

Chapter 5

Deep feedforward networks

Neural networks and deep learning

Chain rule of differentiation

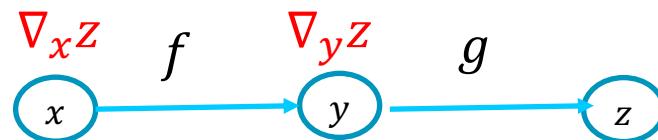
Let $x, y, z \in \mathbf{R}$ be one-dimensional variables and

$$y = f(x)$$
$$z = g(y)$$

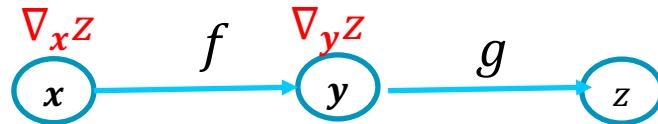
Chain rule of differentiation states that:

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$

$$\nabla_x z = \left(\frac{\partial y}{\partial x} \right) \nabla_y z$$



Chain rule in multidimensions



$\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathbf{R}^n$, $\mathbf{y} = (y_1, y_2, \dots, y_K) \in \mathbf{R}^K$, $z \in \mathbf{R}$, and

$$\begin{aligned}\mathbf{y} &= f(\mathbf{x}) \\ z &= g(\mathbf{y})\end{aligned}$$

Then, the chain rule of differentiation states that:

$$\nabla_{\mathbf{x}} z = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}} z$$

The matrix $\frac{\partial \mathbf{y}}{\partial \mathbf{x}}$ is known as the **Jacobian** of the function f where $\mathbf{y} = f(\mathbf{x})$.

Chain rule in multidimensions

$$\nabla_{\mathbf{x}^Z} = \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right)^T \nabla_{\mathbf{y}^Z}$$

where

$$\nabla_{\mathbf{x}^Z} = \frac{\partial z}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial z}{\partial x_1} \\ \frac{\partial z}{\partial x_2} \\ \vdots \\ \frac{\partial z}{\partial x_n} \end{pmatrix} \text{ and } \nabla_{\mathbf{y}^Z} = \frac{\partial z}{\partial \mathbf{y}} = \begin{pmatrix} \frac{\partial z}{\partial y_1} \\ \frac{\partial z}{\partial y_2} \\ \vdots \\ \frac{\partial z}{\partial y_K} \end{pmatrix},$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_n} \\ \vdots & \vdots & \dots & \vdots \\ \frac{\partial y_K}{\partial x_1} & \frac{\partial y_K}{\partial x_2} & \dots & \frac{\partial y_K}{\partial x_n} \end{pmatrix}$$

Note that differentiation of a scalar by a vector results in a vector and differentiation of a vector by a vector results in a matrix.

Example 1

Let $\mathbf{x} = (x_1, x_2, x_3) \in \mathbf{R}^3$, $\mathbf{y} = (y_1, y_2) \in \mathbf{R}^2$, and $\mathbf{y} = f(\mathbf{x})$ where f is given by

$$\begin{aligned}y_1 &= 5 - 2x_1 + 3x_3 \\y_2 &= x_1 + 5x_2^2 + x_3^3 - 1\end{aligned}$$

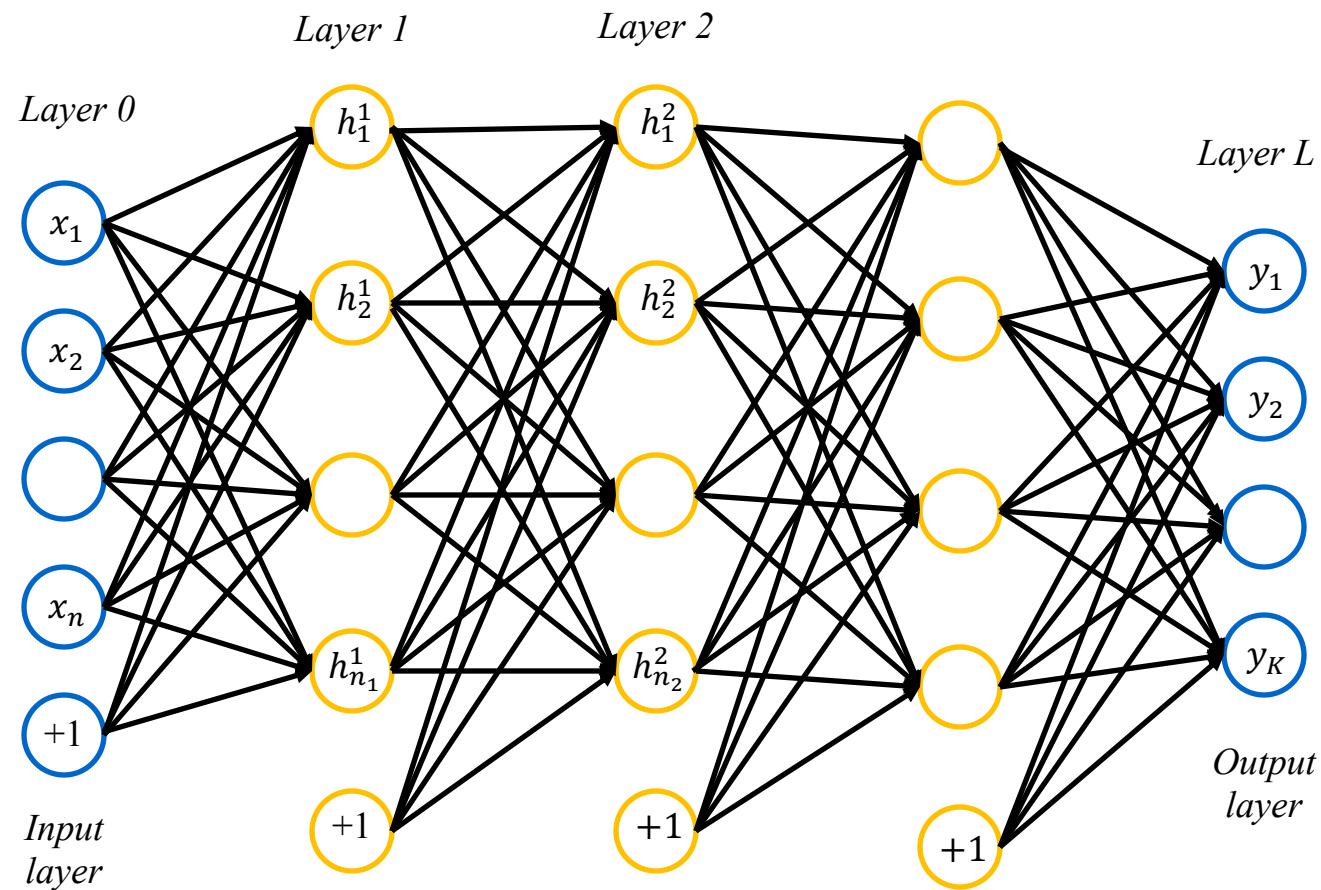
Find the Jacobian of f .

$$\begin{aligned}\frac{\partial y_1}{\partial x_1} &= -2, & \frac{\partial y_1}{\partial x_2} &= 0, & \frac{\partial y_1}{\partial x_3} &= 3 \\ \frac{\partial y_2}{\partial x_1} &= 1, & \frac{\partial y_2}{\partial x_2} &= 10x_2, & \frac{\partial y_2}{\partial x_3} &= 3x_3^2\end{aligned}$$

Jacobian:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \end{pmatrix} = \begin{pmatrix} -2 & 0 & 3 \\ 1 & 10x_2 & 3x_3^2 \end{pmatrix}$$

Deep feedforward networks (FFN)



Deep feedforward networks (FFN)

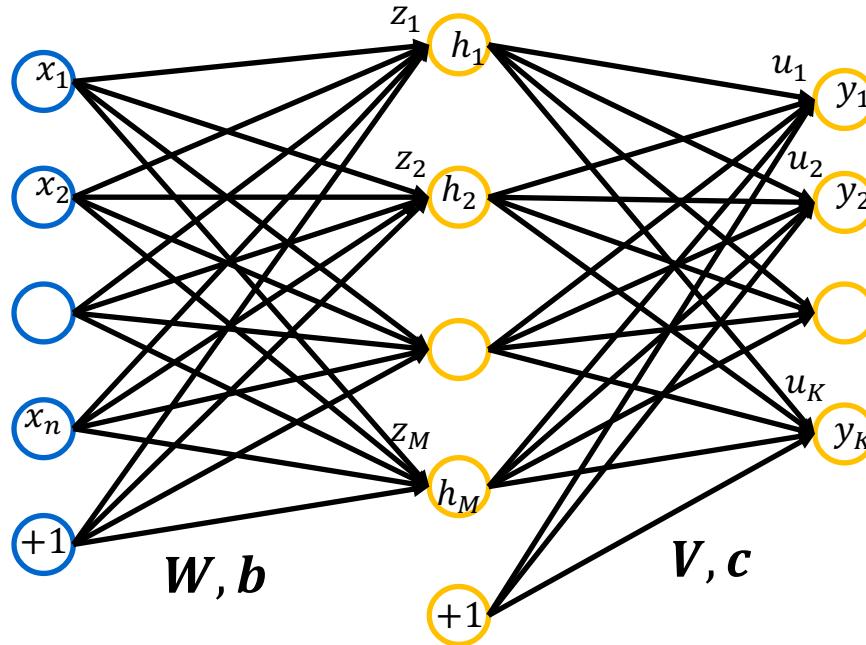
Feedforward networks (FFN) consists of several layers of neurons where activations propagate from input layer to output layer. The layers between the input and output layers are referred to as **hidden layers**.

The number of hidden layers is referred to as the **depth** of the feedforward network. When a network has many hidden layers of neurons, feedforward networks are referred to as **deep neural networks (DNN)**. Learning in deep neural networks is referred to as **deep learning**.

