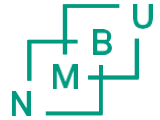# DAT300 – Applied Deep Learning

Math from ANN

# A multi-layer neural network architecture

- Data arrays



**X**

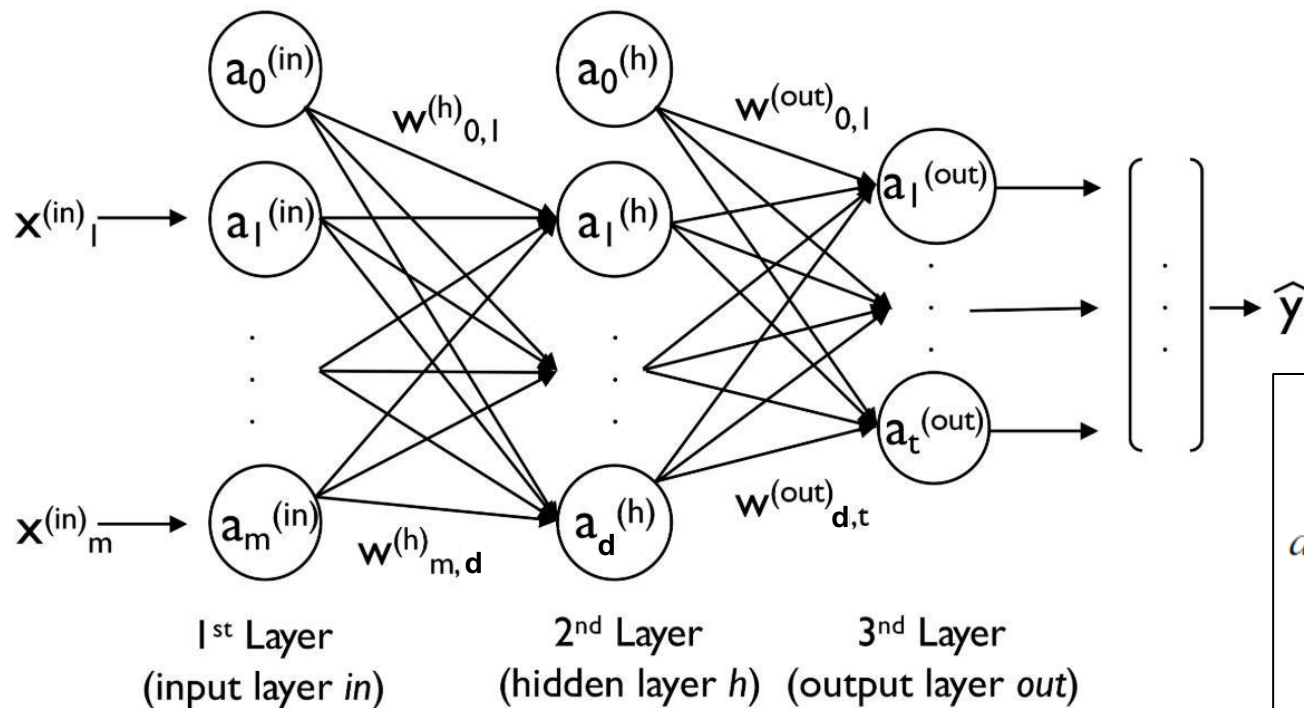**y**

(n x m)     (n x 1)

