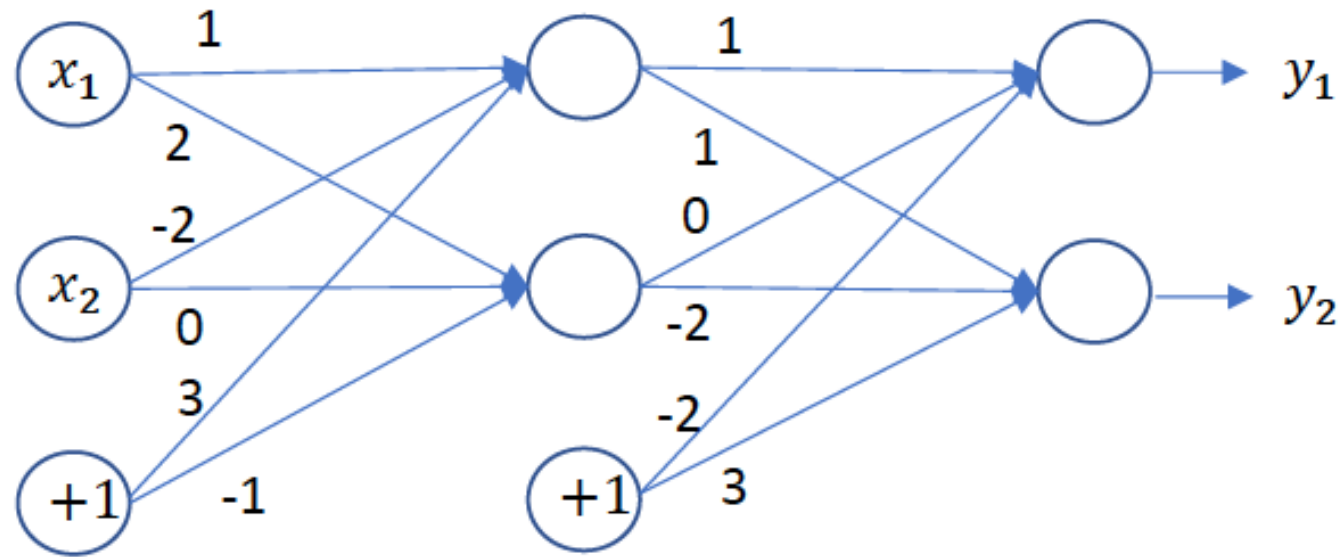


# Deep neural networks

SC4001 – Tutorial 4



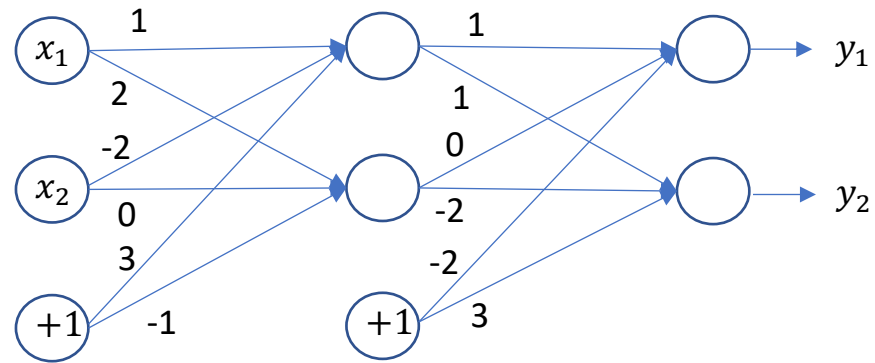
1. The two-layer feedforward perceptron network shown in figure 1 has weights and biases initialized as indicated and receives 2-dimensional inputs  $(x_1, x_2)$ . The network is to respond with  $\mathbf{d}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$  and  $\mathbf{d}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$  for input patterns  $\mathbf{x}_1 = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix}$  and  $\mathbf{x}_2 = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix}$ , respectively.

Analyse a single feedforward and feedback step for gradient decent learning of the two patterns by doing the following:

- (a) Find the weight matrix  $\mathbf{W}$  to the hidden-layer and weight matrix  $\mathbf{V}$  to the output-layer, and the corresponding biases.
- (b) Calculate the synaptic input  $\mathbf{z}$  and output  $\mathbf{h}$  of the hidden-layer, and the synaptic input  $\mathbf{u}$  and output  $\mathbf{y} = (y_1, y_2)$  of the output layer.
- (c) Find the mean square error cost  $J$  between the outputs and targets.
- (d) Calculate the gradients  $\nabla_{\mathbf{u}}J$  and  $\nabla_{\mathbf{z}}J$  at the output-layer and the hidden-layer, respectively.
- (e) Compute the new weights and biases.
- (f) Write a program to continue iterations until convergence and find the final weights and biases.

Assume a learning rate of 0.05.

Repeat above (a) – (f) for stochastic gradient decent learning.



Weight matrix to the hidden layer,  $\mathbf{W} = \begin{pmatrix} 1.0 & 2.0 \\ -2.0 & 0.0 \end{pmatrix}$

Bias vector to the hidden-layer  $\mathbf{b} = \begin{pmatrix} 3.0 \\ -1.0 \end{pmatrix}$

Weight matrix to the output-layer,  $\mathbf{V} = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix}$

Bias vector to the output-layer  $\mathbf{c} = \begin{pmatrix} -2.0 \\ 3.0 \end{pmatrix}$

## GD for 2-layer perceptron network:

Given a training dataset  $(\mathbf{X}, \mathbf{D})$

Set learning parameter  $\alpha$

Initialize  $\mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c}$

Repeat until convergence:

$$\mathbf{Z} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

$$\mathbf{H} = g(\mathbf{Z})$$

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$$

$$\mathbf{Y} = f(\mathbf{U})$$

Forward propagation  
of activation

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U})$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

Backward propagation  
of gradients

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P$$

Weight updating

$$\mathbf{x}_1 = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix} \text{ and } \mathbf{d}_1 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\mathbf{x}_2 = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix} \text{ and } \mathbf{d}_2 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\mathbf{X} = \begin{pmatrix} 1.0 & 3.0 \\ -2.0 & -2.0 \end{pmatrix} \text{ and } \mathbf{D} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

Forward propagation:

Synaptic input to hidden-layer,  $\mathbf{Z} = \mathbf{XW} + \mathbf{B}$

$$= \begin{pmatrix} 1.0 & 3.0 \\ -2.0 & -2.0 \end{pmatrix} \begin{pmatrix} 1.0 & 2.0 \\ -2.0 & 0.0 \end{pmatrix} + \begin{pmatrix} 3.0 & -1.0 \\ 3.0 & -1.0 \end{pmatrix}$$

$$= \begin{pmatrix} -2.0 & 1.0 \\ 5.0 & -5.0 \end{pmatrix}$$

Output of the hidden layer,  $\mathbf{H} = g(\mathbf{Z}) = \frac{1}{1+e^{-Z}} = \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix}$

Synaptic input to output-layer,  $\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$

$$\begin{aligned} &= \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix} \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix} + \begin{pmatrix} -2.0 & 3.0 \\ -2.0 & 3.0 \end{pmatrix} \\ &= \begin{pmatrix} -1.88 & 1.66 \\ -0.99 & 3.98 \end{pmatrix} \end{aligned}$$

Output of the output layer,  $\mathbf{Y} = f(\mathbf{U}) = \frac{1}{1+e^{-U}} = \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix}$

$$\mathbf{D} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

$$\begin{aligned} m.s.e. &= \frac{1}{2} \sum_{p=1}^2 \sum_{k=1}^2 (d_{pk} - y_{pk})^2 \\ &= \frac{1}{2} \left( ((0 - 0.13)^2 + (1 - 0.84)^2) + ((1 - 0.27)^2 + (0 - 0.98)^2) \right) \\ &= 0.77 \end{aligned}$$

**Computing gradients (backward propagation):**

$$f'(\mathbf{U}) = \mathbf{Y} \cdot (\mathbf{1} - \mathbf{Y}) = \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix} \cdot \left( \begin{pmatrix} 1.0 & 1.0 \\ 1.0 & 1.0 \end{pmatrix} - \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix} \right) = \begin{pmatrix} 0.11 & 0.13 \\ 0.20 & 0.02 \end{pmatrix}$$

$$\nabla_{\mathbf{U}} J = -(\mathbf{D} - \mathbf{Y}) \cdot f'(\mathbf{U}) = -\left( \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} - \begin{pmatrix} 0.13 & 0.84 \\ 0.27 & 0.98 \end{pmatrix} \right) \begin{pmatrix} 0.11 & 0.13 \\ 0.20 & 0.02 \end{pmatrix} = \begin{pmatrix} 0.02 & -0.02 \\ -0.14 & 0.02 \end{pmatrix}$$

$$g'(\mathbf{Z}) = \mathbf{H} \cdot (1 - \mathbf{H}) = \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix} \cdot \left( \begin{pmatrix} 1.0 & 1.0 \\ 1.0 & 1.0 \end{pmatrix} - \begin{pmatrix} 0.12 & 0.73 \\ 0.99 & 0.01 \end{pmatrix} \right) = \begin{pmatrix} 0.10 & 0.2 \\ 0.01 & 0.01 \end{pmatrix}$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z}) = \begin{pmatrix} 0.02 & -0.02 \\ -0.14 & 0.02 \end{pmatrix} \begin{pmatrix} 1.0 & 0.0 \\ 1.0 & -2.0 \end{pmatrix} \cdot \begin{pmatrix} 0.10 & 0.2 \\ 0.01 & 0.01 \end{pmatrix} = \begin{pmatrix} -0.001 & 0.01 \\ -0.001 & 0.00 \end{pmatrix}$$



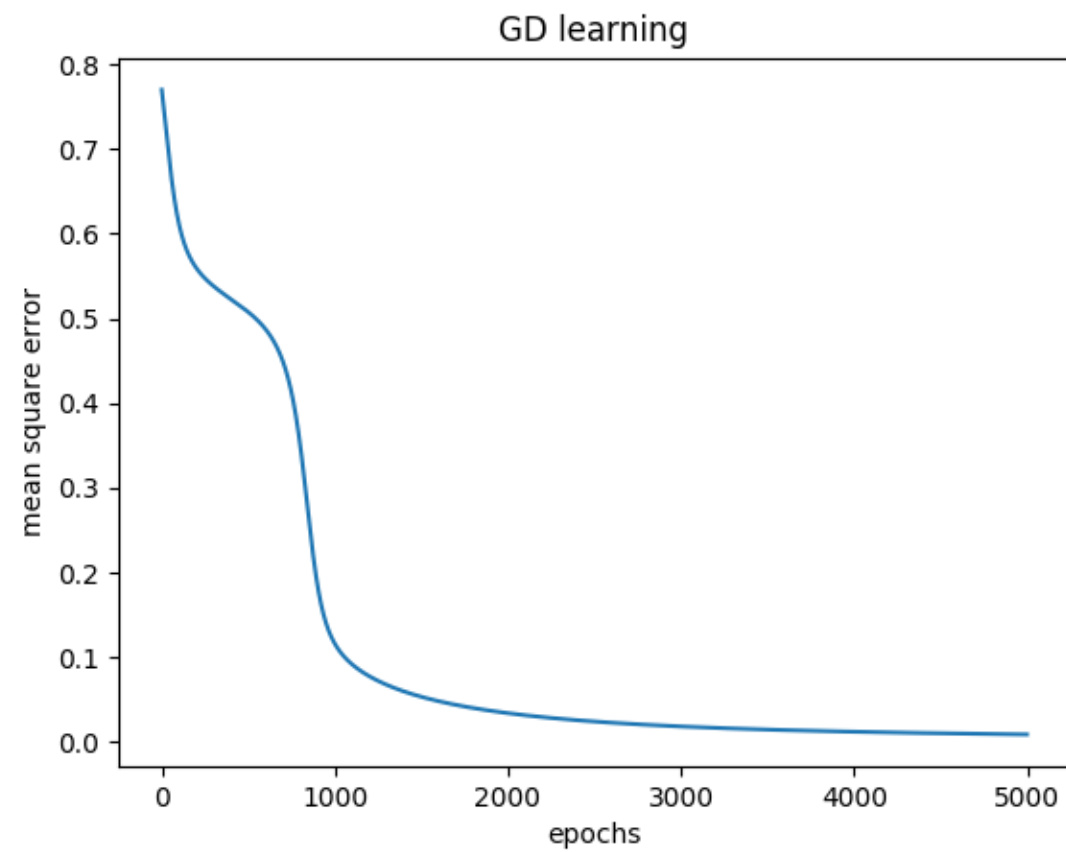
**Updating weights:**

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_U J = \begin{pmatrix} 1.01 & 1.0 \\ 0.0 & -2.0 \end{pmatrix}$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_U J)^T \mathbf{1}_P = \begin{pmatrix} -1.99 \\ 3.00 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_Z J = \begin{pmatrix} 1.0 & 2.0 \\ -2.0 & 0.0 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_Z J)^T \mathbf{1}_P = \begin{pmatrix} 3.0 \\ -1.0 \end{pmatrix}$$



