

# Homework 04

Deepak Sharma

CSCI 630-01 Foundations of Intelligent Systems

February 28, 2018

## 1

Suppose we generate a training set from a decision tree and then apply decision-tree learning to that training set. Will the learning algorithm eventually produce the original decision tree used to generate the sample, if the training set size approaches infinity? Why or why not?

### **Solution:**

Since all possible samples of the dataset were available while constructing two decision trees, hence both decision trees know the true distribution of the data and always give same output for any valid input sample. But these two decision trees generated using same dataset can have different configuration they will be logically equivalent .

A Decision tree perform sequentially multiple logical operations. Here parameters of these operation are values of attributes in the dataset. Since for a input sample the end result may be indepenent of the order of applicable operations, hence, each decision tree with different configuration might be executing such possible logically equivalent sequence of operations for an input sample. In the end each decision tree logically equivalent set of operations and produced the same results. Though, decision tree constructor algorithm use entropy heuristic while deciding the next attribute, but in large dataset with high number of attributes, two attributes can provide same information gain after expansion. Hence the algorithm can pick different attributes for expansion during the two independent executions.

For example, we have a dataset with 2 attributes: grades in FIS and grade in algorithms. All students who work in DPRL lab have scored A grade in both courses. Now one decision tree may first check grade in Algorithm course and then in FIS course for deciding whether he/she will get a chance to work in the lab or not. On the other hand, second decision tree many first check algorithm course's grade and then FIS grade for deciding whether the applicant will get a chance or not.

## 2

Below is a data set for a binary class i cation problem (classes 0 and 1), for which each input has three binary attributes (with values inf0,1g). Each row in the table represents a training instance of the form ((A1;A2;A3),y), where (A1;A2;A3) are the inputs to the decision tree classier, and y is the correct output (i.e. the output for the 'true' function the learned tree

should represent). Use the decision tree learning algorithm (see text, p. 702, or the lecture slides) to learn a decision tree using Information Gain to select the best attribute from available attributes(i.e. attributes that are unused along the current path from the root to the current node in the tree at each recursive call). Show the computations that determine how attributes are selected for the tree.

**Solution:**

Sample	A1	A2	A3	y
x1	1	0	0	0
x2	1	0	1	0
x3	0	1	0	0
x4	1	1	1	1
x5	1	1	0	1

Table 1: data table

### LEVEL=1

Let entropy of the random binary variable y for the provided data in table 2 is  $E(A)$ .

Let  $n$  and  $p$  represent the number of samples with 0 label and 1 label respectively.

Given  $n = 3$ ,  $p = 2$

Let:  $k = p/(n + p) = 2/5$

Using the definition of entropy for a binary distribution [3]:

$$E(A) = B(k) = -1 * ((k) \log_2(k) + (1 - k) \log_2(k))$$

$$E(A) = B(2/5) = -1 * ((2/5) \log_2(2/5) + (3/5) \log_2(3/5)) = 0.972$$

### CASE 1

If we choose attribute **A1** for expansion of the decision tree then

For distribution of **A1 = 1**,  $k = p_{A1=1}/(n_{A1=1} + p_{A1=1}) = 2/4$

$$E(A1 = 1) = B(1/2) = -1 * ((2/4) \log_2(2/4) + (2/4) \log_2(2/4)) = 1 \quad (1)$$

For the distribution of **A1 = 0**,  $k = p/(n + p)0$

$$E(A1 = 0) = B(0) = -1 * ((0) \log_2(0) + (1) \log_2(1))$$

$$B(0) = 0 \quad (2)$$

Total Entropy for the resultant distribution will be a weighted average of the entropy for **A1 = 1** and **A1 = 0** distribution:

$$E(A1) = ((n_{A1=1} + p_{A1=1})/(n + p)) * E(A1 = 1) + ((n_{A1=0} + p_{A1=0})/(n + p)) * E(A1 = 0)$$

using equation 1 and 2:

$$E(A1) = (4/5) * E(A1 = 1) + (1/5) * E(A1 = 0) = (4/5) * 1 + (1/5) * 0 = .80$$

information gain after the expansion  $\mathbf{I}_{A1}$ :

$$\mathbf{I}_{A1} = \mathbf{E}(\mathbf{A}) - \mathbf{E}(\mathbf{A1}) = 0.972 - 0.80 = .172 \quad (3)$$

