

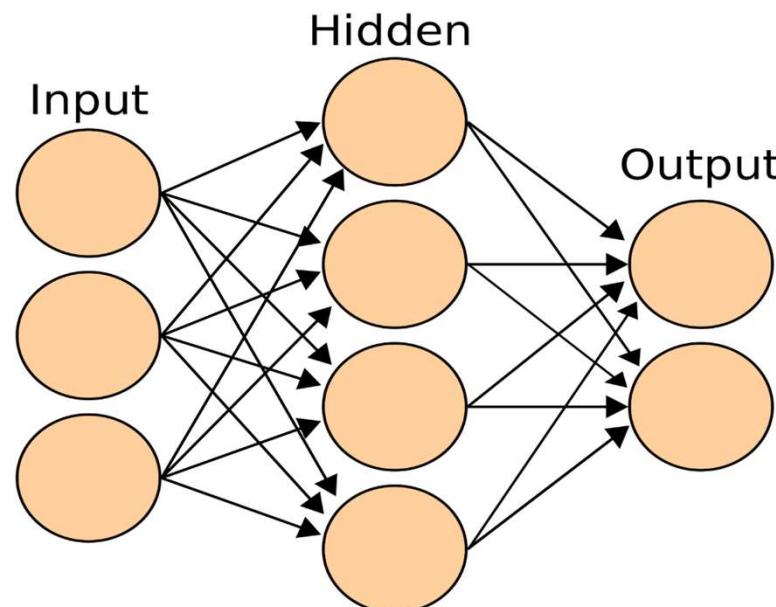
# 6. Multilayer Perceptron Neural Networks

## 6.1 Introduction

In this part, we are to study the multilayer feed-forward network, which is an important class of neural networks. Typically, the network consists of a set of inputs (source nodes) that constitute the input layer, one or more hidden layers of computational nodes, and an output layer of computational nodes. The input signals propagates through the network in a forward direction, on a layer-by-layer basis. These neural networks are commonly referred to as multilayer perceptrons (MLP).

A MLP neural network has three distinctive characteristics as summarized below:

- (1) Smooth activation functions such as sigmoid function are used.
- (2) The network contains one or more hidden layers that are not part of the input or output of the network. The hidden neurons enable the network to learn complex tasks.
- (3) The network exhibits a high degree of connectivity (fully connected).



## **6.2 Preliminary: 1-D gradient descent method**

Suppose we have a 1-D minimization problem:

$$J(x) = (a - bx)^2$$

Where  $a$  and  $b$  are parameters. The goal of the optimization problem is to find a value of  $x$  so that  $J(x)$  is minimized.

To minimize the cost function  $J(x)$ , we take the first-order derivative:

$$\frac{\partial J(x)}{\partial x} = -2b(a - bx)$$

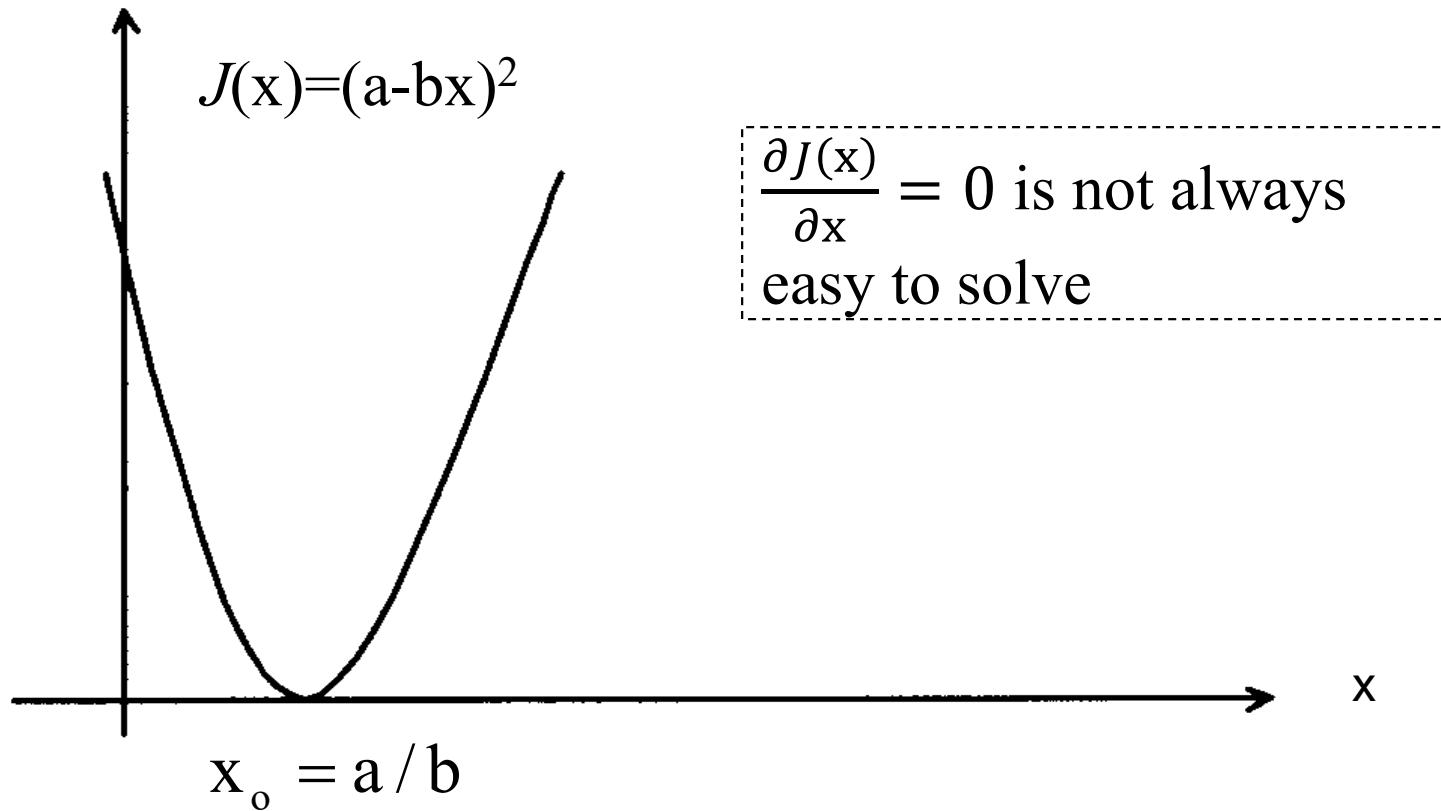
Solve the equation:

$$\frac{\partial J(x)}{\partial x} = 0$$

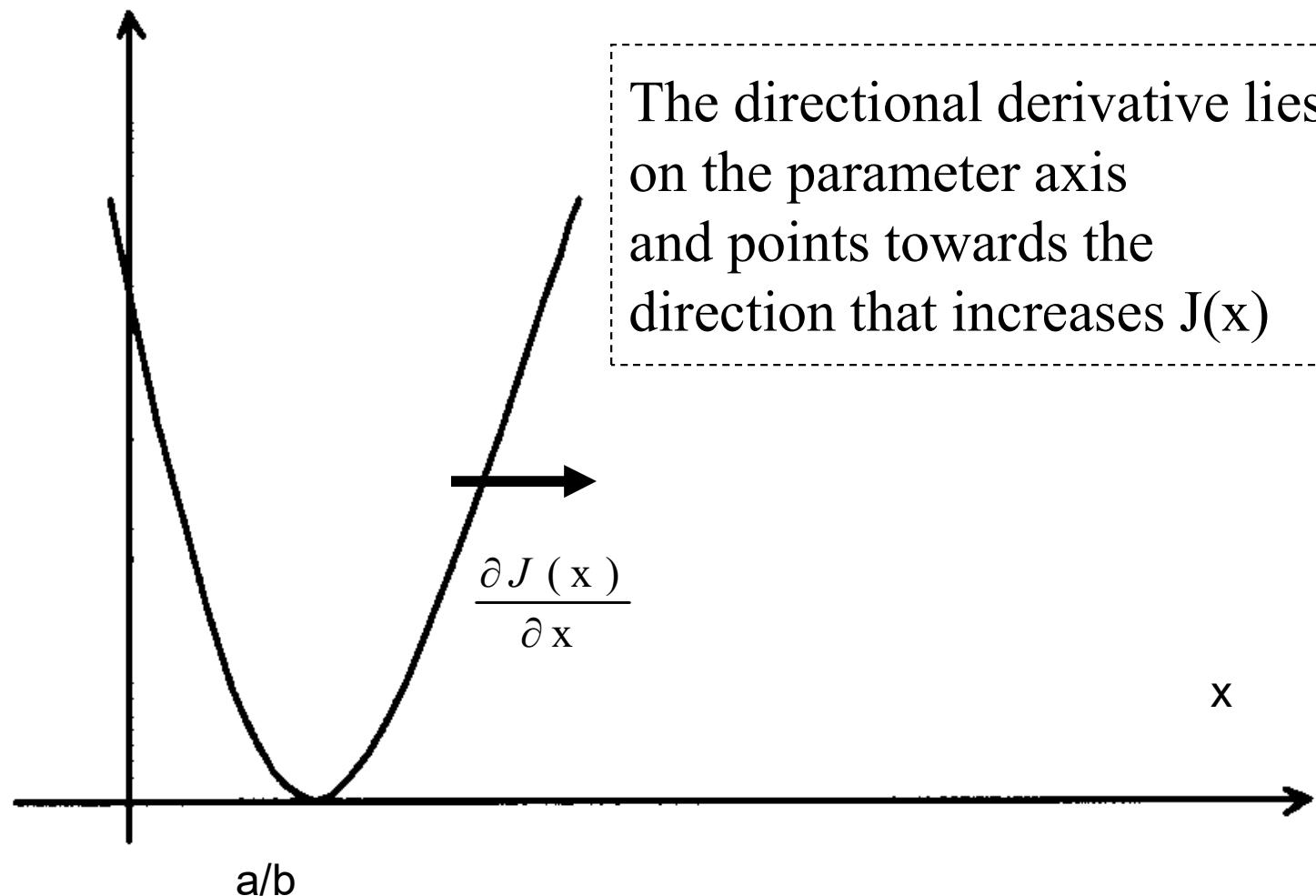
We obtain the optimal solution:

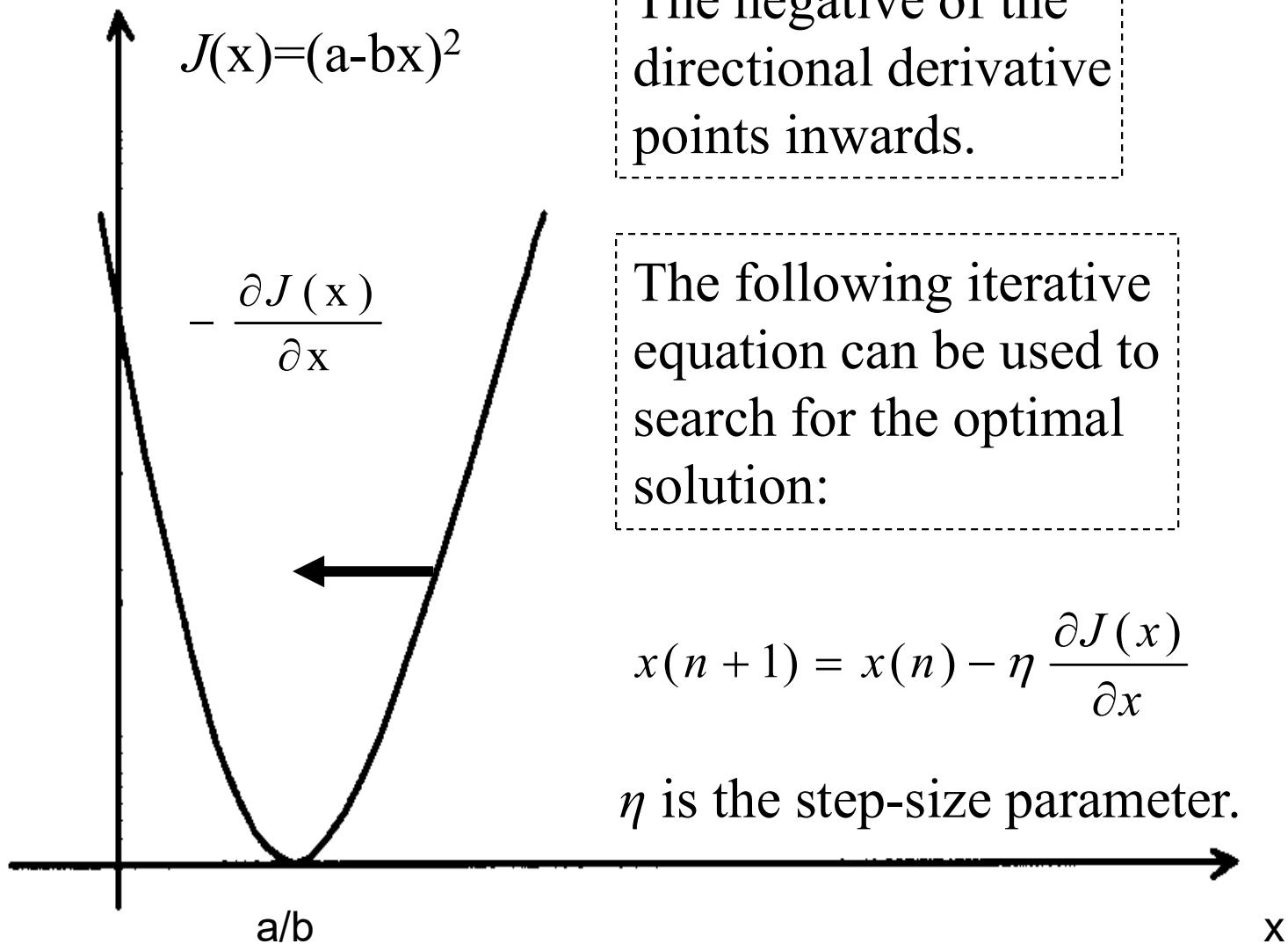
$$x_0 = a / b$$

The following is the graphical illustration:



$$J(x) = (a - bx)^2$$





## 6.3 The MLP neural network learning

Given a set of training (learning) samples:

$$\{\mathbf{x}(1), \mathbf{d}(1)\}, \{\mathbf{x}(2), \mathbf{d}(2)\}, \dots, \{\mathbf{x}(N), \mathbf{d}(N)\}$$

where

$$\mathbf{x}(i) = [x_1(i), x_2(i), \dots, x_n(i)]^T$$

$$\mathbf{d}(i) = [d_1(i), d_2(i), \dots, d_m(i)]^T$$

The objective of learning is to construct a MLP neural network so that the network response approximates the desired response  $\mathbf{d}(i)$  well:

$$\mathbf{y}[\mathbf{x}(i)] \rightarrow \mathbf{d}(i)$$

The error signal of the output neuron  $j$  at iteration  $n$  (i.e. when the  $n$ -th training sample is presented to the network) is defined by:

$$e_j(n) = d_j(n) - y_j(n) \quad (7.1)$$

Where  $d_j(n)$  and  $y_j(n)$  denote desired response and actual response of neuron  $j$  at  $n$ -th iteration.

Then the total instantaneous squared error at iteration  $n$  is obtained by summing the error of all neurons in the output layer:

$$E(n) = \frac{1}{2} \sum_{j \in C} e_j^2(n) \quad (7.2)$$

Where  $C$  is the set of output neurons.

