# Neural Networks - II

**Thang Vu**

**21.11.2025**

Universität Stuttgart

# Outline

- Neuron
- Neural network
- Computing the output
- Training the network
- Backpropagation

**Universität Stuttgart**

# Notation

Output of a neuron:

$$a_i^l$$

layer l

neuron i

Output of one layer:

$a^l$ is a vector

$a_1^{l-1}$

$a_2^{l-1}$

.
.
.

$a_N^{l-1}$

Layer l-1
$N_{l-1}$ nodes

$a_1^l$

$a_2^l$

.
.
.

$a_N^l$

Layer l
$N_l$ nodes

Universität Stuttgart

# Notation



Layer l-1
$N_{l-1}$ nodes

Layer l
$N_l$ nodes

- Weights:

$w_{ij}^l$ → layer $l-1$ to layer $l$

from neuron $j$ (layer $l$ -1)
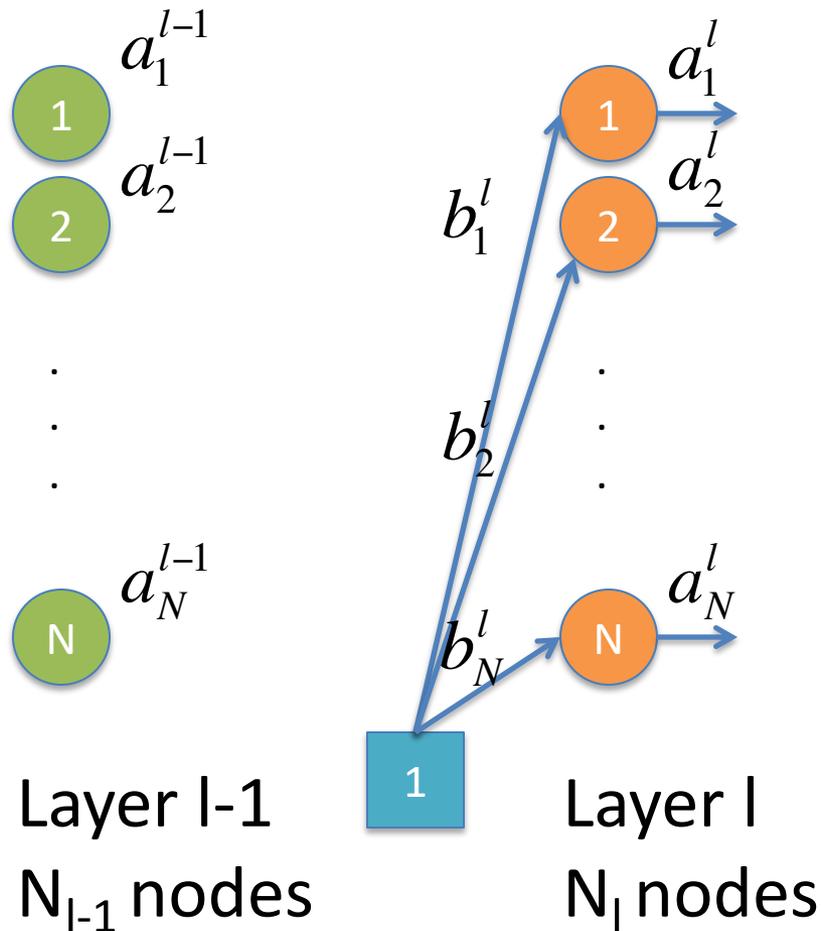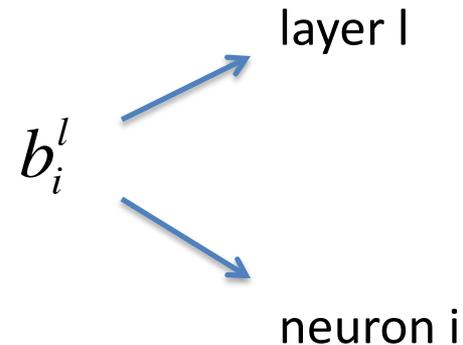to neuron $i$ (layer $l$)

$$W^l = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & & \end{bmatrix}$$

$N_{l-1}$

$N_l$

# Notation



$a_1^{l-1}$

$a_2^{l-1}$

.
.
.

$a_N^{l-1}$

Layer l-1
$N_{l-1}$ nodes

$b_1^l$

$b_2^l$

$b_N^l$

$a_1^l$

$a_2^l$

.
.
.

$a_N^l$

Layer l
$N_l$ nodes

- Biases:

$b_i^l$ → layer l

$b_i^l$ → neuron i

$$b^l = \begin{bmatrix} b_1^l \\ b_2^l \\ \vdots \end{bmatrix}$$ Bias for all the neurons in layer l

Universität Stuttgart

# Notation

$a_1^{l-1}$

$a_2^{l-1}$

$z_i^l$

$a_i^l$

$a_N^{l-1}$

$b_i^l$

Layer l-1
$N_{l-1}$ nodes

Layer l
$N_l$ nodes

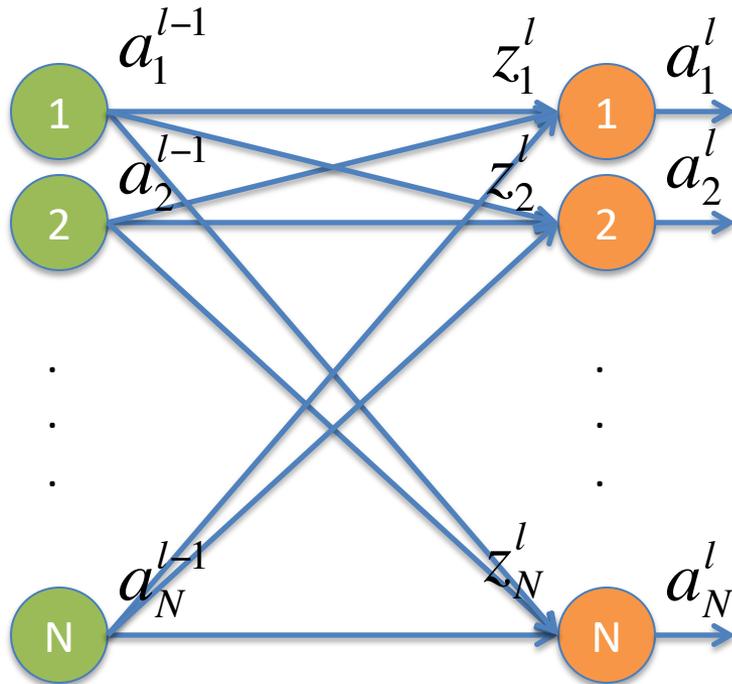$z_i^l$ : input of the activation function for neuron i at layer l

