



Machine Learning

8. week

- Multi Layer Perceptron (MLP) Network
- Back-propagation algorithm
- Stopping Criteria

Umut ORHAN, PhD.

1

1



Multi Layer Perceptron Network

Multi layer perceptron network (MLP) is proposed after Adaline and Perceptron methods have failed producing non-linear solutions to the problem.

MLP is improved in terms of both architectural and training algorithm.

MLP is formed by many neurons each having non-linear activation function.

MLP uses back-propagation learning algorithm in addition to the advantages of Perceptron and Adaline methods.

Umut ORHAN, PhD.

2

2

MLP Structure

General structure of a neuron used in MLP.

General structure of an MLP network.

Umut ORHAN, PhD. 3

3

Activation Function

Activation function is basically a linear or non-linear function that transforms a variable from one dimension into another dimension. Ease of its derivability increases the training speed. Sigmoid, hyperbolic tangent and step functions are the most frequently used activation functions.

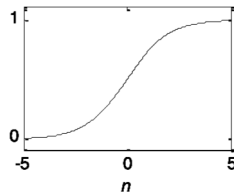
Umut ORHAN, PhD. 4

4

Activation Function

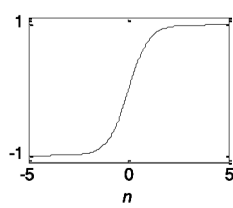
Sigmoid

$$y = \frac{1}{1 + e^{-net}}$$



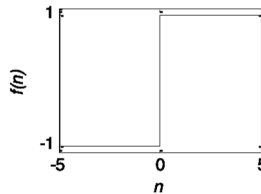
Hyperbolic tangent

$$y = \frac{1 - e^{-2net}}{1 + e^{2net}}$$



Step function

$$y = \begin{cases} 1 & net \geq 0 \\ -1 & net < 0 \end{cases}$$



Umut ORHAN, PhD.

5

5

Backpropagation

All of the weight values are calculated by a method called gradient descent which minimizes the error function. Calculated error is adjusted by the derivation of the activation function while going back from the output end of the neuron to the input end. The applied error value to the derivation of the function is distributed to all of the weights proportional to neurons' input values.

Umut ORHAN, PhD.

6

6



Backpropagation

$$E = \frac{1}{2} \sum_j e_j^2 = \frac{1}{2} \sum_j (d_j - y_j)^2$$

$$\begin{aligned} \Delta w_{ij} &= -\eta \frac{\partial E}{\partial w_{ij}} = -\eta \left[\frac{\partial E}{\partial y_j} \right] \left[\frac{\partial y_j}{\partial net_i} \right] \left[\frac{\partial net_i}{\partial w_{ij}} \right] \\ &= -\eta [-e_j] \left[\frac{\partial y_j}{\partial net_i} \right] [x_i] \end{aligned}$$

If sigmoid is used as activation function;

$$\Delta w_{ij} = -\eta [-e_j] [(1 - y_j) y_j] [x_i] = \eta e_j (1 - y_j) y_j x_i$$

Umut ORHAN, PhD.

7

7



Backpropagation

Prove the statement below for $y = \frac{1 - e^{-2net}}{1 + e^{2net}}$

$$-\eta \frac{\partial E}{\partial w_{ij}} = \eta e_j (1 - y_j^2) x_i$$

Umut ORHAN, PhD.

8

8

Momentum Effect

While system converges through the optimal result step by step using gradient descent, sometimes it gets stuck to one of the local minimum points and the global target cannot be reached. To prevent this, a portion (α) of value change occurred in previous step is used as a momentum effect.

$$\Delta w_{ij}(t+1) = -\eta \frac{\partial E}{\partial w_{ij}} + \alpha \Delta w_{ij}(t)$$

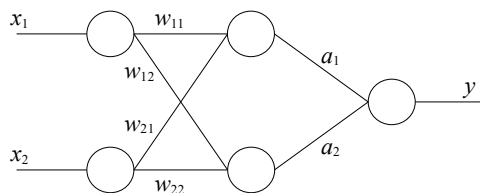
Umut ORHAN, PhD.

9

9

Example

Lets solve XOR problem for an MLP network which has two inputs and one output and has two neurons on its single hidden layer.



$$\eta = 0.5$$

$$\alpha = 0.1$$

$$w_{11} = 1$$

$$w_{12} = -2$$

$$w_{21} = -1$$

$$w_{22} = 2$$

$$a_1 = 3$$

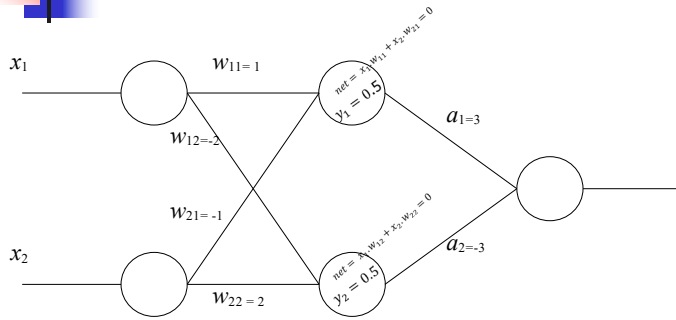
$$a_2 = -3$$

Umut ORHAN, PhD.