**At convergence:**

$$\mathbf{W} = \begin{pmatrix} 0.63 & 0.60 \\ -3.0 & -2.0 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 2.72 \\ -0.74 \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 4.97 & -3.46 \\ 0.25 & -2.37 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} -2.42 \\ 2.56 \end{pmatrix}$$

Predicted values:

$$\mathbf{y}_1 = \begin{pmatrix} 0.08 \\ 0.93 \end{pmatrix} \text{ and } \mathbf{y}_2 = \begin{pmatrix} 0.94 \\ 0.05 \end{pmatrix}$$

$$\text{m.s.e.} = 0.004$$

## SGD learning for 2-layer perceptron network:

Given a training dataset  $\{(x, d)\}$

Set learning parameter  $\alpha$

Initialize  $W, b, V, c$

Repeat until convergence:

For every pattern  $(x, d)$ :

$$z = W^T x + b$$

$$h = g(z)$$

$$u = V^T h + c$$

$$y = f(u)$$

Forward propagation  
of activation

$$\nabla_u J = -(d - y) \cdot f'(z)$$

$$\nabla_z J = V \nabla_u J \cdot g'(z)$$

Backward propagation  
of gradients

$$V \leftarrow V - \alpha h (\nabla_u J)^T$$

$$c \leftarrow c - \alpha \nabla_u J$$

$$W \leftarrow W - \alpha x (\nabla_z J)^T$$

$$b \leftarrow b - \alpha \nabla_z J$$

Weight updating

## Epoch 1:

Apply **first** pattern  $\mathbf{x} = \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix}$  and  $\mathbf{d} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$ :

Synaptic input to the hidden-layer

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b} = \begin{pmatrix} 1.0 & -2.0 \\ 2.0 & 0.0 \end{pmatrix} \begin{pmatrix} 1.0 \\ 3.0 \end{pmatrix} + \begin{pmatrix} 3.0 \\ -1.0 \end{pmatrix} = \begin{pmatrix} -2.0 \\ 1.0 \end{pmatrix}$$

Output of the hidden-layer  $\mathbf{h} = g(\mathbf{z}) = \frac{1}{1+e^{-z}} = \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix}$

Synaptic input to output-layer

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c} = \begin{pmatrix} -1.88 \\ 1.66 \end{pmatrix}$$

Output of the output-layer  $\mathbf{y} = f(\mathbf{u}) = \frac{1}{1+e^{-u}} = \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix}$

$$s.e. = (d_1 - y_1)^2 + (d_2 - y_2)^2 = 0.043$$

### Computing gradients:

$$f'(\mathbf{u}) = \mathbf{y} \cdot (1 - \mathbf{y}) = \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix} \cdot \left( \begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix} - \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix} \right) = \begin{pmatrix} 0.11 \\ 0.13 \end{pmatrix}$$

$$\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) = -\left( \begin{pmatrix} 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 0.13 \\ 0.84 \end{pmatrix} \right) \cdot \begin{pmatrix} 0.12 \\ 0.14 \end{pmatrix} = \begin{pmatrix} 0.02 \\ -0.02 \end{pmatrix}$$

$$g'(\mathbf{z}) = \mathbf{h} \cdot (1 - \mathbf{h}) = \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix} \cdot \left( \begin{pmatrix} 1.0 \\ 1.0 \end{pmatrix} - \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix} \right) = \begin{pmatrix} 0.10 \\ 0.20 \end{pmatrix}$$

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z}) = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix} \begin{pmatrix} 0.02 \\ -0.02 \end{pmatrix} \cdot \begin{pmatrix} 0.11 \\ 0.20 \end{pmatrix} = \begin{pmatrix} -0.001 \\ 0.008 \end{pmatrix}$$

### Updating weights:

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{h} (\nabla_{\mathbf{u}} J)^T = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & -2.0 \end{pmatrix} - 0.2 \begin{pmatrix} 0.12 \\ 0.73 \end{pmatrix} \begin{pmatrix} -0.02 & 0.022 \end{pmatrix} = \begin{pmatrix} 1.0 & 1.0001 \\ 0.00 & -2.0 \end{pmatrix}$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J = \begin{pmatrix} -2.0 \\ 3.0 \end{pmatrix} + 0.2 \begin{pmatrix} 0.02 \\ -0.02 \end{pmatrix} = \begin{pmatrix} -2.00 \\ 3.001 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{z}} J)^T = \begin{pmatrix} 1.0 & 2.0 \\ -2.00 & -0.001 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J = \begin{pmatrix} 3.00 \\ -1.00 \end{pmatrix}$$

Apply **second** pattern  $\mathbf{x} = \begin{pmatrix} -2.0 \\ -2.0 \end{pmatrix}$  and  $\mathbf{d} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ :

Synaptic input to the hidden-layer

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b} = \begin{pmatrix} 5.0 \\ -5.0 \end{pmatrix}$$

Output of the hidden-layer  $\mathbf{h} = g(\mathbf{z}) = \frac{1}{1+e^{-z}} = \begin{pmatrix} 1.0 \\ 0.007 \end{pmatrix}$

Synaptic input to output-layer

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c} = \begin{pmatrix} -0.99 \\ 3.98 \end{pmatrix}$$

Output of the output-layer  $\mathbf{y} = f(\mathbf{u}) = \frac{1}{1+e^{-u}} = \begin{pmatrix} 0.27 \\ 0.98 \end{pmatrix}$

$$s.e. = (d_1 - y_1)^2 + (d_2 - y_2)^2 = 1.5$$

**Computing gradients:**

$$f'(\mathbf{u}) = \mathbf{y} \cdot (1 - \mathbf{y}) = \begin{pmatrix} 0.195 \\ 0.018 \end{pmatrix}$$

$$\nabla_{\mathbf{u}} J = -(\mathbf{d} - \mathbf{y}) \cdot f'(\mathbf{u}) = \begin{pmatrix} -0.14 \\ 0.018 \end{pmatrix}$$

$$g'(\mathbf{z}) = \mathbf{h} \cdot (1 - \mathbf{h}) = \begin{pmatrix} 0.007 \\ 0.007 \end{pmatrix}$$

