# CHALMERS, GÖTEBORGS UNIVERSITET

### EXAM for ARTIFICIAL NEURAL NETWORKS

### COURSE CODES: FFR 135, FIM 720 GU, PhD

Time:	October 26, 2020, at $8^{30} - 12^{30}$
Place:	Zoom
Teachers:	Bernhard Mehlig, 073-420 0988 (mobile)
	Johan Fries, 070-370 1272 (mobile)
Allowed material:	Mathematics Handbook for Science and Engineering
Not allowed:	Any other written material, calculator

Maximum score on this exam: 12 points.

Maximum score for homework problems: 12 points.

To pass the course it is necessary to score at least 5 points on this written exam.

 $\geq$  14 grade 0,  $\geq$  20 grade V0

1. Feature map. The two patterns  $\boldsymbol{x}^{(1)}$  and  $\boldsymbol{x}^{(2)}$  shown in Figure 1(a) are processed by a very simple convolutional network that has one convolution layer with one single 4×4 kernel with ReLU units, zero threshold, weights  $w_{ij}$  as given in Figure 1(b), and stride (1,1). The resulting feature map is fed into a 2×2 max-pooling layer with stride (1,1). Finally there is a fully connected output layer with one output unit  $O^{(\mu)}$  with Heaviside activation function. For both patterns determine the resulting feature map and the output of the max-pooling layer. Determine weights  $W_k$  and a threshold  $\Theta$ so that the network output is  $O^{(1)} = 0$  for input pattern  $\boldsymbol{x}^{(1)}$ , and  $O^{(2)} = 1$ for input pattern  $\boldsymbol{x}^{(2)}$ .



Figure 1: (a) Input patterns  $\boldsymbol{x}^{(1)}$  and  $\boldsymbol{x}^{(2)}$  with 0/1 bits ( $\Box$  corresponds to  $x_i=0$  and  $\blacksquare$  to  $x_i=1$ ). (b) Weights  $w_{ij}$  of a 4×4 kernel of a feature map. The weights are either 0 or 1 ( $\Box$  corresponds to  $w_{ij} = 0$  and  $\blacksquare$  to  $w_{ij} = 1$ ). (Question 1).

2. Hopfield network with hidden units. A Hopfield network with hidden neurons can be used to learn a distribution of input patterns. Consider a Hopfield network with N visible neurons  $v_j$  and M hidden neurons  $h_i$ . The neurons are binary, with values -1 or +1. The network learns by updating the visible neurons according to

$$v_j \leftarrow \operatorname{sgn}\left[b_j^{(v)}\right] \quad \text{with} \quad b_j^{(v)} = \sum_{i=1}^M h_i w_{ij},$$
(1)

and by updating the hidden neurons according to

$$h_i \leftarrow \operatorname{sgn}\left[b_i^{(h)}\right] \quad \text{with} \quad b_i^{(h)} = \sum_{j=1}^N w_{ij} v_j.$$
 (2)

In Equations (1) and (2),  $w_{ij}$  are the elements of a  $M \times N$  weight matrix. Furthermore,  $\operatorname{sgn}[b]$  is the signum function,  $\operatorname{sgn}[b] = -1$  if b < 0 and +1 otherwise. Show that the energy function

$$H = -\sum_{i=1}^{M} \sum_{j=1}^{N} w_{ij} h_i v_j$$
(3)

can not increase upon updating one of the hidden neurons according to eq. (2).

## 3. Backpropagation

Assuming the energy function

$$H = \frac{1}{2} \sum_{i,\mu} (y_i^{(\mu)} - O_i^{(\mu)})^2, \qquad (4)$$

derive the update rule for the weights  $w_{ij}^{(\ell)}$  for  $\ell = 1, 2$  and 3 for the network shown in Figure 2.



Figure 2: Network for Question 3.

