# Feed forward Neural Networks

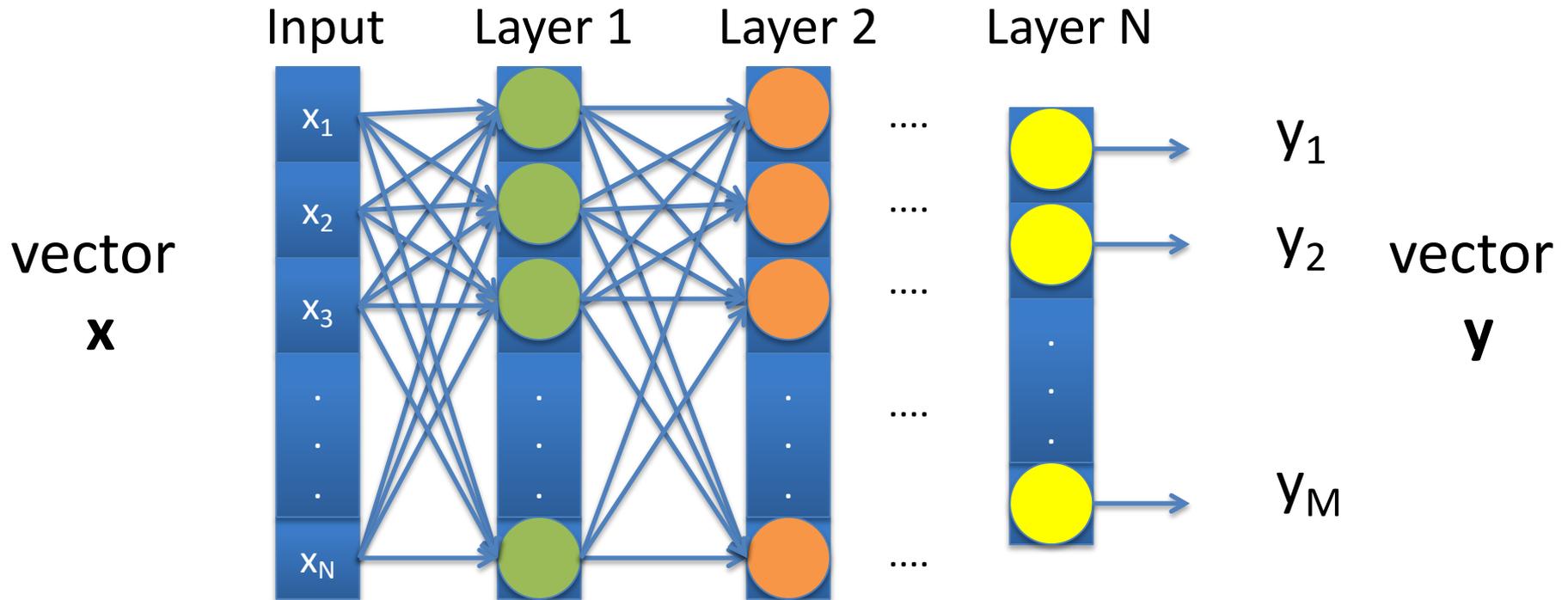**Thang Vu**

**27.11.2025**

**Universität Stuttgart**

# Overview

- Review
  - Backpropagation
  - Cross-entropy loss function

- Activation function
  - Sigmoid
  - Tanh
  - ReLU

- Setting a Network

- Learning: practical usages

Universität Stuttgart

# Computation of the final output



vector **x**

vector **y**

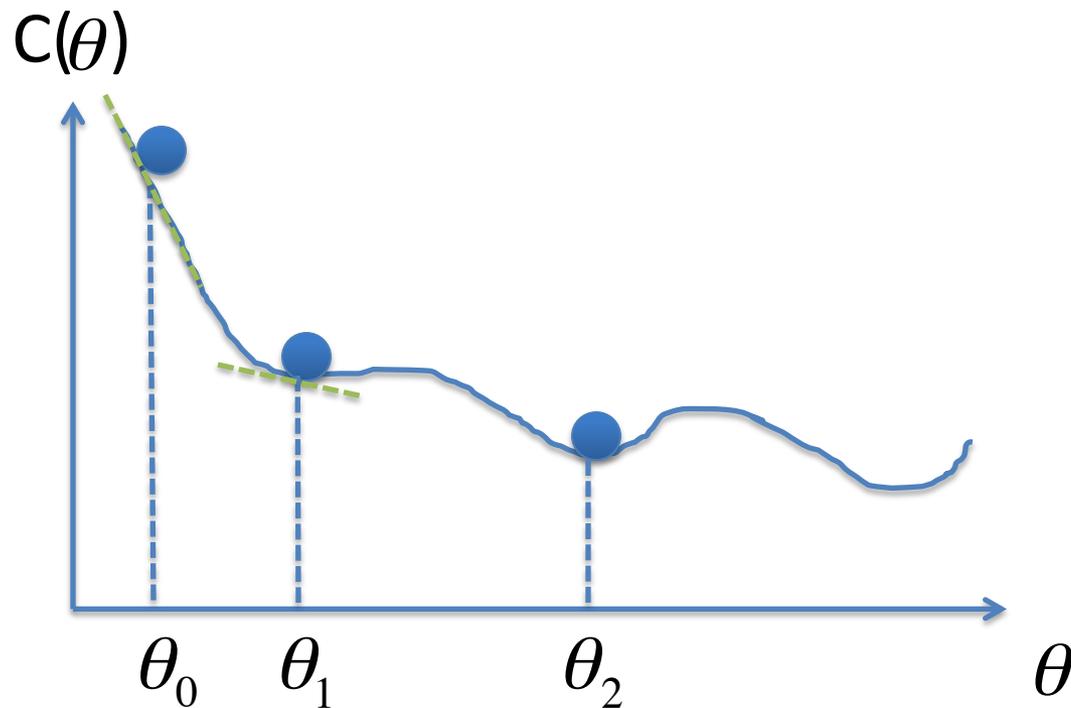$$y = f(x) = \sigma(W^L ... \sigma(W^2 \sigma(W^1 x + b^1) + b^2) ... + b^L)$$

# 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**
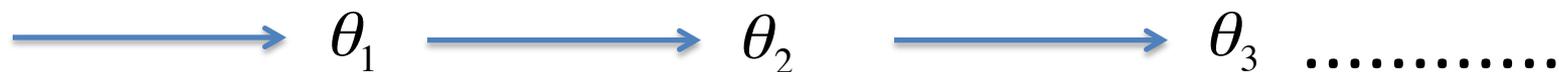
# Gradient descent

- First consider that $\theta$ 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 $dC(\theta_1)/d\theta$
$$\theta_2 \leftarrow \theta_1 - \eta dC(\theta_1)/d\theta$$
- ......

$\theta_0$  $\theta_1$  $\theta_2$  $\theta$

# Gradient descent for Neural Network

- We will do the same thing as presented before

- Starting parameters $\theta_0$

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

- However,

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

- i.e. millions of parameters ☹

- On the other hand, a lot of training data too!!
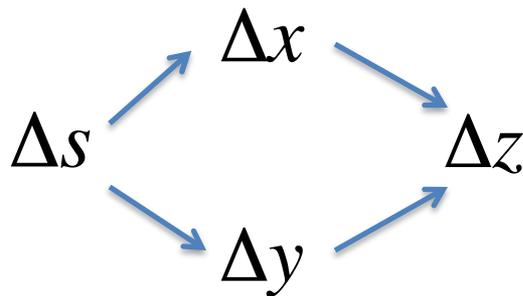
# Backpropagation

Based on the chain rules:

- Case 1:
$$y = g(x)$$
$$z = h(y)$$
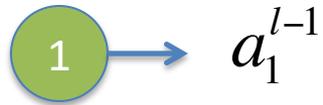$$\Delta x \rightarrow \Delta y \rightarrow \Delta z$$