10

10

Example



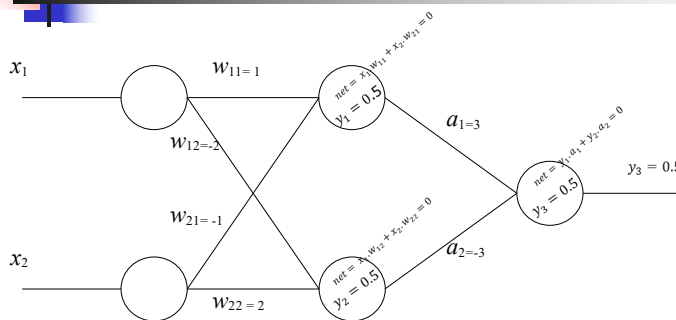
$\eta = 0.5$
 $\alpha = 0.1$
 $w_{11} = 1$
 $w_{12} = -2$
 $w_{21} = -1$
 $w_{22} = 2$
 $a_1 = 3$
 $a_2 = -3$

Umut ORHAN, PhD.

11

11

Example



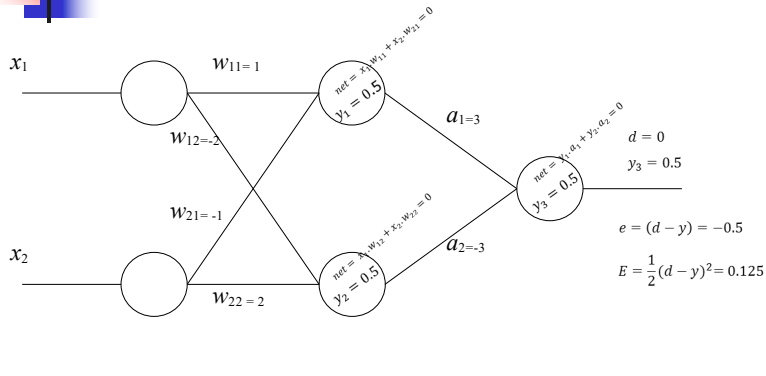
$\eta = 0.5$
 $\alpha = 0.1$
 $w_{11} = 1$
 $w_{12} = -2$
 $w_{21} = -1$
 $w_{22} = 2$
 $a_1 = 3$
 $a_2 = -3$

Umut ORHAN, PhD.

12

12

Example

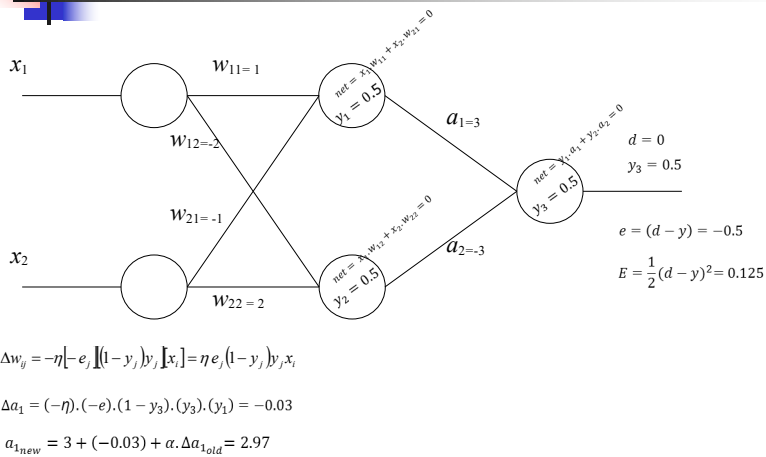


Umut ORHAN, PhD.

13

13

Example

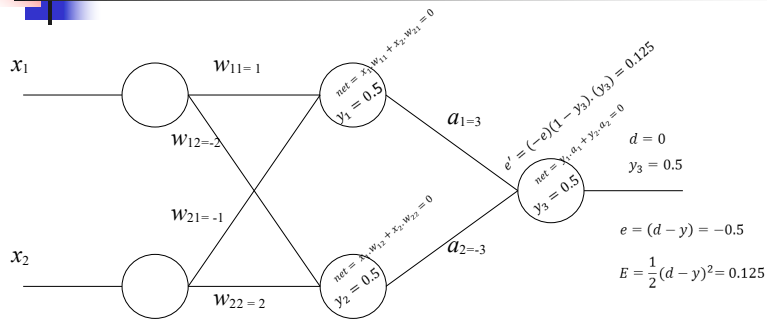


Umut ORHAN, PhD.