$z^l$ : input of the activation function of all the neurons in layer l

$$z_i^l = w_{i1}^l a_1^{l-1} + w_{i2}^l a_2^{l-1} + \ldots + b_i^l$$

or in another form

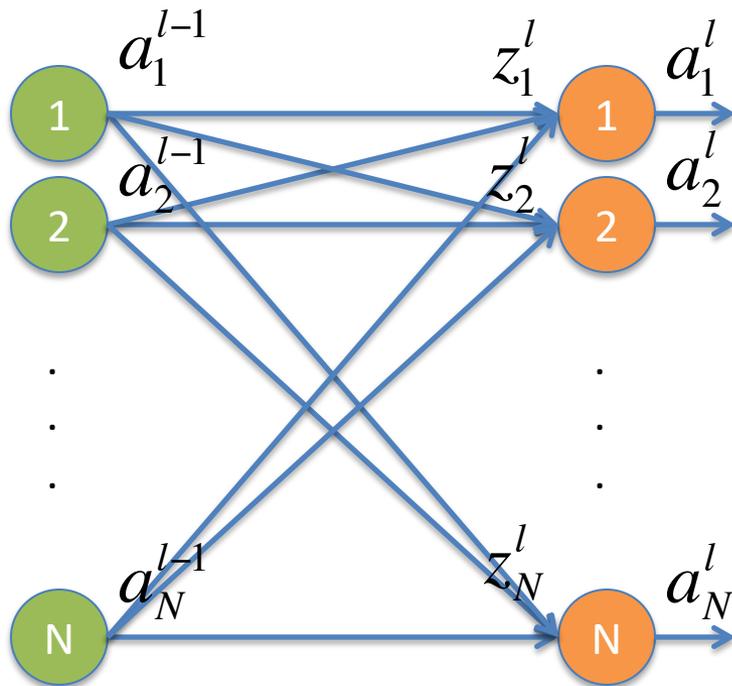$$z_i^l = \sum_{j=1}^{N_{l-1}} w_{ij}^l a_j^{l-1} + b_i^l$$

# Relations between layer outputs



Layer l-1
$N_{l-1}$ nodes

Layer l
$N_l$ nodes

# Relations between layer outputs
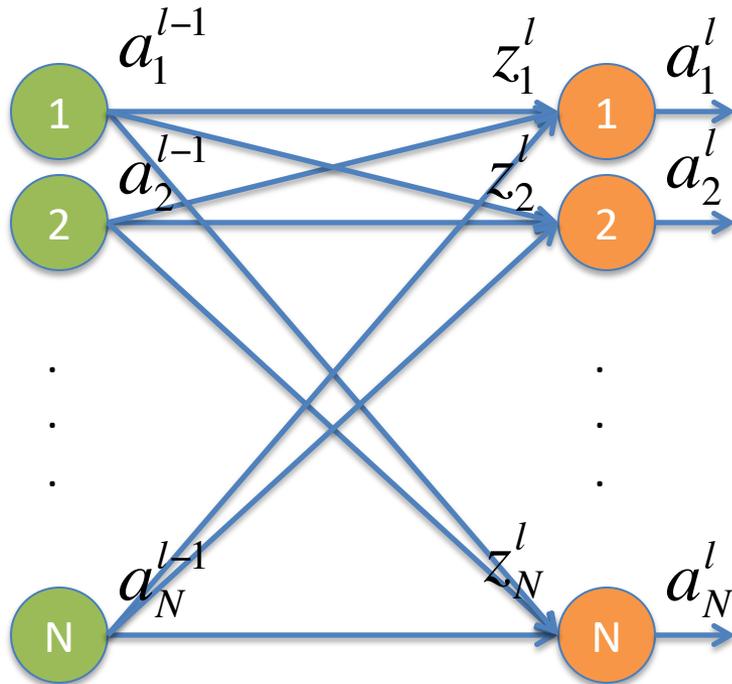


$$z_1^l = w_{11}^l a_1^{l-1} + w_{12}^l a_2^{l-1} + \ldots + b_1^l$$

$$z_2^l = w_{21}^l a_1^{l-1} + w_{22}^l a_2^{l-1} + \ldots + b_2^l$$

.....

$$z_N^l = w_{N1}^l a_1^{l-1} + w_{N2}^l a_2^{l-1} + \ldots + b_N^l$$

$$z^l = W^l a^{l-1} + b^l$$

Layer l-1
$N_{l-1}$ nodes

Layer l
$N_l$ nodes

# Relations between layer outputs



$$a_1^l = \sigma(z_1^l)$$

$$a_2^l = \sigma(z_2^l)$$

....

$$a_N^l = \sigma(z_N^l)$$

$$a^l = \sigma(z^l)$$

Layer l-1
$N_{l-1}$ nodes

Layer l
$N_l$ nodes

Universität Stuttgart

# Relations between layer outputs



$$z^l = W^l a^{l-1} + b^l$$

$$a^l = \sigma(z^l)$$

$$a^l = \sigma(W^l a^{l-1} + b^l)$$

Layer l-1
$N_{l-1}$ nodes

Layer l
$N_l$ nodes

Universität Stuttgart

# Computation of the final output

Input    Layer 1    Layer 2    Layer N

vector **x**

$x_1$

$x_2$

$x_3$

$\cdot$

$\cdot$

$\cdot$

$x_N$

vector **y**

$y_1$

$y_2$

$y_M$

$x$        $a^1$        $a^2$        $a^L$

$$a^1 = \sigma(W^1 x + b^1)$$

$$a^L = \sigma(W^L a^{L-1} + b^L)$$

$$a^2 = \sigma(W^2 a^1 + b^2)$$

Universität Stuttgart

# Computation of the final output



$$y = f(x) = \sigma(W^L ... \sigma(W^2 \sigma(W^1 x + b^1) + b^2) ... + b^L)$$

**Universität Stuttgart**

# Softmax function

Ouput layer L

$z_1^L$

$\hat{y}_1$

$z_2^L$

$\hat{y}_2$

$\cdots$

$\hat{y}_N$

- Outputs are probabilities

$$0 < y_i < 1$$

$$\sum_i y_i = 1$$

- Sofmax function:

$$y_i = \frac{e^{z_i^L}}{\sum_j e^{z_j^L}}$$

$$\frac{\delta y_i}{\delta z_j} = y_j(1 - y_j) \quad i = j$$

$$= -y_i y_j \qquad i \neq j$$

Universität Stuttgart

# Train a Neural Network

**Universität Stuttgart**

# Searching for the best function

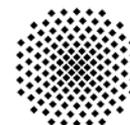- Best function = best parameters

$$y = f(x) = \sigma(W^L ... \sigma(W^2 \sigma(W^1 x + b^1) + b^2) ... + b^L)$$

- The function itself depends on the parameter set

$$f(x) = f(x, \theta)$$

$$\theta = \left\{ W^1, b^1, W^2, b^2, ..., W^L, b^L \right\}$$

- Seach the best function f*

→ Search the best parameter set $\theta$*

Universität Stuttgart

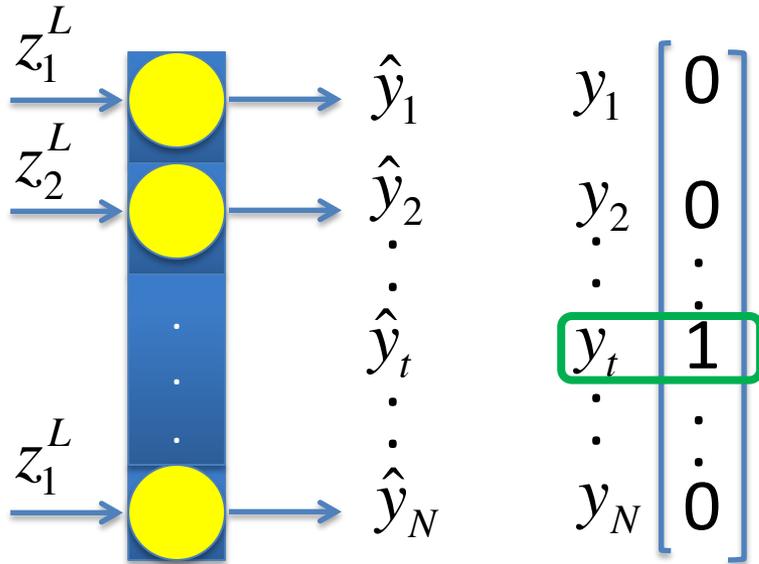# Empirical risk minimization

- Given a finite set of training data
- Empirical risk = average loss on this training data