Assume there are  $N$  training samples, the squared error is obtained by summing  $E(n)$  over all  $n$ :

$$E_{total} = \sum_{n=1}^N E(n) \quad (7.3)$$

$$E_{ave} = \frac{1}{N} \sum_{n=1}^N E(n) \quad (7.4)$$

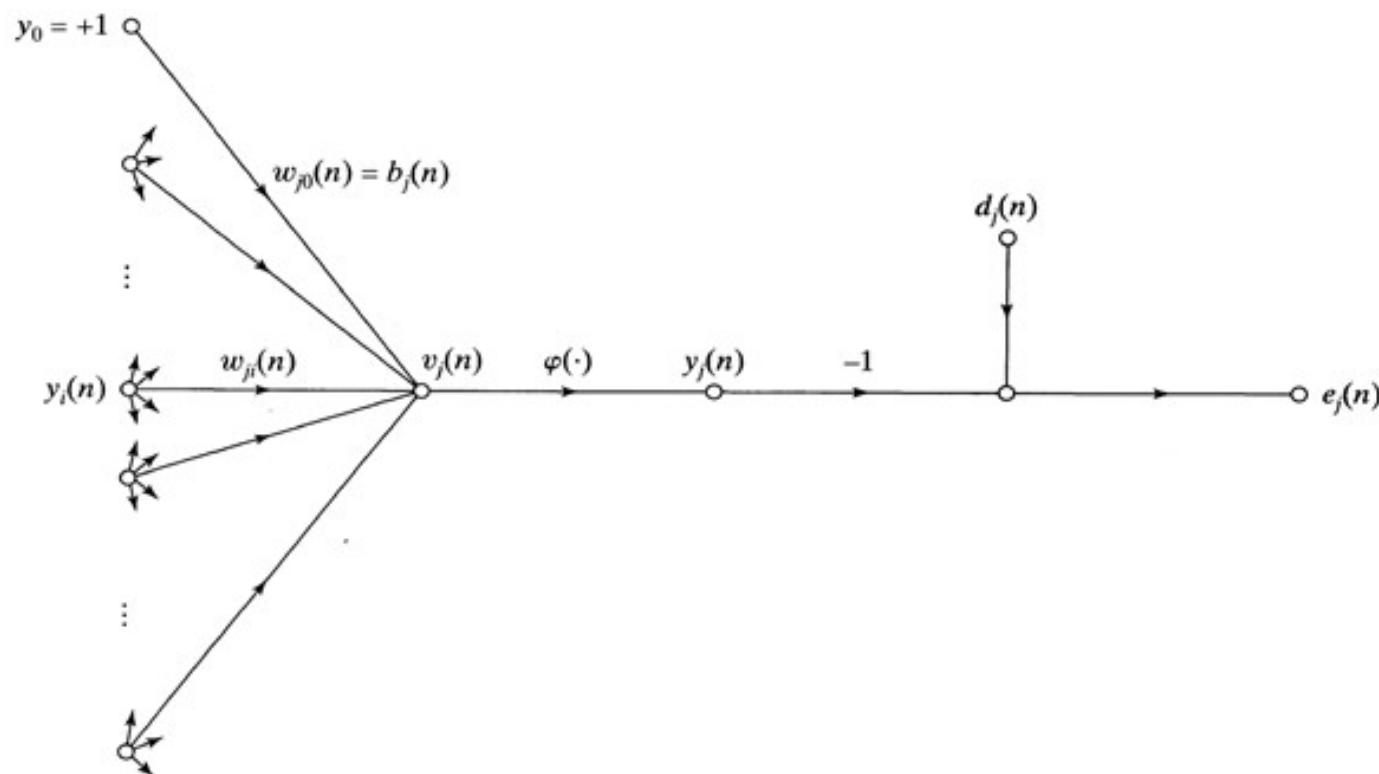
The instantaneous error  $E(n)$  and therefore the average error  $E_{ave}$  is a function of all free parameters, i.e. synaptic weights of the network. For a given training data set,  $E_{ave}$  represents the cost function as a measure of learning performance.

The objective of learning is to adjust the parameters (i.e. the weights) of the network to minimize  $E_{ave}$ . To perform this minimization, we consider a sequential training, in which the weights are updated on a sample-by-sample basis until the entire training samples have been used.

Consider the following diagram, where the neuron  $j$  is fed by a set of signals produced by a layer of neurons to its left. The activation and the output of the neuron  $j$  is therefore:

$$v_j(n) = \sum_{i=0}^m w_{ji}(n) y_i(n) \quad (7.5)$$

$$y_j(n) = \varphi_j[v_j(n)] \quad (7.6)$$



The back-propagation algorithm applies a correction  $\Delta w_{ji}(n)$  to the synaptic weight  $w_{ji}(n)$ , which is proportional to the partial derivative of  $E(n)$  with respect to  $w_{ji}(n)$ .

According to the chain rule of calculus, we may express this derivative (gradient) as:

$$\frac{\partial E(n)}{\partial w_{ji}(n)} = \frac{\partial E(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \frac{\partial v_j(n)}{\partial w_{ji}(n)} \quad (7.7)$$

The partial derivative represents a sensitivity factor that determines the direction of search in weight space for the synaptic weight  $w_{ji}(n)$ .

Differentiating both sides of Eqn (7.2) with respect to  $e_j(n)$ , we get:

$$\frac{\partial E(n)}{\partial e_j(n)} = e_j(n) \quad (7.8)$$

Differentiating both sides of Eqn (7.1) with respect to  $y_j(n)$ , we get:

$$\frac{\partial e_j(n)}{\partial y_j(n)} = -1 \quad (7.9)$$

Differentiating Eqn (7.6) with respect to  $v_j(n)$ , we obtain:

$$\frac{\partial y_j(n)}{\partial v_j(n)} = \varphi'_j[v_j(n)] \quad (7.10)$$

Differentiating Eqn (7.5) with respect to  $w_{ji}(n)$ , we obtain:

$$\frac{\partial v_j(n)}{\partial w_{ji}(n)} = y_i(n) \quad (7.11)$$

Substituting Eqns (7.8)-(7.11) into Eqn (7.7), yields:

$$\frac{\partial E(n)}{\partial w_{ji}(n)} = -e_j(n) \varphi'_j[v_j(n)] y_i(n) \quad (7.12)$$

The correction applied to  $w_{ji}(n)$  is defined by:

$$\Delta w_{ji}(n) = -\eta \frac{\partial E(n)}{\partial w_{ji}(n)} = \eta \delta_j(n) y_i(n) \quad (7.13)$$

Where  $\eta$  is the learning rate parameter of the back-propagation algorithm, and  $\delta_j(n)$  is the local gradient defined by:

$$\begin{aligned} \delta_j(n) &= -\frac{\partial E(n)}{\partial v_j(n)} \\ &= -\frac{\partial E(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \\ &= e_j(n) \varphi'_j[v_j(n)] \end{aligned} \quad (7.14)$$

The above equation shows that the local gradient points to the required changes in synaptic weights. The local gradient is equal to the product of the corresponding error signal for that neuron and the derivative of the associated activation function.

In the above, we note that a key factor involved in the calculation of the weight adjustment  $\Delta w_{ji}(n)$  is the error signal  $e_j(n)$  at the output neuron  $j$ . In this context, we may identify two distinct cases, depending on where in the network neuron  $j$  is located.

### **Case 1: neuron $j$ is at the output layer**

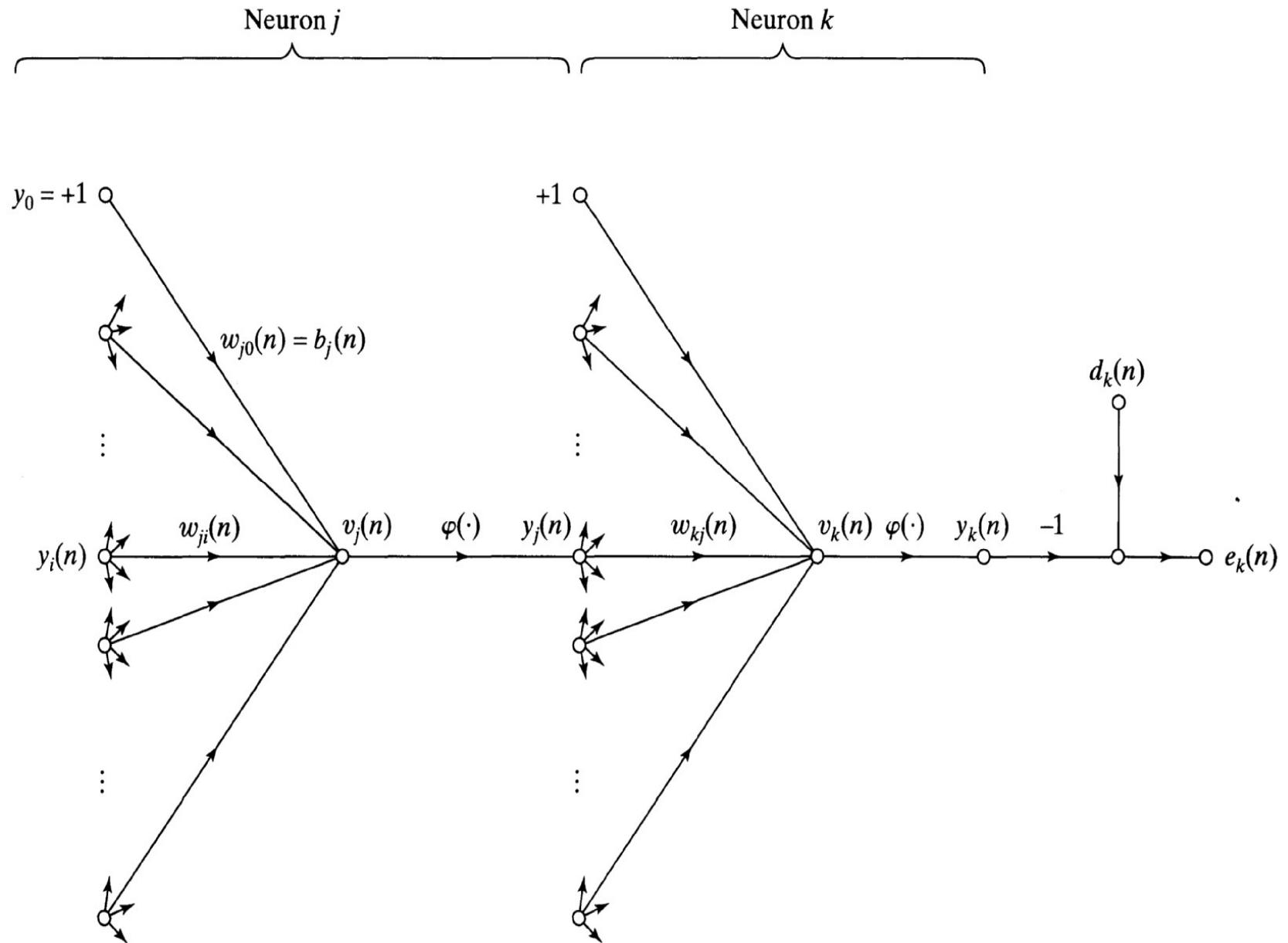
When neuron  $j$  is located in the output layer of the network, it is supplied with a desired response of its own. Thus, it is straightforward to compute the error signal  $e_j(n)$ .

## Case 2: neuron $j$ is at the hidden layer

When neuron  $j$  is located in a hidden layer of the network, there is no supplied desired response for the neuron. Thus, the error signal for a hidden layer neuron has to be determined recursively in terms of the error signals of all neurons to which that hidden layer neuron is directly connected. This makes the back-propagation algorithm very complicated.

The local gradient for a hidden neuron  $j$  is defined as:

$$\begin{aligned}\delta_j(n) &= -\frac{\partial E(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \\ &= -\frac{\partial E(n)}{\partial y_j(n)} \varphi'_j[v_j(n)]\end{aligned}\tag{7.15}$$



To calculate the partial derivative  $\partial E(n)/\partial y_j(n)$ , we may proceed as follows:

$$E(n) = \frac{1}{2} \sum_{k \in C} e_k^2(n) \quad (7.16)$$

Where  $C$  is the set of output neurons.

Differentiating Eqn (7.16) with respect to  $y_j(n)$ , we get:

$$\frac{\partial E(n)}{\partial y_j(n)} = \sum_k e_k \frac{\partial e_k(n)}{\partial y_j(n)} \quad (7.17)$$

Next, we use the chain rule for the partial derivative  $\partial E(n)/\partial y_j(n)$ , and rewrite Eqn (7.17) as:

$$\frac{\partial E(n)}{\partial y_j(n)} = \sum_k e_k \frac{\partial e_k(n)}{\partial v_k(n)} \frac{\partial v_k(n)}{\partial y_j(n)} \quad (7.18)$$

We note that:

$$\begin{aligned} e_k(n) &= d_k(n) - y_k(n) \\ &= d_k(n) - \varphi_k[v_k(n)] \end{aligned} \quad (7.19)$$

Hence, we have:

$$\frac{\partial e_k(n)}{\partial v_k(n)} = -\varphi'_k[v_k(n)] \quad (7.20)$$

We also note from the above diagram that for neuron k, the activation is:

$$v_k(n) = \sum_{j=0}^m w_{kj}(n) y_j(n) \quad (7.21)$$

Differentiating (7.21) with respect to  $y_j(n)$ , yields:

$$\frac{\partial v_k(n)}{\partial y_j(n)} = w_{kj}(n) \quad (7.22)$$

Substituting Eqns (7.20) and (7.22) into Eqn (7.18), yields:

$$\begin{aligned}\frac{\partial E(n)}{\partial y_j(n)} &= - \sum_k e_k(n) \varphi'_k(v_k(n)) w_{kj}(n) \\ &= - \sum_k \delta_k(n) w_{kj}(n)\end{aligned}\quad (7.23)$$

Substituting Eqn (7.23) into Eqn (7.15), we get the back-propagation formula for the local gradient of hidden neuron  $j$ :

$$\delta_j(n) = \varphi'_j[v_j(n)] \sum_k \delta_k(n) w_{kj}(n) \quad (7.24)$$

The factor  $\varphi'_j[v_j(n)]$  involved in the computation of the local gradient in the above equation depends solely on the activation functions associated with hidden neuron  $j$ . The remaining factor  $\delta_k(n)$  requires knowledge of the error signal  $e_k(n)$ , for all neurons that lie in the layer to the immediate right of hidden neuron  $j$ , and  $w_{kj}(n)$  are synaptic weights associated.

Now, we summarize the relations that we have derived for the back-propagation algorithm.

(1) First, the correction  $\Delta w_{ji}(n)$  applied to the synaptic weight connecting neuron  $i$  to neuron  $j$  is defined by the delta rule:

$$\begin{pmatrix} \text{weight} \\ \text{correction} \\ \Delta w_{ji}(k) \end{pmatrix} = \begin{pmatrix} \text{learning} \\ \text{rate} \\ \eta \end{pmatrix} \begin{pmatrix} \text{local} \\ \text{gradient} \\ \delta_j(n) \end{pmatrix} \begin{pmatrix} \text{input} \\ \text{signal} \\ y_i(n) \end{pmatrix}$$

(2) Second, the local gradient depends on whether neuron  $j$  is an output or hidden neuron:

- (i) If neuron  $j$  is hidden, the gradient is computed using Eqn (7.24)
- (ii) If neuron  $j$  is an output neuron, the gradient is computed using Eqn (7.14)

## Two passes of the computation

In the application of the back-propagation algorithm for MLP neural network training, there are two distinct passes: the *forward pass* and the *backward pass*.

### *(1) Forward pass*

In the forward pass, the synaptic weights remain unchanged throughout the network, and the signals of the network are computed on a neuron-by-neuron, layer-by-layer basis.

### *(2) Backward pass*

The backward pass starts from the output layer, passes error signals leftward through the network, layer-by-layer, and recursively computes the local gradient for each neuron and updates the weights accordingly.

## Activation function

The computation of the local gradient and hence weight correction of each neuron of the MLP neural network requires the derivative of the activation function. For this derivative to exist, the activation function needs to be continuous and differentiable. An example of a continuously differentiable nonlinear activation function commonly used in MLP neural network are sigmoid function and hyperbolic tangent function.

### *(1) Sigmoid function*

$$\varphi_j[v_j(n)] = \frac{1}{1 + \exp[-v_j(n)]} \quad (7.25)$$

Where  $v_j(n)$  is the activation signal of neuron  $j$ .