The hidden layers are usually composed of perceptrons or ReLU units and the output layer is

- A linear neuron layer for regression
- A softmax layer for classification

Three-layer FFN



$$\text{Input } \mathbf{x} = (x_1 \quad x_2 \quad \cdots \quad x_n)^T$$

$$\text{Hidden-layer activation } \mathbf{h} = (h_1 \quad h_2 \quad \cdots \quad h_M)^T$$

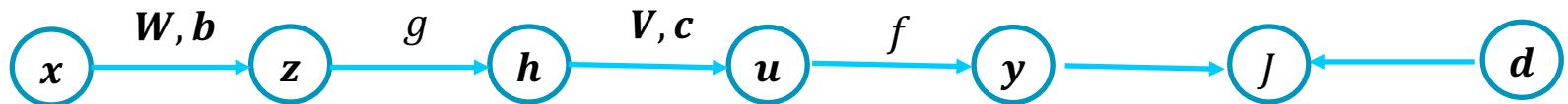
$$\text{Output } \mathbf{y} = (y_1 \quad y_2 \quad \cdots \quad y_K)^T$$

\mathbf{W}, \mathbf{b} – weight and bias of the hidden layer

\mathbf{V}, \mathbf{c} – weight and bias of the output layer

Forward propagation of activations: single pattern

Computational graph of 3-layer FFN for an input pattern (x, d) :



Synaptic input \mathbf{z} to hidden layer:

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

Output \mathbf{h} of the hidden layer:

$$\mathbf{h} = g(\mathbf{z})$$

Synaptic input \mathbf{u} to output layer:

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$$

Output \mathbf{y} of the output layer:

$$\mathbf{y} = f(\mathbf{u})$$

Derivatives

$$\mathbf{y} = f(\mathbf{u})$$

Considering k th neuron:

$$y_k = f(u_k)$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{u}} = \begin{pmatrix} \frac{\partial y_1}{\partial u_1} & \frac{\partial y_1}{\partial u_2} & \dots & \frac{\partial y_1}{\partial u_K} \\ \frac{\partial y_2}{\partial u_1} & \frac{\partial y_2}{\partial u_2} & \dots & \frac{\partial y_2}{\partial u_K} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y_K}{\partial u_1} & \frac{\partial y_K}{\partial u_2} & \dots & \frac{\partial y_K}{\partial u_K} \end{pmatrix} = \begin{pmatrix} f'(u_1) & 0 & \dots & 0 \\ 0 & f'(u_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & f'(u_K) \end{pmatrix} = \text{diag}(f'(\mathbf{u}))$$

That is,

$$\frac{\partial \mathbf{y}}{\partial \mathbf{u}} = \text{diag}(f'(\mathbf{u}))$$

where $\text{diag}(f'(\mathbf{u}))$ is a diagonal matrix composed of derivatives corresponding to individual components of \mathbf{u} in the diagonal.

Derivatives

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$$

Consider k th neuron:

weight vector $\mathbf{v}_k = (v_{k1} \quad v_{k2} \quad \cdots \quad v_{kM})^T$ and bias c_k .

The synaptic input u_k due to \mathbf{h} is given by

$$u_k = v_{k1}h_1 + v_{k2}h_2 + \cdots + v_{kM}h_M + c_k$$
$$\frac{\partial u_k}{\partial h_j} = v_{kj}$$

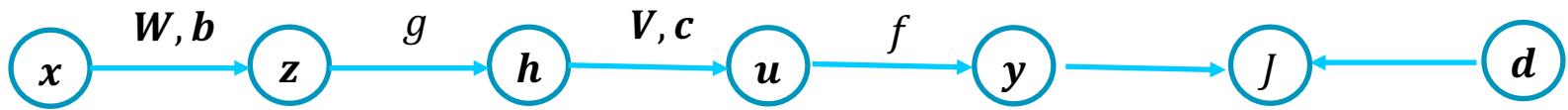
Therefore

$$\frac{\partial \mathbf{u}}{\partial \mathbf{h}} = \begin{pmatrix} \frac{\partial u_1}{\partial h_1} & \frac{\partial u_1}{\partial h_2} & \cdots & \frac{\partial u_1}{\partial h_M} \\ \frac{\partial u_2}{\partial h_1} & \frac{\partial u_2}{\partial h_2} & \cdots & \frac{\partial u_2}{\partial h_M} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial u_K}{\partial h_1} & \frac{\partial u_K}{\partial h_2} & \cdots & \frac{\partial u_K}{\partial h_K} \end{pmatrix} = \begin{pmatrix} v_{11} & v_{12} & \cdots & v_{1M} \\ v_{21} & v_{22} & \cdots & v_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ v_{K1} & v_{K2} & \cdots & v_{KM} \end{pmatrix} = \mathbf{V}^T$$

That is,

$$\frac{\partial \mathbf{u}}{\partial \mathbf{h}} = \mathbf{V}^T$$

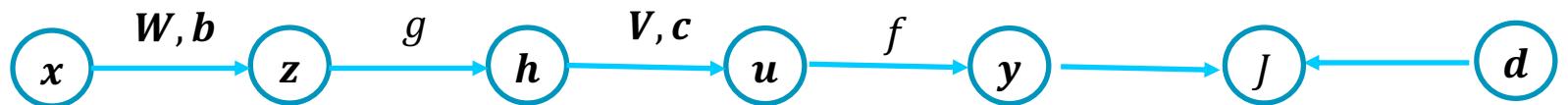
Derivatives



$$\frac{\partial \mathbf{y}}{\partial \mathbf{u}} = \text{diag}(f'(\mathbf{u}))$$

$$\frac{\partial \mathbf{u}}{\partial \mathbf{h}} = \mathbf{V}^T$$

Back-propagation of gradients: single pattern



Considering output layer,

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) & \text{for a linear layer} \\ -(1(\mathbf{k} = \mathbf{d}) - f(\mathbf{u})) & \text{for a softmax layer} \end{cases}$$

From chain rule of differentiation,

$$\nabla_{\mathbf{h}} J = \left(\frac{\partial \mathbf{u}}{\partial \mathbf{h}} \right)^T \nabla_{\mathbf{u}} J = \mathbf{V} \nabla_{\mathbf{u}} J$$

$$\nabla_{\mathbf{z}} J = \left(\frac{\partial \mathbf{h}}{\partial \mathbf{z}} \right)^T \nabla_{\mathbf{h}} J = \text{diag}(g'(\mathbf{z})) \mathbf{V} \nabla_{\mathbf{u}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z}) \quad (\text{A})$$

Back-propagation of gradients: single pattern

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{w}} J \cdot g'(\mathbf{z})$$

Note that gradients at the hidden layer can be derived from the gradients at the output layer.

In other words, the gradients are multiplied by \mathbf{V} and back-propagated from the output layer to the hidden layer. Note that hidden-layer activations are multiplied by \mathbf{V}^T in the forward propagation.

This is known as **back-propagation** (backprop) of gradients.

Proof

For a vector \mathbf{x} :

$$diag(f'(\mathbf{u}))\mathbf{x} = \begin{pmatrix} f'(u_1) & 0 & \cdots & 0 \\ 0 & f'(u_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & f'(u_K) \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{pmatrix} = \begin{pmatrix} f'(u_1)x_1 \\ f'(u_2)x_2 \\ \vdots \\ f'(u_K)x_K \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{pmatrix} \cdot \begin{pmatrix} f'(u_1) \\ f'(u_2) \\ \vdots \\ f'(u_K) \end{pmatrix}$$

That is:

$$diag(f'(\mathbf{u}))\mathbf{x} = \mathbf{x} \cdot f'(\mathbf{u}) = f'(\mathbf{u}) \cdot \mathbf{x}$$

SGD of three-layer FFN

Output layer:

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) & \text{for linear layer} \\ -(1(\mathbf{k} = d) - f(\mathbf{u})) & \text{for softmax layer} \end{cases}$$

Weights of output layer are updated as

$$\begin{aligned}\mathbf{V} &\leftarrow \mathbf{V} - \alpha \mathbf{h} (\nabla_{\mathbf{u}} J)^T \\ \mathbf{c} &\leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J\end{aligned}$$

Hidden layer:

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z})$$

Weights of hidden layer are updated as

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{z}} J)^T \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J\end{aligned}$$

SGD for a three-layer FFN

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

Set learning parameter α

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

Repeat until convergence:

For every pattern (\mathbf{x}, \mathbf{d}) :

$$\mathbf{z} = \mathbf{W}^T \mathbf{x} + \mathbf{b}$$

$$\mathbf{h} = g(\mathbf{z})$$

$$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$$

$$\mathbf{y} = f(\mathbf{u})$$

Forward propagation

$$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) \\ -(1(\mathbf{k} = \mathbf{d}) - f(\mathbf{u})) \end{cases}$$

$$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z})$$

Backward propagation

$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{h} (\nabla_{\mathbf{u}} J)^T$$

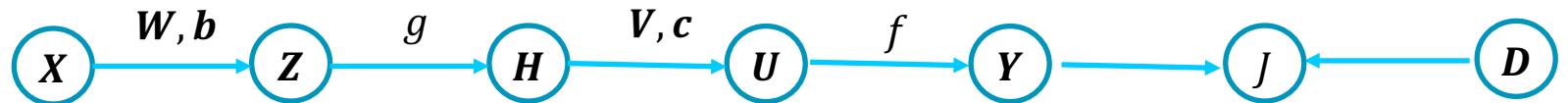
$$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{z}} J)^T$$

$$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J$$

Forward propagation of activations: batch of patterns

Computational graph of 3-layer FFN for an batch of patterns (X, D) :