$$C(\theta) = \frac{1}{|D|} \sum_{(x,y)} c(f(x), y)$$

$$= \frac{1}{|D|} \sum_{(x,y)} c(\theta)$$

**Universität Stuttgart**

# Loss function

Ouput layer L

$z_1^L$   $\hat{y}_1$   $y_1$   0

$z_2^L$   $\hat{y}_2$   $y_2$   0

$\hat{y}_t$   $y_t$   1

$z_1^L$   $\hat{y}_N$   $y_N$   0

- Sofmax function:
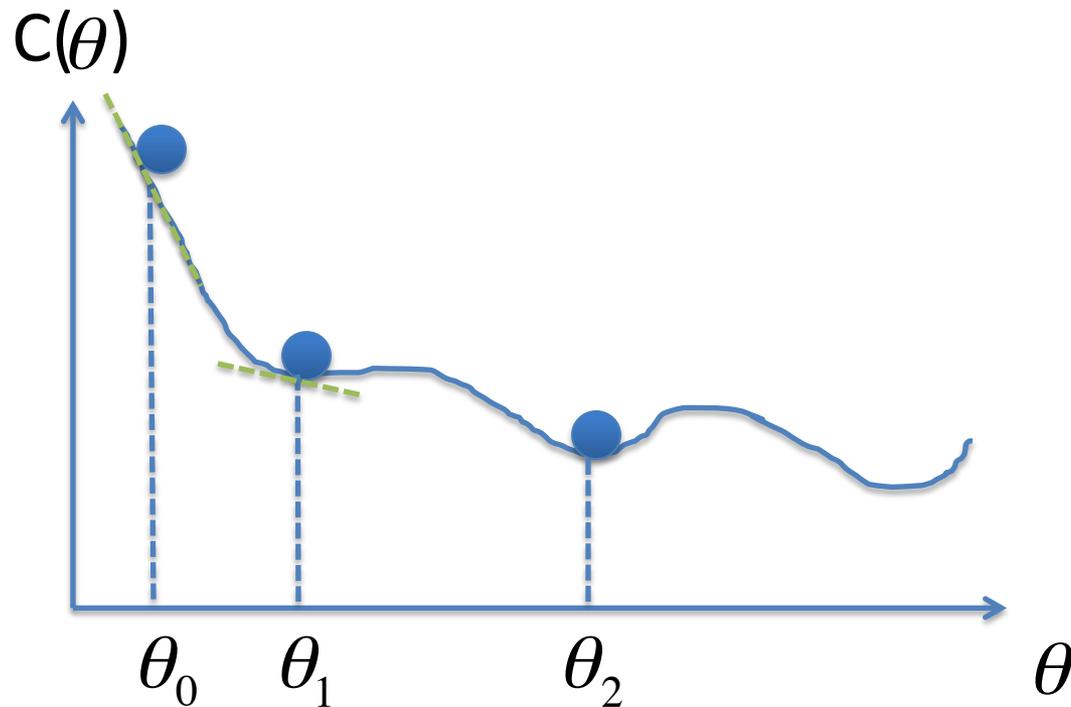
$$y_i = \frac{e^{z_i^L}}{\sum_j e^{z_j^L}}$$

- Mean Squared Error (MSE)

$$c(\theta) = C_x(\theta) = \sum_i (\hat{y}_i - y_i)^2$$

- Cross Entropy (CE):

$$c(\theta) = C_x(\theta) = -\log \hat{y}_t$$

Universität Stuttgart

# Gradient descent

- If $\theta$ has only one variable

C($\theta$)



- Randomly start at $\theta_0$
- Compute dC($\theta_0$)/d$\theta$

$$\theta_1 \leftarrow \theta_0 - \eta dC(\theta_0) / d\theta$$

- Compute dC($\theta_1$)/d$\theta$

$$\theta_2 \leftarrow \theta_1 - \eta dC(\theta_1) / d\theta$$

- ......

Universität Stuttgart

# Gradient descent for Neural Network

- We will do the same thing as presented before

- Starting parameters $\theta_0$

$$\longrightarrow \quad \theta_1 \quad \longrightarrow \quad \theta_2 \quad \longrightarrow \quad \theta_3 \quad \ldots\ldots\ldots$$

- However,

$$\theta = \left\{ W^1, b^1, W^2, b^2, ..., W^L, b^L \right\}$$

- i.e. millions of parameters ☹
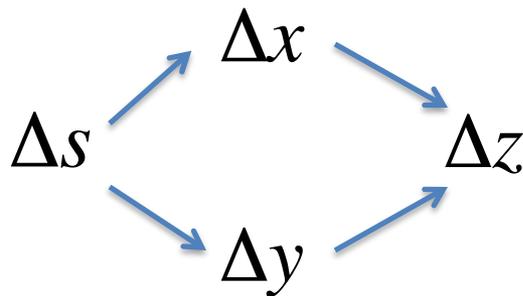
# Backpropagation

Based on the chain rules:

- Case 1:
$$y = g(x)$$
$$z = h(y)$$
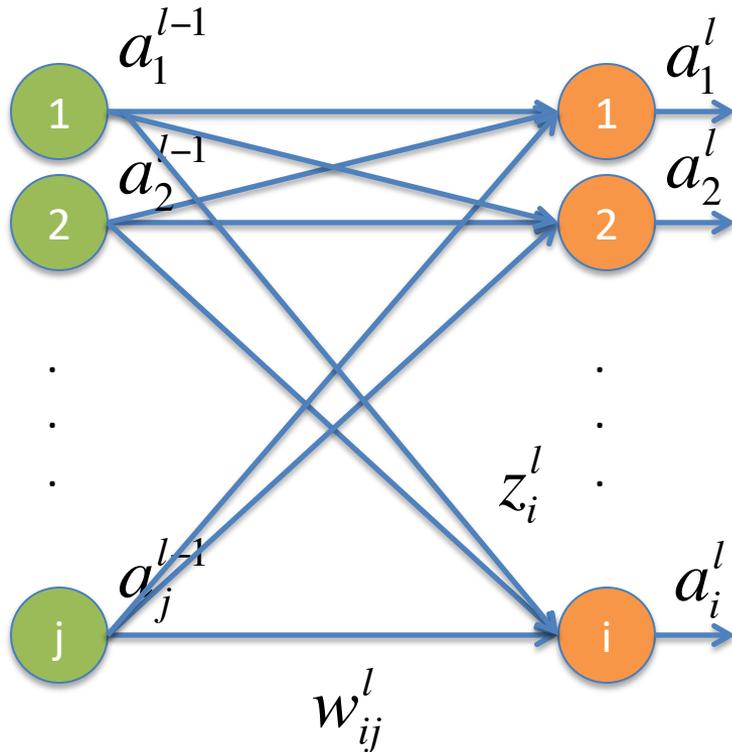$$\Delta x \rightarrow \Delta y \rightarrow \Delta z$$
$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