According to this nonlinearity, the amplitude of the output lies in the range of:

$$0 \leq y_j(n) \leq 1$$

Differentiating Eqn (7.25) with respect to  $v_j(n)$ , we get:

$$\phi'_j[v_j(n)] = \frac{\exp[-v_j(n)]}{[1 + \exp[-v_j(n)]]^2} \quad (7.26)$$

Considering  $y_j(n) = \phi_j[v_j(n)]$ , we may express Eqn (7-26) into:

$$\phi'_j[v_j(n)] = y_j(n)[1 - y_j(n)]$$

For a neuron  $j$  located in the output layer, the local gradient maybe expressed as:

$$\begin{aligned}\delta_j(n) &= e_j(n)\phi'_j[v_j(n)] \\ &= [d_j(n) - y_j(n)]y_j(n)[1 - y_j(n)]\end{aligned}\quad (7.27)$$

On the other hand, for an arbitrary hidden neuron  $j$ , we may express the local gradient as:

$$\begin{aligned}\delta_j(n) &= \phi'_j[v_j(n)]\sum_k \delta_k(n)w_{kj}(n) \\ &= y_j(n)[1 - y_j(n)]\sum_k \delta_k(n)w_{kj}(n)\end{aligned}\quad (7.28)$$

## (2) Hyperbolic tangent function

Hyperbolic tangent function is another commonly used activation function defined by:

$$\varphi_j[v_j(n)] = \tanh[v_j(n)] = \frac{\exp[v_j(n)] - \exp[-v_j(n)]}{\exp[v_j(n)] + \exp[-v_j(n)]}\quad (7.29)$$

The derivative of the hyperbolic tangent function with respect to  $v_j(n)$  is given by:

$$\begin{aligned}
 \phi'_j[v_j(n)] &= \operatorname{sech}^2[v_j(n)] \\
 &= (1 - \tanh^2[bv_j(n)]) \\
 &= [1 - y_j(n)][1 + y_j(n)]
 \end{aligned} \tag{7.30}$$

For a neuron  $j$  located in the output layer, the local gradient is:

$$\begin{aligned}
 \delta_j(n) &= e_j(n)\phi'_j[v_j(n)] \\
 &= [d_j(n) - y_j(n)][1 - y_j(n)][1 + y_j(n)]
 \end{aligned} \tag{7.31}$$

For a neuron  $j$  located in a hidden layer, we have

$$\begin{aligned}
 \delta_j(n) &= \phi'_j[v_j(n)] \sum_k \delta_k(n) w_{kj}(n) \\
 &= [1 - y_j(n)][1 + y_j(n)] \sum_k \delta_k(n) w_{kj}(n)
 \end{aligned}$$

## Learn rate parameter $\eta$

The smaller we make the learning rate parameter  $\eta$ , the smaller the changes to synaptic weights will be from one iteration to the next. If we make the learning rate parameter too large, the resulting large changes in the synaptic weights might make the network learning become unstable. A simple method of increasing learning speed, yet avoiding the danger of instability, is to modify the delta rule by including a momentum term:

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \eta \delta_j(n) y_i(n)$$

Where  $\alpha$  is usually a positive number called the momentum constant. For the convergence of the weight estimation, the momentum constant must be restricted to the range:

$$0 \leq \alpha < 1$$

## Modes of training

In a practical application of the back-propagation algorithm, learning results from the many presentations of a prescribed set of training samples to the MLP. One complete presentation of the entire training set during the learning process is called an *epoch*. The learning process is maintained on an epoch-by-epoch basis, until the synaptic weights of the network stabilize. For a given set of training samples, back-propagation learning may proceed in one of the following three basic ways.

### *(i) Sequential mode*

The sequential mode of the back-propagation learning is also referred to as on-line learning. In this mode of operation, weight updating is performed after the presentation of each training sample, one sample at a time.

*(ii) Batch mode*

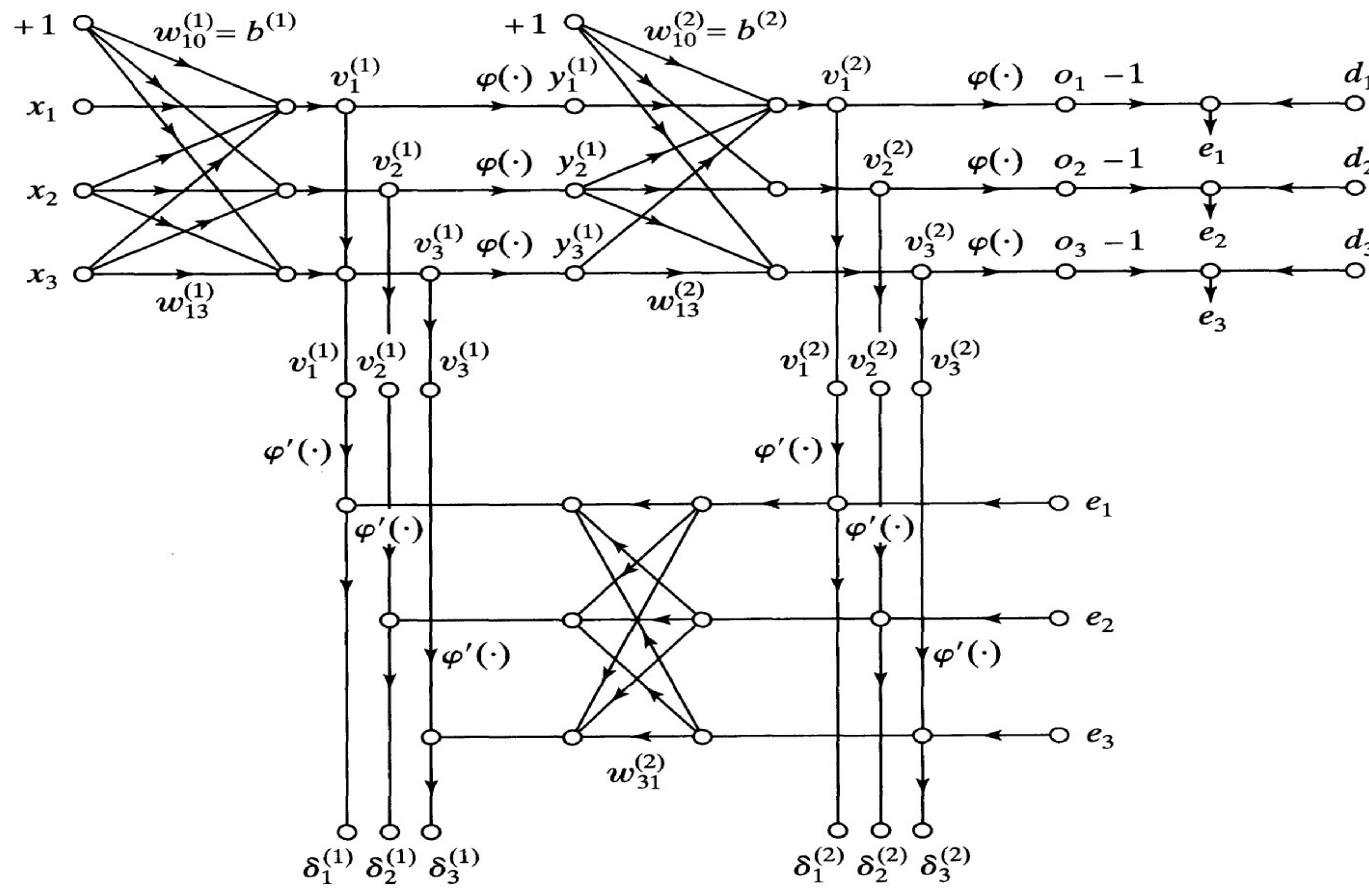
In the batch model, the weight updating is performed when all the samples in the epoch are presented to the network. Details are not described here, since the sequential mode is the commonly used method in MLP neural network training.

*(iii) Mini-batch mode*

Mini-batch splits the training samples into small batches that are used to calculate model error and update model coefficients.

## Summary of the back-propagation (BP) algorithm

The signal flow graph of BP algorithm incorporating both the forward (top part) and backward (lower part) passes in the learning process is shown in the following diagram:



This graph is for the case of 1 hidden layer. Often, the sequential learning is preferred for implementing the back-propagation algorithm. For this mode of operation, the algorithm cycles as follows:

### **BP algorithm (sequential mode)**

#### *(1) Initialization.*

Pick the synaptic weights from a random (uniform or normal) distribution.

#### *(2) Presentation of training samples*

Present the network with an epoch of training samples. For each sample, perform the sequence of forward and backward computations.

### *(2.1) Forward computation*

During forward computation, the weights are fixed. Compute the activation signal and output of each neuron by proceeding through the network, layer by layer.

### *(2.2) Backward computation*

Compute the local gradient of the neurons, and adjust the synaptic weights of network.

### *(3) Iteration*

Iterate the forward and backward computations in steps (2.1)-(2.2) by presenting new epochs of training samples to the network until the stopping criterion is met.

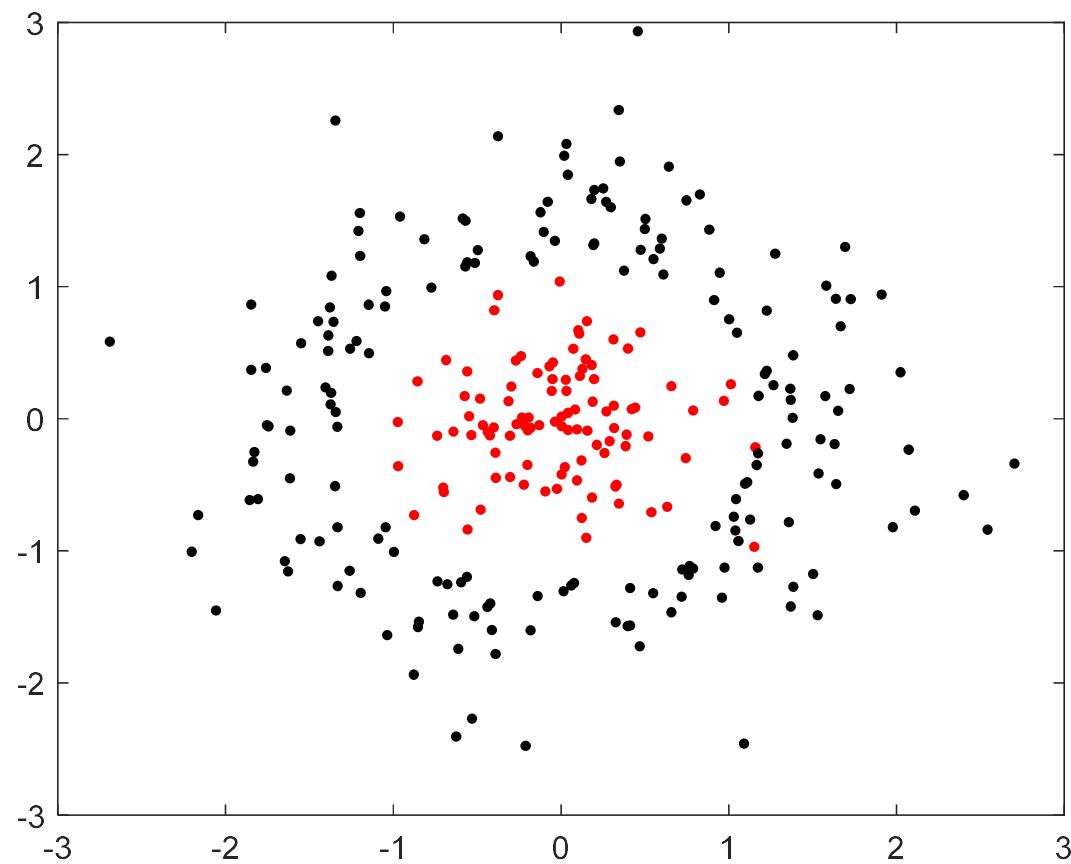
Note, the order of presentation of training samples should be randomized from epoch to epoch.

*The stopping criterion* can be

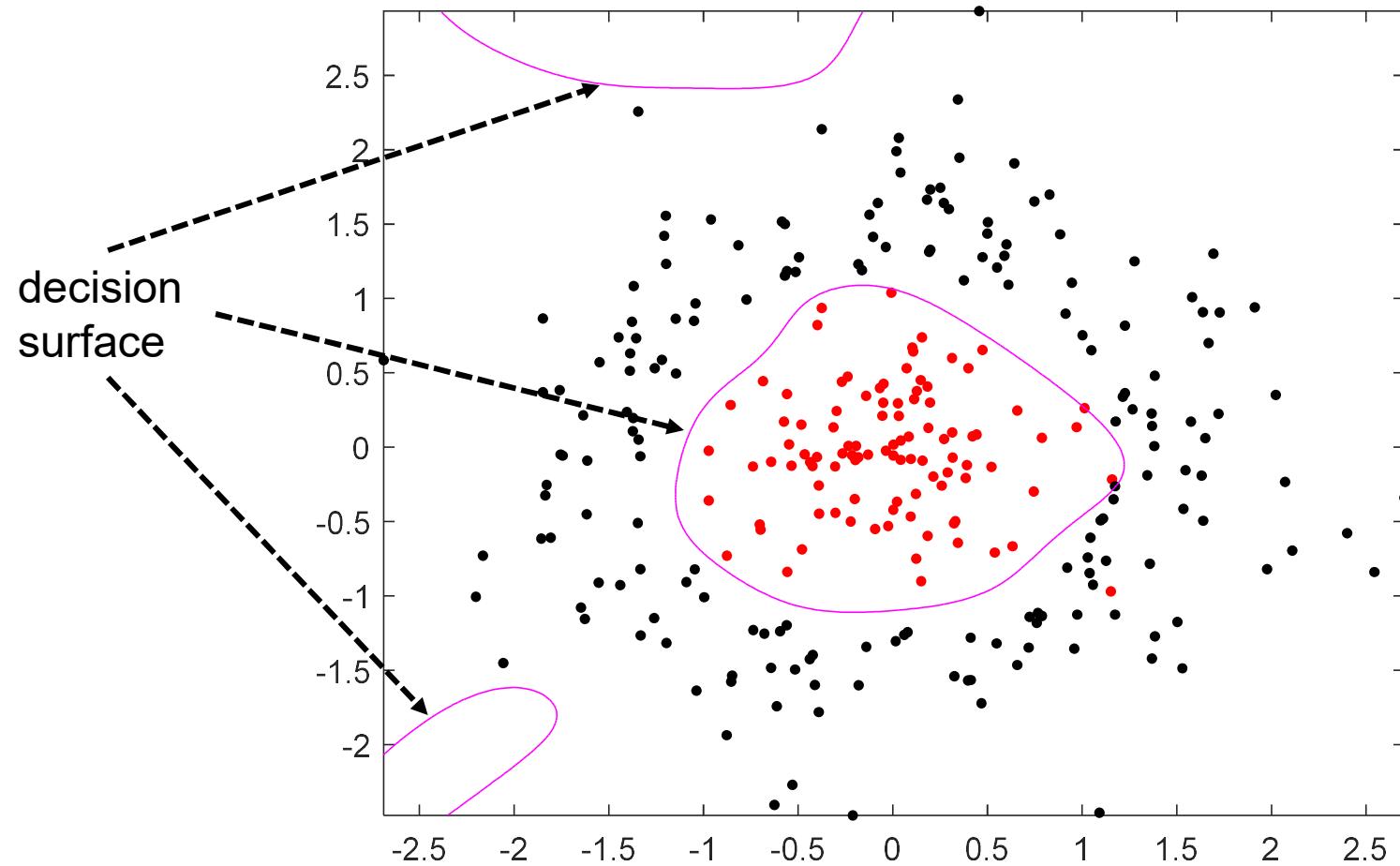
- (i) the number of iterations, or
- (ii) The change of the average error is small enough.