Synaptic input Z to hidden layer:

$$Z = XW + B$$

Output H of the hidden layer:

$$H = g(Z)$$

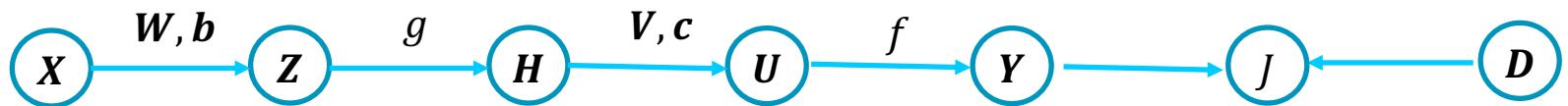
Synaptic input U to output layer:

$$U = HV + C$$

Output Y of the output layer:

$$Y = f(U)$$

Back-propagation of gradients: batch of patterns



$$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) & \text{for linear layer} \\ -(\mathbf{K} - f(\mathbf{U})) & \text{for softmax layer} \end{cases}$$

$$\nabla_{\mathbf{H}} J = \begin{pmatrix} (\nabla_{\mathbf{h}_1} J)^T \\ (\nabla_{\mathbf{h}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{h}_P} J)^T \end{pmatrix} = \begin{pmatrix} (\mathbf{V} \nabla_{\mathbf{u}_1} J)^T \\ (\mathbf{V} \nabla_{\mathbf{u}_2} J)^T \\ \vdots \\ (\mathbf{V} \nabla_{\mathbf{u}_P} J)^T \end{pmatrix} = \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T \mathbf{V}^T \\ (\nabla_{\mathbf{u}_2} J)^T \mathbf{V}^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T \mathbf{V}^T \end{pmatrix} = (\nabla_{\mathbf{U}} J) \mathbf{V}^T$$

$$(XY)^T = Y^T X^T$$

Back-propagation of gradients: batch of patterns

$$\begin{aligned}\nabla_{\mathbf{Z}} J &= \begin{pmatrix} (\nabla_{\mathbf{z}_1} J)^T \\ (\nabla_{\mathbf{z}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{z}_P} J)^T \end{pmatrix} = \begin{pmatrix} \left(V \nabla_{\mathbf{u}_1} J \cdot g'(\mathbf{z}_1) \right)^T \\ \left(V \nabla_{\mathbf{u}_2} J \cdot g'(\mathbf{z}_2) \right)^T \\ \vdots \\ \left(V \nabla_{\mathbf{u}_P} J \cdot g'(\mathbf{z}_P) \right)^T \end{pmatrix} && \text{from (A)} \\ &= \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T V^T \cdot (g'(\mathbf{z}_1))^T \\ (\nabla_{\mathbf{u}_2} J)^T V^T \cdot (g'(\mathbf{z}_2))^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T V^T \cdot (g'(\mathbf{z}_P))^T \end{pmatrix} && (XY)^T = Y^T X^T \\ &= \begin{pmatrix} (\nabla_{\mathbf{u}_1} J)^T \\ (\nabla_{\mathbf{u}_2} J)^T \\ \vdots \\ (\nabla_{\mathbf{u}_P} J)^T \end{pmatrix} V^T \cdot \begin{pmatrix} (g'(\mathbf{z}_1))^T \\ (g'(\mathbf{z}_2))^T \\ \vdots \\ (g'(\mathbf{z}_P))^T \end{pmatrix} \\ &= (\nabla_{\mathbf{u}} J) V^T \cdot g'(\mathbf{Z})\end{aligned}$$

Back-propagation of gradients: batch of patterns

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

GD of three-layer FFN

Output layer:

$$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(\mathbf{K} - f(\mathbf{U})) \end{cases}$$

Learning equations:

$$\begin{aligned}\mathbf{V} &\leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J \\ \mathbf{c} &\leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P\end{aligned}$$

Hidden layer:

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

Learning equations

$$\begin{aligned}\mathbf{W} &\leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J \\ \mathbf{b} &\leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P\end{aligned}$$

GD for a three-layer FFN

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

Set learning parameter α

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

$$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(\mathbf{K} - f(\mathbf{U})) \end{cases}$$

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

Backward propagation

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

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

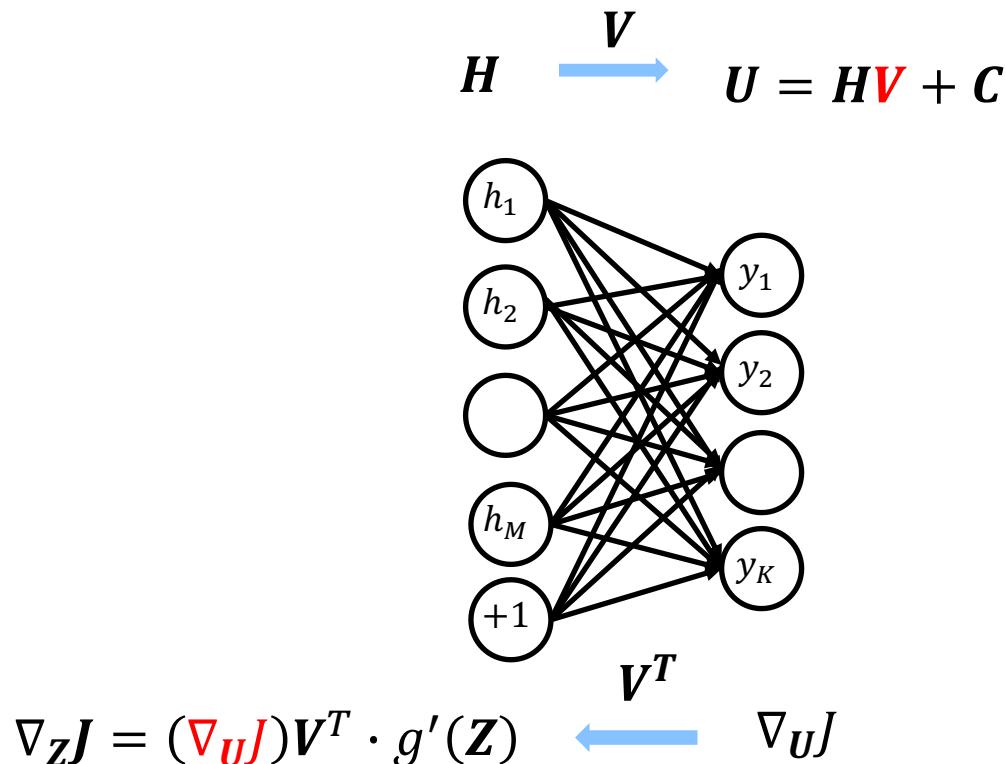
$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J$$

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

Learning in three-layer FFN

GD	SGD
(X, \mathbf{D})	(x, \mathbf{d})
$\mathbf{Z} = X\mathbf{W} + \mathbf{B}$	$\mathbf{z} = \mathbf{W}^T x + \mathbf{b}$
$\mathbf{H} = g(\mathbf{Z})$	$\mathbf{h} = g(\mathbf{z})$
$\mathbf{U} = \mathbf{H}\mathbf{V} + \mathbf{C}$	$\mathbf{u} = \mathbf{V}^T \mathbf{h} + \mathbf{c}$
$\mathbf{Y} = f(\mathbf{U})$	$\mathbf{y} = f(\mathbf{u})$
$\nabla_{\mathbf{U}} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(\mathbf{K} - f(\mathbf{U})) \end{cases}$	$\nabla_{\mathbf{u}} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) \\ (1(\mathbf{k} = d) - f(\mathbf{u})) \end{cases}$
$\nabla_{\mathbf{Z}} J = (\nabla_{\mathbf{U}} J) \mathbf{V}^T \cdot g'(\mathbf{Z})$	$\nabla_{\mathbf{z}} J = \mathbf{V} \nabla_{\mathbf{u}} J \cdot g'(\mathbf{z})$
$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J$	$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{x} (\nabla_{\mathbf{z}} J)^T$
$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P$	$\mathbf{b} \leftarrow \mathbf{b} - \alpha \nabla_{\mathbf{z}} J$
$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J$	$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{h} (\nabla_{\mathbf{u}} J)^T$
$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P$	$\mathbf{c} \leftarrow \mathbf{c} - \alpha \nabla_{\mathbf{u}} J$

Back-propagation



The error can be seen as propagating to the hidden layer from the output layer and therefore learning in feedforward networks is known as the *back-propagation* algorithm

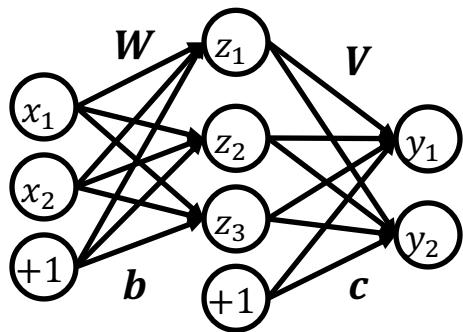
Example 2

Design a three-layer FFN, using gradient descent to perform the following mapping. Use a learning factor = 0.05 and three perceptrons in the hidden-layer.

Inputs $x = (x_1, x_2)$	Outputs $y = (y_1, y_2)$
(0.77, 0.02)	(0.44, -0.42)
(0.63, 0.75)	(0.84, 0.43)
(0.50, 0.22)	(0.09, -0.72)
(0.20, 0.76)	(-0.25, 0.35)
(0.17, 0.09)	(-0.12, -0.13)
(0.69, 0.95)	(0.24, 0.03)
(0.00, 0.51)	(0.30, 0.20)
(0.81, 0.61)	(0.61, 0.04)

Example 2

$$X = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \\ 0.81 & 0.61 \end{pmatrix} \text{ and } D = \begin{pmatrix} 0.44 & -0.42 \\ 0.84 & 0.43 \\ 0.09 & -0.72 \\ -0.25 & 0.35 \\ -0.12 & -0.13 \\ 0.24 & 0.03 \\ 0.30 & 0.20 \\ 0.61 & 0.04 \end{pmatrix}$$



Output layer is a linear neuron layer

Hidden layer is a sigmoidal layer

Initialized:

$$W = \begin{pmatrix} -3.97 & 1.10 & 0.42 \\ 2.79 & -2.64 & 3.13 \end{pmatrix}, b = \begin{pmatrix} 0.0 \\ 0.0 \\ 0.0 \end{pmatrix}, V = \begin{pmatrix} 3.58 & -1.58 \\ -3.58 & -1.75 \\ -3.38 & 2.88 \end{pmatrix}, c = \begin{pmatrix} 0.0 \\ 0.0 \end{pmatrix}$$

Example 2

Iteration 1:

$$\mathbf{Z} = \mathbf{XW} + \mathbf{B} = \begin{pmatrix} 0.77 & 0.02 \\ 0.63 & 0.75 \\ 0.50 & 0.22 \\ 0.20 & 0.76 \\ 0.17 & 0.09 \\ 0.69 & 0.95 \\ 0.00 & 0.51 \\ 0.81 & 0.61 \end{pmatrix} \begin{pmatrix} -3.97 & 1.10 & 0.42 \\ 2.79 & -2.64 & 3.13 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} -3.00 & 0.8 & 0.39 \\ -0.42 & -1.27 & 2.61 \\ -1.35 & -0.04 & 0.91 \\ 1.34 & -1.79 & 2.46 \\ -0.42 & -0.05 & 0.35 \\ -0.05 & -1.76 & 3.27 \\ 1.42 & -1.34 & 1.60 \\ -1.51 & -0.72 & 2.25 \end{pmatrix}$$

$$\mathbf{H} = g(\mathbf{Z}) = \frac{1}{1 + e^{-\mathbf{z}}} = \begin{pmatrix} 0.05 & 0.69 & 0.60 \\ 0.40 & 0.22 & 0.93 \\ 0.21 & 0.49 & 0.71 \\ 0.79 & 0.14 & 0.92 \\ 0.40 & 0.49 & 0.59 \\ 0.49 & 0.15 & 0.96 \\ 0.80 & 0.21 & 0.83 \\ 0.18 & 0.33 & 0.91 \end{pmatrix}$$