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

- 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**

# Notation



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

Layer l
$N_l$ nodes

- Output of a neuron:

$$a_i^l \nearrow \text{layer l}$$
$$\searrow \text{neuron i}$$

- Output of one layer:

$a^l$ is a vector

Universität Stuttgart

# Notation



- Weights:

$$w_{ij}^l \nearrow \text{layer } l-1 \text{ to layer } l$$

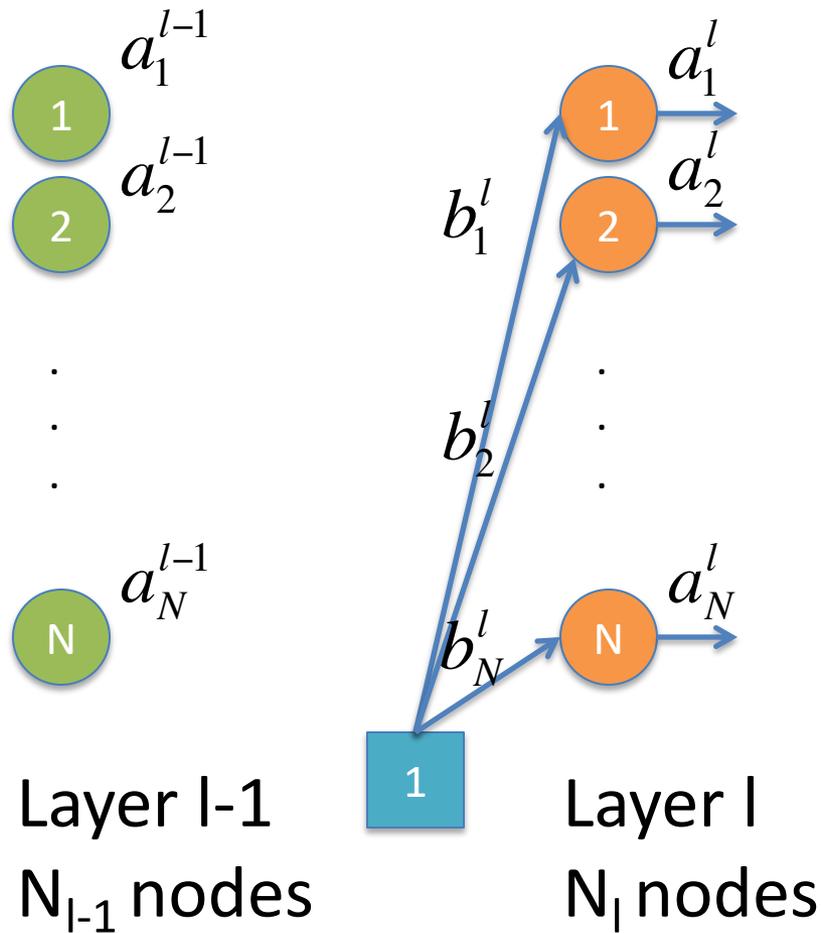$$\searrow \text{from neuron } j \text{ (layer } l\text{ -1)}$$
$$\text{to neuron } i \text{ (layer } l\text{)}$$