4. XNOR function. The Boolean XNOR function takes two binary inputs. For the inputs [-1, -1] and [1, 1] the function evaluates to +1, for the other two to -1. Encode the XNOR function as weights  $w_{ij}$  in a Hopfield net with three neurons by storing the patterns  $\boldsymbol{x}^{(1)} = [-1, -1, 1], \, \boldsymbol{x}^{(2)} = [1, 1, 1],$  $\boldsymbol{x}^{(3)} = [-1, 1, -1], \text{ and } \boldsymbol{x}^{(4)} = [1, -1, -1]$  using Hebb's rule:

$$w_{ij} = \frac{1}{3} \sum_{\mu=1}^{4} x_i^{(\mu)} x_j^{(\mu)}$$
 where  $i, j = 1, \dots, 3.$  (5)

The update rule for bit  $S_i$  is

$$S_i \leftarrow \operatorname{sgn}\left[\sum_{j=1}^3 w_{ij}S_j\right],$$
 (6)

where  $\operatorname{sgn}[b]$  is the signum function,  $\operatorname{sgn}[b] = -1$  if b < 0 and +1 otherwise.

- (a) What is the weight matrix that you obtain? Feed the stored patterns to the net, and test whether they are stable under synchronous updating.
- (b) Use the weight matrix to compute the energy function,

$$H = -\frac{1}{2} \sum_{ij} w_{ij} s_i s_j. \tag{7}$$

Use the fact that the elements  $s_i$  only take values  $\pm 1$ .

- (c) Based on your answers to the previous parts, conclude with one or two sentences whether the network is useful for recognising the XNOR function.
- (d) What would be the difference if one tried to store just patterns 1, 2 and 3, and not all 4 patterns?

5. Gradient descent and momentum. Consider the given energy function H as a function of a single weight w as shown in Figure 3. Use the following gradient-descent update rule:

$$\delta w_{n+1} = -\eta \frac{\partial H}{\partial w} + \alpha \, \delta w_n. \tag{8}$$

Here  $\eta$  is the learning rate, and  $\alpha$  is the momentum parameter. The weight at time step n+1 is then given by  $w_{n+1} = w_n + \delta w_n$ . Assume that the system is initially at point A. The slope of the segment AB in Figure 3 is -s and the slope of the segment BC is 0. The slope at point A is defined to be -sand that at point B to be 0. The system starts at time step 1, and assume that  $\delta w_0 = 0$ . Assume that  $\eta s = 1/2$ .

- (a) At which time step n does the system reach point B for  $\alpha = 0$ ?
- (b) Repeat the previous calculation for the case  $\alpha = 1/2$ . You should find that the final equation you obtain for the number of time steps n involves a linear term in n, and an exponential term in n. Plot the linear and exponential functions schematically with n on the x-axis. In this plot, mark the value of n where the two functions intersect, thus obtaining the value of n at which the system reaches point B.
- (c) Which of the two cases:  $\alpha = 0$  and  $\alpha = 1/2$  reaches point B faster? Use the results of the previous two parts to justify your answer.
- (d) What is the fate of the two systems  $\alpha = 0$  and  $\alpha = 1/2$  once they cross point B?



Figure 3: Energy as a function of weight for Question 5.

6. Linear activation function Consider using a linear activation function g(b) = b in a fully connected simple perceptron with one output unit. Fed with a training pattern  $\mathbf{x}^{(\mu)}$ , the output  $O^{(\mu)}$  is given by

$$O^{(\mu)} = \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}^{(\mu)} - \boldsymbol{\theta}.$$
<sup>(9)</sup>

Here  $\boldsymbol{w}$  is a column vector of weights, and  $\boldsymbol{\theta}$  is a scalar threshold. There are p training patterns,  $\mu = 1, \ldots, p$ . Their target outputs are denoted by  $t^{(\mu)}$ . For the perceptron considered, the energy function

$$H = \frac{1}{2} \sum_{\mu=1}^{p} \left( O^{(\mu)} - t^{(\mu)} \right)^2 \tag{10}$$

has only one minimum, and it can be found analytically. In the following, you will derive the threshold  $\theta$  at the minimum.