Example 2

$$Y = HV + C = \begin{pmatrix} 0.05 & 0.69 & 0.60 \\ 0.40 & 0.22 & 0.93 \\ 0.21 & 0.49 & 0.71 \\ 0.79 & 0.14 & 0.92 \\ 0.40 & 0.49 & 0.59 \\ 0.49 & 0.15 & 0.96 \\ 0.80 & 0.21 & 0.83 \\ 0.18 & 0.33 & 0.91 \end{pmatrix} \begin{pmatrix} 2.38 & -4.2 \\ 1.17 & 2.18 \\ -0.01 & -2.41 \end{pmatrix} + \begin{pmatrix} 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \\ 0.0 & 0.0 \end{pmatrix} = \begin{pmatrix} -4.32 & 0.44 \\ -2.52 & 1.67 \\ -3.43 & 0.87 \\ -0.79 & 1.15 \\ -2.32 & 0.21 \\ -2.04 & 1.75 \\ -0.67 & 0.76 \\ -3.59 & 1.75 \end{pmatrix}$$

$$\nabla_U J = -(\mathbf{D} - Y) = - \left(\begin{pmatrix} 0.44 & -0.42 \\ 0.84 & 0.43 \\ 0.09 & -0.72 \\ -0.25 & 0.35 \\ -0.12 & -0.13 \\ 0.24 & 0.03 \\ 0.30 & 0.20 \\ 0.61 & 0.04 \end{pmatrix} - \begin{pmatrix} -4.32 & 0.44 \\ -2.52 & 1.67 \\ -3.43 & 0.87 \\ -0.79 & 1.15 \\ -2.32 & 0.21 \\ -2.04 & 1.75 \\ -0.67 & 0.76 \\ -3.59 & 1.75 \end{pmatrix} \right) = \begin{pmatrix} -4.76 & 0.85 \\ -3.35 & 1.24 \\ -3.52 & 1.59 \\ -0.54 & 0.80 \\ -2.20 & 0.34 \\ -2.28 & 1.72 \\ -0.98 & 0.56 \\ -4.20 & 1.70 \end{pmatrix}$$

Example 2

$$g'(\mathbf{Z}) = \mathbf{H} \cdot (\mathbf{1} - \mathbf{H}) = \begin{pmatrix} 0.04 & 0.21 & 0.24 \\ 0.24 & 0.17 & 0.06 \\ 0.16 & 0.25 & 0.20 \\ 0.16 & 0.12 & 0.07 \\ 0.24 & 0.25 & 0.24 \\ 0.25 & 0.13 & 0.04 \\ 0.16 & 0.16 & 0.14 \\ 0.15 & 0.22 & 0.09 \end{pmatrix}$$

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

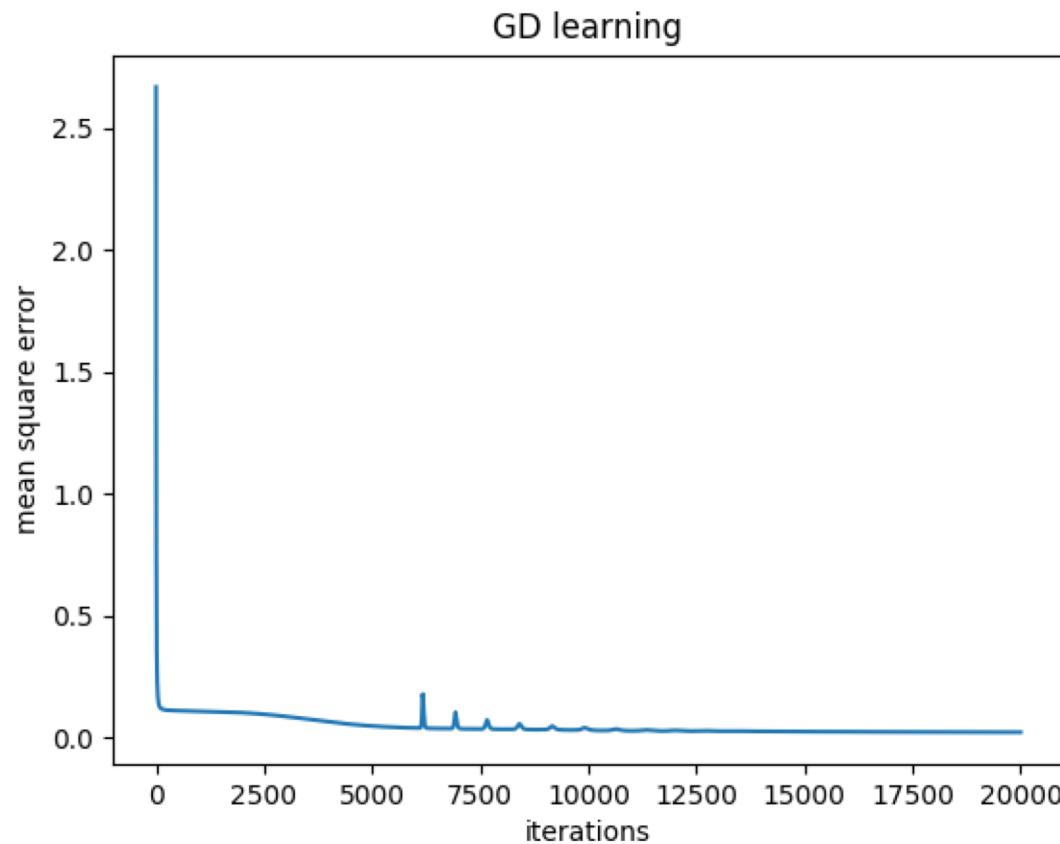
$$= \begin{pmatrix} -4.76 & 0.85 \\ -3.35 & 1.24 \\ -3.52 & 1.59 \\ -0.54 & 0.80 \\ -2.20 & 0.34 \\ -2.28 & 1.72 \\ -0.98 & 0.56 \\ -4.20 & 1.70 \end{pmatrix} \begin{pmatrix} 3.58 & -1.58 \\ -3.58 & -1.75 \\ -3.38 & 2.88 \end{pmatrix}^T \cdot \begin{pmatrix} 0.04 & 0.21 & 0.24 \\ 0.24 & 0.17 & 0.06 \\ 0.16 & 0.25 & 0.20 \\ 0.16 & 0.12 & 0.07 \\ 0.24 & 0.25 & 0.24 \\ 0.25 & 0.13 & 0.04 \\ 0.16 & 0.16 & 0.14 \\ 0.15 & 0.22 & 0.09 \end{pmatrix}$$

Example 2

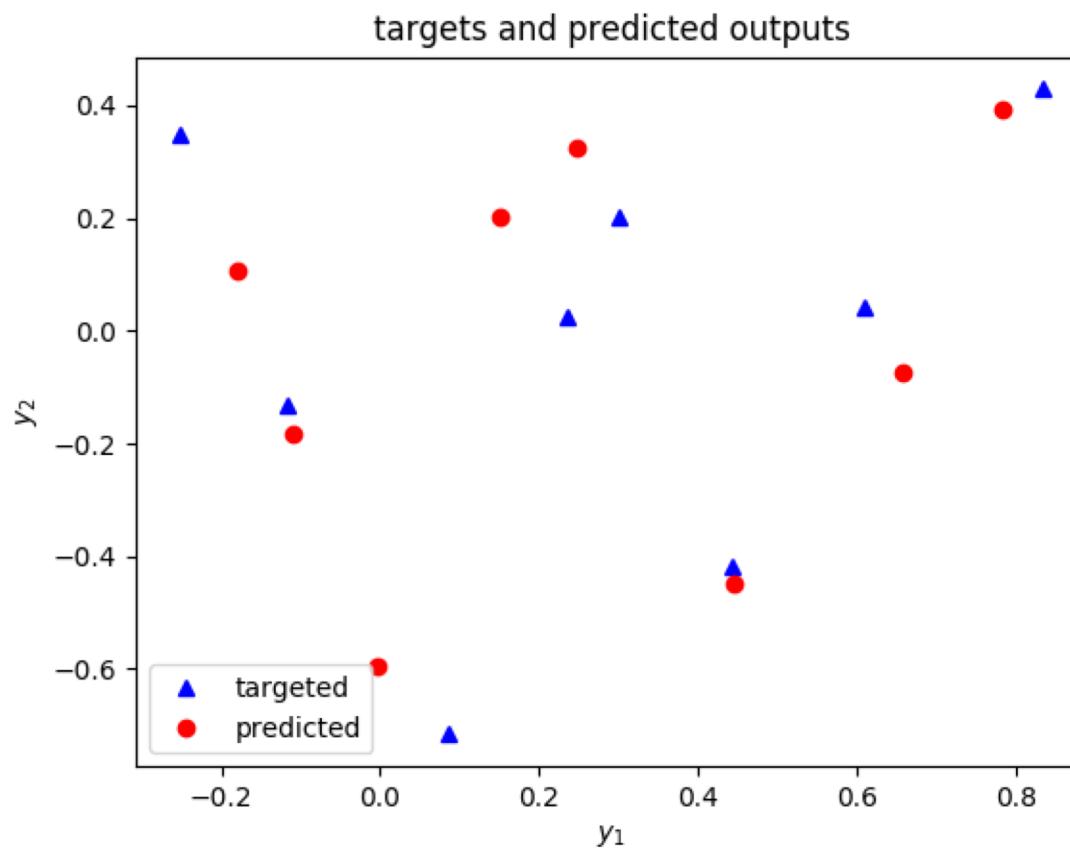
$$\mathbf{V} \leftarrow \mathbf{V} - \alpha \mathbf{H}^T \nabla_{\mathbf{U}} J = \begin{pmatrix} 3.89 & -1.74 \\ -3.14 & -1.89 \\ -2.53 & 2.51 \end{pmatrix}$$
$$\mathbf{c} \leftarrow \mathbf{c} - \alpha (\nabla_{\mathbf{U}} J)^T \mathbf{1}_P = \begin{pmatrix} 1.09 \\ -0.44 \end{pmatrix}$$

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \mathbf{X}^T \nabla_{\mathbf{Z}} J = \begin{pmatrix} -3.55 & 0.72 & 0.03 \\ 3.21 & -2.87 & 2.94 \\ 0.76 \end{pmatrix}$$
$$\mathbf{b} \leftarrow \mathbf{b} - \alpha (\nabla_{\mathbf{Z}} J)^T \mathbf{1}_P = \begin{pmatrix} -0.66 \\ -0.70 \end{pmatrix}$$

Example 2



Example 2



The universal approximation theorem

Kolmogorov Theorem:

Given any function $\phi: I^n \rightarrow \mathbf{R}^K$ where I is the closed unit interval $[0,1]$, an arbitrary function ϕ can be implemented exactly by a multilayer perceptron (MLP) network with one hidden layer, n input nodes, $2n+1$ hidden layer neurons and m output layer neurons.