14

14

Example



$\eta = 0.5$
 $\alpha = 0.1$
 $w_{11} = 1$
 $w_{12} = -2$
 $w_{21} = -1$
 $w_{22} = 2$
 $a_1 = 3$
 $a_2 = -3$

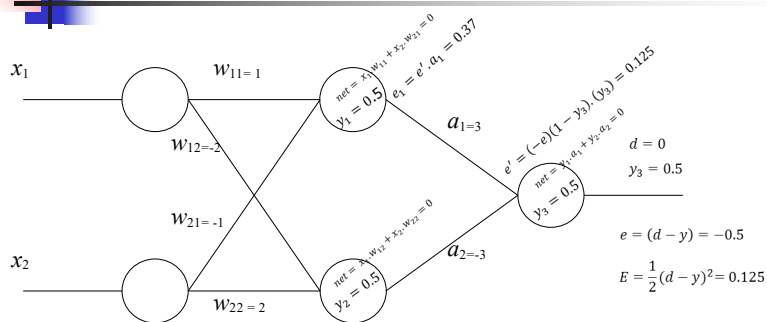
$\Delta w_{ij} = -\eta [-e_j] [1 - y_j] y_j [x_i] = \eta e_j (1 - y_j) y_j x_i$
 $\Delta a_1 = (-\eta) \cdot (-e) \cdot (1 - y_3) \cdot (y_3) \cdot (y_1) = -0.03$
 $a_{1_{new}} = 3 + (-0.03) + \alpha \cdot \Delta a_{1_{old}} = 2.97$

Umut ORHAN, PhD.

15

15

Example



$\eta = 0.5$
 $\alpha = 0.1$
 $w_{11} = 1$
 $w_{12} = -2$
 $w_{21} = -1$
 $w_{22} = 2$
 $a_1 = 3$
 $a_2 = -3$

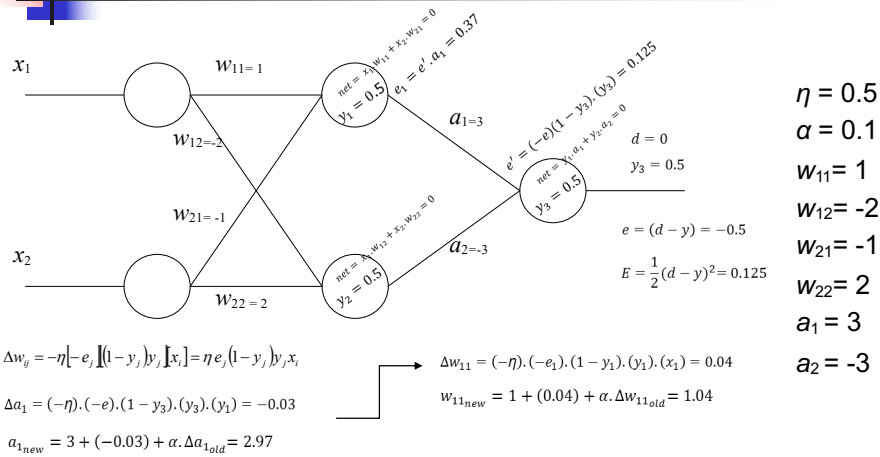
$\Delta w_{ij} = -\eta [-e_j] [1 - y_j] y_j [x_i] = \eta e_j (1 - y_j) y_j x_i$
 $\Delta a_1 = (-\eta) \cdot (-e) \cdot (1 - y_3) \cdot (y_3) \cdot (y_1) = -0.03$
 $a_{1_{new}} = 3 + (-0.03) + \alpha \cdot \Delta a_{1_{old}} = 2.97$

Umut ORHAN, PhD.

16

16

Example



Umut ORHAN, PhD.

17

17

Stopping Criteria

Without a stopping criteria to determine when to stop the training, overfitting occurs. The training set, which is a subset of the entire set, is divided into two and while one part of it is used for updating the weight values, the other part is used for validation to calculate the success of the training. If success drops rapidly, training is stopped.

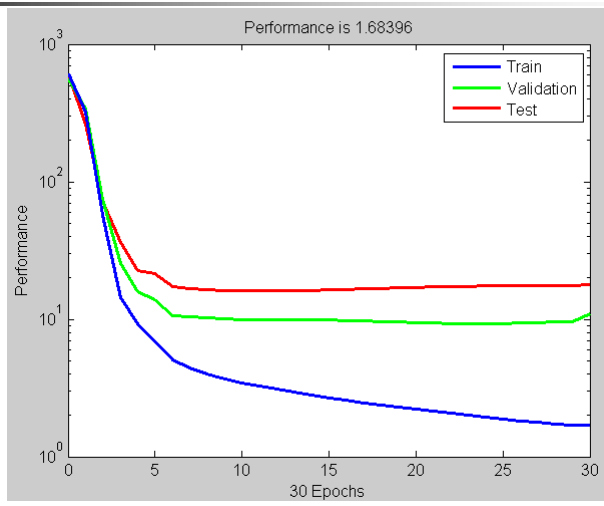
Umut ORHAN, PhD.

18

18



Stopping Criteria



Umut ORHAN, PhD.

19

19