- Indexing

$$i = 1, \ldots, n \qquad \text{sample / instance index}$$

$$j = 0, \ldots, m \qquad \text{feature / instance index}$$

# A multi-layer neural network architecture

$1^{st}$ Layer (input layer $in$)   $2^{nd}$ Layer (hidden layer $h$)   $3^{nd}$ Layer (output layer $out$)
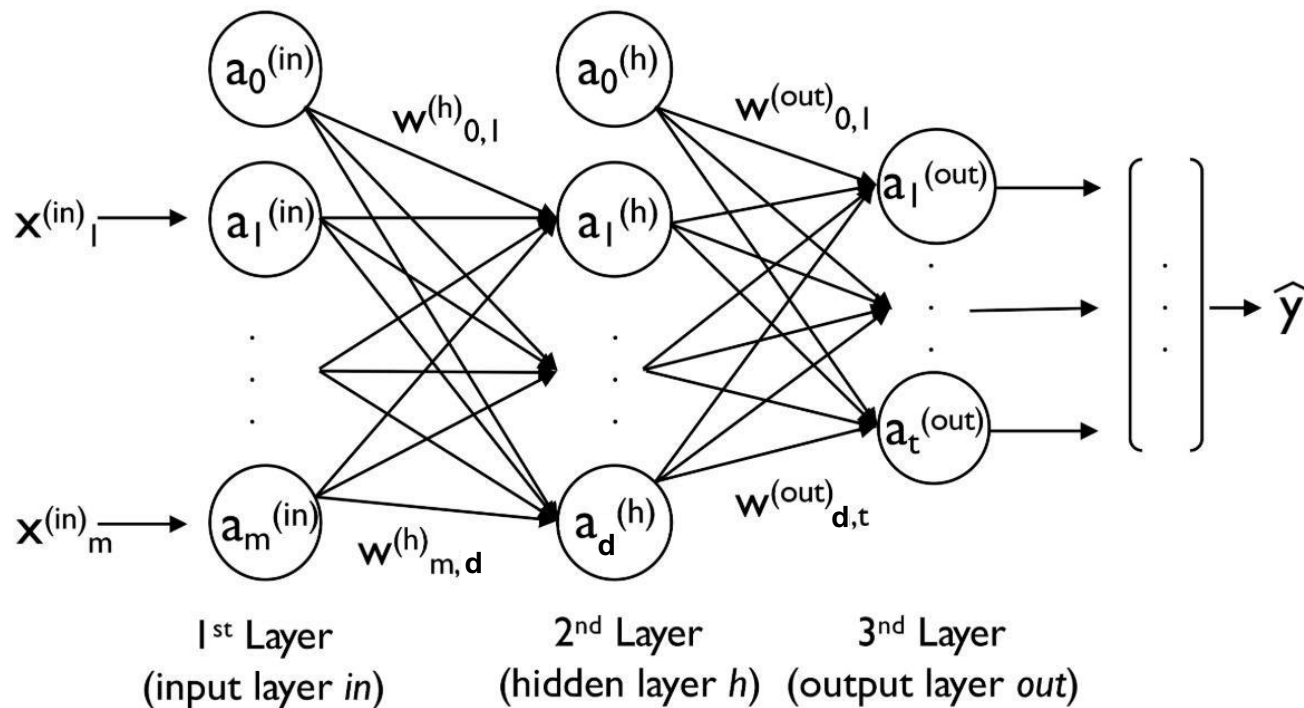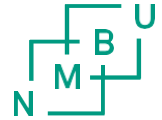
$$a^{(in)} = \begin{bmatrix} a_0^{(in)} \\ a_1^{(in)} \\ \vdots \\ a_m^{(in)} \end{bmatrix} = \begin{bmatrix} 1 \\ x_1^{(in)} \\ \vdots \\ x_m^{(in)} \end{bmatrix}$$