The theorem does not refer to the algorithms or the amount of training data needed for universal approximation. Recent researches find that due to the limited amount of training data and limitations of learning algorithms, deeper neural networks and a large number of training data are required in order to learn an arbitrary complex function.

Normalization of inputs

If inputs have similar variations, better approximation of inputs is achieved. Mainly, there are two approaches to normalization of inputs.

Suppose i th input $x_i \in [x_{i,min}, x_{i,max}]$ and has a mean μ_i and a standard deviation σ_i .

If \tilde{x}_i denotes the normalized input.

1. Scale the inputs such that $x_i \in [0, 1]$:

$$\tilde{x}_i = \frac{x_i - x_{i,min}}{x_{i,max} - x_{i,min}}$$

2. Normalize the input to have standard normal distributions $x_i \sim N(0,1)$:

$$\tilde{x}_i = \frac{x_i - \mu_i}{\sigma_i}$$

Output-layer activation function

If a sigmoid is used as the activation function of output neurons, the output neuron activation function should be scaled and translated to match the space of output.

If output $y \in [y_{min}, y_{max}]$, the output-layer activation function:

$$f(u) = \frac{a}{1 + e^{-u}} + b$$

where $a = y_{max} - y_{min}$

$$b = y_{min}$$

If a linear output layer is used, it is useful to normalize the outputs to have zero mean and unit standard deviations.

Boston housing dataset

<https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html>

14 variables

CRIM - per capita crime rate by town

ZN - proportion of residential land zoned for lots over 25,000 sq.ft.

INDUS - proportion of non-retail business acres per town.

CHAS - Charles River dummy variable (1 if tract bounds river; 0 otherwise)

NOX - nitric oxides concentration (parts per 10 million)

RM - average number of rooms per dwelling

AGE - proportion of owner-occupied units built prior to 1940

DIS - weighted distances to five Boston employment centres

RAD - index of accessibility to radial highways

TAX - full-value property-tax rate per \$10,000

PTRATIO - pupil-teacher ratio by town

B - $1000(Bk - 0.63)^2$ where Bk is the proportion of blacks by town

LSTAT - % lower status of the population

MEDV - Median value of owner-occupied homes in \$1000's

506 data points

Try to predict MEDV by using other 13 variables

Example 3: Predicting housing prices in Boston

Thirteen input variables

One output variable

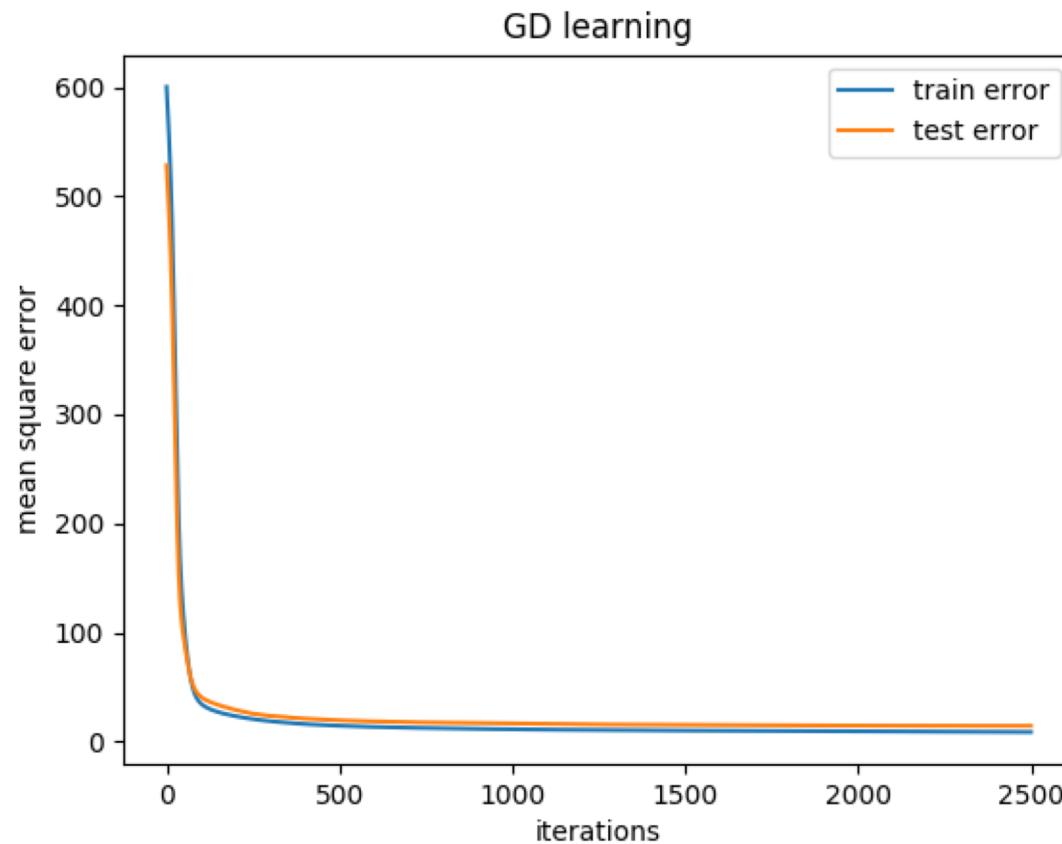
We use FFN with one hidden layer

Hidden number of neurons = 10

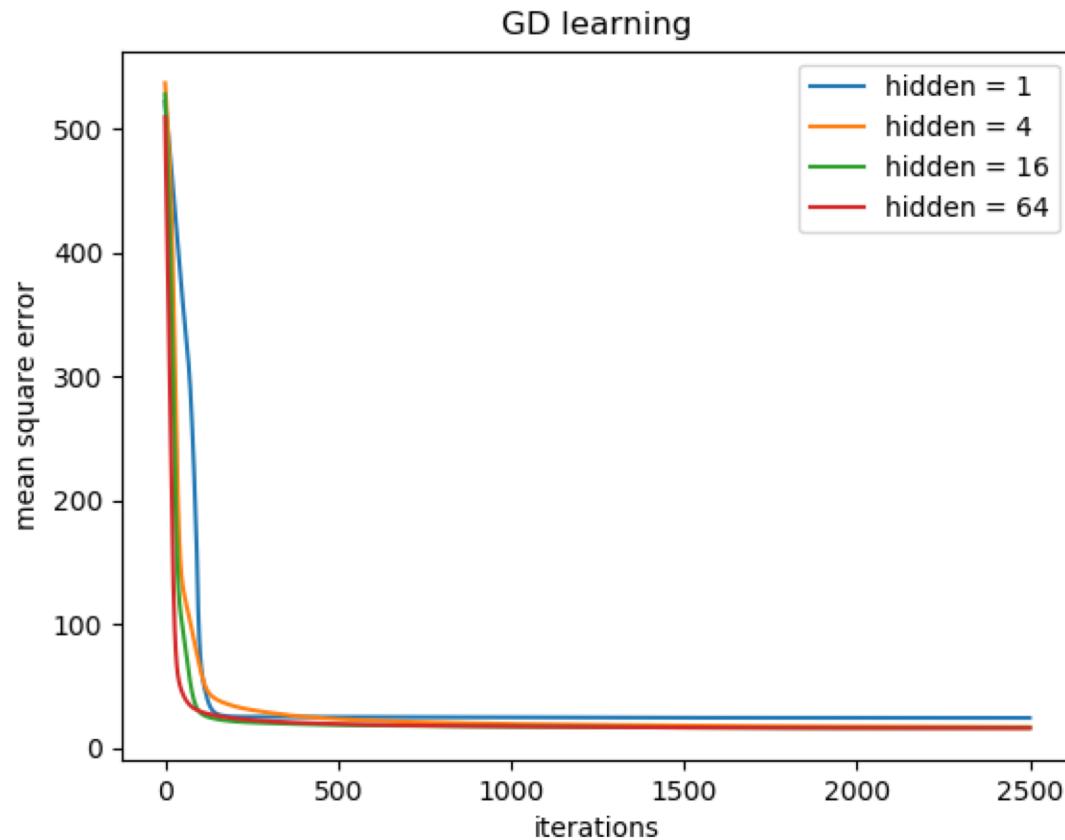
Network size: [13, 10, 1]

We divide the dataset into train and test partitions: [0.8, 0.2]