### CASE 2

If we decide attribute  $\mathbf{A2}$  for expansion of the decision tree then for the distribution of  $\mathbf{A2} = 1$ ,  $\mathbf{k} = \mathbf{p}_{A2=1}/(\mathbf{p}_{A2=1} + \mathbf{n}_{A2=1}) = 2/3$

$$\mathbf{E}(\mathbf{A2} = 1) = \mathbf{B}(2/3) = -1 * ((2/3) \log_2(2/3) + (1/3) \log_2(1/3)) = .918 \quad (4)$$

For distribution of  $\mathbf{A2} = 0$ ,  $\mathbf{k} = \mathbf{p}_{A2=0}/(\mathbf{p}_{A2=0} + \mathbf{n}_{A2=0}) = 0$

Using equation 2:

$$\mathbf{E}(\mathbf{A2} = 0) = \mathbf{B}(0) = 0$$

Total Entropy for the resultant distribution will be a weighted average of the entropy for  $\mathbf{A2} = 1$  and  $\mathbf{A2} = 0$  distribution:

$$\mathbf{E}(\mathbf{A2}) = ((\mathbf{n}_{A2=1} + \mathbf{p}_{A2=1})/(\mathbf{n} + \mathbf{p})) * \mathbf{E}(\mathbf{A2} = 1) + ((\mathbf{n}_{A2=0} + \mathbf{p}_{A2=0})/(\mathbf{n} + \mathbf{p})) * \mathbf{E}(\mathbf{A2} = 0)$$

using equation 4 and 1:

$$\mathbf{E}(\mathbf{A2}) = (3/5) * \mathbf{E}(\mathbf{A2} = 1) + (2/5) * \mathbf{E}(\mathbf{A2} = 0) = .546$$

Information gain after the expansion  $\mathbf{I}_{A2}$ :

$$\mathbf{I}_{A2} = \mathbf{E}(\mathbf{A}) - \mathbf{E}(\mathbf{A2}) = 0.972 - 0.546 = .426 \quad (5)$$

### CASE 3

If we decide attribute  $\mathbf{A3}$  for expansion of the decision tree then for distribution of  $\mathbf{A3} = 1$ :  $\mathbf{k} = \mathbf{p}_{A3=1}/(\mathbf{p}_{A3=1} + \mathbf{n}_{A3=1}) = 1/2$

Using equation 1:

$$\mathbf{E}(\mathbf{A3} = 1) = \mathbf{B}(1/2) = 1$$

For distribution of  $\mathbf{A3} = 0$ :  $\mathbf{k} = 1/3$

Using equation 4:

$$\mathbf{E}(\mathbf{A3} = 0) = \mathbf{B}(1/3) = \mathbf{B}(2/3) = .918 \quad (6)$$

Total Entropy for the resultant distribution will be a weighted average of the entropy for  $\mathbf{A3} = 1$  and  $\mathbf{A3} = 0$  distribution:

$$\mathbf{E}(\mathbf{A3}) = ((\mathbf{n}_{A3=1} + \mathbf{p}_{A3=1})/(\mathbf{n} + \mathbf{p})) * \mathbf{E}(\mathbf{A3} = 1) + ((\mathbf{n}_{A3=0} + \mathbf{p}_{A3=0})/(\mathbf{n} + \mathbf{p})) * \mathbf{E}(\mathbf{A3} = 0)$$

$$\mathbf{E}(\mathbf{A3}) = (2/5) * \mathbf{B}(1/2) + (3/5) * \mathbf{B}(1/3) = (2/5) * 1 + (3/5) * .918 = .946$$

information gain after the expansion  $\mathbf{I}_{A3}$ :

$$\mathbf{I}_{A3} = \mathbf{E}(\mathbf{A}) - \mathbf{E}(\mathbf{A3}) = 0.972 - 0.946 = .026 \quad (7)$$

Using equation 3, 5 and 7 we can say:

$$\mathbf{I}_{A2} > \mathbf{I}_{A1} > \mathbf{I}_{A3}$$

Sample	A1	A3	y
x3	0	0	0
x4	1	1	1
x5	1	0	1

Table 2: data distribution for A2=1

Sample	A1	A3	y
x1	1	0	0
x2	1	1	0

Table 3: data distribution for A2=0

Therefore, we will decide to expend decision tree with A2 attribute. In the resultant distribution has been shown in table 2 and 3 for A2=1 and A2=0 values respectively. Since the resultant distribution for A2=0 (table 3) belongs to only y=0, hence this distribution is a base case for decision tree algorithm and it can not be expended any further.