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z}) = \begin{pmatrix} -0.0008 \\ -0.0002 \end{pmatrix}$$

**Updating weights:**

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{h} (\nabla_{\mathbf{u}} J)^T = \begin{pmatrix} 1.007 & 0.99 \\ 0.0 & -2.0 \end{pmatrix}$$

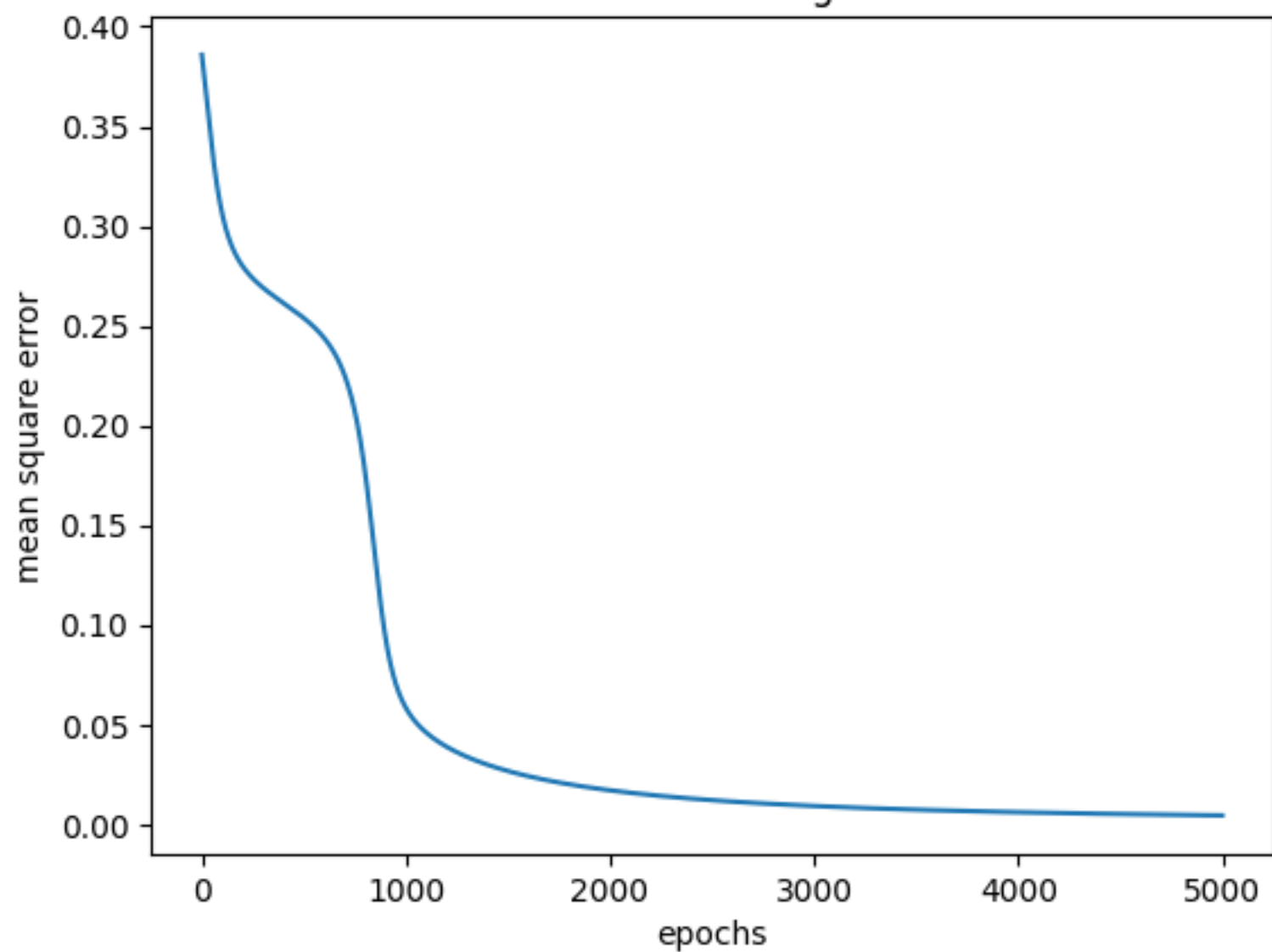
$$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J = \begin{pmatrix} -1.99 \\ 3.0 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{z}} J)^T = \begin{pmatrix} 0.999 & 1.99 \\ -1.99 & 0.00 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J = \begin{pmatrix} 3.00 \\ -1.00 \end{pmatrix}$$



SGD learning



At convergence:

$$\mathbf{W} = \begin{pmatrix} 0.63 & 0.60 \\ -3.0 & -2.0 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 2.72 \\ -0.74 \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 4.97 & -3.46 \\ 0.25 & -2.37 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} -2.42 \\ 2.56 \end{pmatrix}$$

Predicted values:

$$\mathbf{y}_1 = \begin{pmatrix} 0.08 \\ 0.93 \end{pmatrix} \text{ and } \mathbf{y}_2 = \begin{pmatrix} 0.94 \\ 0.05 \end{pmatrix}$$

$$\text{m.s.e.} = 0.004$$

2. A feedforward neural network with one hidden layer to perform the following classification:

class	inputs
A	(1.0, 1.0), (0.0, 1.0)
B	(3.0, 4.0), (2.0, 2.0)
C	(2.0, -2.0), (-2.0, -3.0)

The network has a hidden layer consisting of three perceptrons and a softmax output layer.

Show one iteration of gradient descent learning and plot learning curves until convergence at a learning rate  $\alpha = 0.1$ .

Initialize the weights  $\mathbf{W}$  and biases  $\mathbf{b}$  to the hidden layer, and the weights  $\mathbf{V}$  and biases  $\mathbf{c}$  to the output layer as follows:

$$\mathbf{W} = \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$
$$\mathbf{V} = \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

Determine the weights and biases at convergence.