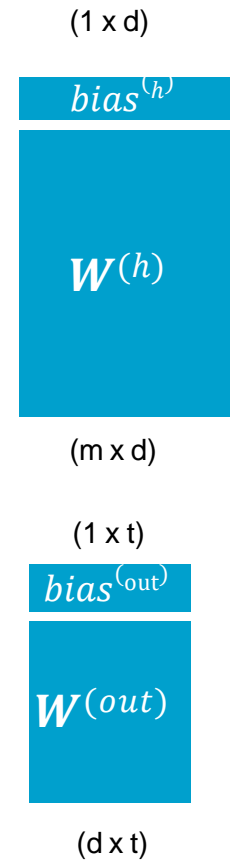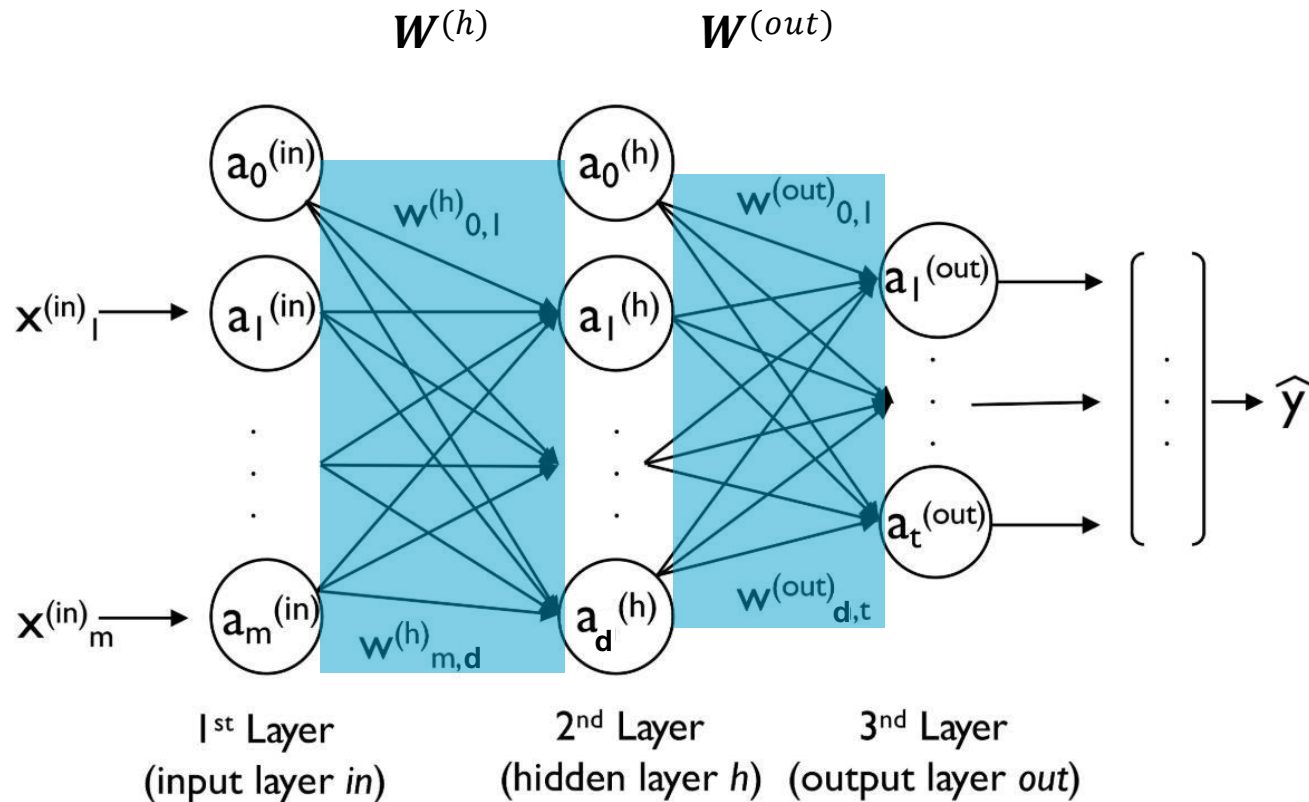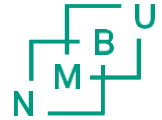(m x 1)     (m x 1)
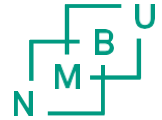+ bias     + bias

# A multi-layer neural network architecture



Example of model for three classes with t=3

$$0 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, 1 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, 2 = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