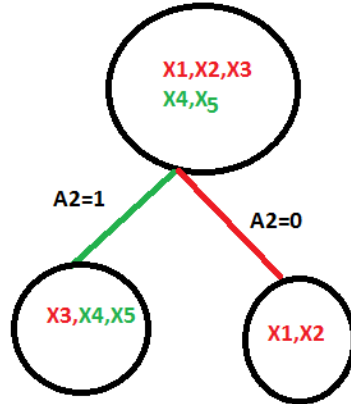


Figure 1: After expending the decision tree with attribute A2

## LEVEL=2

For the data distribution of table 2 entropy is:

$$E(A) = B(2/3) = .91$$

## CASE=1

Expending the decision tree for the attribute A1:

For distribution of **A1 = 1**:  $k = (p_{A1=1} / (n_{A1=1} + p_{A1=1})) = 2/2$

Using equation 2 :

$$E(A1 = 1) = B(1) = B(0) = 0$$

For distribution of  $\mathbf{A1} = 0$ :  $k(p_{\mathbf{A1}=0}/(n_{\mathbf{A1}=0} + p_{\mathbf{A1}=0})) = 0$   
using equation 2

$$\mathbf{E}(\mathbf{A1} = 0) = \mathbf{B}(0) = 0$$

Total Entropy for the resultant distribution will be a weighted average of the entropy for  $\mathbf{A1} = 1$  and  $\mathbf{A1} = 0$  distribution:

$$\mathbf{E}(\mathbf{A1}) = ((n_{\mathbf{A1}=1} + p_{\mathbf{A1}=1})/(n + p)) * \mathbf{E}(\mathbf{A1} = 1) + ((n_{\mathbf{A1}=0} + p_{\mathbf{A1}=0})/(n + p)) * \mathbf{E}(\mathbf{A1} = 0)$$

$$\mathbf{E}(\mathbf{A1}) = (2/3) * \mathbf{E}(\mathbf{A1} = 1) + (1/3) * \mathbf{E}(\mathbf{A1} = 0) = 0$$

Information gain after the expansion  $\mathbf{I_{A1}}$ :

$$\mathbf{I_{A1}} = \mathbf{E}(\mathbf{A}) - \mathbf{E}(\mathbf{A1}) = .91 - 0.0 = .91 \quad (8)$$

## CASE 2

If we select to expend binary tree using A3 attribute then for distribution of  $\mathbf{A3} = 1$ :  
 $k = (p/(n + p)) = 1/1$

Using equation 2 :

$$\mathbf{E}(\mathbf{A3} = 1) = \mathbf{B}(1) = 0$$

For distribution of  $\mathbf{A3} = 0$ :  $k = p/(n + p) = 1/2$   
using equation 1

$$\mathbf{E}(\mathbf{A3} = 0) = \mathbf{B}(1/2) = 1$$

Total Entropy for the resultant distribution will be a weighted average of the entropy for  $\mathbf{A3} = 1$  and  $\mathbf{A3} = 0$  distribution:

$$\mathbf{E}(\mathbf{A1}) = ((n_{\mathbf{A3}=1} + p_{\mathbf{A3}=1})/(n + p)) * \mathbf{E}(\mathbf{A3} = 1) + ((n_{\mathbf{A3}=0} + p_{\mathbf{A3}=0})/(n + p)) * \mathbf{E}(\mathbf{A3} = 0)$$

$$\mathbf{E}(\mathbf{A3}) = (1/3) * \mathbf{E}(\mathbf{A3} = 1) + (2/3) * \mathbf{E}(\mathbf{A3} = 0) = 2/3$$

Information gain after the expansion  $\mathbf{I_{A1}}$ :

$$\mathbf{I_{A3}} = \mathbf{E}(\mathbf{A}) - \mathbf{E}(\mathbf{A3}).91 - 0.67 = .24 \quad (9)$$

Using equation 8 and 9 we can say:

$$\mathbf{I_{A1}} > \mathbf{I_{A3}}$$

Hence we will select attribute A1 for the expansion of the binary tree.

After expansion, in the resultant distribution for A1=1, all the samples have value 1 for target variable y. Similarly, for A1=0 all the samples have value 0 for target variable y. Therefore we can not expend the tree any further as both these distribution are the base case of the algorithm.