Find the class labels predicted by the trained network for patterns:

$$\mathbf{x}_1 = \begin{pmatrix} 2.5 \\ 1.5 \end{pmatrix} \text{ and } \mathbf{x}_2 = \begin{pmatrix} -1.5 \\ 0.5 \end{pmatrix}$$

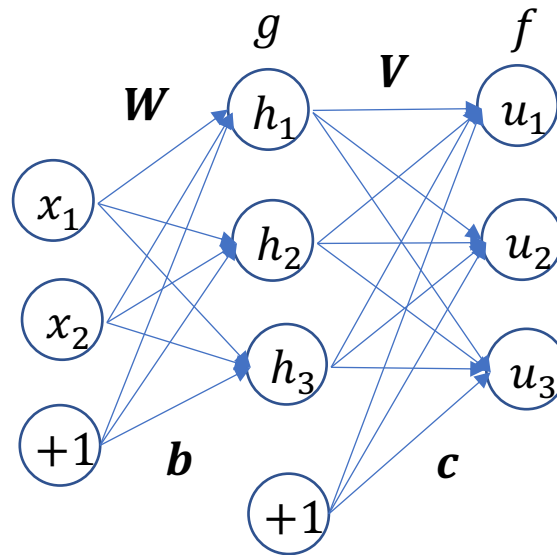
## Training examples (patterns):

class	inputs	Target Label
A	(1.0, 1.0), (0.0, 1.0)	1
B	(3.0, 4.0), (2.0, 2.0)	2
C	(2.0, -2.0), (-2.0, -3.0)	3

Feedforward network :

Perceptron hidden layer with 3 neurons

Softmax output layer with 3 neurons



$$g(\mathbf{Z}) = \frac{1}{1 + e^{-\mathbf{Z}}}$$

$$f(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{\mathbf{U}_k}}$$

### Training inputs and targets

$$\mathbf{X} = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & 1.0 \\ 3.0 & 4.0 \\ 2.0 & 2.0 \\ 2.0 & -2.0 \\ -2.0 & -3.0 \end{pmatrix} \text{ and } \mathbf{D} = \begin{pmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

### Targets as a one hot matrix:

$$\mathbf{K} = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$$

### Initial weights and biases:

To the hidden layer,

$$\mathbf{W} = \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix}, \mathbf{b} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

To the output-layer

$$\mathbf{V} = \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}$$

**Learning factor  $\alpha = 0.1$**

### Activation functions:

Hidden layer is a continuous perceptron layer:  $g(\mathbf{Z}) = \frac{1}{1+e^{-\mathbf{Z}}}$

Output layer is a softmax layer:  $f(\mathbf{U}) = \frac{e^U}{\sum_{k=1}^K e^{U_k}}$

## GD for the feedforward network

Given a training dataset  $(\mathbf{X}, \mathbf{D})$

Set learning parameter  $\alpha$

Initialize  $\mathbf{W}, \mathbf{b}, \mathbf{V}, \mathbf{c}$

Repeat until convergence:

$$\mathbf{Z} = \mathbf{X}\mathbf{W} + \mathbf{B}$$

$$\mathbf{H} = g(\mathbf{Z})$$

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$$

$$\mathbf{Y} = \arg \max_k f(\mathbf{U})$$

Forward propagation  
of activation

$$\nabla_{\mathbf{U}} J = -(\mathbf{K} - f(\mathbf{U}))$$

$$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$$

Backward propagation  
of gradients

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P$$

Weight updating

Synaptic input to hidden-layer,

$$\mathbf{Z} = \mathbf{XW} + \mathbf{B} = \begin{pmatrix} 1.0 & 1.0 \\ 0.0 & 1.0 \\ 3.0 & 4.0 \\ 2.0 & 2.0 \\ 2.0 & -2.0 \\ -2.0 & -3.0 \end{pmatrix} \begin{pmatrix} -0.10 & 0.97 & 0.18 \\ -0.70 & 0.38 & 0.93 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} -0.80 & 1.35 & 1.10 \\ -0.70 & 0.38 & 0.93 \\ -3.08 & 4.44 & 4.23 \\ -1.59 & 2.70 & 2.21 \\ 1.20 & 1.18 & -1.50 \\ 2.29 & -3.08 & -3.13 \end{pmatrix}$$

$$\text{Output of the hidden layer, } \mathbf{H} = g(\mathbf{Z}) = \frac{1}{1+e^{-\mathbf{Z}}} = \begin{pmatrix} 0.31 & 0.79 & 0.75 \\ 0.33 & 0.59 & 0.72 \\ 0.04 & 0.99 & 0.99 \\ 0.17 & 0.94 & 0.90 \\ 0.77 & 0.77 & 0.18 \\ 0.91 & 0.04 & 0.04 \end{pmatrix}$$



Synaptic input to output-layer,

$$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C} = \begin{pmatrix} 0.31 & 0.79 & 0.75 \\ 0.33 & 0.59 & 0.72 \\ 0.04 & 0.99 & 0.99 \\ 0.17 & 0.94 & 0.90 \\ 0.77 & 0.77 & 0.18 \\ 0.91 & 0.04 & 0.04 \end{pmatrix} \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} 1.15 & 0.40 & -0.34 \\ 1.01 & 0.46 & -0.31 \\ 1.10 & 0.51 & -0.30 \\ 1.16 & 0.47 & -0.33 \\ 1.43 & -0.09 & -0.48 \\ 0.96 & 0.11 & -0.36 \end{pmatrix}$$

$$\text{Output layer activation } f(\mathbf{U}) = \frac{e^U}{\sum_{k=1}^K e^{U_k}} = \begin{pmatrix} 0.59 & 0.28 & 0.13 \\ 0.54 & 0.31 & 0.15 \\ 0.56 & 0.31 & 0.14 \\ 0.58 & 0.29 & 0.13 \\ 0.73 & 0.16 & 0.11 \\ 0.59 & 0.25 & 0.16 \end{pmatrix}$$

$$\text{Output } \mathbf{Y} = \arg \max_k f(\mathbf{U}) = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\mathbf{D} = \begin{pmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}, \quad \mathbf{Y} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}, \quad f(\mathbf{U}) = \begin{pmatrix} 0.59 & 0.28 & 0.13 \\ 0.54 & 0.31 & 0.15 \\ 0.56 & 0.31 & 0.14 \\ 0.58 & 0.29 & 0.13 \\ 0.73 & 0.16 & 0.11 \\ 0.59 & 0.25 & 0.16 \end{pmatrix}$$