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

$$\overbrace{\qquad\qquad}^{N_{l-1}}$$

**Universität Stuttgart**

# Notation



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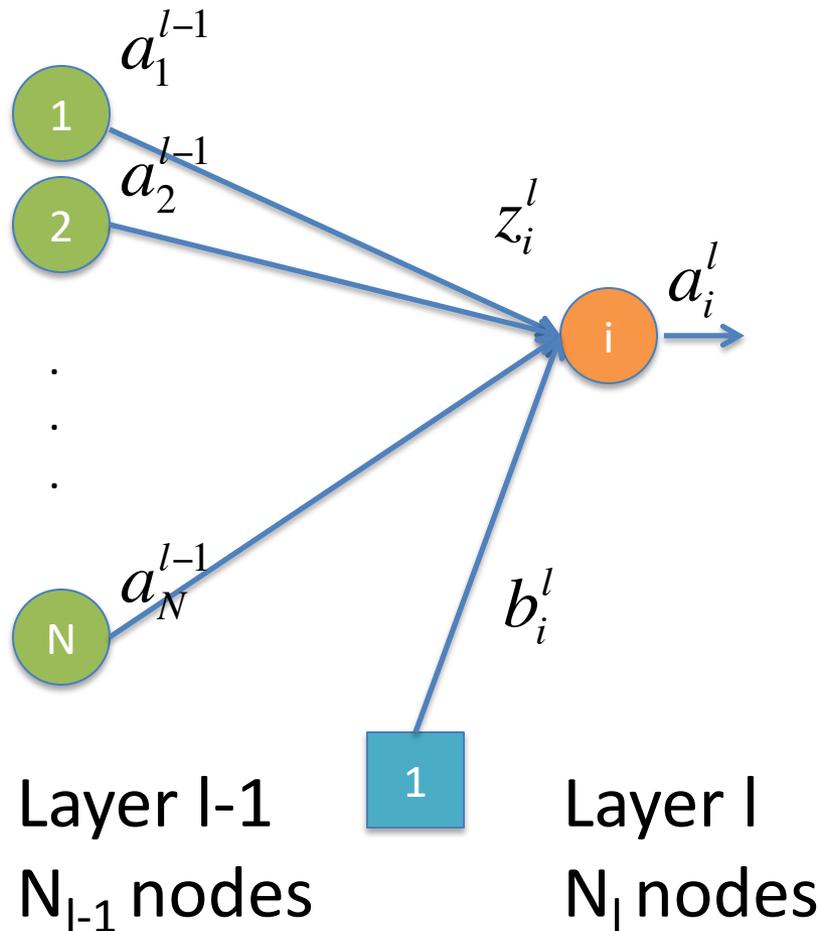$N_{l-1}$ nodes

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

# Notation



$a_1^{l-1}$

$a_2^{l-1}$

$z_i^l$

$a_i^l$

1

2

i

.
.
.

N

$a_N^{l-1}$

$b_i^l$

1

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

Layer l
$N_l$ nodes

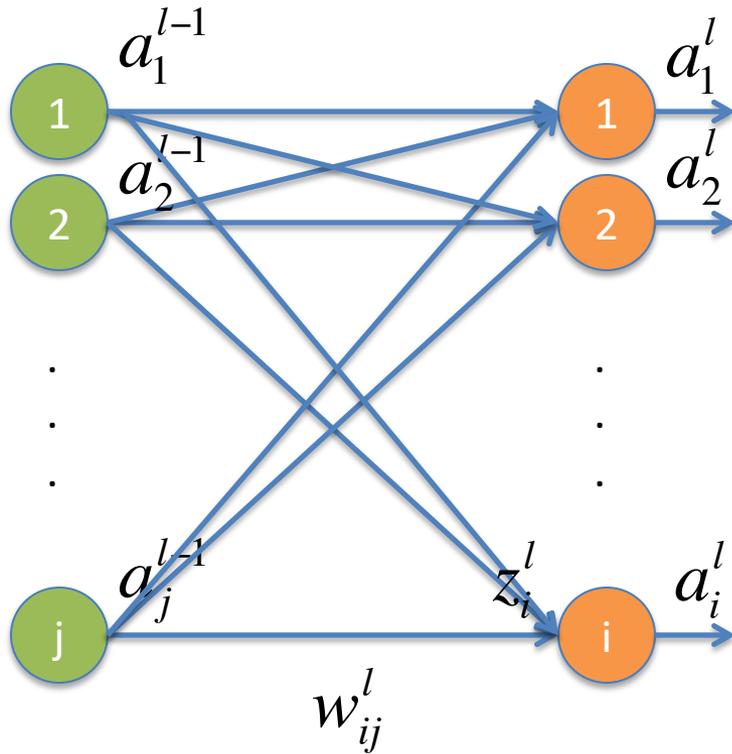$z_i^l$ : input of the activation function for neuron i at layer l