a) Start by showing that the minimum implies

$$\mathbb{G}\boldsymbol{w} = \boldsymbol{\alpha} + \theta\boldsymbol{\beta} \tag{11a}$$

$$\boldsymbol{\beta}^{\mathsf{T}} \boldsymbol{w} = \boldsymbol{\theta} + \boldsymbol{\gamma} \tag{11b}$$

with

$$\mathbb{G} = \langle \boldsymbol{x}\boldsymbol{x}^{\mathsf{T}} \rangle, \quad \boldsymbol{\alpha} = \langle t\boldsymbol{x} \rangle, \quad \boldsymbol{\beta} = \langle \boldsymbol{x} \rangle \quad \text{and} \quad \boldsymbol{\gamma} = \langle t \rangle, \quad (12)$$

where  $\langle \dots \rangle$  denotes an average over the training patterns.

b) Assume that  $\mathbb{G}$  is invertible, with inverse  $\mathbb{G}^{-1}$ . Furthermore, assume that  $\boldsymbol{\beta}^{\mathsf{T}}\mathbb{G}^{-1}\boldsymbol{\beta} \neq 1$  and solve eqs. (11) for  $\theta$ .

c) If, in a fully connected multi-layer perceptron, one uses a linear activation function g(b) = b, it holds that

$$\boldsymbol{V}^{(\mu,\ell)} = \boldsymbol{w}^{(\ell)} \boldsymbol{V}^{(\mu,\ell-1)} - \boldsymbol{\theta}^{(\ell)}$$
  
=  $\left[ \boldsymbol{w}^{(\ell)} \boldsymbol{w}^{(\ell-1)} \right] \boldsymbol{V}^{(\mu,l-2)} - \left[ \boldsymbol{w}^{(l)} \boldsymbol{\theta}^{(\ell-1)} + \boldsymbol{\theta}^{(\ell)} \right].$  (13)

Here,  $\boldsymbol{V}^{(\mu,\ell)}$  is the  $\mu^{\text{th}}$  neuron in the  $\ell^{\text{th}}$  hidden layer. Furthermore,  $\boldsymbol{w}^{(\ell)}$  and  $\boldsymbol{\theta}^{(\ell)}$  are the weight matrix and the shold vector for the neurons in the  $\ell^{\text{th}}$  hidden layer. Write at most three sentences where you, based on eq. (13), argue that a non-linear activation function is essential for a multi-layer perceptron.

#### CHALMERS, GÖTEBORGS UNIVERSITET

#### SOLUTIONS FOR EXAM for ARTIFICIAL NEURAL NETWORKS

#### COURSE CODES: FFR 135, FIM 720 GU, PhD

1. Feature map. The two patterns  $\boldsymbol{x}^{(1)}$  and  $\boldsymbol{x}^{(2)}$  shown in Figure 1(a) are processed by a very simple convolutional network that has one convolution layer with one single 4×4 kernel with ReLU units, zero threshold, weights  $w_{ij}$  as given in Figure 1(b), and stride (1,1). The resulting feature map is fed into a 2×2 max-pooling layer with stride (1,1). Finally there is a fully connected output layer with one output unit  $O^{(\mu)}$  with Heaviside activation function. For both patterns determine the resulting feature map and the output of the max-pooling layer. Determine weights  $W_k$  and a threshold  $\Theta$ so that the network output is  $O^{(1)} = 0$  for input pattern  $\boldsymbol{x}^{(1)}$ , and  $O^{(2)} = 1$ for input pattern  $\boldsymbol{x}^{(2)}$ .



Figure 1: (a) Input patterns  $\boldsymbol{x}^{(1)}$  and  $\boldsymbol{x}^{(2)}$  with 0/1 bits ( $\Box$  corresponds to  $x_i=0$  and  $\blacksquare$  to  $x_i=1$ ). (b) Weights  $w_{ij}$  of a 4×4 kernel of a feature map. The weights are either 0 or 1 ( $\Box$  corresponds to  $w_{ij} = 0$  and  $\blacksquare$  to  $w_{ij} = 1$ ). (Question 4).