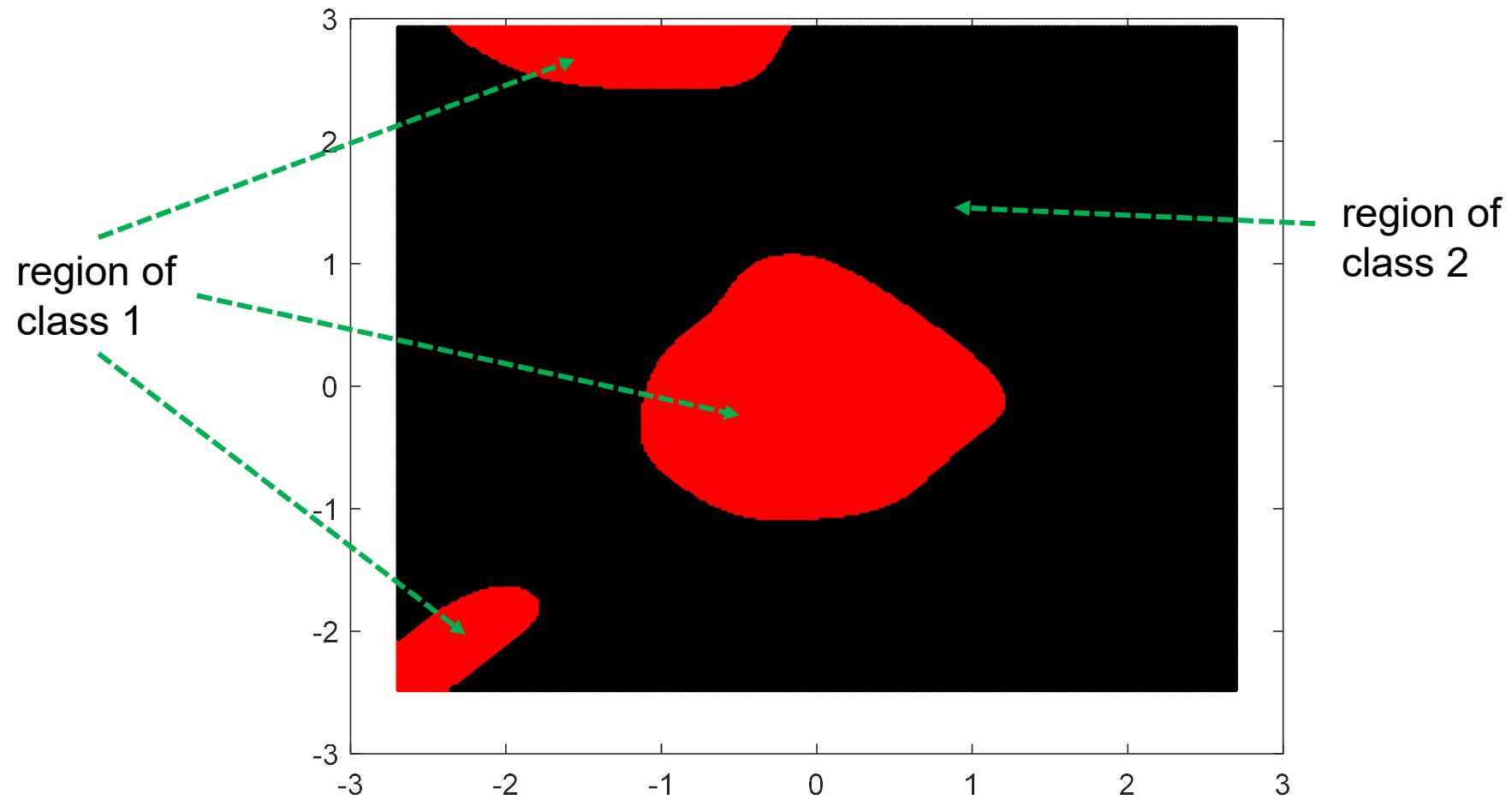
## 6.4 Discussions on MLP neural networks

*Observation 1: The neural network learned is affected by the initial values of the weights*

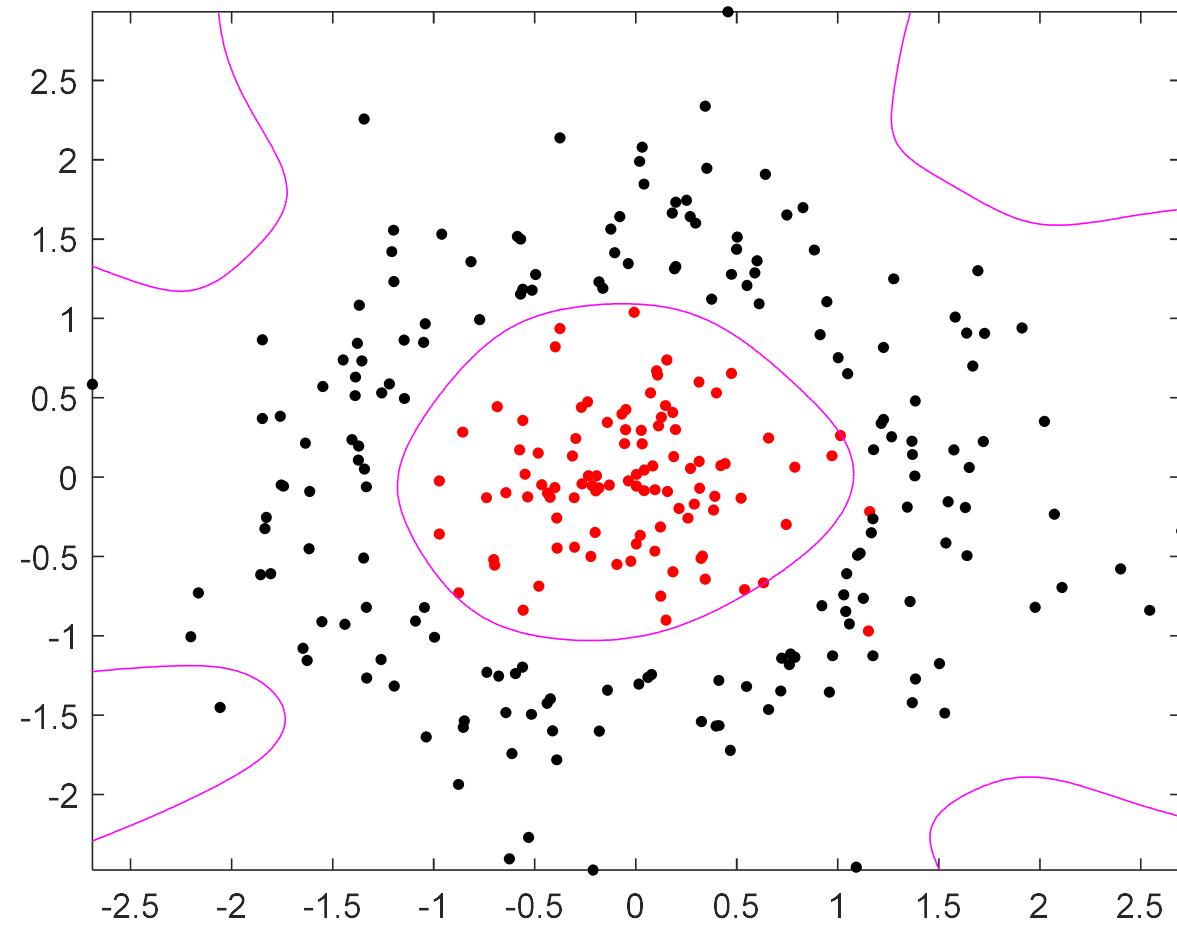


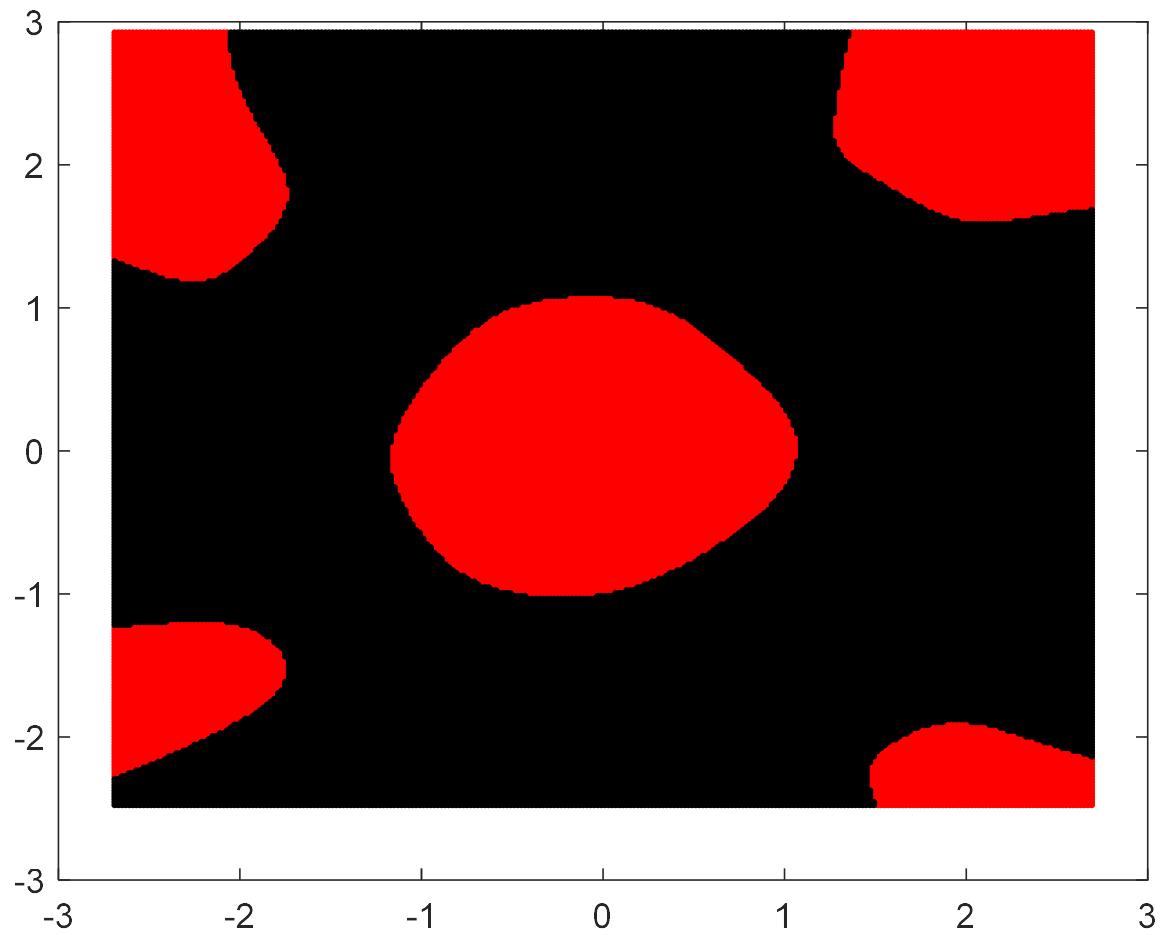
# 1<sup>st</sup> run: one hidden layer with 20 neurons



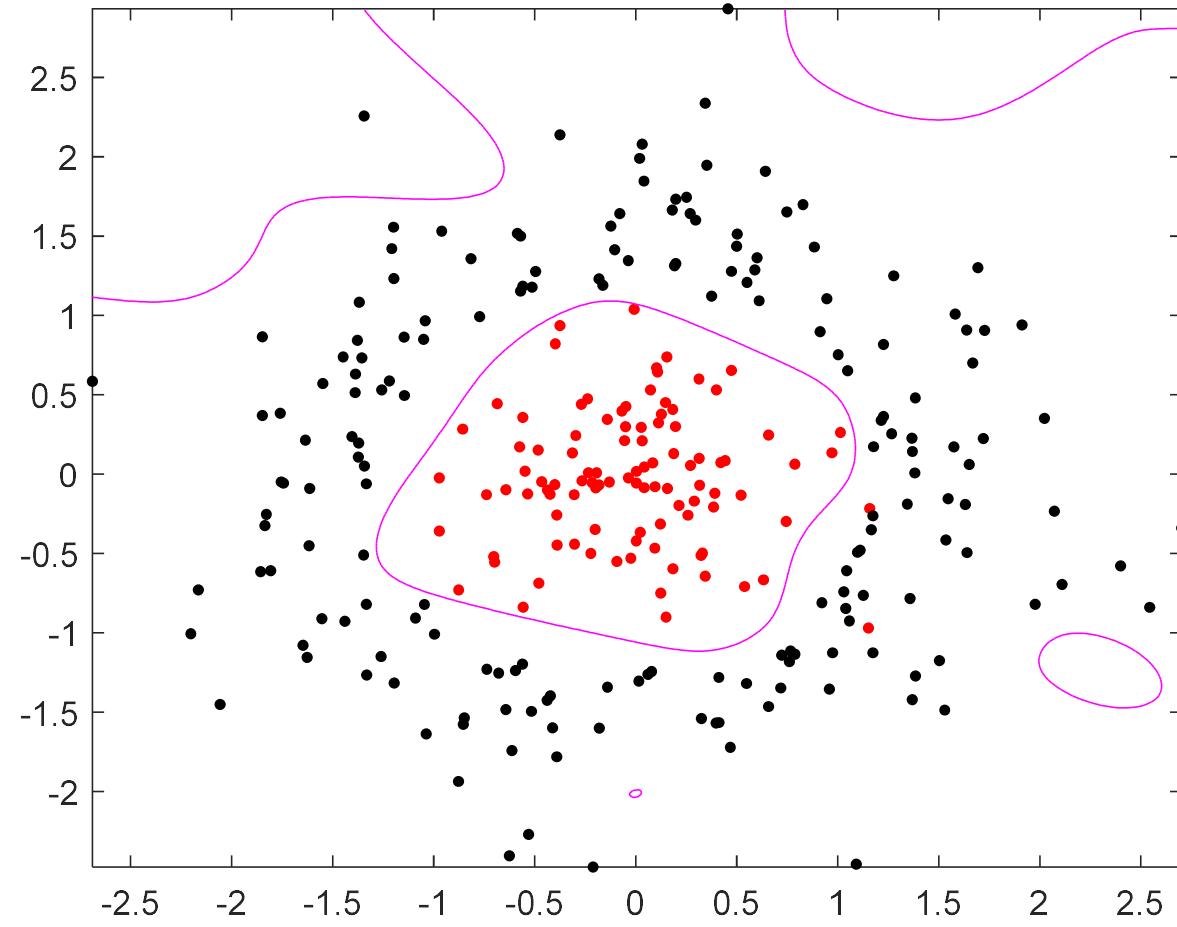


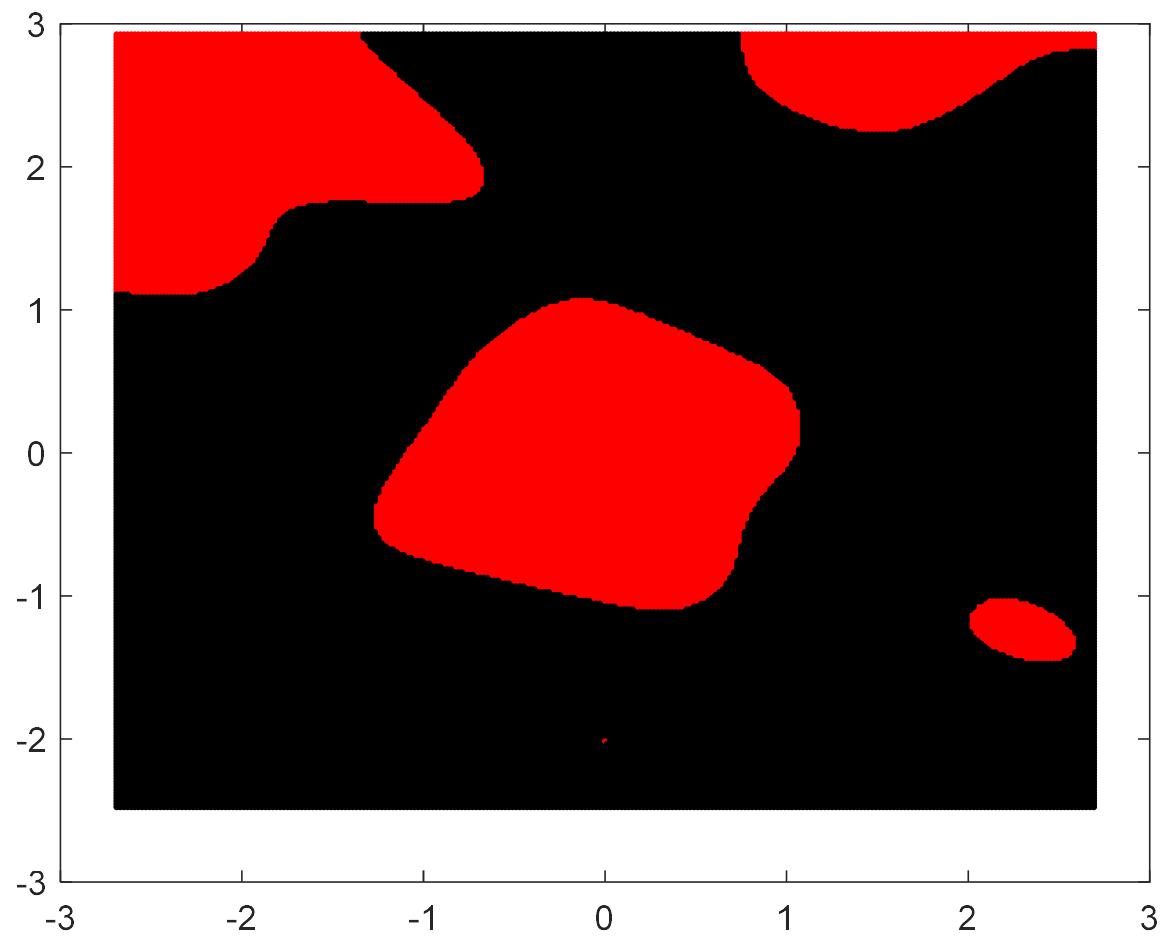
2<sup>nd</sup> run: one hidden layer with 20 neurons