$z^l$ : 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

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

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



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

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

Universität Stuttgart

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

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

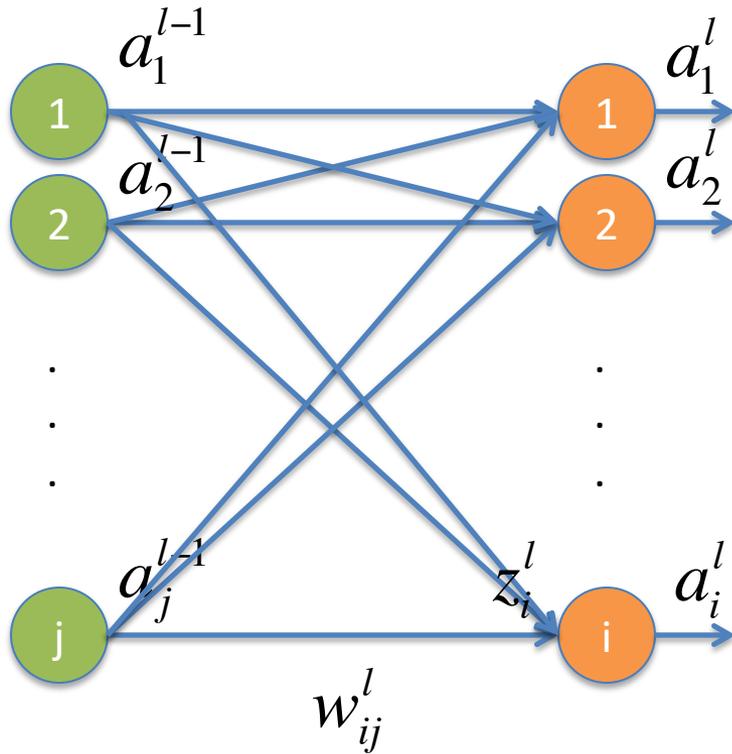$$\frac{\partial z_i^1}{\partial w_{ij}^1} = x_j$$
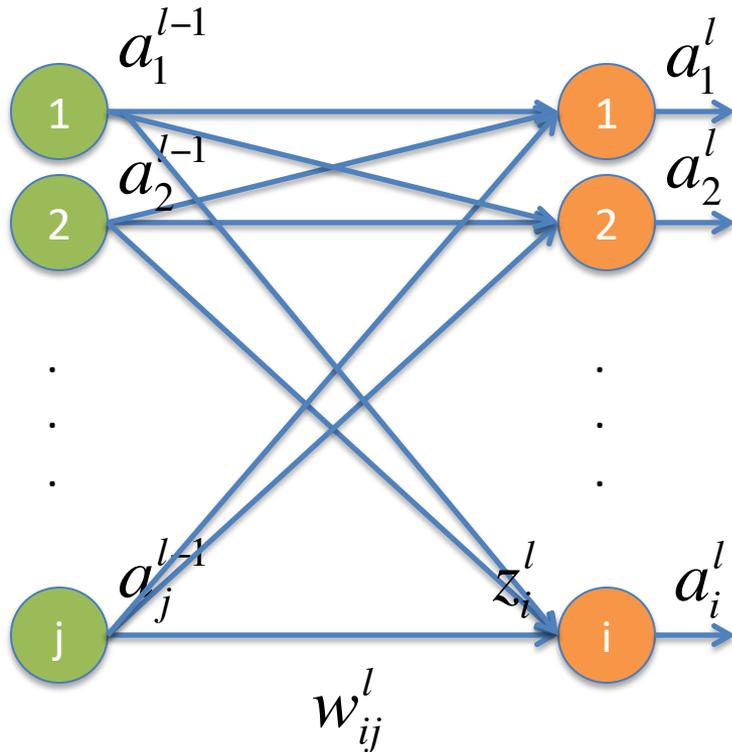
Universität Stuttgart