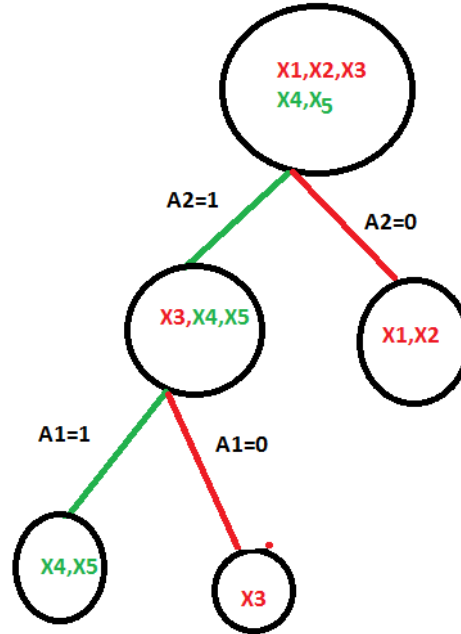


Figure 2: After expanding the decision tree with attribute A2 and A1

### 3

Provide a neural network that computes the XOR function for two binary inputs (i.e. with values in  $\{0,1\}$ ). Make sure to provide the explicit network architecture (including activation functions), and all weights in the network. Note: you should not train your network using a program to answer this question set the weights of your network manually.

**Solution**

Weight	Value	Weight	Value	Weight	Value
$W_{14}$	1	$W_{15}$	1	$W_{24}$	1
$W_{25}$	1	$W_{34}$	-.5	$W_{35}$	-1.5
$W_{47}$	1	$W_{57}$	-1	$W_{67}$	0

Table 4: Weights for the Neural network executing XOR function

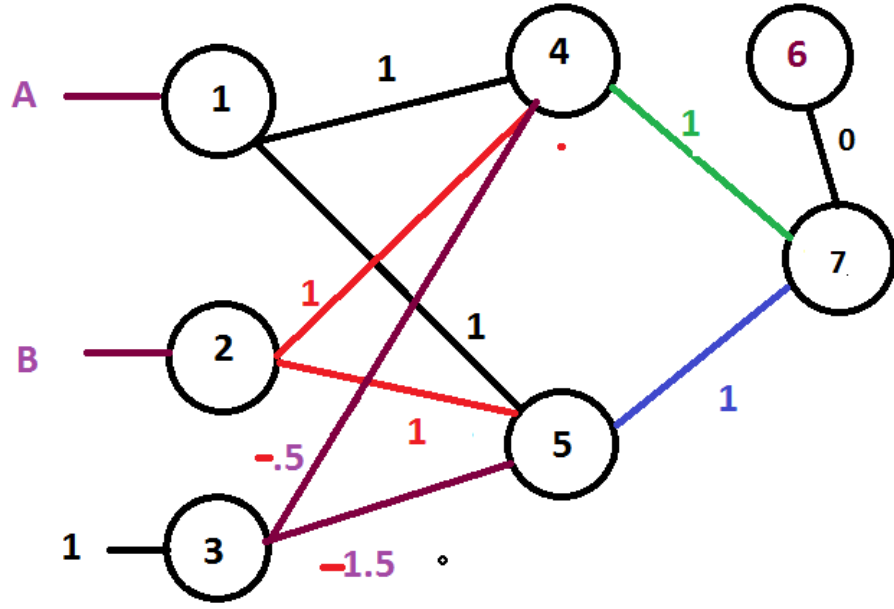


Figure 3: Neural Network with network weights for performing XOR operation. A and B are binary inputs where A and B  $\in \{0, 1\}$

### Activation function

$$g(X) = \begin{cases} 0, & \text{if } X < 0 \\ 1, & \text{otherwise} \end{cases} \quad (10)$$

Activation value of neuron, bias and input units:

$a_1(input)$	$a_2(input)$	$a_4$	$a_5$	$a_7(output)$
0	0	0	0	0
0	1	1	0	1
1	0	1	0	1
1	1	1	1	0

Table 5: Inputs and Activation of the neurons in neural network exhibiting XOR functionality

### Calculation

For input  $a_1=0, a_2=0$

$$a_4 = g(a_1 * W_{14} + a_2 * W_{24} + a_3 * W_{34}) = g(.1 + .1 + 1 * (-.5)) = g(-.5) = 0$$

$$a_5 = g(a_1 * W_{15} + a_2 * W_{25} + a_3 * W_{35}) = g(0 * 1 + 0 * 1 + 1 * (-1.5)) = g(-1.5) = 0$$