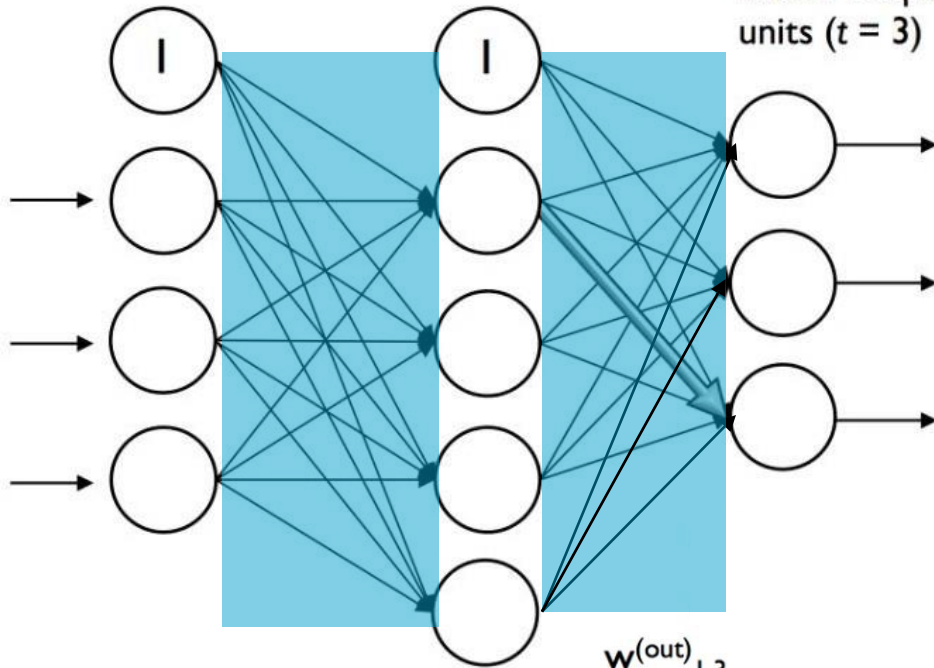# A multi-layer neural network architecture

# A 3-4-3 multi-layer perceptron



Input layer with 3 input units plus bias unit ($m = 3+1$)

Hidden layer with 4 hidden units plus bias unit ($d = 4+1$)

Output layer with 3 output units ($t = 3$)

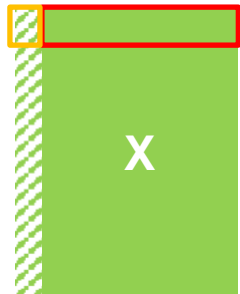Number of layers: $L = 3$
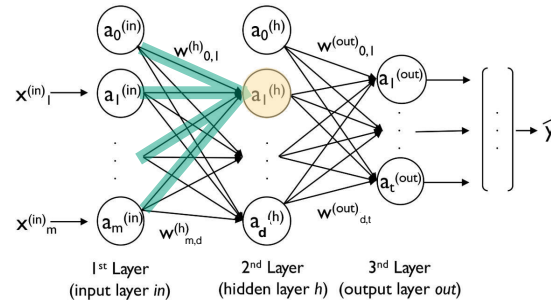
$W^{(out)}_{1,3}$ connects 1st non-bias neuron in the 2nd layer (hidden layer $h$) to the 3rd unit in the 3rd layer (output layer $out$)

(1 x 4)

$bias^{(h)}$    4 weights

$W^{(h)}$    12 weights

(3 x 4)    **Total: 16 weights**

(1 x 3)

$bias^{(out)}$    3 weights

$W^{(out)}$    12 weights

(4 x 3)    **Total: 15 weights**

Need to train 16 + 15 = 31 weights in NN model

(n x m)

$$z_1^{(h)} = a_0^{(in)} w_{0,1}^{(h)} + a_1^{(in)} w_{1,1}^{(h)} + \cdots + a_m^{(in)} w_{m,1}^{(h)}$$

(1 x 1)          (1 x 1)   (1 x 1)     (1 x 1)   (1 x 1)          (1 x 1)   (1 x 1)

Computations for one sample (row) $x_i$ in $X$ for one neuron in hidden layer

$$a_1^{(h)} = \phi\left(z_1^{(h)}\right)$$

(1 x 1)              (1 x 1)

# A 3-4-3 multi-layer perceptron

**X**

(n x m)

m: features
d: neurons
# Samples: 1

$$z^{(h)} = a^{(in)} W^{(h)}$$

(1 x d)    (1 x m)    (m x d)

Computations
for one sample
(row) $x_i$ in $X$
for all neurons
in hidden layer

$$a^{(h)} = \phi\left(z^{(h)}\right)$$

(1 x d)         (1 x d)

# A 3-4-3 multi-layer perceptron



**X**

(n x m)

m: features
d: neurons
# Samples: n

$$Z^{(h)} = A^{(in)}W^{(h)}$$

(n x d)      (n x m)      (m x d)

$$A^{(h)} = \phi\left(Z^{(h)}\right)$$

(n x d)      (n x d)