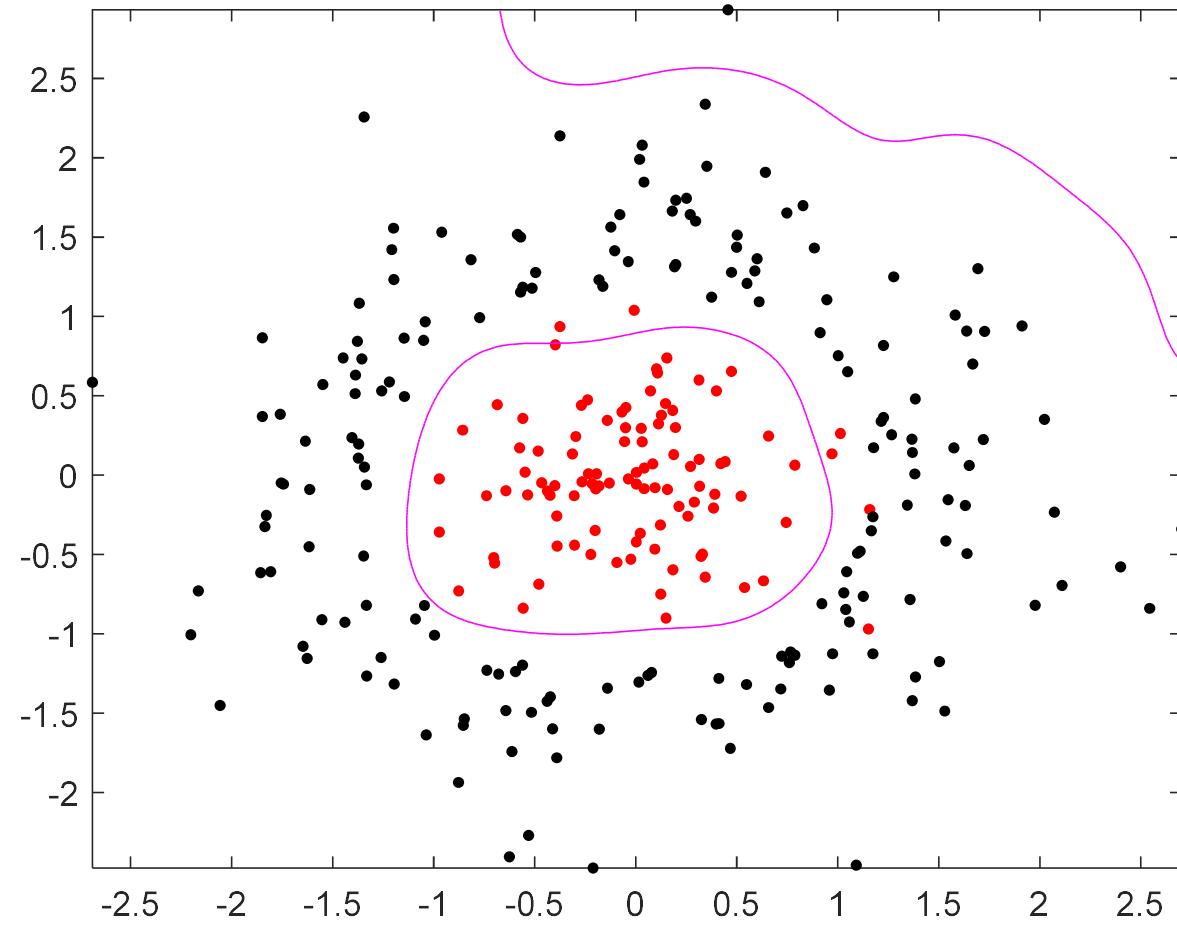


### 3<sup>rd</sup> run: one hidden layer with 20 neurons



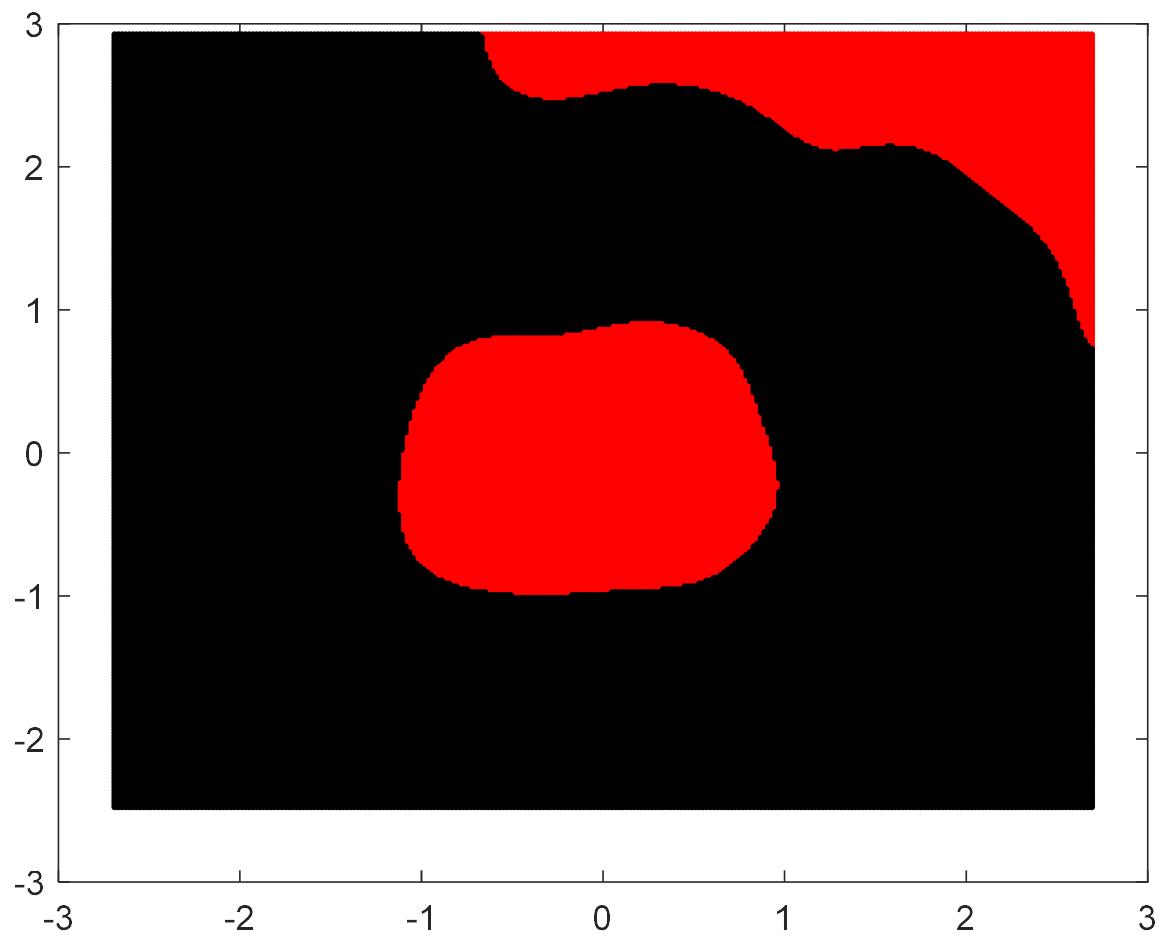


## 4<sup>th</sup> run: one hidden layer with 20 neurons



## Discussion 1:

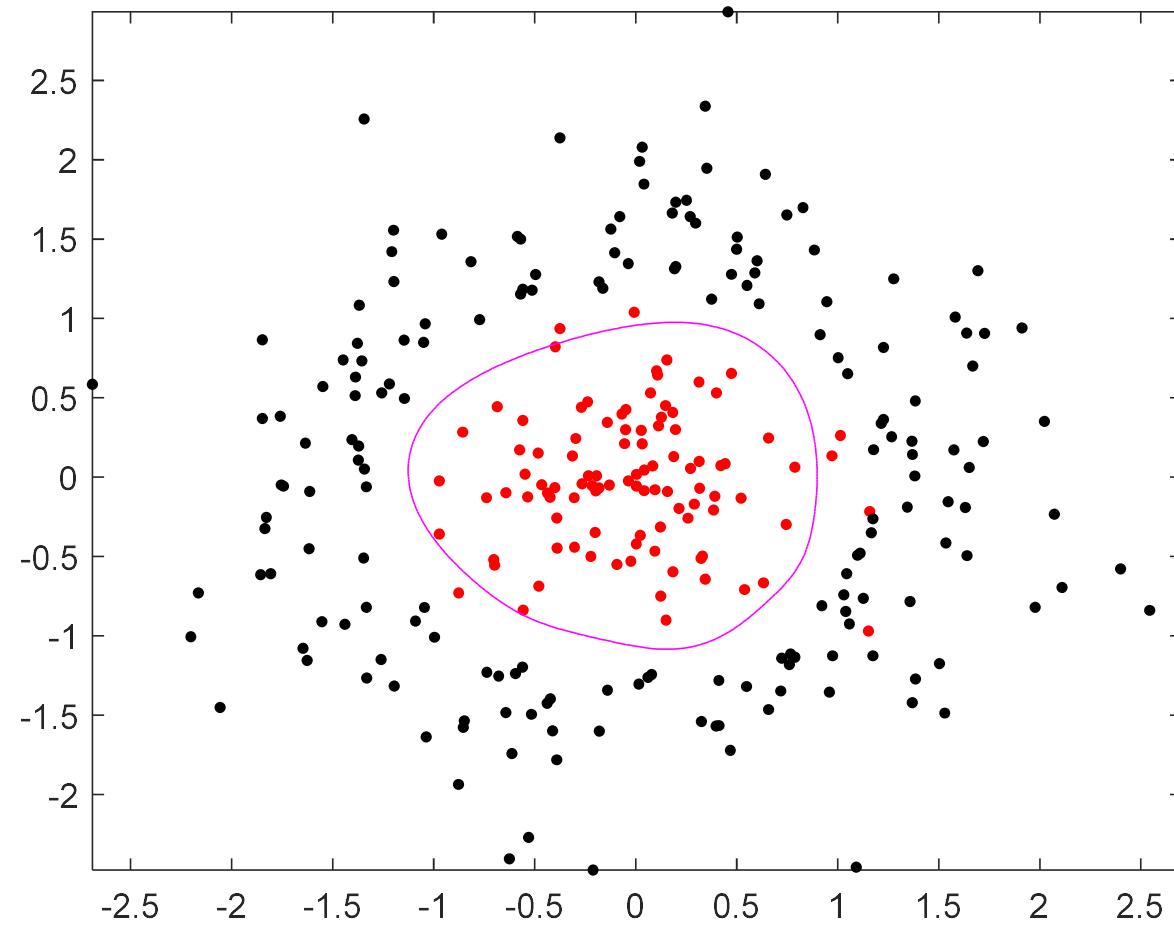
Why 4 different runs give different results?



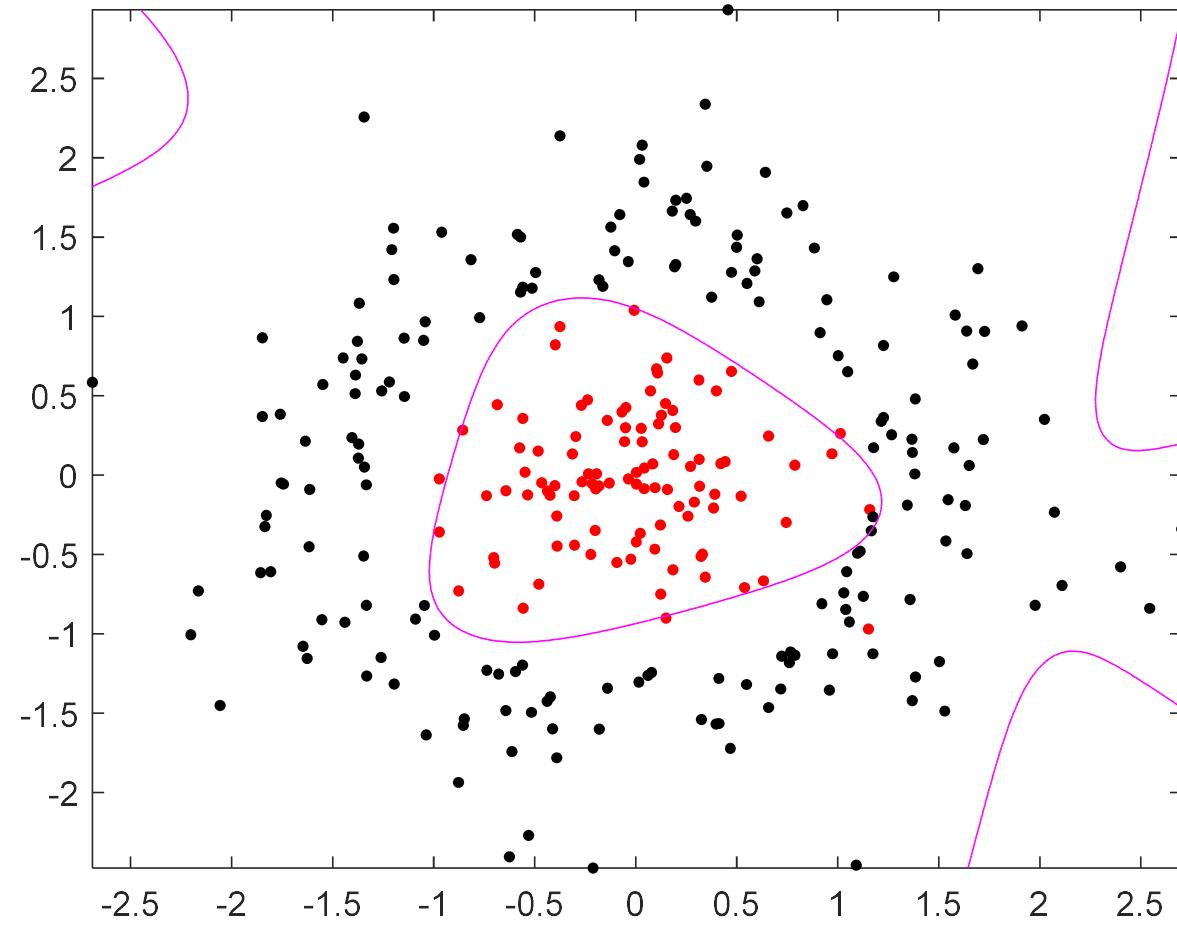
*Observation 2: The more hidden layer neurons are used, the more complex the decision surface could be.*

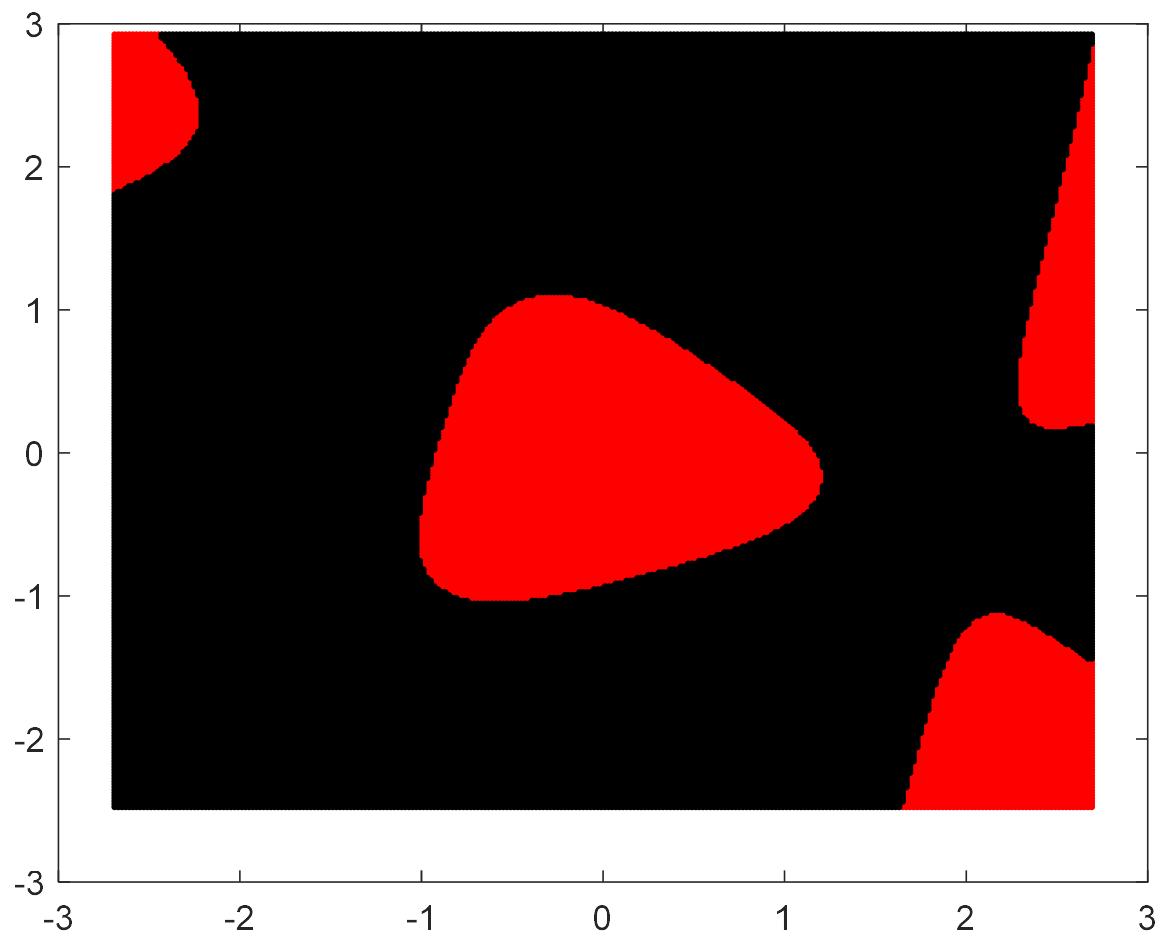
Next, we set the number of neurons to 5 and 200 respectively, and repeat the experiments.