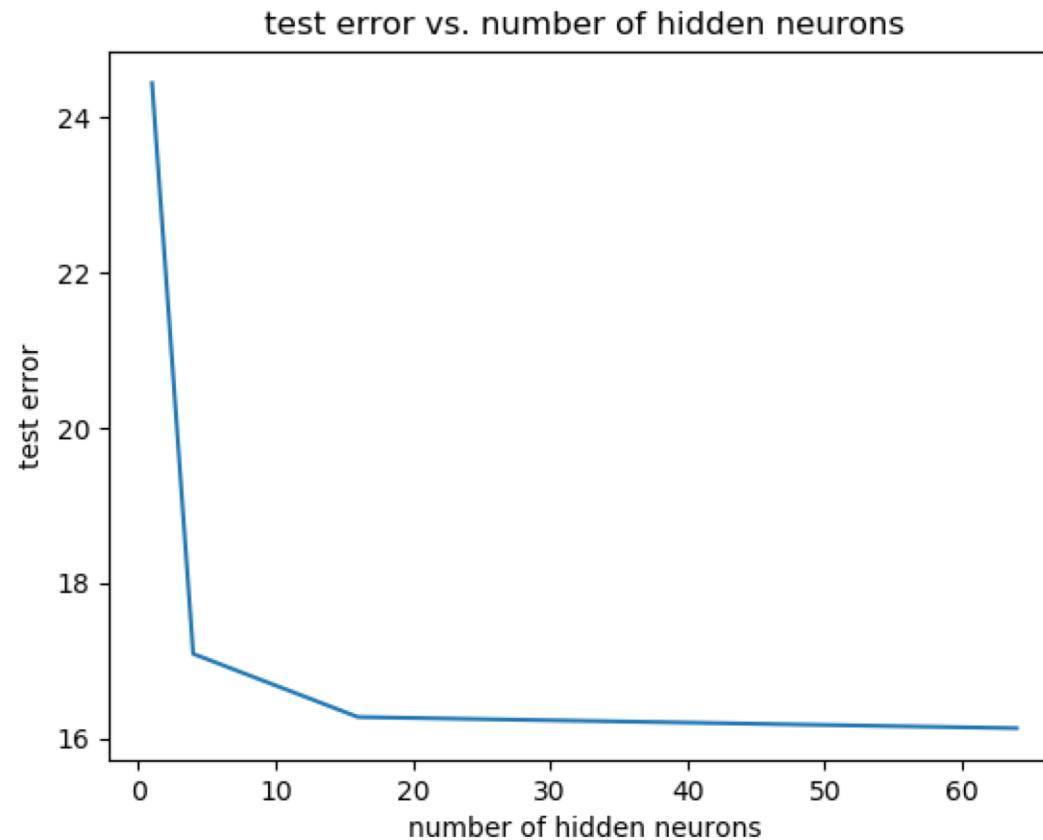
Example 3



Example 3: Number of hidden neurons



Example 3

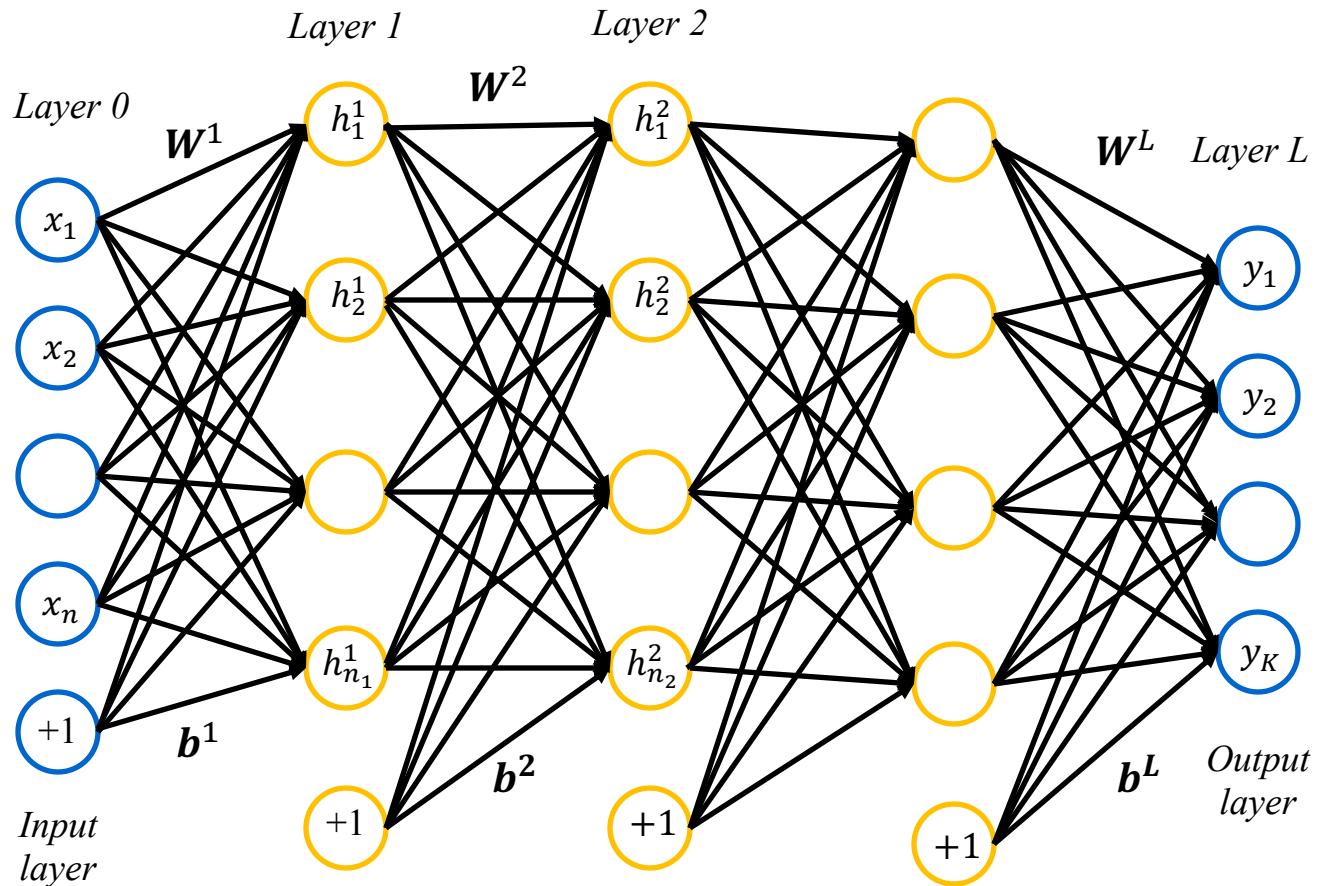


Width: the number of hidden neurons

The number of parameters of the network increases with the number of hidden units. Therefore, the network attempts to remember the training patterns with increasing number of parameters. In other words, the network aims at minimizing the training error at the expense of its generalization ability on unseen data.

As the number of hidden units increases, the test error decreases initially but tends to increase at some point. The optimal number of hidden units is often determined empirically (that is, by trial and error).

Deep neural networks (DNN)



Depth = $L + 1$

DNN

Input layer $l = 0$:

Width = n

Input x, X

Hidden layers $l = 1, 2, \dots, L - 1$

Width: n_l

Weight matrix \mathbf{W}^l , bias vector \mathbf{b}^l

Synaptic input $\mathbf{u}^l, \mathbf{U}^l$

Activation function f^l

Output $\mathbf{h}^l, \mathbf{H}^l$

Output layer $l = L$

Width: K

Synaptic input $\mathbf{u}^L, \mathbf{U}^L$

Activation function f^L

Output y, Y

Desired output d, D

Forward propagation of activation in DNN: single pattern

Input (x, d)

$$\mathbf{u}^1 = \mathbf{W}^{1T} x + \mathbf{b}^1$$

For layers $l = 1, 2, \dots, L - 1$:

$$\mathbf{h}^l = f^l(\mathbf{u}^l)$$

$$\mathbf{u}^{l+1} = \mathbf{W}^{l+1T} \mathbf{h}^l + \mathbf{b}^{l+1}$$

$$\mathbf{y} = f^L(\mathbf{u}^L)$$

Back-propagation of gradients in DNN: single pattern

if $l = L$:

$$\nabla_{\mathbf{u}^l} J = \begin{cases} -(d - y) & \text{for linear layer} \\ -(1(k = d) - f^L(\mathbf{u}^L)) & \text{for softmax layer} \end{cases}$$

else:

$$\nabla_{\mathbf{u}^l} J = \mathbf{W}^{l+1} (\nabla_{\mathbf{u}^{l+1}} J) \cdot f^{l'}(\mathbf{u}^l) \quad \text{from (A)}$$

$$\nabla_{\mathbf{W}^l} J = \mathbf{h}^{l-1} (\nabla_{\mathbf{u}^l} J)^T$$

$$\nabla_{\mathbf{b}^l} J = \nabla_{\mathbf{u}^l} J$$

Stochastic backpropagation for DNN

For layers $l = L, L - 1, \dots, 2$:

If $l = L$:

$$\nabla_{\mathbf{u}^l} J = \begin{cases} -(\mathbf{d} - \mathbf{y}) \\ -(1(\mathbf{k} = d) - f^L(\mathbf{u}^L)) \end{cases}$$

Else:

$$\nabla_{\mathbf{u}^l} J = (\mathbf{W}^{l+1} \nabla_{\mathbf{u}^{l+1}} J) \cdot f^{l'}(\mathbf{u}^l)$$

$$\nabla_{\mathbf{W}^l} J = \mathbf{h}^{l-1} (\nabla_{\mathbf{u}^l} J)^T$$

$$\nabla_{\mathbf{b}^l} J = \nabla_{\mathbf{u}^l} J$$

$$\nabla_{\mathbf{u}^1} J = (\mathbf{W}^2 \nabla_{\mathbf{u}^2} J) \cdot f^{1'}(\mathbf{u}^1)$$

$$\nabla_{\mathbf{W}^1} J = \mathbf{x} (\nabla_{\mathbf{u}^1} J)^T$$

$$\nabla_{\mathbf{b}^1} J = \nabla_{\mathbf{u}^1} J$$

Forward propagation of activation in DNN: batch of patterns

Input (\mathbf{X}, \mathbf{D})

$$\mathbf{U}^1 = \mathbf{X}\mathbf{W}^1 + \mathbf{B}^1$$

For layers $l = 1, 2, \dots, L - 1$:

$$\mathbf{H}^l = f^l(\mathbf{U}^l)$$

$$\mathbf{U}^{l+1} = \mathbf{H}^l \mathbf{W}^{l+1} + \mathbf{B}^{l+1}$$

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

Back-propagation of gradients in DNN: batch of patterns

If $l = L$:

$$\nabla_{\mathbf{U}^l} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(\mathbf{K} - f^L(\mathbf{U}^L)) \end{cases}$$

Else:

$$\nabla_{\mathbf{U}^l} J = (\nabla_{\mathbf{U}^{l+1}} J) \mathbf{W}^{l+1T} \cdot f^{l'}(\mathbf{U}^l)$$

$$\nabla_{\mathbf{W}^l} J = \mathbf{H}^{l-1T} (\nabla_{\mathbf{U}^l} J)$$

$$\nabla_{\mathbf{b}^l} J = (\nabla_{\mathbf{U}^l} J)^T \mathbf{1}_P$$

Back-propagation for DNN

For layers $l = L, L - 1, \dots, 2$:

If $l = L$:

$$\nabla_{\mathbf{U}^l} J = \begin{cases} -(\mathbf{D} - \mathbf{Y}) \\ -(\mathbf{K} - f^L(\mathbf{U}^L)) \end{cases}$$