Computations
for all $n$ samples
(rows) $x_i$ in $X$
for all neurons in
hidden layer

# A 3-4-3 multi-layer perceptron

**X**

(n x m)

t: labels
d: neurons
# Samples: n

Computations
for all $n$ samples
(rows) $x_i$ in $X$
for all neurons in
output layer

$$Z^{(out)} = A^{(h)} W^{(out)}$$

(n x t)      (n x d)      (d x t)

$$A^{(out)} = \phi\left(Z^{(out)}\right)$$

(n x t)      (n x t)

One sample

$$\delta^{(out)} = a^{(out)} - y$$

(1 x t)    (1 x t)    (1 x t)

**3** Compute the loss gradient:

$$\frac{\partial}{\partial w_{i,j}^{(out)}} J(\boldsymbol{W}) = a_j^{(h)} \delta_i^{(out)}$$

**2** Error term of the output layer:

$$\boldsymbol{\delta}^{(out)} = \boldsymbol{a}^{(out)} - \boldsymbol{y}$$

Input $\boldsymbol{x}$

**1** Output $\widehat{y}$ ← Target $y$

**5** Compute the loss gradient:

$$\frac{\partial}{\partial w_{i,j}^{(h)}} J(\boldsymbol{W}) = a_j^{(in)} \delta_i^{(h)}$$

**4** Error term of the hidden layer:

$$\boldsymbol{\delta}^{(h)} = \boldsymbol{\delta}^{(out)} \left( \boldsymbol{W}^{(out)} \right)^{\top} \odot \frac{\partial \phi(\boldsymbol{z}^{(h)})}{\partial \boldsymbol{z}^{(h)}}$$

One sample

t: labels
d: neurons
# Samples: 1

$$\delta^{(out)} = a^{(out)} - y$$

(1 x t)     (1 x t)     (1 x t)

$$\delta^{(h)} = \delta^{(out)} \left(W^{(out)}\right)^T \odot \frac{\partial \phi\left(z^{(h)}\right)}{\partial z^{(h)}}$$

(1 x d)   (1 x t)     (t x d)           (1 x d)

$$\frac{\partial \phi(z)}{\partial z} = \left(a^{(h)} \odot \left(1 - a^{(h)}\right)\right)$$

(1 x d)     (1 x d)     (1 x d)

Derivative of
activation function

Derivative of *sigmoid*
activation function

# Backpropagation
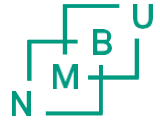
$n$ samples

t: labels
d: neurons
# Samples: n

$$\Delta^{(out)} = A^{(out)} - Y$$

(n x t)    (n x t)    (n x t)

$$\Delta^{(h)} = \Delta^{(out)}\left(W^{(out)}\right)^T \odot \frac{\partial \phi(Z^{(h)})}{\partial Z^{(h)}} \qquad \frac{\partial \phi(Z^{(h)})}{\partial Z^{(h)}} = \left(A^{(h)} \odot (1 - A^{(h)})\right)$$

(n x d)    (n x t)    (t x d)         (n x d)                (n x d)              (n x d)           (n x d)

Derivative of
activation function

Derivative of _sigmoid_
activation function

31

Norwegian University of Life Sciences

# Vanishing gradients problem

# Vanishing gradients problem



| Activation Function | Equation | Example | 1D Graph |
|---|---|---|---|
| Linear | $\phi(z) = z$ | Adaline, linear regression | |
| Unit Step (Heaviside Function) | $\phi(z) = \begin{cases} 0 & z < 0 \\ 0.5 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| Sign (signum) | $\phi(z) = \begin{cases} -1 & z < 0 \\ 0 & z = 0 \\ 1 & z > 0 \end{cases}$ | Perceptron variant | |
| Piece-wise Linear | $\phi(z) = \begin{cases} 0 & z \leq -\frac{1}{2} \\ z + \frac{1}{2} & -\frac{1}{2} \leq z \leq \frac{1}{2} \\ 1 & z \geq \frac{1}{2} \end{cases}$ | Support vector machine | |
| Logistic (sigmoid) | $\phi(z) = \dfrac{1}{1 + e^{-z}}$ | Logistic regression, Multilayer NN | |
| Hyperbolic Tangent (tanh) | $\phi(z) = \dfrac{e^{z} - e^{-z}}{e^{z} + e^{-z}}$ | Multilayer NN, RNNs | |
| ReLU | $\phi(z) = \begin{cases} 0 & z < 0 \\ z & z > 0 \end{cases}$ | Multilayer NN, CNNs | |