- Case 2:
$$x = g(s) \quad y = h(s) \quad z = k(x,y)$$



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

Universität Stuttgart

# Compute $\dfrac{\partial C}{\partial w_{ij}^l}$



$$\Delta w_{ij}^l \rightarrow \Delta z_i^l \ldots \rightarrow \Delta C$$

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \frac{\partial C}{\partial z_i^l}$$
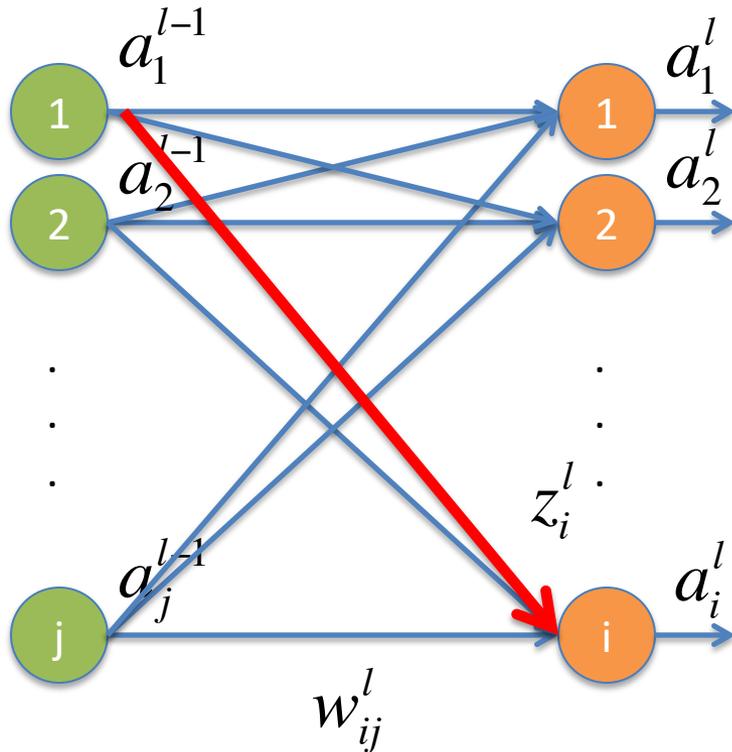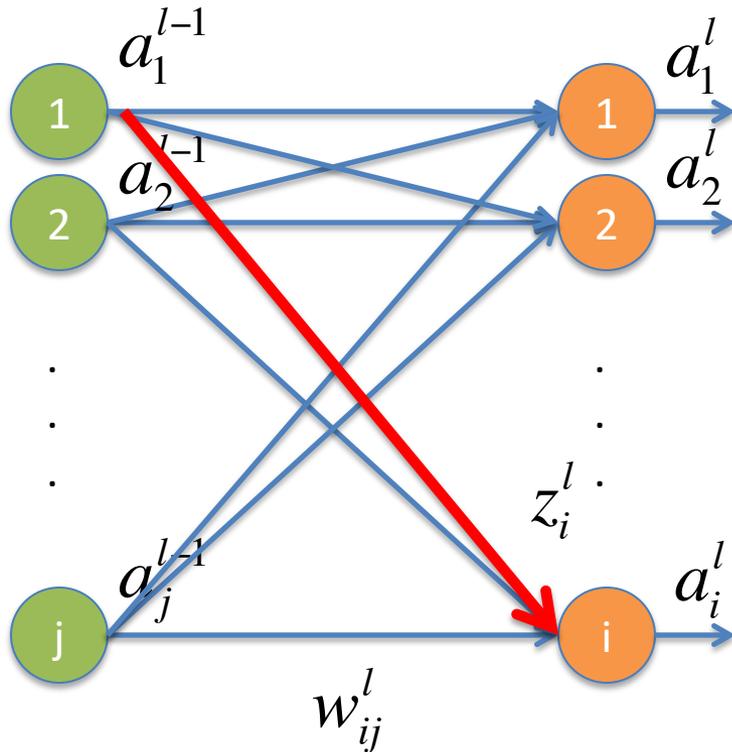
Universität Stuttgart

# Compute $\dfrac{\partial C}{\partial w_{ij}^l}$



$$\Delta w_{ij}^l \rightarrow \Delta z_i^l \ldots \rightarrow \Delta C$$

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l}\frac{\partial C}{\partial z_i^l}$$
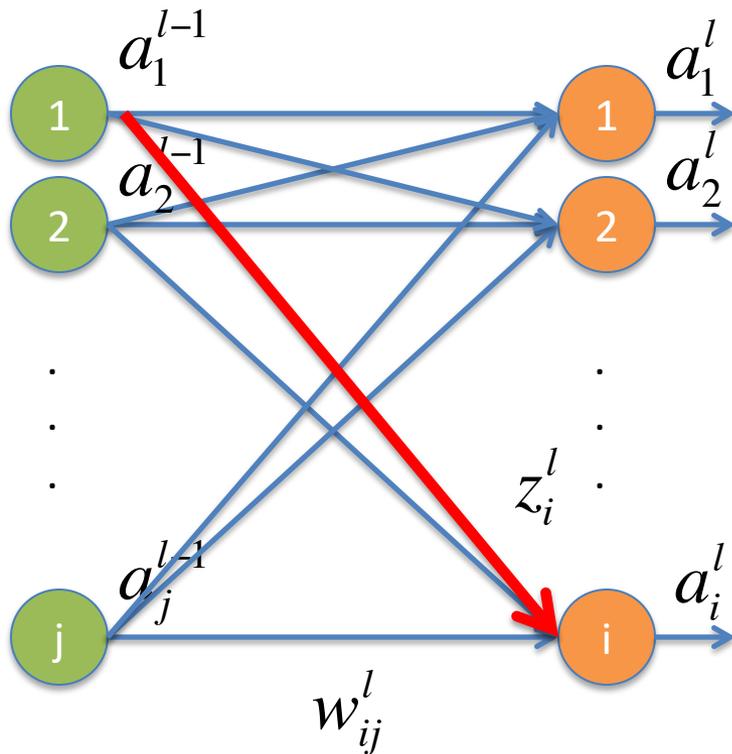
Universität Stuttgart

# Compute - $\frac{\partial C}{\partial w_{ij}^l}$ First term



$$\Delta w_{ij}^l \rightarrow \Delta z_i^l \ldots \rightarrow \Delta C$$

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \frac{\partial C}{\partial z_i^l}$$