## Solution to "1. Feature map"

Input to feature map of pattern  $\boldsymbol{x}^{(1)}$ :

$$\begin{bmatrix} 8 & 6 \\ -2 & -6 \\ -6 & -2 \\ 6 & 8 \end{bmatrix}.$$
 (1)

Input to feature map of pattern  $x^{(2)}$ :

$$\begin{bmatrix} -2 & -2\\ 0 & 0\\ 0 & 0\\ 0 & 0 \end{bmatrix}.$$
 (2)

Feature map of pattern  $\boldsymbol{x}^{(1)}$ :

$$\begin{bmatrix} 8 & 6 \\ 0 & 0 \\ 0 & 0 \\ 6 & 8 \end{bmatrix} .$$
(3)

Feature map of pattern  $x^{(2)}$ :

$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}.$$
 (4)

Max-pooling layer of pattern  $\boldsymbol{x}^{(1)}$ :

$$\begin{bmatrix} 8\\0\\8 \end{bmatrix}.$$
 (5)

Max-pooling layer of pattern  $\boldsymbol{x}^{(2)}$ :

$$\begin{bmatrix} 0\\0\\0 \end{bmatrix}.$$
 (6)

With  $W_k = -\delta_{k1}$  and  $\theta = -4$  we have

$$\sum_{k=1}^{3} W_k \begin{bmatrix} 8\\0\\8 \end{bmatrix}_k - T = -4 \tag{7}$$

and

$$\sum_{k=1}^{3} W_k \begin{bmatrix} 0\\0\\0 \end{bmatrix}_k - T = 4.$$
 (8)

Applying the Heaviside activation function results in the requested outputs.

2. Hopfield network with hidden units A Hopfield network with hidden neurons can be used to learn a distribution of input patterns. Consider a Hopfield network with N visible neurons  $v_j$  and M hidden neurons  $h_i$ . The neurons are binary, with values -1 or +1. The network learns by updating the visible neurons according to

$$v_j \leftarrow \operatorname{sgn}\left[b_j^{(v)}\right] \quad \text{with} \quad b_j^{(v)} = \sum_{i=1}^M h_i w_{ij},$$
(9)

and by updating the hidden neurons according to

$$h_i \leftarrow \operatorname{sgn}\left[b_i^{(h)}\right] \quad \text{with} \quad b_i^{(h)} = \sum_{j=1}^N w_{ij} v_j.$$
 (10)

In Equations (9) and (10),  $w_{ij}$  are the elements of a  $M \times N$  weight matrix. Furthermore,  $\operatorname{sgn}[b]$  is the signum function,  $\operatorname{sgn}[b] = -1$  if b < 0 and +1 otherwise. Show that the energy function

$$H = -\sum_{i=1}^{M} \sum_{j=1}^{N} w_{ij} h_i v_j$$
(11)

can not increase upon updating one of the hidden neurons according to eq. (10).

## Solution to "2. Hopfield network with hidden units"

Denote the the value of hidden neuron i after the update by  $h'_i$ . Suppose that the  $k^{\text{th}}$  hidden neuron changes sign. We then have:

$$h_i' = h_i - 2h_i \delta_{ik},\tag{12}$$

The energy after the update is

$$H' = -\sum_{i=1}^{M} \sum_{j=1}^{N} w_{ij} h'_i v_j$$
(13)

$$= -\sum_{j=1}^{N} v_j \sum_{i=1}^{M} w_{ij} (h_i - 2h_i \delta_{ik})$$
(14)

$$= -\sum_{j=1}^{N} v_j \sum_{i=1}^{M} w_{ij} (h_i - 2h_i \delta_{ik})$$
(15)

$$= -\sum_{j=1}^{N} v_j \left[ \sum_{i=1}^{M} w_{ij} h_i - 2 \sum_{i=1}^{M} w_{ij} h_i \delta_{ik} \right]$$
(16)

$$= -\sum_{j=1}^{N} v_j \left[ \sum_{i=1}^{M} w_{ij} h_i - 2w_{kj} h_k \right]$$
(17)

$$= -\sum_{j=1}^{N} \sum_{i=1}^{M} w_{ij} h_i v_j + 2h_k \sum_{j=1}^{N} w_{kj} v_j$$
(18)

$$=H+2h_k b_k^{(h)}.$$
(19)

If the  $k^{\text{th}}$  hidden neuron change sign, then  $h_k b_k^{(h)} < 0$ .