# Vanishing gradients problem

$$\delta^{(h)} = \delta^{(out)} \left( W^{(out)} \right)^{\top} \odot \frac{\partial \phi(z^{(h)})}{\partial z^{(h)}}$$
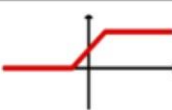
Derivative of
activation function

Page 419, 3rd edition

Derivative of *sigmoid*
activation function

$$\frac{\partial \phi(z)}{\partial z} = (a \odot (1 - a))$$

$$a = \phi(z) = \frac{1}{1 + e^{-z}}$$

# Vanishing gradients problem

Derivative of *sigmoid* activation function

$$\frac{\partial \phi(z)}{\partial z} = (a \odot (1-a))$$



$$\delta^{(h)} = \delta^{(out)} \left( W^{(out)} \right)^{\top} \odot \frac{\partial \phi(z^{(h)})}{\partial z^{(h)}}$$

$$\delta^{(h)} = \delta^{(out)} \left( W^{(out)} \right)^{\top} \odot \left( a^{(h)} \odot \left(1 - a^{(h)}\right) \right)$$

| $a$ | $1 - a$ | $a \odot (1-a)$ |
|-----|---------|-----------------|
| 0.1 | 0.9 | 0.09 |
| 0.2 | 0.8 | 0.16 |
| 0.3 | 0.7 | 0.21 |
| 0.4 | 0.6 | 0.24 |
| 0.5 | 0.5 | 0.25 |
| 0.6 | 0.4 | 0.24 |
| 0.7 | 0.3 | 0.21 |
| 0.8 | 0.2 | 0.16 |
| 0.9 | 0.1 | 0.09 |

# Vanishing gradients problem

$$\delta^{(h)} = \delta^{(out)} \left( \mathbf{W}^{(out)} \right)^{\top} \odot \frac{\partial \phi(\mathbf{z}^{(h)})}{\partial \mathbf{z}^{(h)}}$$

Derivative of
activation function

Derivative of *sigmoid* activation function

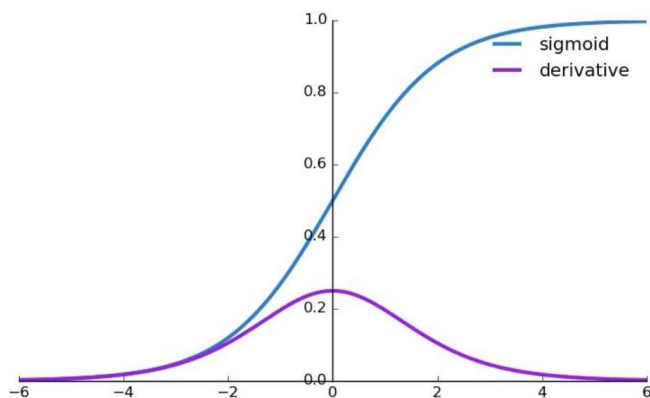$$\frac{\partial \phi(z)}{\partial z} = (a \odot (1 - a))$$