For input  $\mathbf{a}_1=0, \mathbf{a}_2=1$

$$\mathbf{a}_4 = \mathbf{g}(1 * 1 + .1 + 1 * (-.5)) = \mathbf{g}(.5) = 1$$

$$\mathbf{a}_5 = \mathbf{g}(0 * 1 + 1 * 1 + 1 * (-1.5)) = \mathbf{g}(-0.5) = 0$$

For input  $\mathbf{a}_1=1, \mathbf{a}_2=0$

$$\mathbf{a}_4 = \mathbf{g}(1 * 1 + .1 + 1 * (-.5)) = \mathbf{g}(.5) = 1$$

$$\mathbf{a}_5 = \mathbf{g}(1 * 1 + 0 * 1 + 1 * (-1.5)) = \mathbf{g}(-0.5) = 0$$

For input  $\mathbf{a}_1=1, \mathbf{a}_2=1$

$$\mathbf{a}_4 = \mathbf{g}(1 * 1 + .1 * 1 + 1 * (-.5)) = \mathbf{g}(1.5) = 1$$

$$\mathbf{a}_5 = \mathbf{g}(1 * 1 + 1 * 1 + 1 * (-1.5)) = \mathbf{g}(0.5) = 1$$

For input  $\mathbf{a}_4=0, \mathbf{a}_5=0$

$$\mathbf{a}_7 = \mathbf{g}(\mathbf{a}_4 * \mathbf{W}_{45} + \mathbf{a}_5 * \mathbf{W}_{57} + \mathbf{a}_6 * \mathbf{W}_{67}) = \mathbf{g}(0 * 1 + 0 * 1 + 1 * (0)) = \mathbf{g}(0) = 0$$

For input  $\mathbf{a}_4=1, \mathbf{a}_5=0$

$$\mathbf{a}_7 = \mathbf{g}(1 * 1 + 0 * (-1) + 1 * (0)) = \mathbf{g}(1) = 1$$

For input  $\mathbf{a}_4=1, \mathbf{a}_5=1$

$$\mathbf{a}_7 = \mathbf{g}(1 * 1 + 1 * (-1) + 1 * (0)) = \mathbf{g}(0) = 0$$

## 4

**Hyper Parameters during experiment:**

**Learning Rate:.1**

**Epochs:1000**

**Outputs:**

**Linear dataset final SSE: 011176**

**Non Linear dataset final SSE 0.14259**



	correct	incorrect
Class 0	20	0
Class 1	16	0

Table 6: Number of samples predicted Correct and Incorrect using the weight gained after training on iris dataset(Linearly separable) [2]

	correct	incorrect
Class 0	59	41
Class 1	93	7

Table 7: Number of samples predicted Correct and Incorrect using the weight gained after training with artificially generated dataset(not linearly separable) [1]

## Results:

After performing logistic regression using batch gradient decent optimizer on linearly separable and non-separable dataset, we have found that after each epoch the Sum of Square Error has reduced. Initially rate of SSE reduction was fast but latter the process become slow as shown in figure 4 and 5. After a number of epochs, the SSE was reducing very less with each epoch. The rate of reduction for SSE also depends on the learning rate, with high learning rate SSE reduces early. But the final value of SSE after 10000 epoch was grater than value(.1425) achieved by a low learning rates(.001, value=.1420 non linear dataset). Also, if we perform experiment with very high learning rate (grater than 2 for non linear dataset) than the SSE was not reducing with each epoch. Running the algorithm for high number of epochs might results in overfitting on the test dataset [3].

The suggested one neuron network was not able to separate non-linear distribution as the underneath hypothesis which gradient decent algorithm was optimizing, was a linear hypothesis of form.

$$\mathbf{y} = \mathbf{w}_1 * \mathbf{x}_1 + \mathbf{w}_2 * \mathbf{x}_2 + \mathbf{bias}$$

by finding a right combination of  $\mathbf{w}_1$  ,  $\mathbf{w}_2$  and bias. Since the given data distribution can not be defined by any linear equation the gradient decent algorithm was not able to reduce the loss under certain value.