## 3. Backpropagation

Assuming the energy function

$$H = \frac{1}{2} \sum_{i,\mu} (y_i^{(\mu)} - O_i^{(\mu)})^2, \qquad (20)$$

derive the update rule for the weights  $w_{ij}^{(L)}$  for L = 1, 2 and 3 for the network shown in Figure 2.



Figure 2: Network for Question 4.

# Solution to "3. Backpropagation" See course book.

4. XOR function. The Boolean XOR function takes two binary inputs. For the inputs [-1, -1] and [1, 1] the function evaluates to -1, for the other two inputs it evaluates to +1. Encode the XOR function as weights  $w_{ij}$  in a Hopfield net with three neurons by storing the patterns  $\boldsymbol{x}^{(1)} = [-1, -1, -1]$ ,  $\boldsymbol{x}^{(2)} = [1, 1, -1], \, \boldsymbol{x}^{(3)} = [-1, 1, 1], \text{ and } \boldsymbol{x}^{(4)} = [1, -1, 1]$  using Hebb's rule:

$$w_{ij} = \frac{1}{3} \sum_{\mu=1}^{4} x_i^{(\mu)} x_j^{(\mu)}$$
 where  $i, j = 1, \dots, 3.$  (21)

The update rule for bit  $S_i$  is

$$S_i \leftarrow \operatorname{sgn}\left[\sum_{j=1}^3 w_{ij}S_j\right],$$
 (22)

where  $\operatorname{sgn}[b]$  is the signum function,  $\operatorname{sgn}[b] = -1$  if b < 0 and +1 otherwise. Feed the stored patterns to the net, and test whether they are stable under synchronous updating. Conclude with one or two sentences whether the network is useful for recognising the XOR function.

# Solution to "4. XOR function"

$$3\mathbb{W} = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 1 & -1 \\ 1 & 1 & -1 \\ -1 & -1 & 1 \end{bmatrix}$$
(23)

$$+ \begin{bmatrix} 1 & -1 & -1 \\ -1 & 1 & 1 \\ -1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & -1 & 1 \\ -1 & 1 & -1 \\ 1 & -1 & 1 \end{bmatrix}$$
(24)  
 
$$\begin{bmatrix} 4 & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 4 \end{bmatrix}.$$
(25)

The weight matrix is proportional to the identity matrix. All patterns are therefore stable, and the network does not recognise the XOR function.

#### 5. Gradient descent and momentum

Consider the given energy function  $\mathcal{H}$  as a function of weight w as shown in Fig. 4. Use the following gradient descent update rule,

$$\delta w_{n+1} = -\eta \frac{\partial \mathcal{H}}{\partial w} + \alpha \, \delta w_n. \tag{26}$$

Assume that the system is initially at point A, and that  $\eta s = 1/2$ . The slope of the segment AB in Fig. 4 is -s and the slope of the segment BC is 0. The system starts at time step 1, and assume that  $\delta w_0 = 0$ .

- 1. Find the number of time steps required to travel from point A to point B for  $\alpha = 0$ .
- 2. Repeat the previous calculation for the case  $\alpha = 1/2$ , and graphically find the solution of the final equation you obtain.
- 3. Indicate the results of the previous two parts on the same graph. Which of the two cases:  $\alpha = 0$  and  $\alpha = 1/2$  converges faster?
- 4. What is the fate of the two systems  $\alpha = 0$  and  $\alpha = 1/2$  once they cross point B?



Figure 3: Energy as a function of weight for problem: Gradient descent and momentum.