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

Universität Stuttgart

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

Output layer L

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

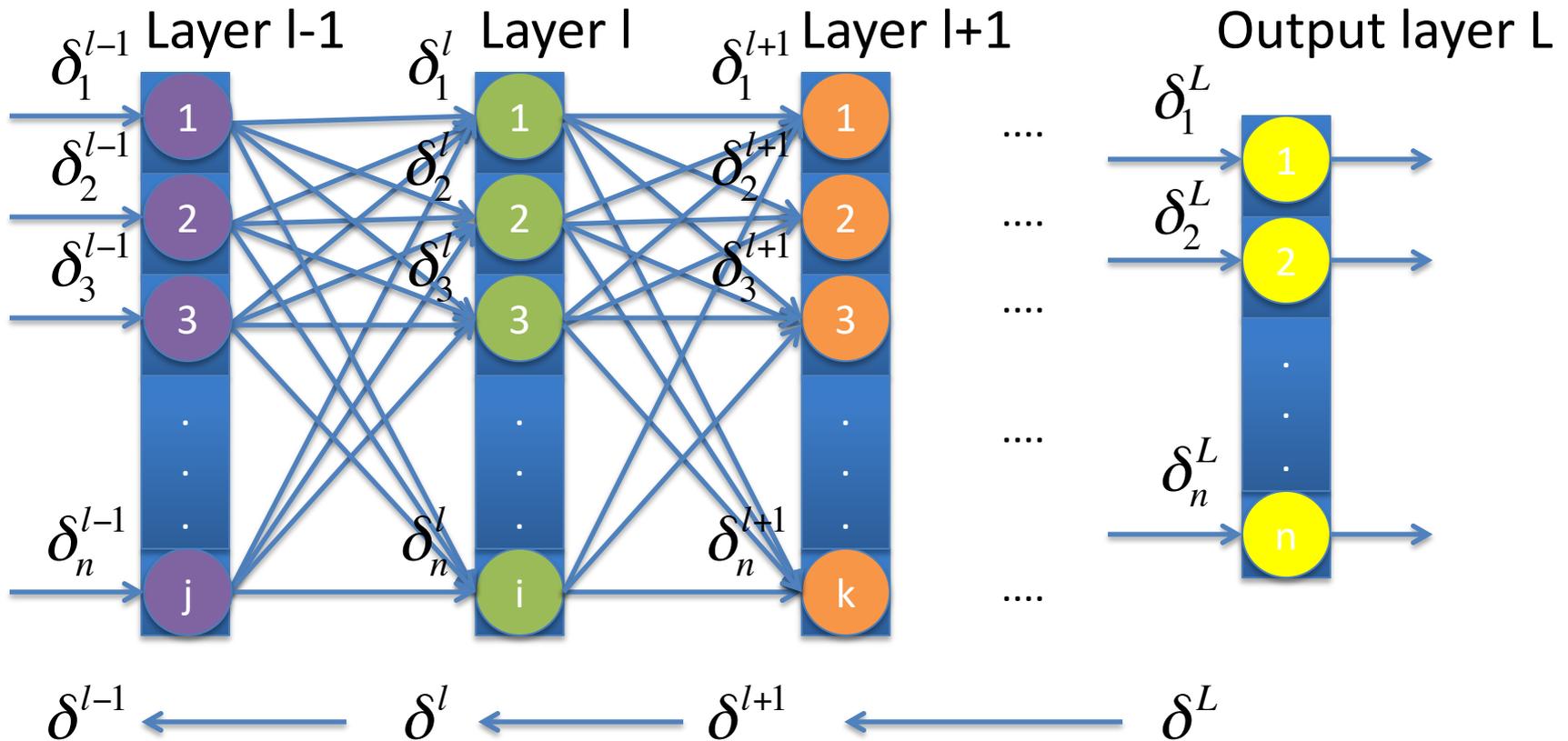$$\delta_n^L = \frac{\partial C}{\partial z_n^L} = \frac{\partial a_n^L}{\partial z_n^L} \frac{\partial C}{\partial a_n^L}$$

$\delta_1^L$

1

$\delta_2^L$

2

.