$$\text{Classification error} = \sum 1(\mathbf{D} \neq \mathbf{Y}) = 4$$

$$\begin{aligned} \text{Entropy } J &= -\sum_{p=1}^P \log \left( f(u_{pd_p}) \right) \\ &= -(\log(0.59) + \log(0.54) + \log(0.31) + \log(0.29) + \log(0.11) + \log(0.16)) \\ &= 7.63 \end{aligned}$$

$$\nabla_U J = -(\mathbf{K} - f(\mathbf{U})) = -\left(\begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 0.59 & 0.28 & 0.13 \\ 0.54 & 0.31 & 0.15 \\ 0.56 & 0.31 & 0.14 \\ 0.58 & 0.29 & 0.13 \\ 0.73 & 0.16 & 0.11 \\ 0.59 & 0.25 & 0.16 \end{pmatrix}\right) = \begin{pmatrix} -0.41 & 0.28 & 0.13 \\ -0.46 & 0.31 & 0.15 \\ 0.56 & -0.69 & 0.14 \\ 0.58 & -0.71 & 0.13 \\ 0.73 & 0.16 & -0.89 \\ 0.59 & 0.25 & -0.84 \end{pmatrix}$$

$$g'(\mathbf{Z}) = \mathbf{H} \cdot (\mathbf{1} - \mathbf{H}) = \begin{pmatrix} 0.21 & 0.16 & 0.19 \\ 0.22 & 0.24 & 0.20 \\ 0.04 & 0.01 & 0.01 \\ 0.14 & 0.06 & 0.09 \\ 0.18 & 0.18 & 0.15 \\ 0.08 & 0.04 & 0.04 \end{pmatrix}$$

$$\nabla_Z J = (\nabla_U J) \mathbf{V}^T \cdot g'(\mathbf{Z}) = \begin{pmatrix} -0.41 & 0.28 & 0.13 \\ -0.46 & 0.31 & 0.15 \\ 0.56 & -0.69 & 0.14 \\ 0.58 & -0.71 & 0.13 \\ 0.73 & 0.16 & -0.89 \\ 0.59 & 0.25 & -0.84 \end{pmatrix} \begin{pmatrix} 1.01 & 0.09 & -0.39 \\ 0.79 & -0.45 & -0.22 \\ 0.28 & 0.96 & -0.07 \end{pmatrix}^T \cdot \begin{pmatrix} 0.21 & 0.16 & 0.19 \\ 0.22 & 0.24 & 0.20 \\ 0.04 & 0.01 & 0.01 \\ 0.14 & 0.06 & 0.09 \\ 0.18 & 0.18 & 0.15 \\ 0.08 & 0.04 & 0.04 \end{pmatrix} = \begin{pmatrix} -0.09 & -0.08 & 0.03 \\ -0.11 & -0.13 & 0.03 \\ 0.02 & 0.01 & -0.01 \\ 0.07 & 0.04 & -0.05 \\ 0.20 & 0.13 & 0.06 \\ 0.08 & 0.02 & 0.02 \end{pmatrix}$$