Else:

$$\nabla_{\mathbf{U}^l} J = (\nabla_{\mathbf{U}^{l+1}} J) \mathbf{W}^{l+1T} \cdot f^{l'}(\mathbf{U}^l)$$

$$\nabla_{\mathbf{W}^l} J = \mathbf{H}^{l-1T} (\nabla_{\mathbf{U}^l} J)$$

$$\nabla_{\mathbf{b}^l} J = (\nabla_{\mathbf{U}^l} J)^T \mathbf{1}_P$$

$$\nabla_{\mathbf{U}^1} J = (\nabla_{\mathbf{U}^2} J) \mathbf{W}^{2T} \cdot f^{1'}(\mathbf{U}^1)$$

$$\nabla_{\mathbf{W}^1} J = \mathbf{X}^T (\nabla_{\mathbf{U}^1} J)$$

$$\nabla_{\mathbf{b}^1} J = (\nabla_{\mathbf{U}^1} J)^T \mathbf{1}_P$$

Example 4

Boston housing data:

<https://www.cs.toronto.edu/~delve/data/boston/bostonDetail.html>

Predicting housing price from other 13 variables

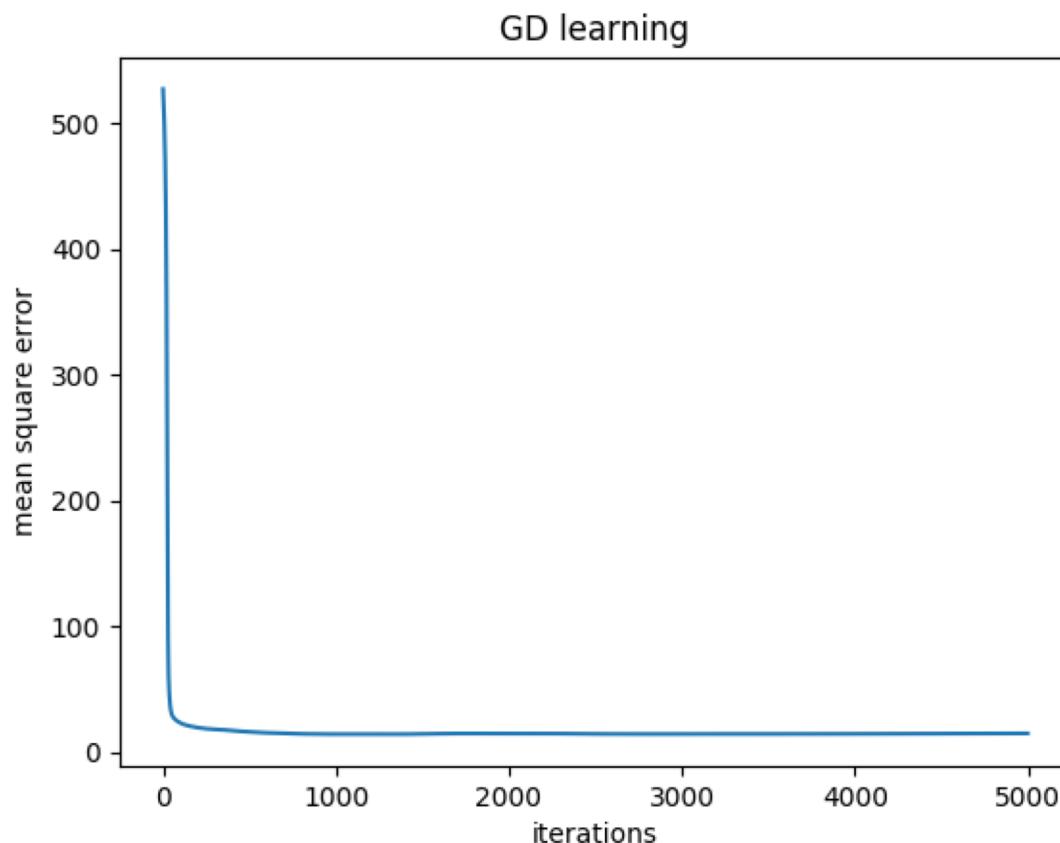
506 data points

Train and test data split: [0.8, 0.2]

DNN with two hidden layers:

[13, 5, 5, 1]

Example 4



Example 4: Varying the number of layers

Architectures:

One hidden layer: [13, 5, 1]

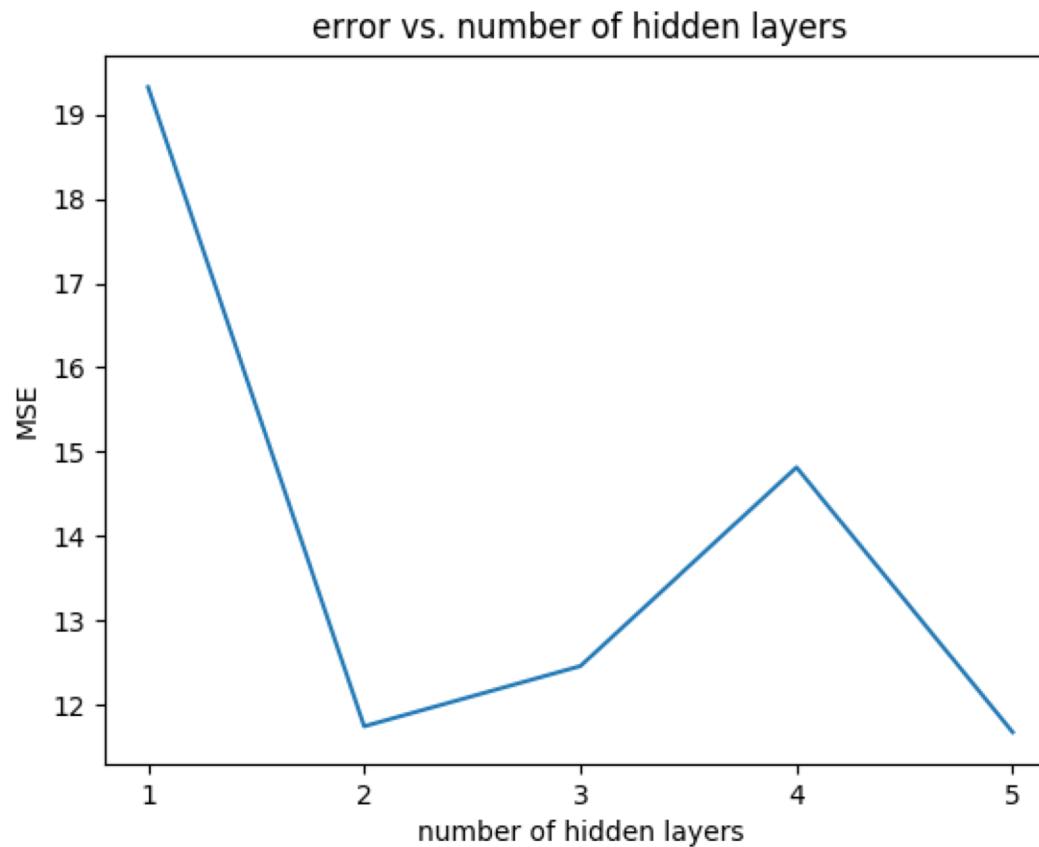
Two hidden layers: [13, 5, 5, 1]

Three hidden layers: [13, 5, 5, 5, 1]

Four hidden layers: [13, 5, 5, 5, 5, 1]

Five hidden layers: [13, 5, 5, 5, 5, 5, 1]

Example 4



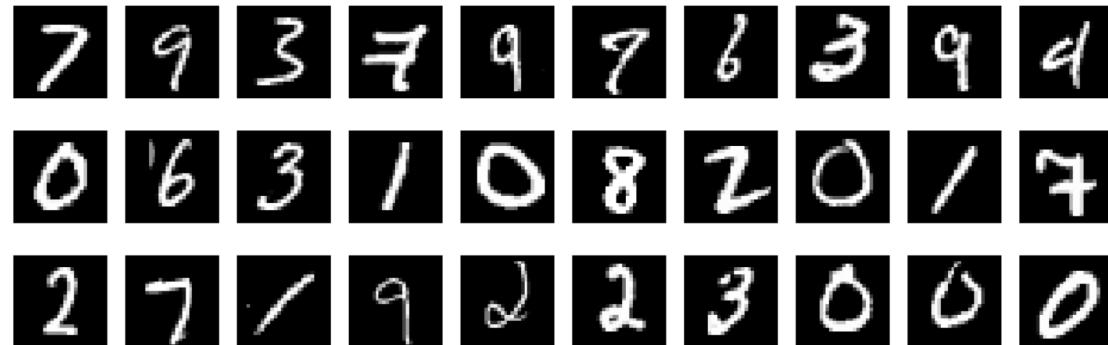
Depth: number of hidden layers

The deep networks extract features at different levels of complexity for regression or classification. However, the depth or the number of layers that you can have for the networks depend on the number of training patterns available. The deep networks have more parameters (weights and biases) to learn, so need more data to train.

The optimal number of layers is determined usually through experiments. The optimal architecture minimizes the error (training, test, and validation).

MNIST images

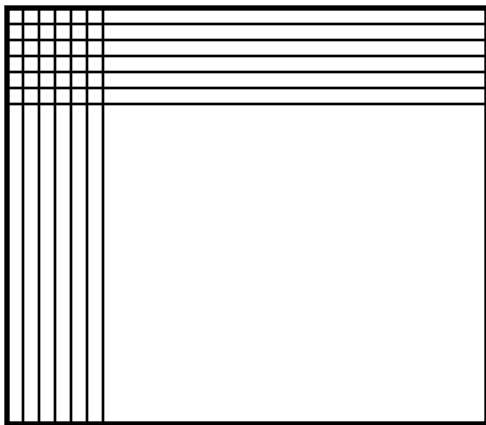
The MNIST database of gray level images of handwritten digits:
<http://yann.lecun.com/exdb/mnist/>



Each image is 28x28 size.
Intensities are in the range [0, 255].

Training set: 60,000
Test set: 10,000

MNIST images



An image is divided into rows and columns and defined by its pixels.

Size of the image = rows x columns pixels

Pixels of grey-level image are assigned intensity values: For example, integer values between 0 and 255 assigned as intensities (grey-values) for pixels with 0 representing ‘black’ and 255 representing ‘white’.

Color images has three color channels: red, green, and blue. A pixel in a color image is a vector (r, g, b) denoting intensity in red, green, and blue channels.