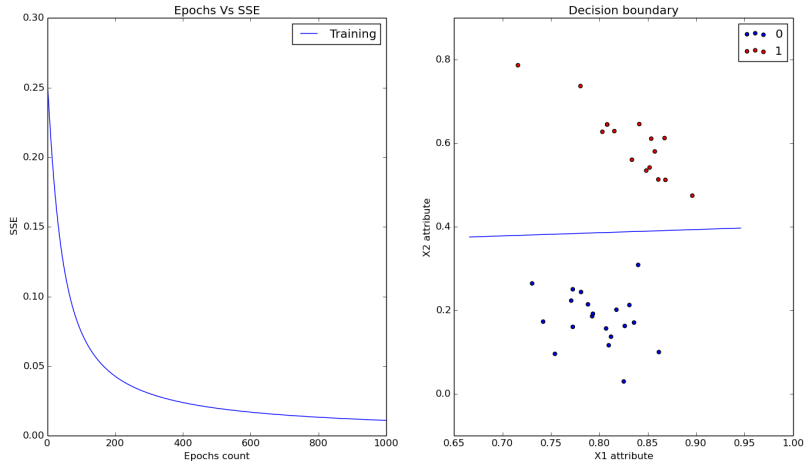


Figure 4: SSE vs Epoch curve and Decision Boundary for a linearly separable dataset [2]

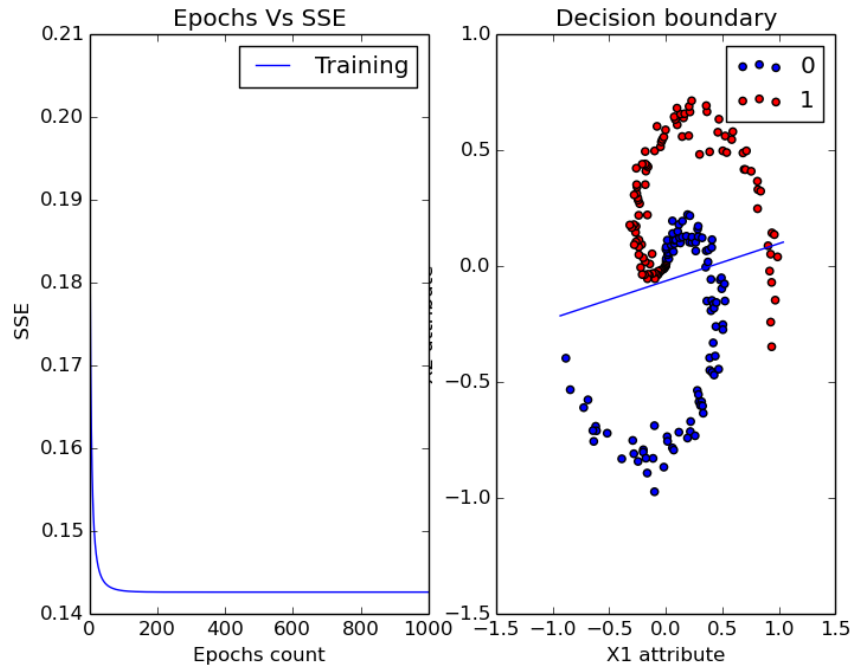


Figure 5: SSE vs Epoch curve and Decision Boundary for a not linearly separable dataset (artificially generated [1])

## 5

For the neural network shown in figure 6, compute the updated weights for the remaining edges in the network. Then compute the output values at nodes ( $v_1; v_2$ ) using the updated weights for the same training sample. The network weights are provided in table .

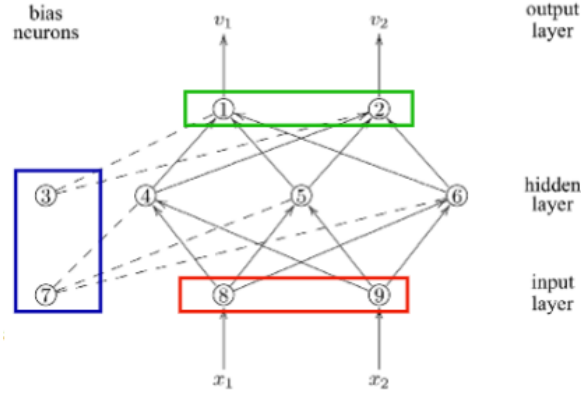


Figure 6: Given MLP

Weight	Value	Weight	Value	Weight	Value	Weight	Value
$W_{31}$	.43	$W_{41}$	.05	$W_{51}$	.71	$W_{61}$	.75
$W_{32}$	.63	$W_{42}$	.57	$W_{52}$	.96	$W_{62}$	.74
$W_{74}$	.55	$W_{84}$	.82	$W_{94}$	.96		
$W_{75}$	.26	$W_{85}$	.67	$W_{95}$	.06		
$W_{76}$	.60	$W_{86}$	1.0	$W_{96}$	.36		

Table 8: Weights for the given Neural network

Node	Activation	Node	Activation	Node	Activation
$v_1$	.8488	$v_2$	.9267	$v_3$	1
$v_4$	.7738	$v_5$	.8235	$v_6$	.9038
$v_7$	1	$v_8$	2	$v_9$	-1