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - First term



- If l > 1:

$$z_i^l = \sum_j w_{ij}^l a_j^{l-1} + b_i^l$$

$$\frac{\partial z_i^l}{\partial w_{ij}^l} = a_j^{l-1}$$
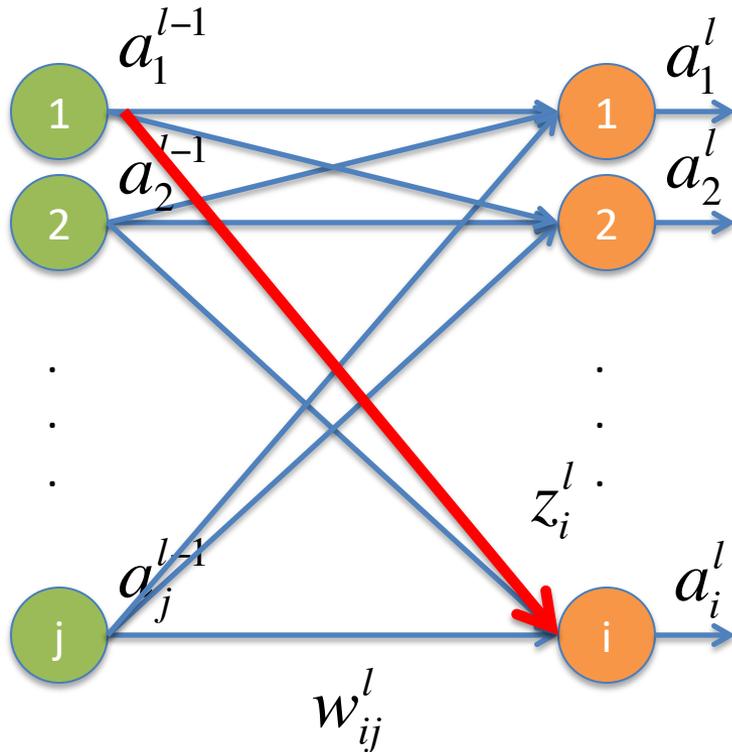
- If l = 1:

$$z_i^l = \sum_j w_{ij}^1 x_j + b_i^1$$

$$\frac{\partial z_i^1}{\partial w_{ij}^1} = x_j$$

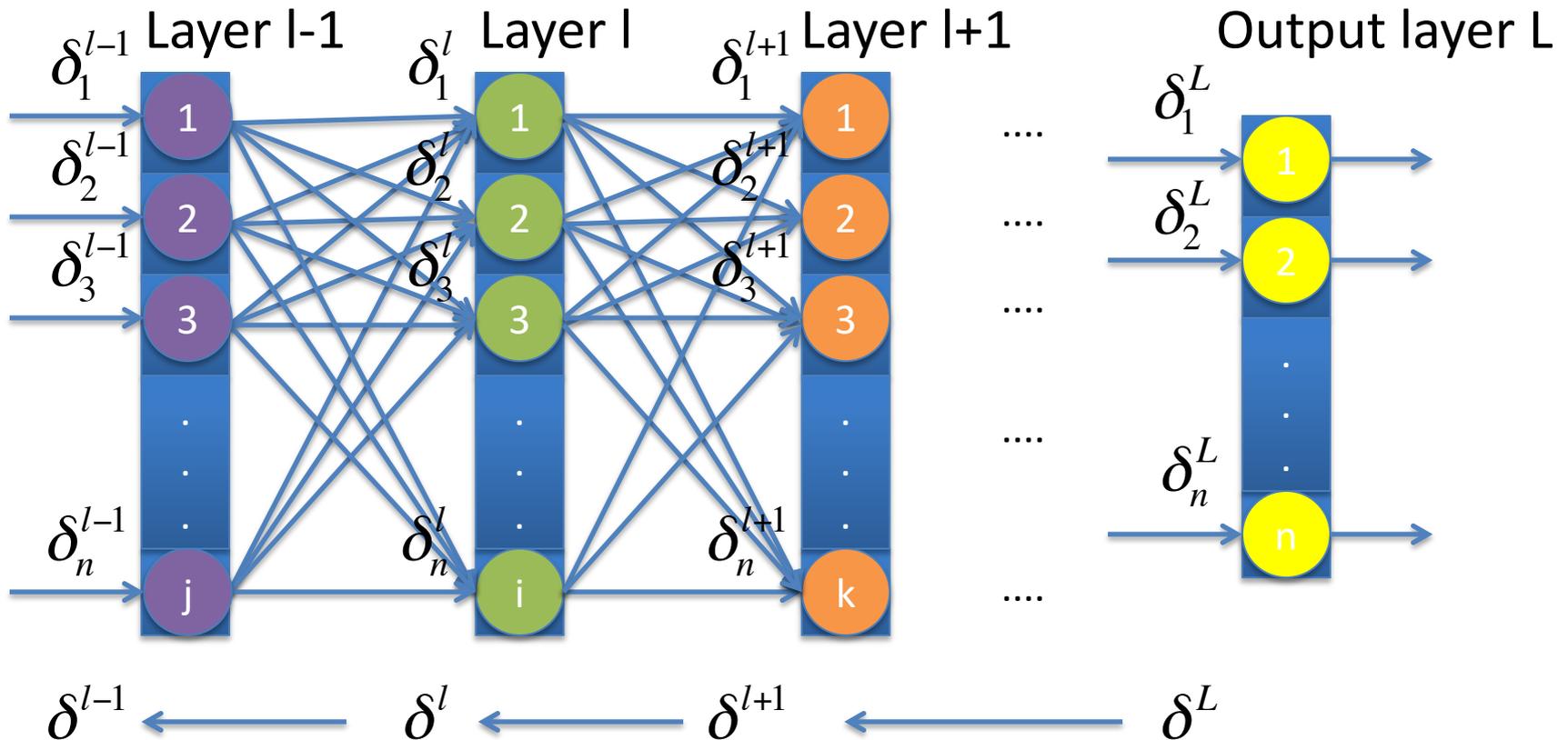# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term



$$\Delta w_{ij}^l \to \Delta z_i^l \ldots \to \Delta C$$

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \frac{\partial C}{\partial z_i^l}$$

Universität Stuttgart

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term

Universität Stuttgart

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term

Output layer L

$$\Delta z_n^L \rightarrow \Delta a_n^L = \Delta y_n \rightarrow \Delta C$$

$$\delta_n^L = \frac{\partial C}{\partial z_n^L} = \frac{\partial a_n^L}{\partial z_n^L} \frac{\partial C}{\partial a_n^L}$$

$\delta_1^L$

$\delta_2^L$

$\delta_n^L$

$$\sigma'(z_n^L)$$

Depending on
The loss function

$\delta^L$

$$\delta^L = \sigma'(z^L)\nabla C(y)$$    Elementwise multiplication

Universität Stuttgart

# Second term – Last layer

- If the last layer uses softmax and cross-entropy loss is used

$$c(\theta) = C_x(\theta) = -\log \hat{y}_t$$

$$\delta_i^L = \frac{\partial a_i^L}{\partial z_i^L} \frac{\partial C}{\partial a_i^L}$$

$$\frac{\partial C}{\partial a_i^L} = -\frac{1}{\hat{y}_t}$$

Ground truth

$$
\begin{array}{cc}
y_1 & 0 \\
y_2 & 0 \\
\vdots & \vdots \\
y_t & 1 \\
\vdots & \vdots \\
\vdots & \vdots \\
y_N & 0 \\
\end{array}
$$