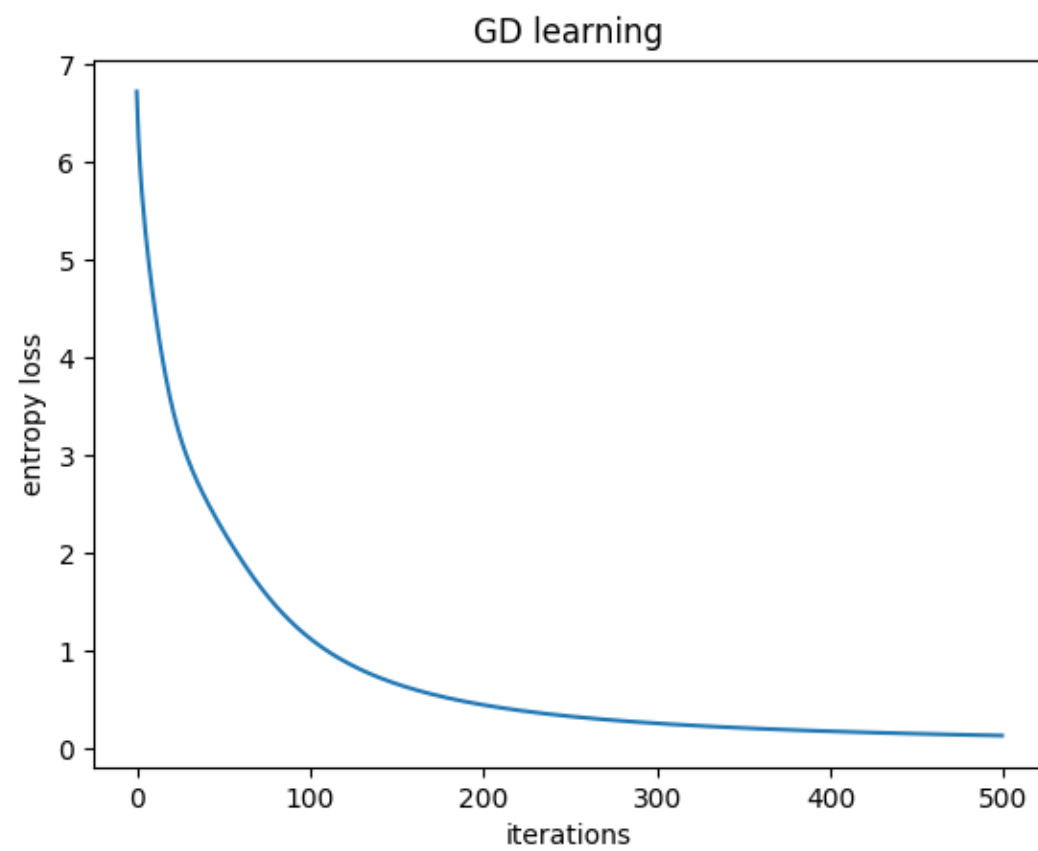
Learning rate  $\alpha = 0.1$

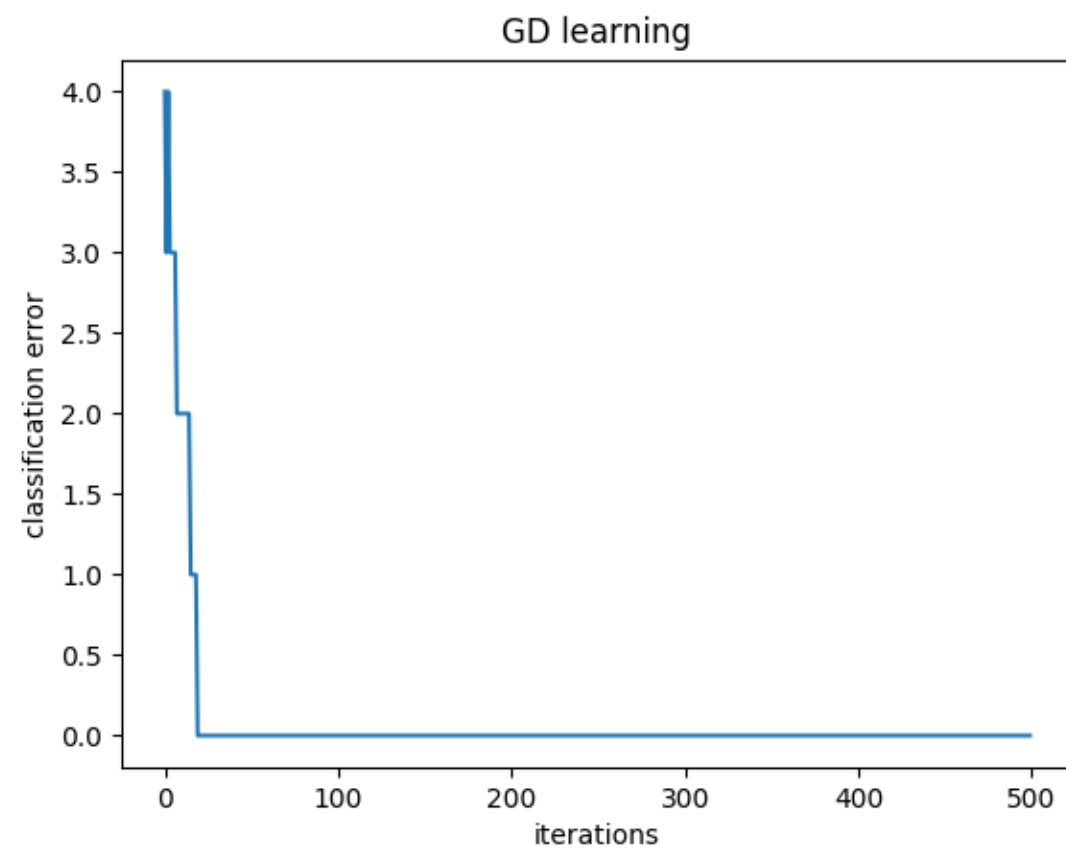
$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_U J = \begin{pmatrix} 0.92 & 0.05 & -0.26 \\ 0.68 & -0.36 & -0.19 \\ 0.22 & 1.05 & -0.10 \end{pmatrix}$$

$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_U J)^T \mathbf{1}_P = \begin{pmatrix} -0.16 \\ 0.04 \\ 0.12 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_Z J = \begin{pmatrix} -0.13 & 0.95 & 0.18 \\ -0.63 & 0.42 & 0.95 \end{pmatrix}$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_Z J)^T \mathbf{1}_P = \begin{pmatrix} -0.02 \\ 0.00 \\ -0.01 \end{pmatrix}$$





At convergence:

$$\mathbf{V} = \begin{pmatrix} 2.93 & -5.33 & 3.12 \\ 2.80 & 1.20 & -3.87 \\ 0.09 & 4.55 & -3.47 \end{pmatrix}, \quad \mathbf{c} = \begin{pmatrix} -1.94 \\ -0.06 \\ 2.01 \end{pmatrix}$$

$$\mathbf{W} = \begin{pmatrix} -1.81 & 0.32 & 0.08 \\ -1.40 & 2.92 & 1.91 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 4.36 \\ 0.73 \\ -1.71 \end{pmatrix}$$

$$\mathbf{Y} = \begin{pmatrix} 1 \\ 1 \\ 2 \\ 2 \\ 3 \\ 3 \end{pmatrix}$$

Entropy = 0.138

Error = 0

Testing patterns:

$$\mathbf{x}_1 = \begin{pmatrix} 2.5 \\ 1.5 \end{pmatrix} \text{ and } \mathbf{x}_2 = \begin{pmatrix} -1.5 \\ 0.5 \end{pmatrix}$$

In batch mode

$$\mathbf{X} = \begin{pmatrix} 2.5 & 1.5 \\ -1.5 & 0.5 \end{pmatrix}$$

Forward propagation of activations

$$\mathbf{Z} = \mathbf{XW} + \mathbf{B} = \begin{pmatrix} -2.26 & 5.91 & 1.36 \\ 6.35 & 1.72 & -0.91 \end{pmatrix}$$

$$\mathbf{H} = f(\mathbf{Z}) = \begin{pmatrix} 0.09 & 1.0 & 0.8 \\ 1.0 & 0.85 & 0.29 \end{pmatrix}$$

$$\mathbf{U} = \mathbf{HV} + \mathbf{C} = \begin{pmatrix} 1.2 & 4.25 & -4.33 \\ 3.38 & -3.006 & 0.82 \end{pmatrix}$$

$$g(\mathbf{U}) = \frac{e^{\mathbf{U}}}{\sum_{k=1}^K e^{U_k}} = \begin{pmatrix} 0.05 & 0.95 & 0.0 \\ 0.93 & 0.0 & 0.07 \end{pmatrix}$$

$$\mathbf{Y} = \arg \max_k g(\mathbf{U}) = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

The class labels:  $\mathbf{x}_1 \rightarrow \text{class } B, \mathbf{x}_2 \rightarrow \text{class } A$