#### Solution to "5. Gradient descent and momentum"

1 and 2: We calculate the total change in weight at time step  $n, \Delta w_n = \sum_{i=1}^n \delta w_i$ , equate  $\Delta w_n$  to L and solve for n. Proceed by solving for  $\delta w_n$ . Iterating the equation for  $\delta w$  we find,

$$\delta w_{i+1} = \sum_{j=0}^{i} \eta s \, \alpha^j + \alpha^{i+1} \delta w_0, \qquad (27)$$

$$=\eta s \, \frac{1-\alpha^{i+1}}{1-\alpha}.\tag{28}$$

Next compute  $\Delta w_n$ ,

$$\Delta w_n = \sum_{i=1}^n \delta w_i,\tag{29}$$

$$= \eta s \sum_{i=1}^{n} \frac{1 - \alpha^{i+1}}{1 - \alpha}, \tag{30}$$

$$=\frac{\eta s}{1-\alpha}\left(n-\alpha\frac{1-\alpha^n}{1-\alpha}\right).$$
(31)

Thus using  $\eta s = 1/2$  we obtain, for  $\alpha = 0$ ,  $\Delta w_n(\alpha = 0) = n/2$ , and for  $\alpha = 1/2$ ,  $\Delta w_n(\alpha = 1/2) = n - 1 + 2^{-m}$ . Equating  $\Delta w = L$  we obtain,

$$n_{\alpha=0} = 2L,\tag{32}$$

$$n_{\alpha=1/2} - 1 + 2^{-n_{\alpha=1/2}} = L.$$
(33)

graphing the above equations, we see that  $n_{\alpha=1/2} < n_{\alpha=0}$ , thus,  $\alpha = 1/2$  converges faster.



Figure 4: Graphical solution of problem : gradient descent and momentum.

After crossing point B,  $\delta w(\alpha = 0) = 0$  so that this system stays stationary, however  $\delta w_{\alpha=1/2} > 0$  so that this system keeps on moving.

#### 6. Linear activation function

Consider using a linear activation function g(b) = b in a fully connected simple perceptron with one output unit. Fed with a training pattern  $\boldsymbol{x}^{(\mu)}$ , the output  $O^{(\mu)}$  is given by

$$O^{(\mu)} = \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}^{(\mu)} - \boldsymbol{\theta}. \tag{34}$$

Here  $\boldsymbol{w}$  is a column vector of weights, and  $\boldsymbol{\theta}$  is a scalar threshold. There are p training patterns,  $\mu = 1, \ldots, p$ . Their target outputs are denoted by  $t^{(\mu)}$ . For the perceptron concidered, the energy function

$$H = \frac{1}{2} \sum_{\mu=1}^{p} \left( O^{(\mu)} - t^{(\mu)} \right)^2 \tag{35}$$

has only one minimum, and it can be found analytically. In the following, you will derive the threshold  $\theta$  at the minimum.

a) Start by showing that the minimum implies

$$\mathbb{G}\boldsymbol{w} = \boldsymbol{\alpha} + \theta\boldsymbol{\beta} \tag{36a}$$

$$\boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{w} = \boldsymbol{\theta} + \boldsymbol{\gamma} \tag{36b}$$

with

$$\mathbb{G} = \langle \boldsymbol{x}\boldsymbol{x}^{\mathsf{T}} \rangle, \quad \boldsymbol{\alpha} = \langle t\boldsymbol{x} \rangle, \quad \boldsymbol{\beta} = \langle \boldsymbol{x} \rangle \quad \text{and} \quad \boldsymbol{\gamma} = \langle t \rangle, \quad (37)$$

where  $\langle \dots \rangle$  denotes an average over the training patterns.

b) Assume that  $\mathbb{G}$  is invertible, with inverse  $\mathbb{G}^{-1}$ . Furthermore, assume that  $\boldsymbol{\beta}^{\mathsf{T}}\mathbb{G}^{-1}\boldsymbol{\beta} \neq 1$  and solve eqs. (36) for  $\theta$ .

c) If, in a fully connected multi-layer perceptron, one uses a linear activation function g(b) = b, it holds that

$$\boldsymbol{V}^{(\mu,\ell)} = \boldsymbol{w}^{(\ell)} \boldsymbol{V}^{(\mu,\ell-1)} - \boldsymbol{\theta}^{(\ell)}$$
  