Universität Stuttgart

# Second term – Last layer

- If the last layer uses softmax and cross-entropy loss is used

$$\delta_i^L = \frac{\partial a_i^L}{\partial z_i^L} \frac{\partial C}{\partial a_i^L}$$

$$y_i = \frac{e^{z_i^L}}{\sum_j e^{z_j^L}}$$

$$\frac{\delta y_i}{\delta z_j} = y_j(1 - y_j) \quad i = j$$

$$= -y_i y_j \quad i \neq j$$

$$\frac{\partial y_t^L}{\partial z_i^L} = y_t(1 - y_t) \text{ if } i = t$$

$$= -y_t y_i \text{ if } i \neq t$$

Ground truth

$$
\begin{array}{c}
y_1 \\
y_2 \\
\vdots \\
y_t \\
\vdots \\
y_N
\end{array}
\begin{bmatrix}
0 \\
0 \\
\vdots \\
1 \\
\vdots \\
0
\end{bmatrix}
$$

Universität Stuttgart

# Second term – Last layer

- If the last layer uses softmax and cross-entropy loss is used

$$\delta_i^L = \left(\frac{\partial a_i^L}{\partial z_i^L}\right)\frac{\partial C}{\partial a_i^L}$$

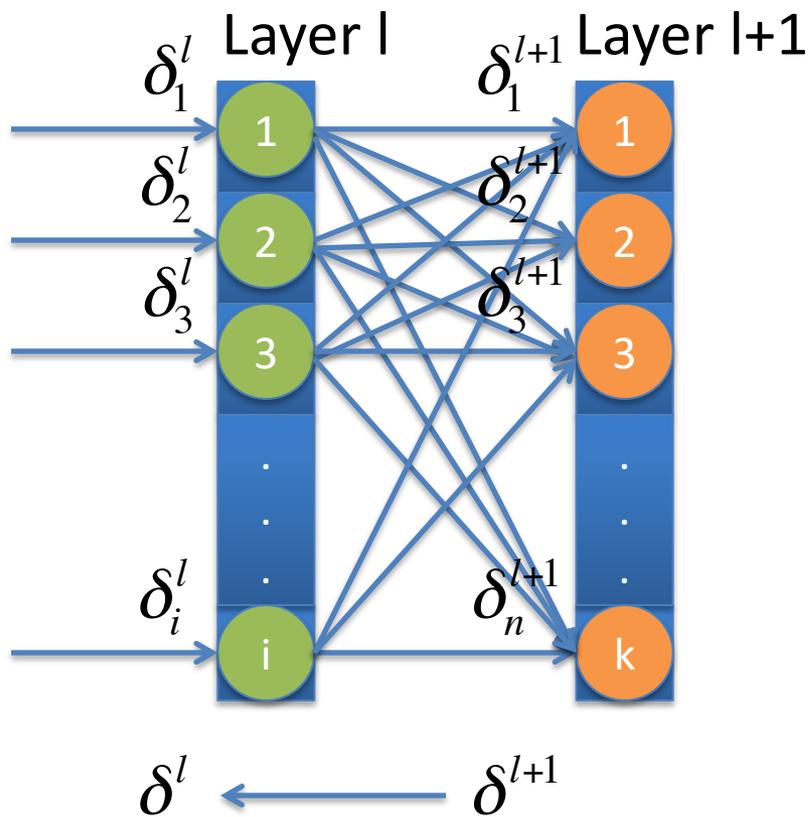$$\frac{\partial y_t^L}{\partial z_i^L} = y_t(1 - y_t) \text{ if } i = t$$

$$= -y_t y_i \text{ if } i \neq t$$

$$\frac{\partial C}{\partial a_i^L} = -\frac{1}{\hat{y}_t}$$

$$\delta_i^L = \hat{y}_t(1 - \hat{y}_t)\left(-\frac{1}{\hat{y}_t}\right) = \hat{y}_t - 1 \quad \text{if } i = t$$

$$= -\hat{y}_t \hat{y}_i \left(-\frac{1}{\hat{y}_t}\right) = \hat{y}_i \qquad \text{if } i \neq t$$

Ground truth

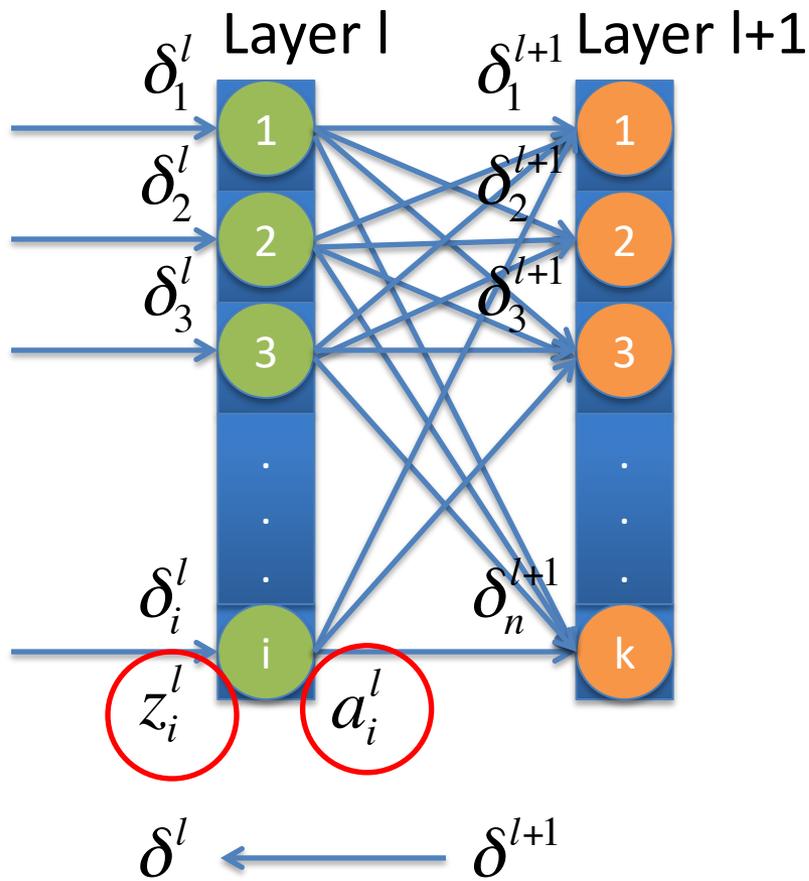$$\begin{array}{cc} y_1 & 0 \\ y_2 & 0 \\ \vdots & \vdots \\ y_t & 1 \\ \vdots & \vdots \\ y_N & 0 \end{array}$$