Table 9: Activation of neurons and bias units in given Neural Network

**Updated Weights After Backpropagation:**

**Activation values generated using updated weights:**

Weight	Value	Weight	Value	Weight	Value	Weight	Value
$\mathbf{W}_{31}$	.4191	$\mathbf{W}_{41}$	.0416	$\mathbf{W}_{51}$	.6910	$\mathbf{W}_{61}$	.7402
$\mathbf{W}_{32}$	.6305	$\mathbf{W}_{42}$	.5704	$\mathbf{W}_{52}$	.9604	$\mathbf{W}_{62}$	.7405
$\mathbf{W}_{74}$	.54995	$\mathbf{W}_{84}$	.8199	$\mathbf{W}_{94}$	.9601		
$\mathbf{W}_{75}$	.25896	$\mathbf{W}_{85}$	.6679	$\mathbf{W}_{95}$	.0610		
$\mathbf{W}_{76}$	.5993	$\mathbf{W}_{86}$	.9986	$\mathbf{W}_{96}$	.3607		

Table 10: Updated Weights after backpropagation

Node	Activation	Node	Activation	Node	Activation
$\mathbf{v}_1$	.8440	$\mathbf{v}_2$	.9263	$\mathbf{v}_3$	1
$\mathbf{v}_4$	.7738	$\mathbf{v}_5$	.8226	$\mathbf{v}_6$	.9034
$\mathbf{v}_7$	1	$\mathbf{v}_8$	2	$\mathbf{v}_9$	-1

Table 11: New Activation values produced by neurons and bias units using updated weights in forward propagation

## Calculations:

During the backpropagation using the error value in prediction, we can calculate  $\delta$  for all layers for updating the previous weights. [3]

For output layer the value of  $\delta$  will be:

$$\delta_1 = (\mathbf{v}_1 - \text{label}_1) * \mathbf{v}_1 * (1 - \mathbf{v}_1) = (.8488 - 0) * .8488 * (1 - .8488) = .1089$$

$$\delta_2 = (\mathbf{v}_2 - \text{label}_2) * \mathbf{v}_2 * (1 - \mathbf{v}_2) = (.9267 - 1) * .9267 * (1 - .9267) = -.0050$$

Using the  $\delta$  values of the output layer, we can calculate  $\delta$  for previous hidden layer.

$$\delta_3 = (\delta_1 * \mathbf{W}_{31} + \delta_2 * \mathbf{W}_{32}) * \mathbf{v}_3 * (1 - \mathbf{v}_3) = 0$$

$$\delta_4 = (\delta_1 * \mathbf{W}_{41} + \delta_2 * \mathbf{W}_{42}) * \mathbf{v}_4 * (1 - \mathbf{v}_4) = .0005$$

$$\delta_5 = (\delta_1 * \mathbf{W}_{51} + \delta_2 * \mathbf{W}_{52}) * \mathbf{v}_5 * (1 - \mathbf{v}_5) = .01038218978$$

$$\delta_6 = (\delta_1 * \mathbf{W}_{61} + \delta_2 * \mathbf{W}_{62}) * \mathbf{v}_6 * (1 - \mathbf{v}_6) = .00677958004$$

Using the  $\delta$  values of the second hidden layer, we can calculate  $\delta$  for the first hidden layer.

$$\delta_7 = (\delta_4 * \mathbf{W}_{74} + \delta_5 * \mathbf{W}_{75} + \delta_6 * \mathbf{W}_{76}) * \mathbf{v}_7 * (1 - \mathbf{v}_7) = 0$$

$$\delta_8 = (\delta_4 * \mathbf{W}_{84} + \delta_5 * \mathbf{W}_{85} + \delta_6 * \mathbf{W}_{86}) * \mathbf{v}_8 * (1 - \mathbf{v}_8) = -.028356$$

$$\delta_9 = (\delta_4 * \mathbf{W}_{94} + \delta_5 * \mathbf{W}_{95} + \delta_6 * \mathbf{W}_{96}) * \mathbf{v}_9 * (1 - \mathbf{v}_9) = -.007104$$

Using learning rate  $\eta=.1$ , respective  $\delta$  values and inputs to the neuron we can update respective weights.

