bias – gradient descent equation
gradient_first_layer_weights_bias[N] = 0
for j = 1 to N:
    gradient_first_layer_weights_bias[j] = 1 * (d[i]-y) * (1 - (tanh((w_first_layer[j] *
        x[i][j]) + w_first_layer_bias[j] )2 ) * w_second_layer[j]

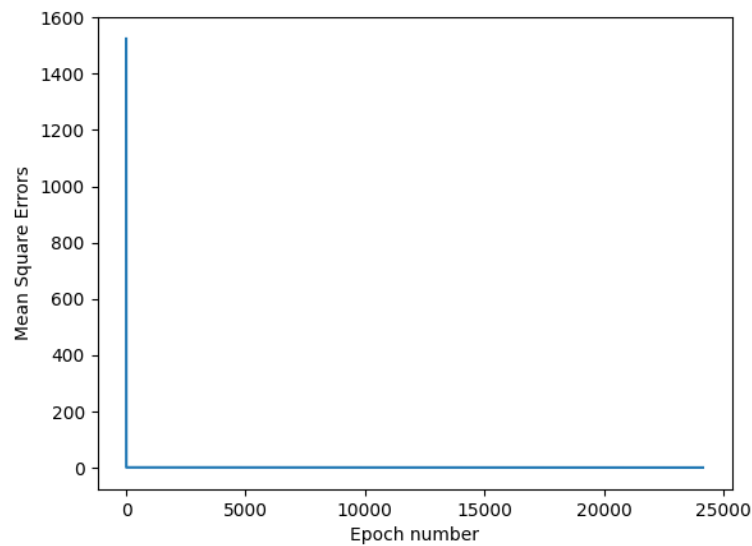
for j = 1 to N:
    w_first_layer_bias[j] = w_first_layer_bias[j] + η*gradient_first_layer_weights_bias[j]

// Updating w_second_layer – gradient descent equation
gradient_second_layer_weights[N] = 0
for j = 1 to N:
    gradient_second_layer_weights[j] = (d[i]-y) * (tanh((w_first_layer[j] * x[i][j]) +
        w_first_layer_bias[j] )

for j = 1 to N:
    w_second_layer[j] = w_second_layer[j] + η*gradient_second_layer_weights[j]

// Updating w_second_layer_bias – gradient descent equation
gradient_second_layer_weights_bias = (d[i]-y) * 1
w_second_layer_bias[j] = w_second_layer_bias[j] +
    η*gradient_second_layer_weights[j]
```

Graph for the number of epochs vs the MSE in the backpropagation algorithm

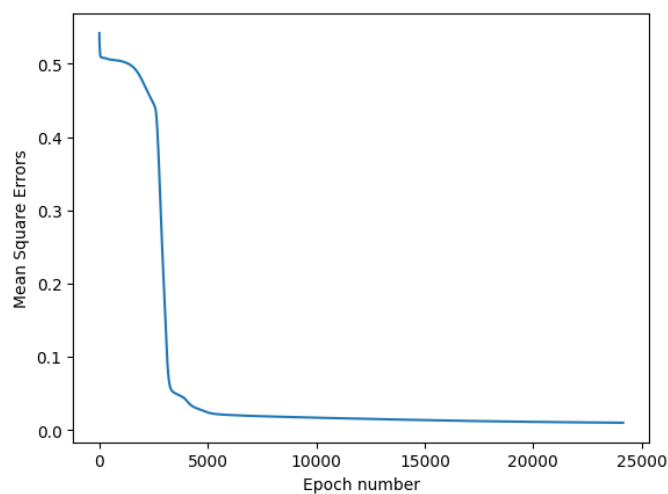


Here, approx 24147 epochs were used to converge MSE to 0.01. I can see that for initial epochs the MSE is very large and then it reduces.

Below is the array for MSE obtained :

```
1524.6108229296, 0.542018697263829, 0.5303619663186013, 0.5327218338390594  
....., 0.010000610234480069, 0.010000397445203679, 0.010000184677296412,  
0.009999971930748348]
```

As for the 0th epoch the MSE is very large compared to other epochs, the graph looks like a straight line but when zoomed in (after removing the 0th epoch) we can see that it is not a straight line.



Plot of desired output and actual output :  $(x_i, d_i), i = 1, \dots, n.$  and  $(x, f(x, w_0)), x \in [0, 1]$