1<sup>st</sup> run: one hidden layer with 5 neurons

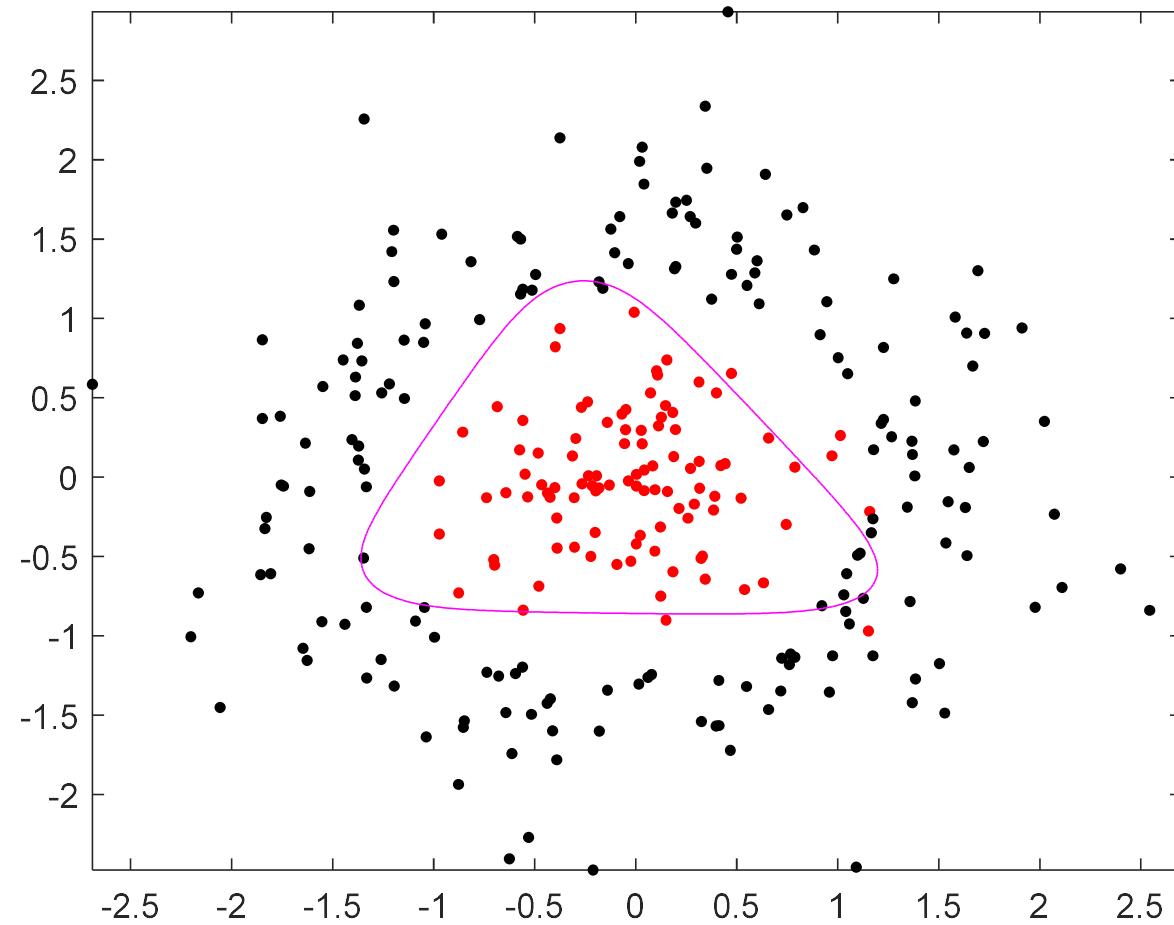


## 2<sup>nd</sup> run: one hidden layer with 5 neurons

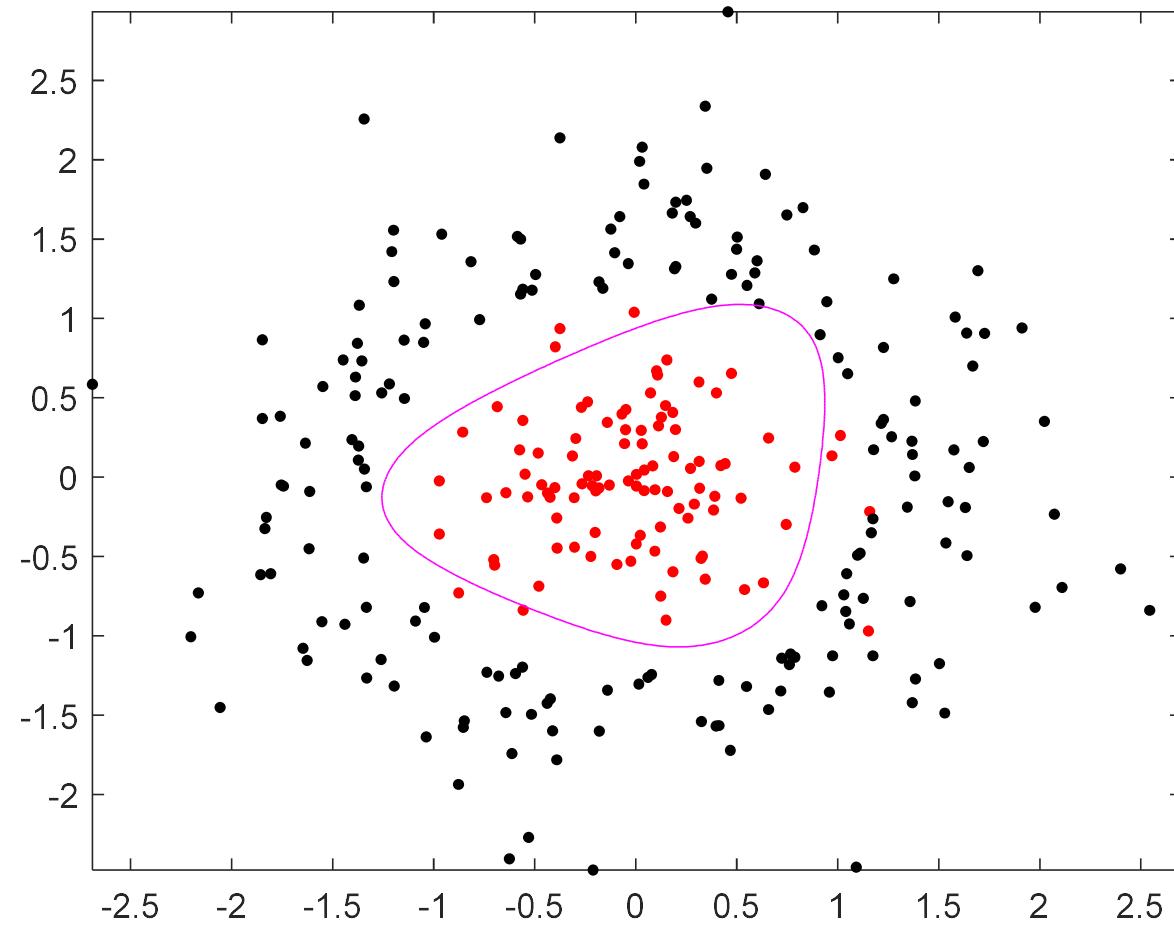




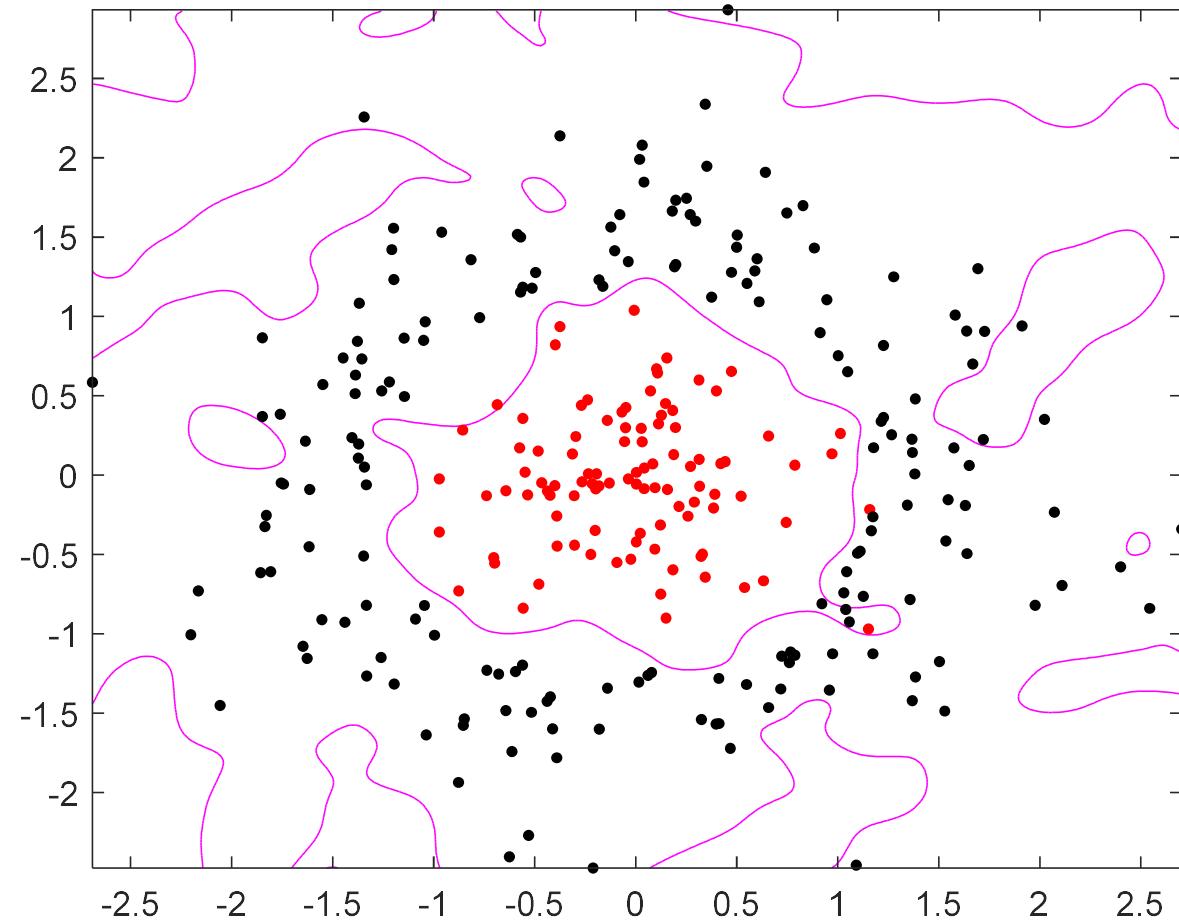
### 3<sup>st</sup> run: one hidden layer with 5 neurons