$$\mathbf{W}_{31} = \mathbf{W}_{31} - \eta * \delta_1 * \mathbf{v}_3 = .43 - .1 * .1089 = 0.41911$$

$$\begin{aligned}
W_{32} &= W_{32} - \eta * \delta_1 * v_3 = .63 - .1 * (-.0050) * 1 = 0.6305 \\
W_{41} &= W_{41} - \eta * \delta_1 * v_4 = .0500 - .1 * .1089 * .7738 = .041573318 \\
W_{42} &= W_{42} - \eta * \delta_2 * v_4 = .5700 + .1 * .005 * .7738 = .5703869 \\
W_{51} &= W_{51} - \eta * \delta_1 * v_5 = .70 - .1 * .1089 * .8235 = .691032085 \\
W_{52} &= W_{52} - \eta * \delta_2 * v_5 = .96 + .1 * .005 * .8235 = .96041175 \\
W_{61} &= W_{61} - \eta * \delta_1 * v_6 = .75 - .1 * .1089 * .9038 = .740157618 \\
W_{62} &= W_{62} - \eta * \delta_2 * v_6 = .74 + .1 * .005 * .9038 = .7404519 \\
W_{74} &= W_{74} - \eta * \delta_4 * v_7 = .55 - .1 * .0005 * 1 = .54995 \\
W_{75} &= W_{75} - \eta * \delta_5 * v_7 = .26 - .1 * .0103822 = .2589617 \\
W_{76} &= W_{76} - \eta * \delta_6 * v_7 = .60 - .1 * .00677958 * 1 = .599322042 \\
W_{84} &= W_{84} - \eta * \delta_4 * v_8 = .82 - .1 * .0005 * 2 = 0.8199 \\
W_{85} &= W_{85} - \eta * \delta_5 * v_8 = .67 - .1 * .01038 * 2 = .667924 \\
W_{86} &= W_{86} - \eta * \delta_6 * v_8 = 1 - .1 * .006779 * 2 = .9986442 \\
\\ 
W_{94} &= W_{94} - \eta * \delta_4 * v_9 = .96 - .1 * .0005 * (-1) = .96005 \\
W_{95} &= W_{95} - \eta * \delta_5 * v_9 = .06 - .1 * .01038 * (-1) = .061038 \\
W_{96} &= W_{96} - \eta * \delta_6 * v_9 = .36 - .1 * .006779 * (-1) = .3606779
\end{aligned}$$

Forward propagation with updated weights:

$$\begin{aligned}
v_4 &= g(v_8 * W_{84} + v_9 * W_{94} + v_7 * W_{74}) = g(2 * .8199 - 1 * .96005 + 1 * .54995) \\
v_4 &= g(1.2297) = .77376606294 \\
\\ 
v_5 &= g(v_8 * W_{85} + v_9 * W_{95} + v_7 * W_{75}) = g(2 * .667924 - 1 * .061038 + 1 * .2589617) \\
v_5 &= g(1.5337717) = .82255748887 \\
\\ 
v_6 &= g(v_8 * W_{86} + v_9 * W_{96} + v_7 * W_{76}) = g(2 * .9986442 - 1 * .36106779 + 1 * .599322042) \\
v_6 &= g(2.235542652) = .9033961575 \\
\\ 
v_1 &= g(v_4 * W_{41} + v_5 * W_{51} + v_6 * W_{61} + v_3 * W_{31})
\end{aligned}$$

$$\begin{aligned} \mathbf{v}_1 &= g(.77376606294 * .041573318 + .82255748887 * .691032085 \\ &\quad + .90345530776 * .740157618 + 1 * 0.41911) \\ \mathbf{v}_1 &= g(1.68839096772) = .84401243926 \end{aligned}$$

$$\begin{aligned} \mathbf{v}_2 &= g(\mathbf{v}_4 * \mathbf{W}_{42} + \mathbf{v}_5 * \mathbf{W}_{52} + \mathbf{v}_6 * \mathbf{W}_{62} + \mathbf{v}_3 * \mathbf{W}_{32}) \\ \mathbf{v}_2 &= g(.77376606294 * .5703869 + .82255748887 * .961032085 \\ &\quad + .90345530776 * .7404519 + 1 * .6305) \\ \mathbf{v}_2 &= g(2.53131536372) = .92630819223 \end{aligned}$$

## References

- [1] karpathy. cs231 training a softmax linear classifier, 2016. URL <http://cs231n.github.io/neural-networks-case-study/>.
- [2] M. Lichman. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml/datasets/Iris>.
- [3] Stuart J. Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach*. Pearson Education, 2 edition, 2003. ISBN 0137903952.