=  $\left[ \boldsymbol{w}^{(\ell)} \boldsymbol{w}^{(\ell-1)} \right] \boldsymbol{V}^{(\mu,l-2)} - \left[ \boldsymbol{w}^{(l)} \boldsymbol{\theta}^{(\ell-1)} + \boldsymbol{\theta}^{(\ell)} \right].$  (38)

Here,  $\boldsymbol{V}^{(\mu,\ell)}$  is the  $\mu^{\text{th}}$  neuron in the  $\ell^{\text{th}}$  hidden layer. Furthermore,  $\boldsymbol{w}^{(\ell)}$  and  $\boldsymbol{\theta}^{(\ell)}$  are the weight matrix and the shold vector for the neurons in the  $\ell^{\text{th}}$  hidden layer. Write at most three sentences where you, based on eq. (38), argue that a non-linear activation function is essential for a multi-layer perceptron.

# Solution to "6. Linear activation function" a)

$$\frac{\partial H}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{\mu=1}^p \left( O^{(\mu)} - t^{(\mu)} \right)^2 \tag{39}$$

$$=\sum_{\substack{\mu=1\\p}}^{p} \left( O^{(\mu)} - t^{(\mu)} \right) \frac{\partial O^{(\mu)}}{\partial w_{i}}$$
(40)

$$=\sum_{\mu=1}^{p} \left( O^{(\mu)} - t^{(\mu)} \right) x_{i}^{\mu}$$
(41)

$$=\sum_{\mu=1}^{p} \left(\sum_{j=1}^{N} w_j x_j^{(\mu)} - \theta - t^{(\mu)}\right) x_i^{\mu}$$
(42)

$$=\sum_{\mu=1}^{p} \left(\sum_{j=1}^{N} w_{j} x_{j}^{(\mu)}\right) x_{i}^{\mu} + \sum_{\mu=1}^{p} \left(-\theta\right) x_{i}^{\mu} + \sum_{\mu=1}^{p} \left(-t^{(\mu)}\right) x_{i}^{\mu} \qquad (43)$$

$$=\sum_{\mu=1}^{p}\sum_{j=1}^{N}w_{j}x_{j}^{(\mu)}x_{i}^{\mu}-\sum_{\mu=1}^{p}\theta x_{i}^{\mu}-\sum_{\mu=1}^{p}t^{(\mu)}x_{i}^{\mu}$$
(44)

$$=\sum_{j=1}^{N}\sum_{\mu=1}^{p}w_{j}x_{j}^{(\mu)}x_{i}^{\mu}-\theta\sum_{\mu=1}^{p}x_{i}^{\mu}-\sum_{\mu=1}^{p}t^{(\mu)}x_{i}^{\mu}$$
(45)

$$=\sum_{\substack{j=1\\N}}^{N} w_j \sum_{\mu=1}^{p} x_j^{(\mu)} x_i^{\mu} - \theta \sum_{\mu=1}^{p} x_i^{\mu} - \sum_{\mu=1}^{p} t^{(\mu)} x_i^{\mu}$$
(46)

$$=\sum_{j=1}^{N} w_j p G_{ji} - \theta p \beta_i - p \alpha_i \tag{47}$$

$$= p\left(\sum_{j=1}^{N} G_{ij}w_j - \theta\beta_i - \alpha_i\right)$$
(48)

$$\frac{\partial H}{\partial w_i} = 0 \Rightarrow \mathbb{G}\boldsymbol{w} = \boldsymbol{\alpha} + \theta\boldsymbol{\beta}$$
(50)

$$\frac{\partial H}{\partial \theta} = \frac{\partial}{\partial \theta} \frac{1}{2} \sum_{\mu=1}^{p} \left( O^{(\mu)} - t^{(\mu)} \right)^2 \tag{51}$$

$$=\sum_{\mu=1}^{p} \left( O^{(\mu)} - t^{(\mu)} \right) \frac{\partial O^{(\mu)}}{\partial \theta}$$
(52)