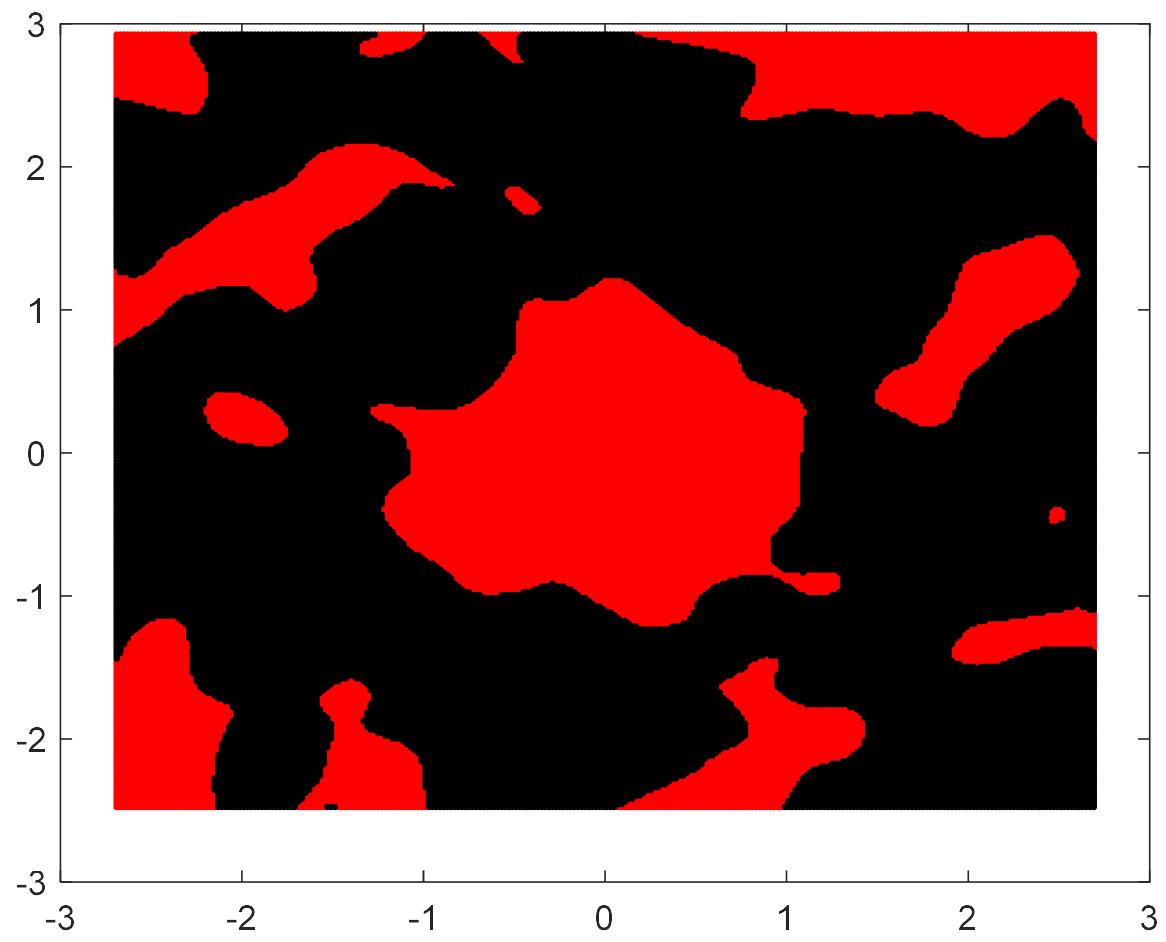


## 4<sup>st</sup> run: one hidden layer with 5 neurons

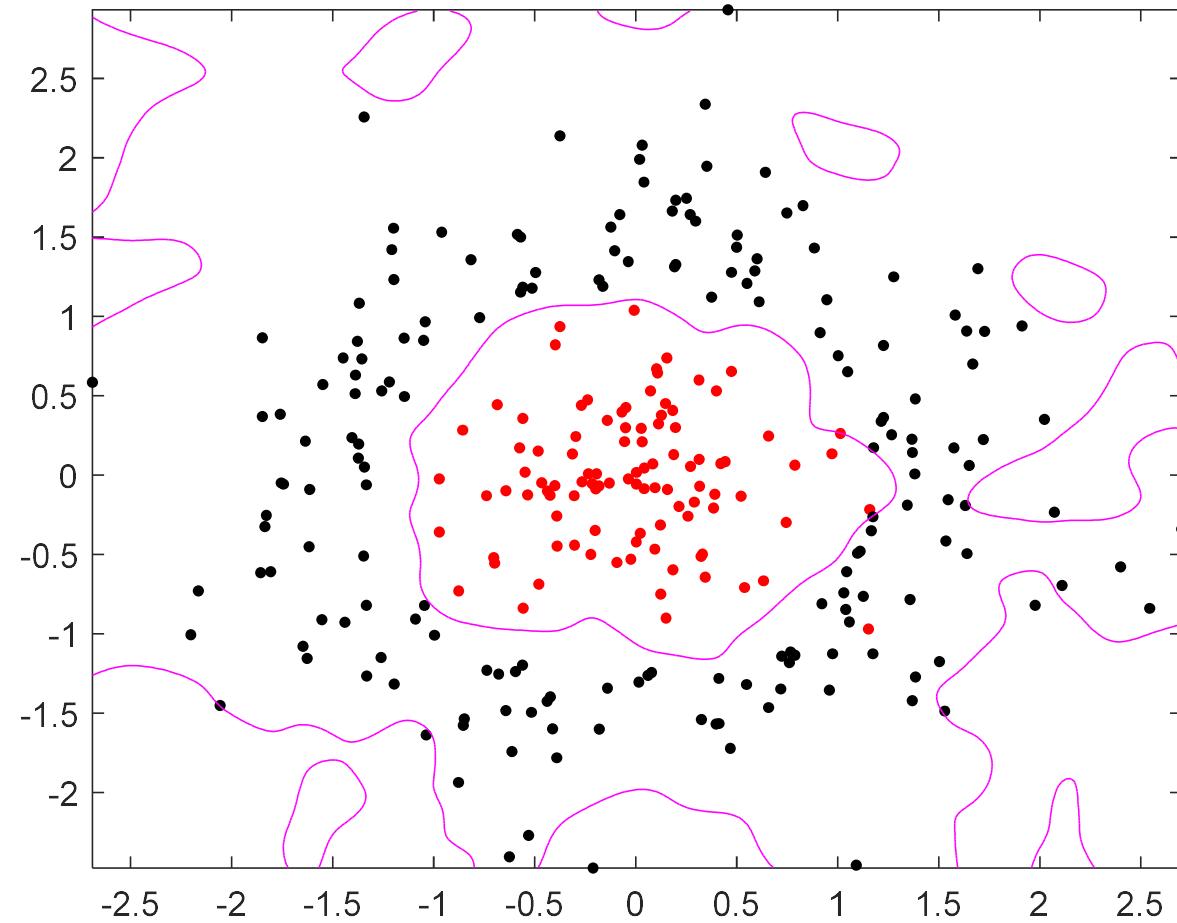


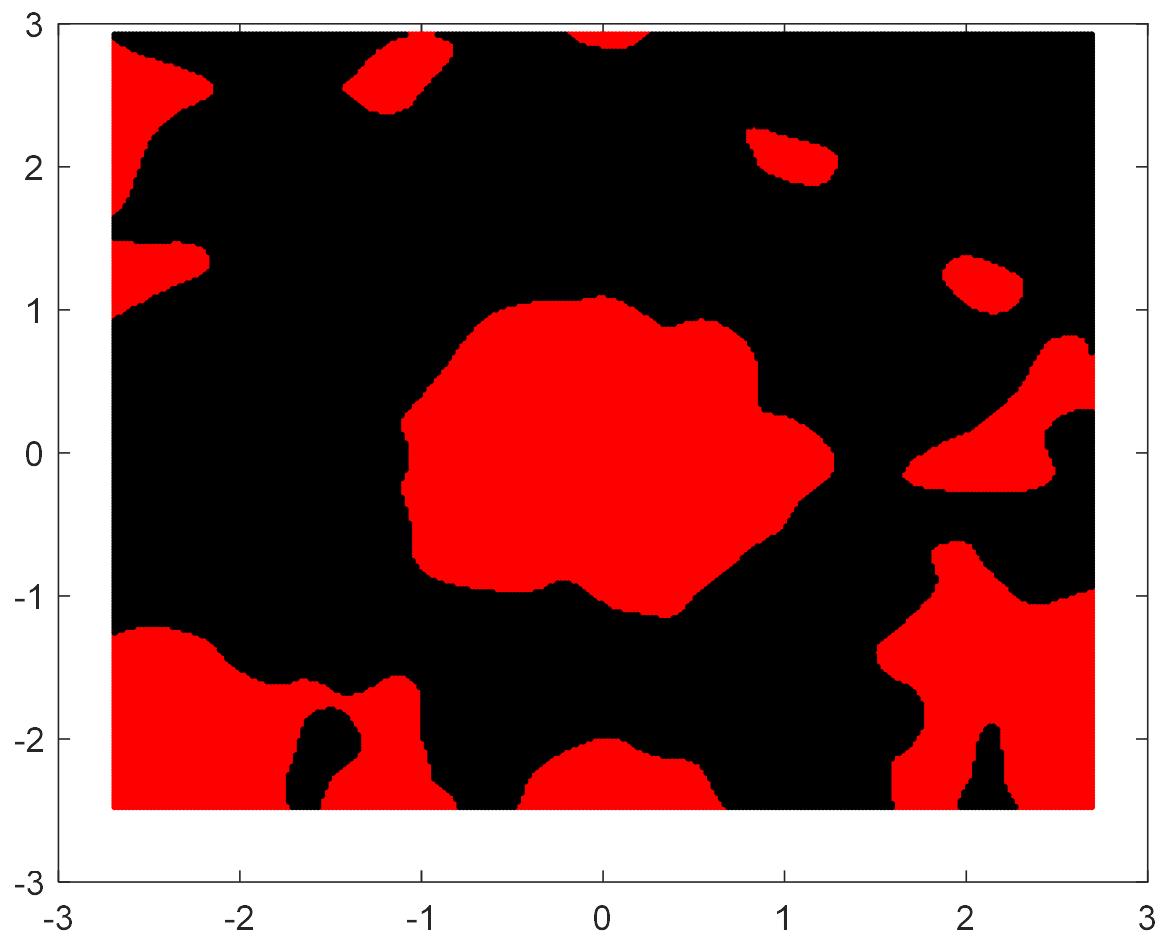
# 1<sup>st</sup> run: one hidden layer with 200 neurons





2<sup>nd</sup> run: one hidden layer with 200 neurons





## Discussion 2

What lessons do you learn from these results?

*Observation 3: The gradient vanishing problem of MLP if more than one hidden layers are used.*

We show this problem using the example of the MLP neural network for recognition of handwritten digits:

|   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 4 | 1 | 9 | 2 | 1 | 3 | 1 | 4 | 3 |
| 5 | 3 | 6 | 1 | 7 | 2 | 8 | 6 | 9 | 4 |
| 0 | 9 | 1 | 1 | 2 | 4 | 3 | 2 | 7 | 3 |
| 8 | 6 | 9 | 0 | 5 | 6 | 0 | 7 | 6 | 1 |
| 8 | 7 | 9 | 3 | 9 | 8 | 5 | 9 | 3 | 3 |
| 0 | 7 | 4 | 9 | 8 | 0 | 9 | 4 | 1 | 4 |
| 4 | 6 | 0 | 4 | 5 | 6 | 1 | 0 | 0 | 1 |
| 7 | 1 | 6 | 3 | 8 | 2 | 1 | 1 | 7 | 9 |
| 0 | 2 | 6 | 7 | 8 | 3 | 9 | 0 | 4 | 6 |
| 7 | 4 | 6 | 8 | 0 | 7 | 8 | 3 | 1 | 5 |

Each digit is an image. If the image is with 28\*28 pixels, then each image can be represented by a 784-dimensional vector.

Next, we design the MLP neural network architecture, including

- (i) The number of neurons in the input layer,
- (ii) The number of hidden layers and number of neurons in each layer
- (iii) The number of neurons in the output layer

The *number of neurons in the input layer* should be the same as the dimension of input vector. Therefore 784 neurons are used in the input layer.

The *number of neurons in the output layer* is usually set to the number of classes. There are 10 classes corresponding to the 10 digits. Therefore 10 neurons are used in the output layer.

Assuming that 30 neurons are used in the hidden layer. Next, we will train three MLP neural networks with one, two, three and four hidden layers using MNIST data.

MNIST data comes in two parts. The first part contains 50,000 training images. These images are scanned handwriting samples from 250 people including US Census Bureau employees and high school students. The second part of the MNIST data set is 10,000 testing images from a different set of 250 people.

Each training image has an associated label (e.g. 6 for image digit 6). This label should be coded into 10-dimensional vector. For example, label 1 and label 6 can be coded as:

$$D = [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1]^T$$

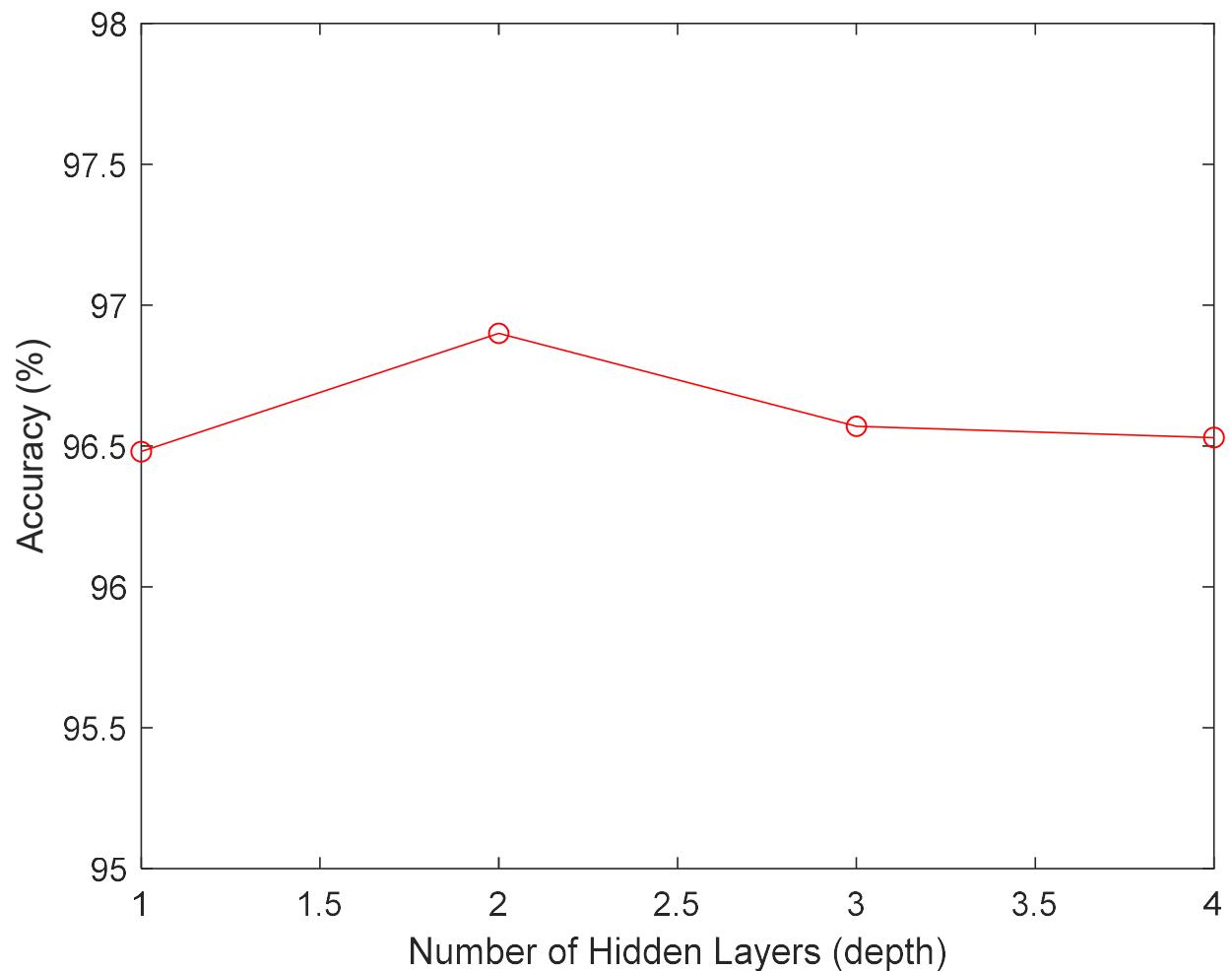
$$D = [0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0 \ 0 \ 0 \ 0]^T$$

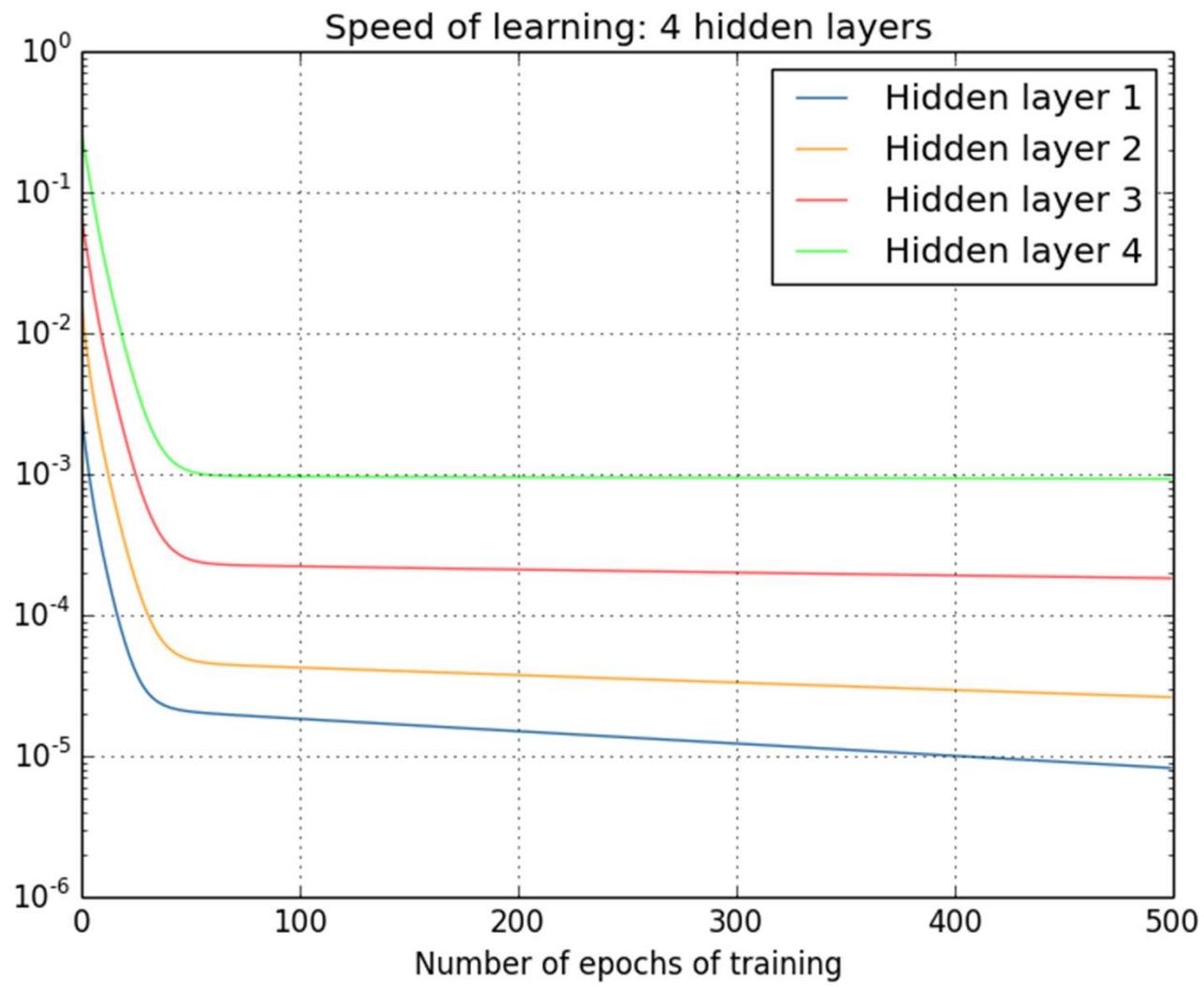
When one hidden layer is used, the performance on the test data is 96.4%.