Example 5: Classification of MNIST images

No of inputs $n = 28 \times 28 = 784$

Inputs were normalized to $[0.0, 1.0]$

Use a 4-layer FFN

- Hidden-layer-1 is a sigmoid layer
- Hidden-layer-2 is sigmoid layer
- Output-layer is a softmax layer

Input-layer size $n = 784$

Hidden-layer-1 size $n_1 = 625$

Hidden-layer-2 size $n_2 = 100$

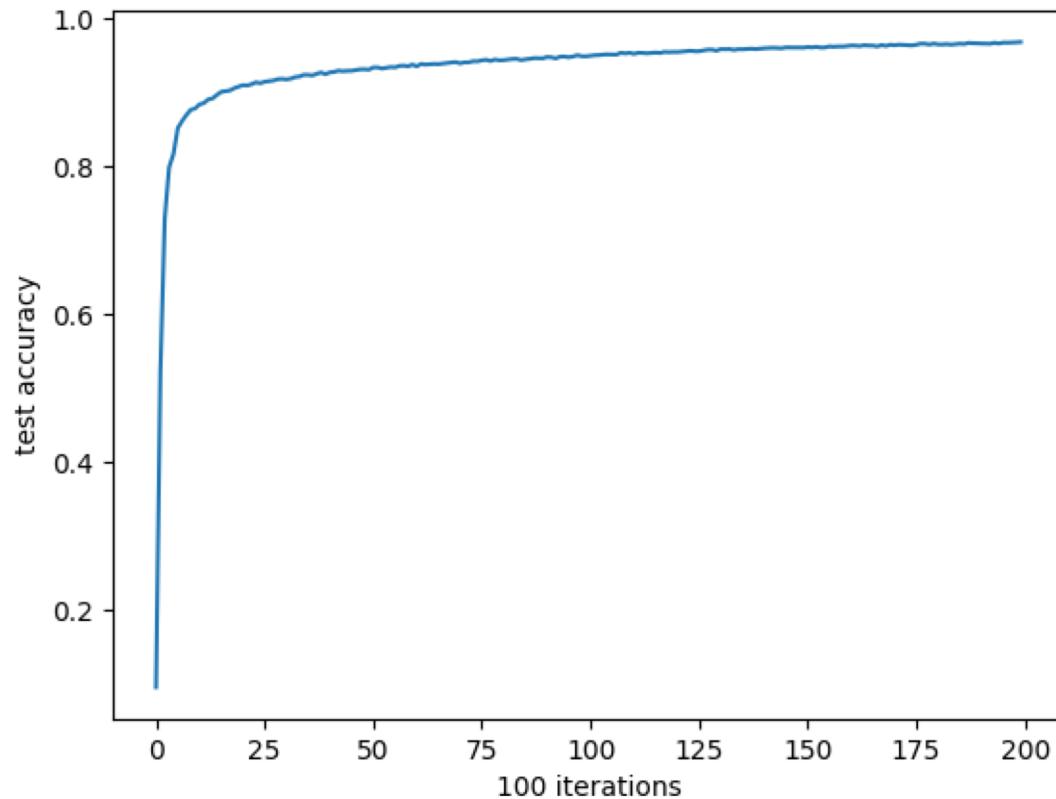
Output-layer size $K = 10$

Training:

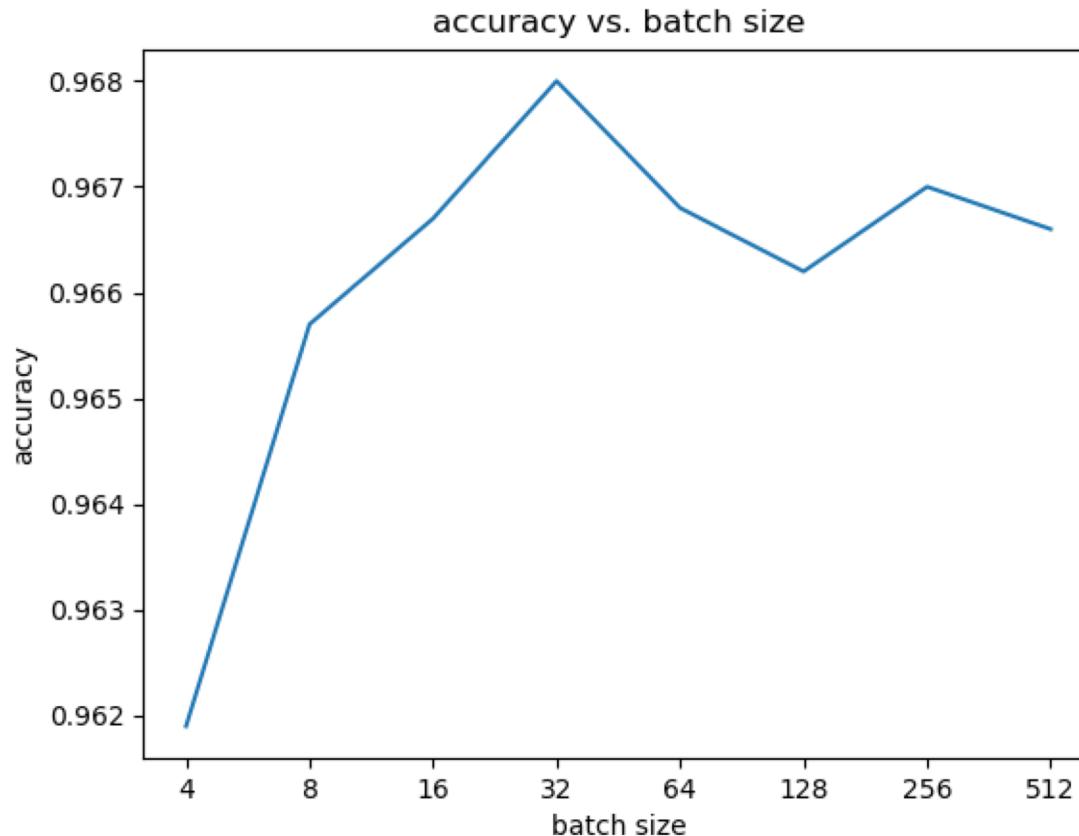
Batch size = 64

Learning rate $\alpha = 0.01$

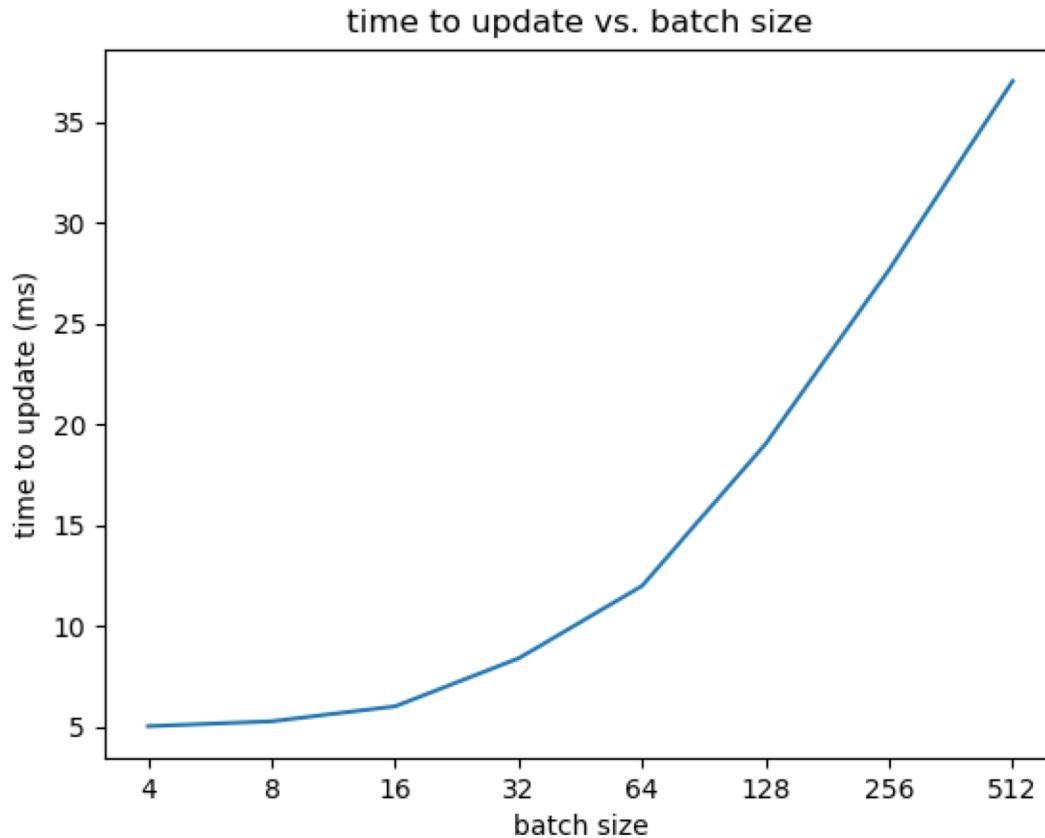
Example 5



Example 5: batch size



Example 5



Mini-batch SGD

In practice, gradient descent is performed on *mini-batch* updates of gradients within a *batch* or *block* of data of size B . In mini-batch SGD, the data is divided into blocks and the gradients are averaged for on each block in an epoch.

$B = 1$: stochastic (online) gradient descent

$B = P$ (size of training data): (batch) gradient descent

$1 < B < P \rightarrow$ mini-batch stochastic gradient descent

When B increases, more add-multiply operations per second, taking advantages of parallelism and matrix computations. On the other hand, as B increases, the number of computations per update (of weights, biases) increases.

Therefore, the curve of the time for weight update against batch size usually take a U-shape curve. There exists an optimal value of B – that depends on the sizes of the caches as well.

Batch size for mini-batch SGD

For SGD, it is desirable to randomly sample the patterns from training data in each epoch. In order to efficiently sample blocks, the training patterns are shuffled at the beginning of every training epoch and then blocks are sequentially fetched from memory.

Typical batch sizes: 16, 32, 64, 128, and 256.

The batch size is dependent on the size of caches of CPU and GPUs.

Summary

- FFN with one hidden layer
- Backpropagation for FFN with one hidden layer:

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

- Backpropagation for deep FFN (DNN)
- Training deep neural networks (GD and SGD):
 - Forward propagation of activation
 - Backpropagation of gradients
 - Updating weights
- Parameters: depth, width, and batch size