Universität Stuttgart

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term



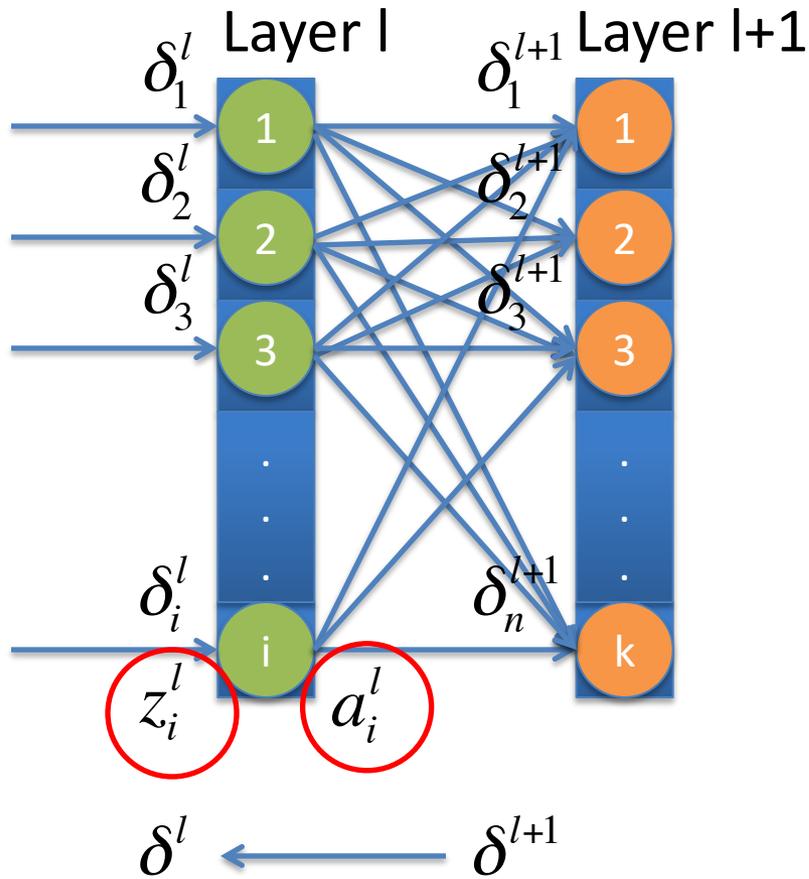$$\delta_i^l = \frac{\partial C}{\partial z_i^l}$$

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term



Layer l     Layer l+1

$$\delta_i^l = \frac{\partial C}{\partial z_i^l}$$

$$\Delta z_i^l \rightarrow \Delta a_i^l$$

$\delta^l \longleftarrow \delta^{l+1}$

Universität Stuttgart

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term



Layer l   Layer l+1

$\delta_1^l$   $\delta_1^{l+1}$
$\delta_2^l$   $\delta_2^{l+1}$
$\delta_3^l$   $\delta_3^{l+1}$
$\delta_i^l$   $\delta_n^{l+1}$

$z_i^l$   $a_i^l$

$\delta^l \longleftarrow \delta^{l+1}$

$$\delta_i^l = \frac{\partial C}{\partial z_i^l}$$

$$\Delta z_i^l \longrightarrow \Delta a_i^l \begin{array}{c} \Delta z_i^{l+1} \\ \Delta z_2^{l+1} \\ \Delta z_k^{l+1} \end{array} \Delta C$$
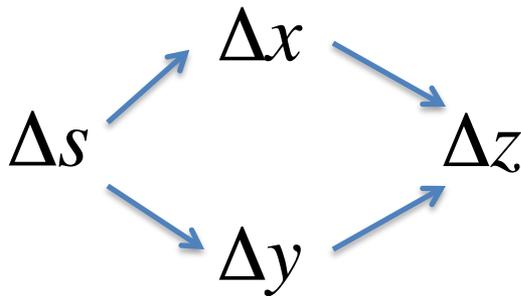
# Chain Rules

- Case 1:

$$y = g(x)$$

$$z = h(y)$$

$$\Delta x \rightarrow \Delta y \rightarrow \Delta z$$

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$
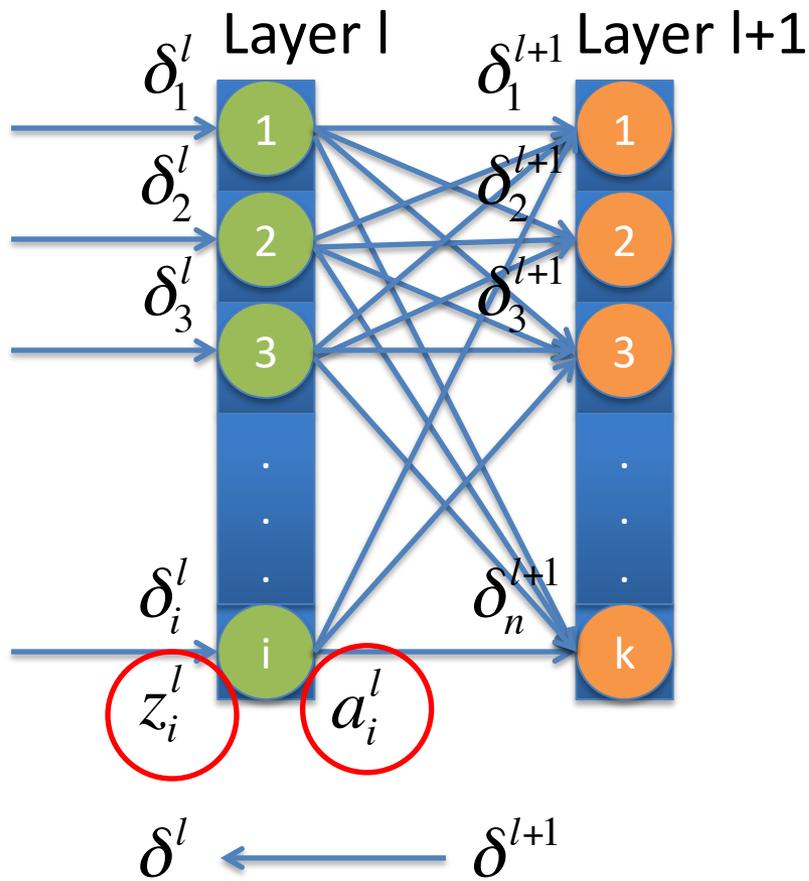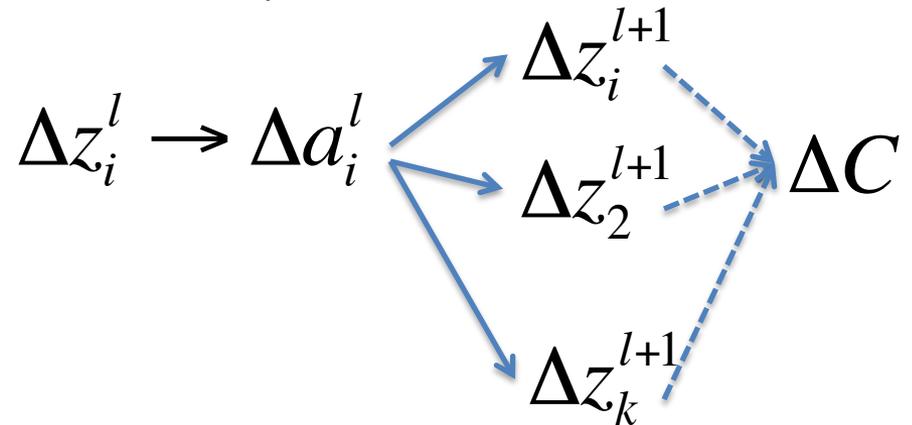
- Case 2:

$$x = g(s) \quad y = h(s) \quad z = k(x,y)$$



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

Universität Stuttgart

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term

Layer l     Layer l+1

$\delta_1^l$   $\delta_1^{l+1}$

$\delta_2^l$   $\delta_2^{l+1}$

$\delta_3^l$   $\delta_3^{l+1}$

$\delta_i^l$   $\delta_n^{l+1}$

$z_i^l$   $a_i^l$

$\delta^l \longleftarrow \delta^{l+1}$

$$\delta_i^l = \frac{\partial C}{\partial z_i^l}$$

$$\Delta z_i^l \rightarrow \Delta a_i^l \begin{array}{c} \Delta z_i^{l+1} \\ \Delta z_2^{l+1} \\ \Delta z_k^{l+1} \end{array} \Delta C$$

$$\delta_i^l = \frac{\partial C}{\partial z_i^l} = \frac{\partial a_i^l}{\partial z_i^l} \sum_k \frac{\partial z_k^{l+1}}{\partial a_i^l} \frac{\partial C}{\partial z_k^{l+1}}$$