Now, let's add another hidden layer, also with 30 neurons, and retrain the neural network with a structure of [784 30 30 10], which gives an improved classification accuracy, 96.90%. This result is encouraging: a little more depth is helping.

Let's add another 30-neuron hidden layer. Now the network has 3 hidden layers. But the result drops back down to 96.57 percent, close to our original shallow network.

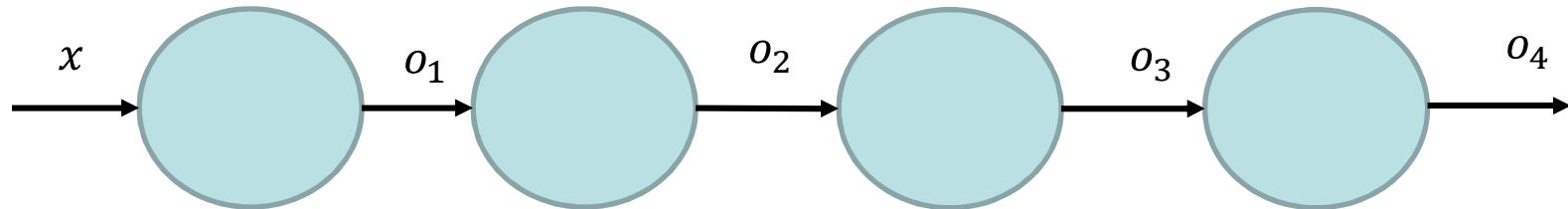
Let's add another 30-neuron hidden layer. Now the network has 4 hidden layers. The result obtained is 96.53 percent. This is not a statistically significant drop, but it is not encouraging, either. The results of MLP neural networks with different number of hidden layer are summarized below:





We have here an important observation: the gradient tends to get smaller as we move backward through the hidden layers. This means that neurons in the earlier layers learn much more slowly than neurons in later layers. And while we've seen this in just a single example, there are fundamental reasons why this happens in MLP neural networks. The phenomenon is known as the *vanishing gradient problem*.

To get insight into why the vanishing gradient problem occurs, let's consider the simple neural network with just a single neuron in each layer. Here's a network with three hidden layers:



where

$$o_j = \varphi(v_j) = \frac{1}{1 + \exp(-v_j)}$$

$$v_j = w_j o_{j-1} + b_j$$

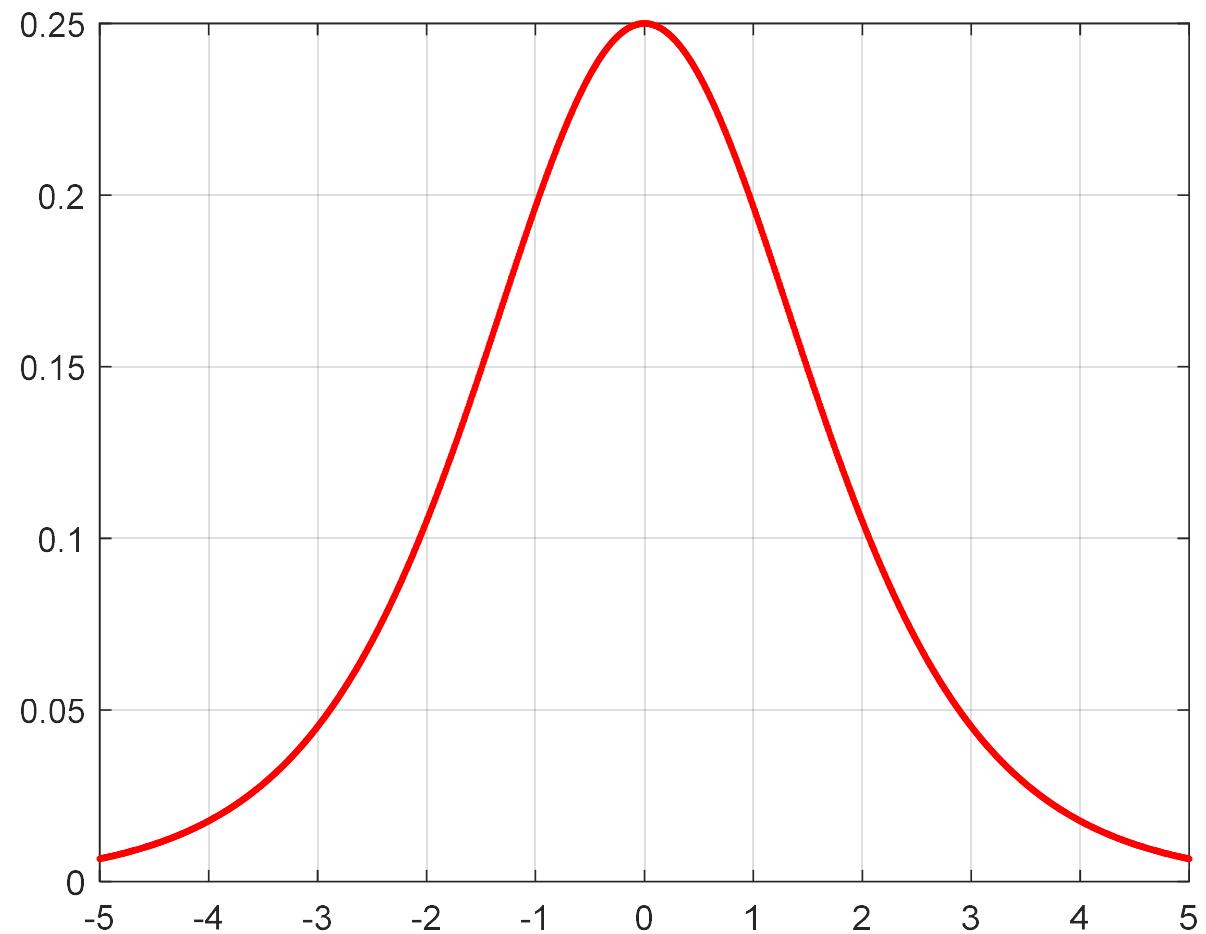
Assuming that error function is given by:

$$J = \frac{1}{2}[o_4 - d]^2$$

The gradient of  $J$  with respect to weight in the first hidden layer is given by:

$$\begin{aligned}\frac{\partial J}{\partial w_1} &= \frac{\partial J}{\partial o_4} \times \frac{\partial o_4}{\partial o_3} \times \frac{\partial o_3}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial w_1} \\ &= \frac{\partial J}{\partial o_4} \times \frac{\partial o_4}{\partial v_4} \times \frac{\partial v_4}{\partial o_3} \times \frac{\partial o_3}{\partial v_3} \times \frac{\partial v_3}{\partial o_2} \times \frac{\partial o_2}{\partial v_2} \times \frac{\partial v_2}{\partial o_1} \times \frac{\partial o_1}{\partial v_1} \times \frac{\partial v_1}{\partial w_1} \\ &= [d - o_4] \times \varphi'(v_4) \times w_4 \times \varphi'(v_3) \times w_3 \times \varphi'(v_2) \times w_2 \times \varphi'(v_1) \times x\end{aligned}$$

$$\varphi(v_j) = \frac{\exp(-v_j)}{[1 + \exp(-v_j)]^2}$$



Obviously,

$$\varphi'(v_j) \leq \frac{1}{4}$$

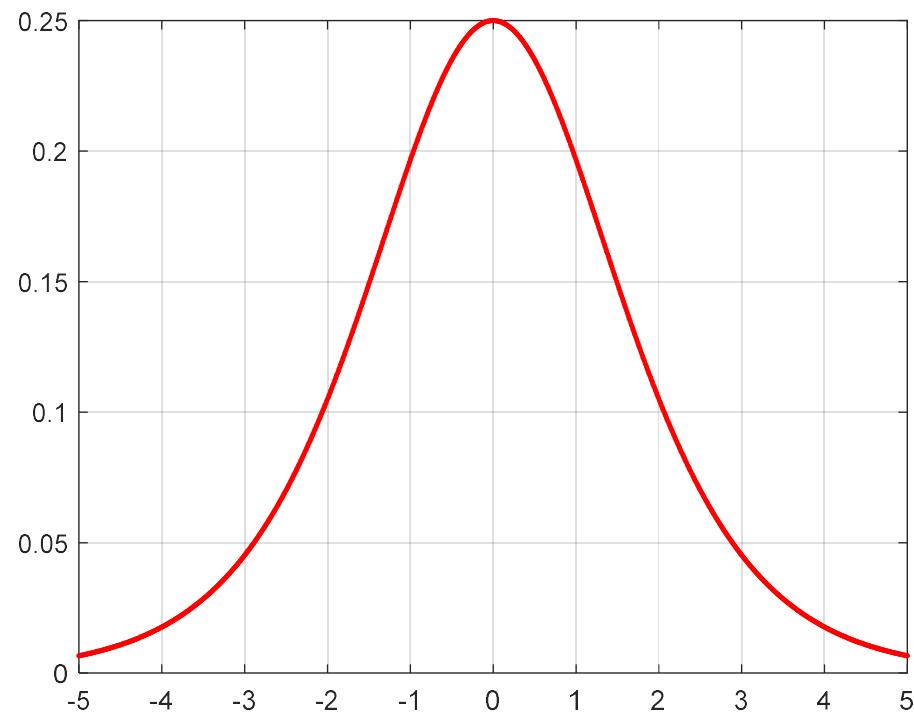
In addition, if  $w_j$  is from a normal distribution with zero mean, and unit standard deviation, most likely we have:

$$|w_j| < 1$$

Hence: 
$$\left| \frac{\partial J}{\partial w_1} \right| < |d - o_4| \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} \times \frac{1}{4} |x|$$

As can be seen, the gradient of the cost function  $J$  to  $w_{j-1}$  will typically be one quarter of the gradient to  $w_j$ , which means the learning speed at the earlier layers will be slower than later layers. This is the essential origin of the vanishing gradient problem.

Of course, this is an informal argument, not a rigorous proof. There are several possible escape clauses. In particular, we might wonder whether the weights will grow during training. Indeed,  $|w_j|$  can be greater than 1, but this will increase the activation signal  $v_j$ , and eventually  $|\varphi'(v_j)w_j|$  may not be necessarily increased as shown below.



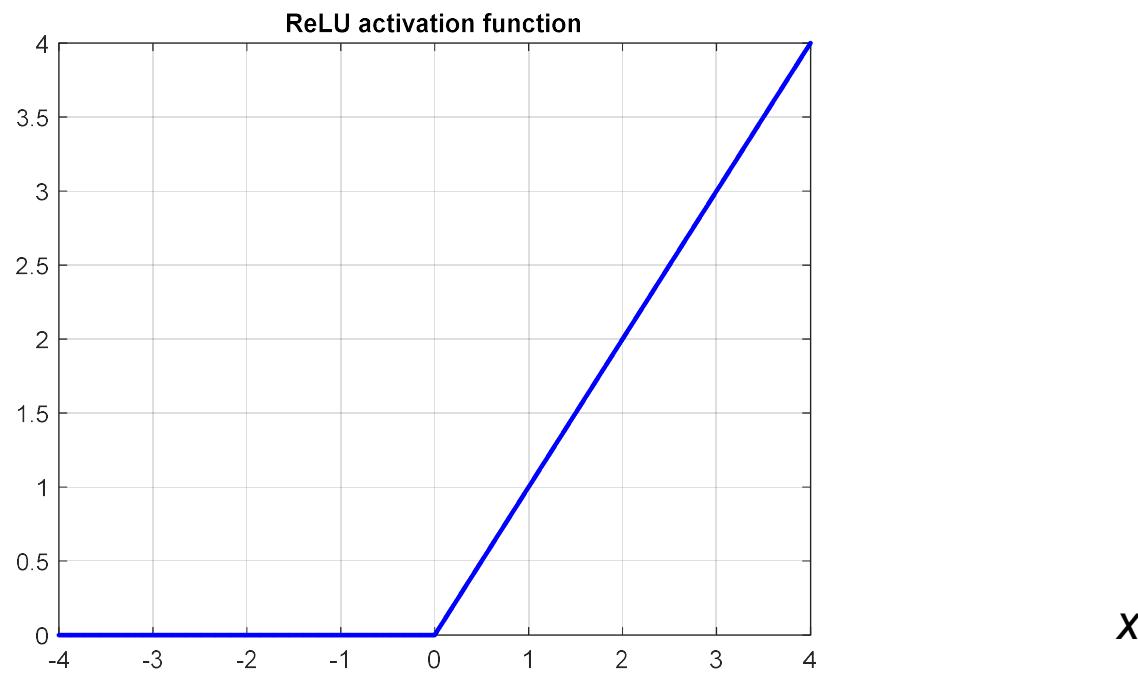
## Discussion 3

How to solve the gradient vanishing problem?

To solve the gradient vanishing problem, ReLU activation function and its variants can be used.

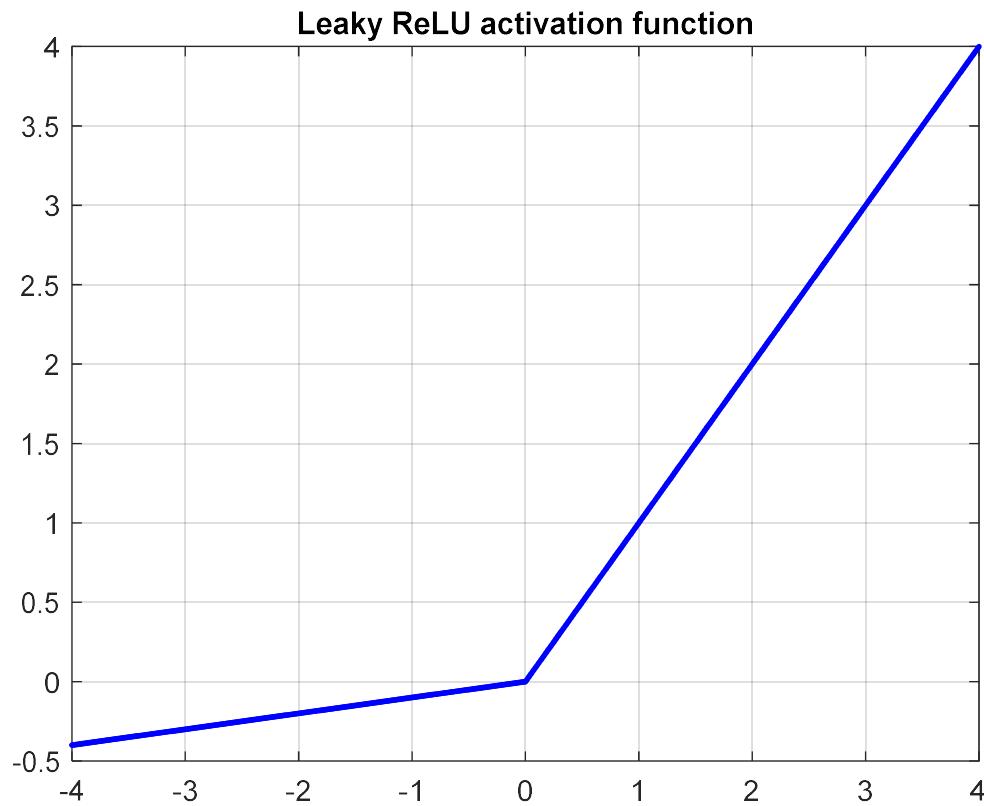
*Rectified Linear Unit (ReLU) activation function*

$$f(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$



## *Leaky/Parametric ReLU*

$$f(x) = \begin{cases} x & x \geq 0 \\ ax & x < 0 \end{cases}$$



## *Exponential Linear Unit (ELU) activation function*

$$f(x) = \begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