$$a = \phi(z) = \frac{1}{1 + e^{-z}}$$

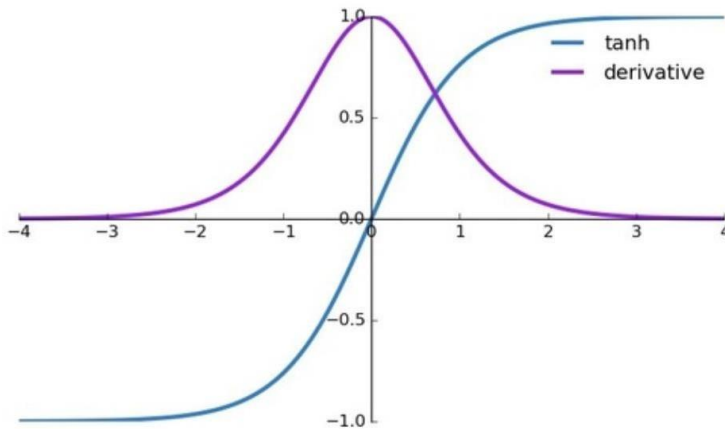Derivative of *tanh* activation function

$$\frac{\partial \phi(z)}{\partial z} = (1 - a^2)$$

$$a = \phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

# Vanishing gradients problem

<span style="color:red">Derivative of <u>tanh</u> activation function</span>

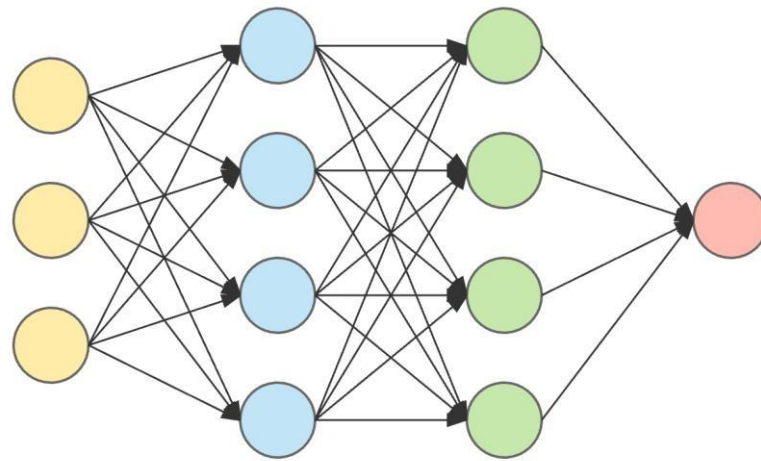$$\frac{\partial \phi(z)}{\partial z} = 1 - a^2$$



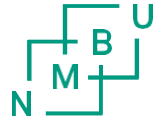| $a$ | $1-a^2$ |
| --- | --- |
| - 0.9 | 0.19 |
| - 0.5 | 0.75 |
| - 0.2 | 0.96 |
| - 0.1 | 0.99 |
| 0.0 | 1.00 |
| 0.1 | 0.99 |
| 0.2 | 0.96 |
| 0.5 | 0.75 |
| 0.9 | 0.19 |

$$\delta^{(h)} = \delta^{(out)} \left(W^{(out)}\right)^{\top} \odot \frac{\partial \phi(z^{(h)})}{\partial z^{(h)}}$$

$$\delta^{(h)} = \delta^{(out)} \left(W^{(out)}\right)^{\top} \odot \left(1 - a^{(h)^2}\right)$$

# Backpropagation



input layer  hidden layer 1  hidden layer 2  output layer

$$\boldsymbol{\delta}^{(out)} = \boldsymbol{a}^{(out)} - \boldsymbol{y}$$

$$\boldsymbol{\delta}^{(h_2)} = \boldsymbol{\delta}^{(out)} (\boldsymbol{W}^{(out)})^T \odot \frac{\partial \phi(\boldsymbol{z}^{(h_2)})}{\partial \boldsymbol{z}^{(h_2)}}$$

$$\boldsymbol{\delta}^{(h_1)} = \boldsymbol{\delta}^{(h_2)} (\boldsymbol{W}^{(h_2)})^T \odot \frac{\partial \phi(\boldsymbol{z}^{(h_1)})}{\partial \boldsymbol{z}^{(h_1)}}$$

$$\boldsymbol{\delta}^{(h_1)} = \boldsymbol{\delta}^{(out)} (\boldsymbol{W}^{(out)})^T \odot \frac{\partial \phi(\boldsymbol{z}^{(h_2)})}{\partial \boldsymbol{z}^{(h_2)}} (\boldsymbol{W}^{(h_2)})^T \odot \frac{\partial \phi(\boldsymbol{z}^{(h_1)})}{\partial \boldsymbol{z}^{(h_1)}}$$