3. Design a feedforward neural network consisting of two-hidden layers to approximate the following function:

$$\phi(x, y) = 0.8x^2 - y^3 + 2.5xy$$

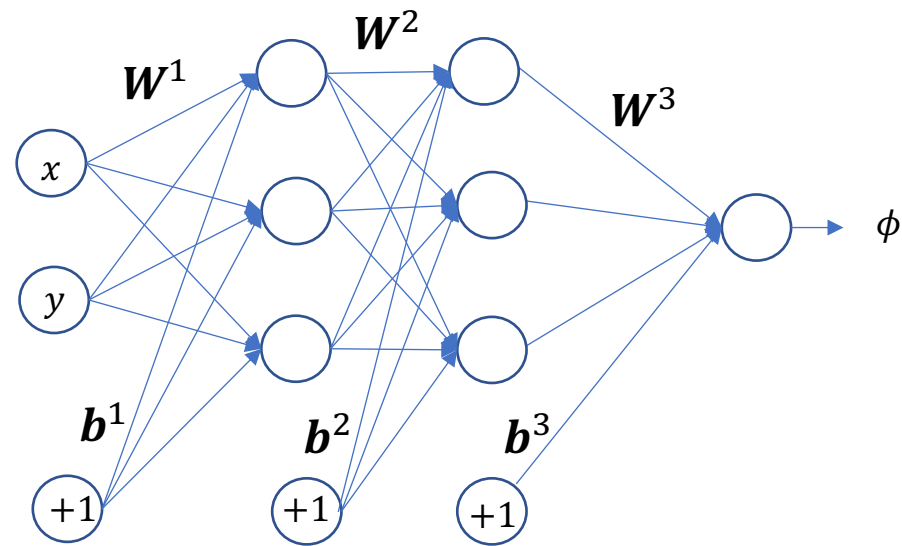
for  $-1.0 \leq x, y \leq 1.0$ .

Use three ReLU neurons at each hidden layer and a linear neuron at the output layer.

- (a) Divide the input space equally into square regions of size  $0.25 \times 0.25$  and use grid points as data to learn the function  $\phi$ .
- (b) Train the network using gradient decent learning at learning rate  $\alpha = 0.01$  and plot the learning curve (mean square error vs. iterations) and the predicted data points.
- (c) Compare the learning curves when learning the function at learning rates  $\alpha = 0.005, 0.01, 0.05$ , and  $0.1$ .

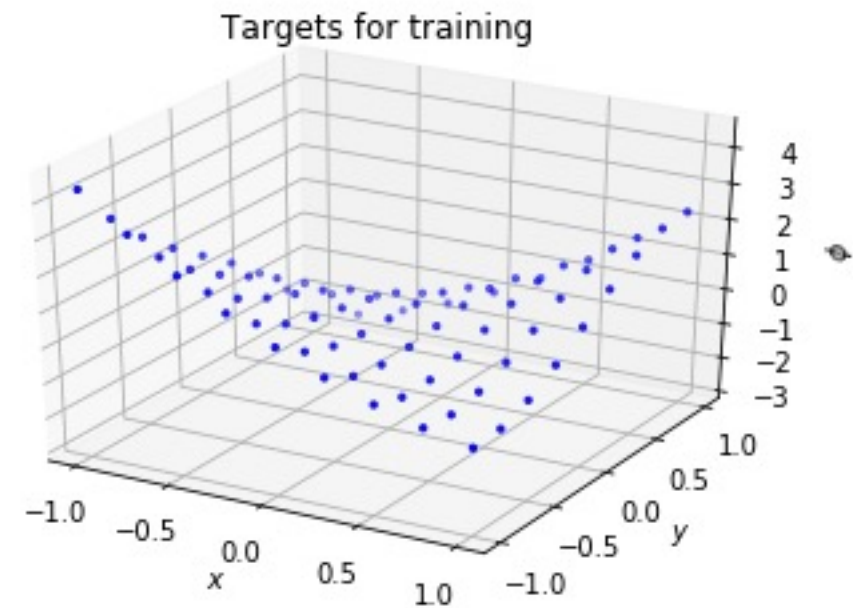
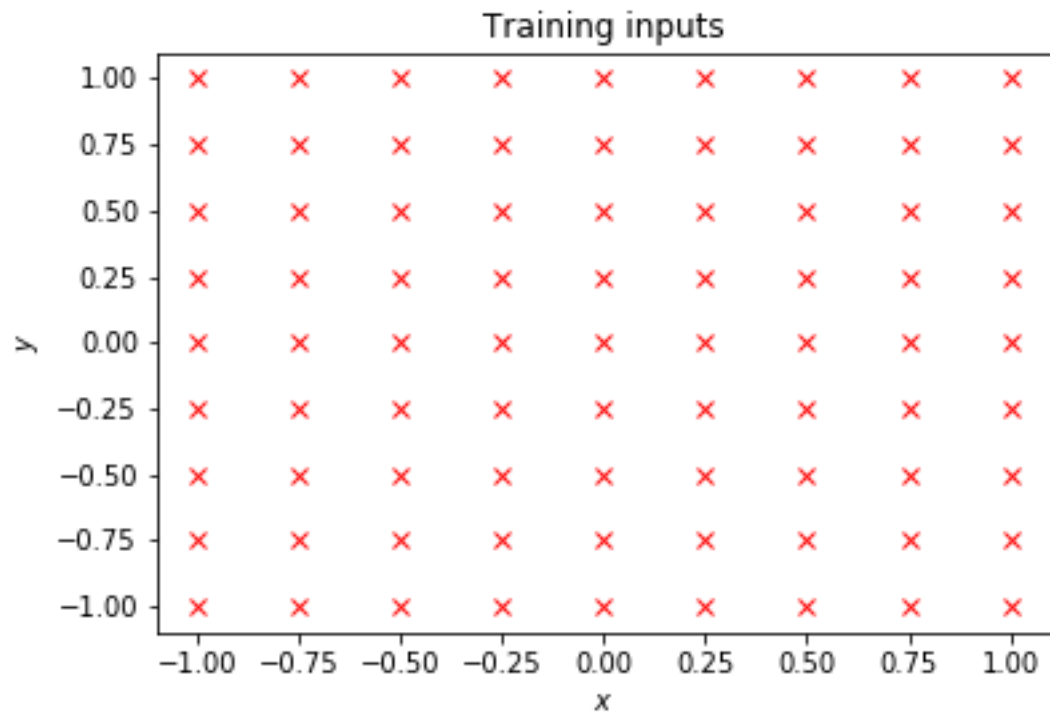
$$\phi(x, y) = 0.8x^2 - x^3 + 2.5xy \quad \text{for } -1.0 \leq x, y \leq 1.0$$

Feedforward neural network with two hidden layers:



$$\phi(x, y) = 0.8x_1^2 - y^3 + 2.5xy \quad \text{for } -1.0 \leq x, y \leq 1.0$$

Data points in a grid of size 0.25x0.25:



```
class FFN(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.relu_stack = nn.Sequential(  
            nn.Linear(2, 10),  
            nn.ReLU(),  
            nn.Linear(10, 5),  
            nn.ReLU(),  
            nn.Linear(5, 1),  
        )  
  
    def forward(self, x):  
        logits = self.relu_stack(x)  
        return logits
```