$$=\sum_{\mu=1}^{p} \left( O^{(\mu)} - t^{(\mu)} \right) x_{i}^{\mu}$$
(53)

$$=\sum_{\mu=1}^{p} \left(\sum_{j=1}^{N} w_j x_j^{(\mu)} - \theta - t^{(\mu)}\right) (-1)$$
(54)

$$= -\sum_{\mu=1}^{p} \left( \sum_{j=1}^{N} w_j x_j^{(\mu)} \right) - \sum_{\mu=1}^{p} (-\theta) - \sum_{\mu=1}^{p} \left( -t^{(\mu)} \right)$$
(55)

$$= -\sum_{\mu=1}^{p} \sum_{j=1}^{N} w_j x_j^{(\mu)} + \sum_{\mu=1}^{p} \theta + \sum_{\mu=1}^{p} t^{(\mu)}$$
(56)

$$= -\sum_{j=1}^{N} w_j \sum_{\mu=1}^{p} x_j^{(\mu)} + p\theta + p\gamma$$
(57)

$$= -p\sum_{j=1}^{N} w_j \beta_j + p\theta + pc \tag{58}$$

(59)

$$\frac{\partial H}{\partial \theta} = 0 \Rightarrow \boldsymbol{w}^{\mathsf{T}} \boldsymbol{\beta} = \theta + \gamma.$$
(60)

b) The first equation gives:

$$\boldsymbol{w} = \mathbb{G}^{-1}\boldsymbol{\alpha} + \theta \mathbb{G}^{-1}\boldsymbol{\beta}.$$
 (61)

Insert into the second, and use that  $\boldsymbol{w}^{\mathsf{T}}\boldsymbol{\beta} = \boldsymbol{\beta}^{\mathsf{T}}\boldsymbol{w}$ :

$$\boldsymbol{\beta}^{\mathsf{T}} \left[ \mathbb{G}^{-1} \boldsymbol{\alpha} + \boldsymbol{\theta} \mathbb{G}^{-1} \boldsymbol{\beta} \right] = \boldsymbol{\theta} + \boldsymbol{\gamma}$$
(62)

$$\beta \left[ \mathbb{G} \quad \boldsymbol{\alpha} + \theta \mathbb{G} \quad \boldsymbol{\beta} \right] = \theta + \gamma$$

$$\Rightarrow \boldsymbol{\beta}^{\mathsf{T}} \mathbb{G}^{-1} \boldsymbol{\alpha} + \theta \boldsymbol{\beta}^{\mathsf{T}} \mathbb{G}^{-1} \boldsymbol{\beta} = \theta + \gamma$$

$$(62)$$

$$\Rightarrow \theta \left[ \boldsymbol{\beta}^{\mathsf{T}} \mathbb{G}^{-1} \boldsymbol{\beta} - 1 \right] = \gamma - \boldsymbol{\beta}^{\mathsf{T}} \mathbb{G}^{-1} \boldsymbol{\alpha} \tag{64}$$

$$\Rightarrow \theta = \frac{\gamma - \beta^{\mathsf{T}} \mathbb{G}^{-1} \boldsymbol{\alpha}}{\beta^{\mathsf{T}} \mathbb{G}^{-1} \boldsymbol{\beta} - 1}.$$
 (65)

c) The equation can be written as

$$\boldsymbol{V}^{(\mu,\ell)} = \boldsymbol{W}\boldsymbol{V}^{(\mu,\ell-2)} - \boldsymbol{\Theta},\tag{66}$$

where

$$\boldsymbol{W} = \boldsymbol{w}^{(\ell)} \boldsymbol{w}^{(\ell-1)}, \tag{67}$$

and

$$\boldsymbol{\Theta} = \boldsymbol{w}^{(l)} \boldsymbol{\theta}^{(\ell-1)} + \boldsymbol{\theta}^{(\ell)}.$$
 (68)

The two layers can tehrefore be collapsed into one single layer, and with a linear activation function in all layers the whole perceptron collapses into a simple perceptron with linear activation function. Such a perceptron can only solve linearly separable problems.