# Compute $\dfrac{\partial C}{\partial w_{ij}^l}$ - Second term

Layer l

$\delta_1^l$
$\delta_2^l$
$\delta_3^l$
$\delta_i^l$

$\delta_1^{l+1}$
$\delta_2^{l+1}$
$\delta_3^{l+1}$
$\delta_n^{l+1}$

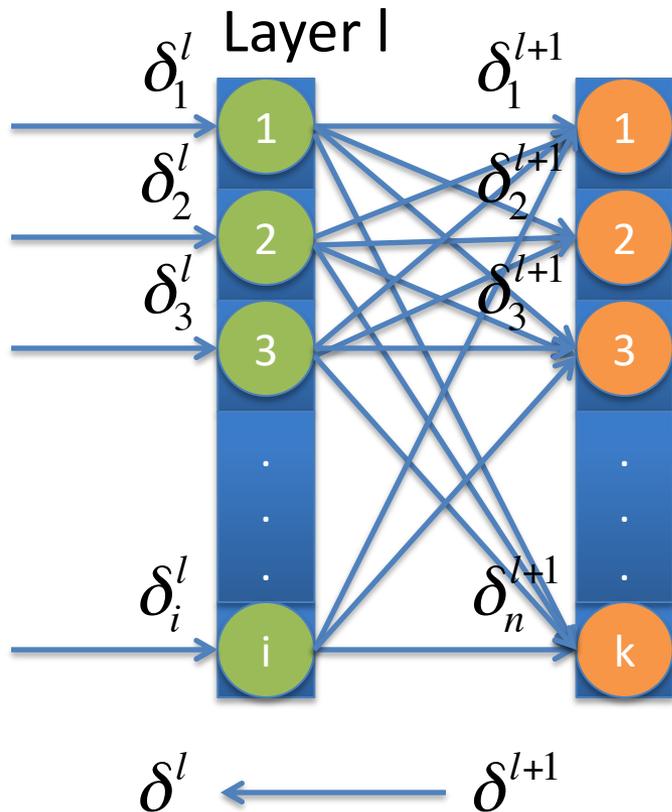$\delta^l \longleftarrow \delta^{l+1}$

$$\delta_i^l = \frac{\partial C}{\partial z_i^l} = \frac{\partial a_i^l}{\partial z_i^l} \sum_k \frac{\partial z_k^{l+1}}{\partial a_i^l} \frac{\partial C}{\partial z_k^{l+1}}$$

$$\sigma'(z_i^l) \qquad \delta_k^{l+1}$$

$$z_k^{l+1} = \sum_i w_{ki}^{l+1} a_i^l + b_k^{l+1}$$

$$\delta_i^l = \sigma'(z_i^l) \sum_k w_{ki}^{l+1} \delta_k^{l+1}$$

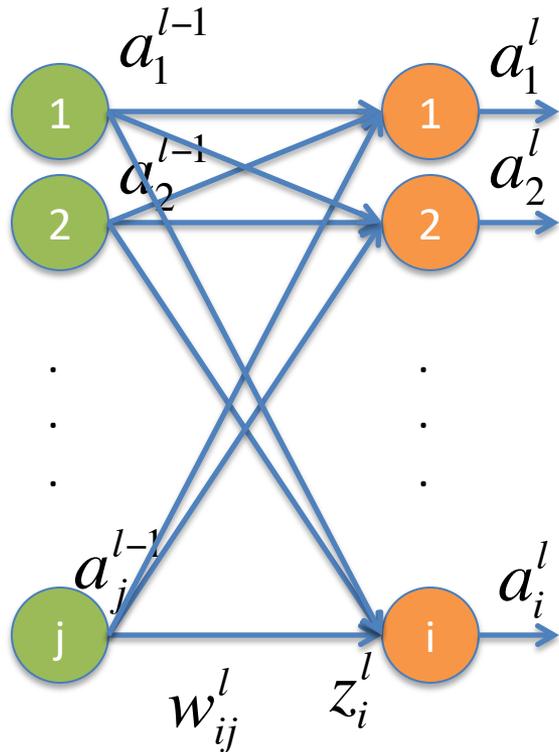Universität Stuttgart

# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term



$$\delta_i^l = \sigma'(z_i^l) \sum_k w_{ki}^{l+1} \delta_k^{l+1}$$

$$\delta^l = \sigma'(z^l).(W^{l+1})^T \delta^{l+1}$$

Universität Stuttgart

# Compute $\dfrac{\partial C}{\partial w_{ij}^l}$

$$\frac{\partial C}{\partial w_{ij}^l} = \frac{\partial z_i^l}{\partial w_{ij}^l} \frac{\partial C}{\partial z_i^l}$$



$$\begin{cases} a_j^{l-1} & l > 1 \\ x_j & l = 1 \end{cases}$$

$$\delta_i^l$$

**Forward pass**
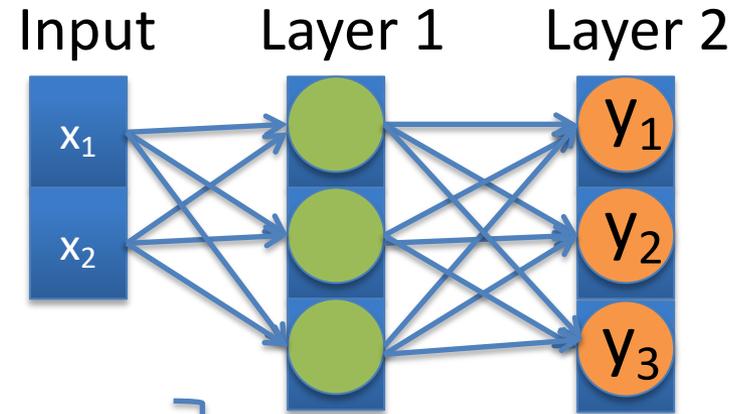
$$z^l = W^l a^{l-1} + b^l$$

$$a^l = \sigma(z^l)$$

**Backward pass**

$$\delta^L = \sigma'(z^L) \nabla C(y)$$

$$\delta^l = \sigma'(z^l)(W^{l+1})^T \delta^{l+1}$$

Universität Stuttgart

# Exercise

- Given the following network
- Layer 1 uses ReLU
- Layer 2 uses Softmax

Input     Layer 1     Layer 2

$x_1$

$x_2$

$y_1$

$y_2$

$y_3$

$$x = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \qquad W^1 = \begin{bmatrix} 1 & 1 \\ -1 & 1 \\ 1 & -1 \end{bmatrix} \qquad W^2 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Ground truth:

$$y = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

$$b^1 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \qquad b^2 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

**Universität Stuttgart**

# Exercise

- Compute the cross entropy loss

- Compute $\dfrac{\delta C}{\delta w_{21}^1}$

Universität Stuttgart

# Live Voting

Universität Stuttgart

# Thanks for listening!

**Universität Stuttgart**