.

$\delta_n^L$

.

n

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

Depending on
The loss function

$\delta^L$

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

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

Layer l        Layer l+1

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

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

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

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

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

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

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

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

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

$\delta_2^l$

$\delta_3^l$

$\delta_2^{l+1}$

$\delta_3^{l+1}$

$\delta_i^l$

$\delta_n^{l+1}$

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

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

Layer l

$\delta_1^l$

$\delta_2^l$

$\delta_3^l$

$\delta_i^l$

$z_i^l$  $a_i^l$

Layer l+1

$\delta_1^{l+1}$

$\delta_2^{l+1}$

$\delta_3^{l+1}$

$\delta_n^{l+1}$

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

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

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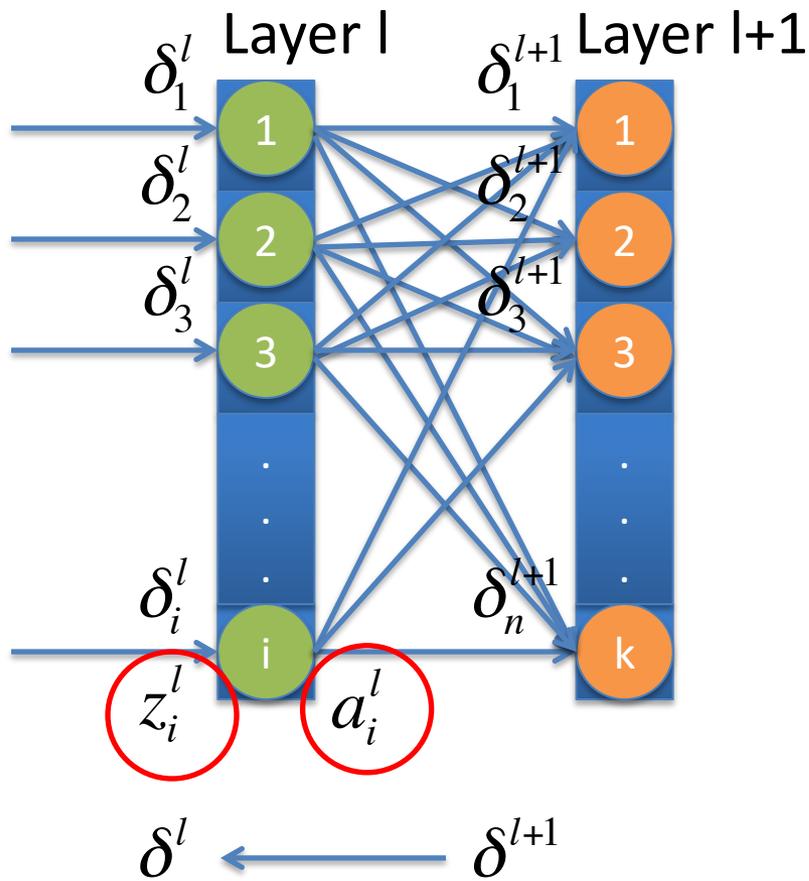
Universität Stuttgart

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



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

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

Universität Stuttgart

# Chain Rules

- Case 1:

$$y = g(x)$$

$$z = h(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

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

- Case 2:

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

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

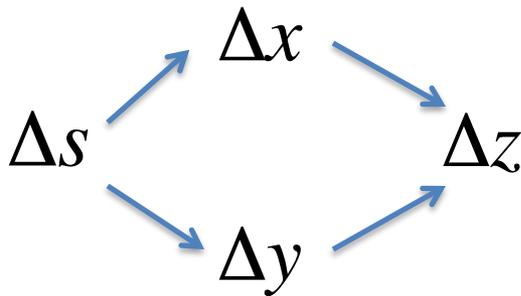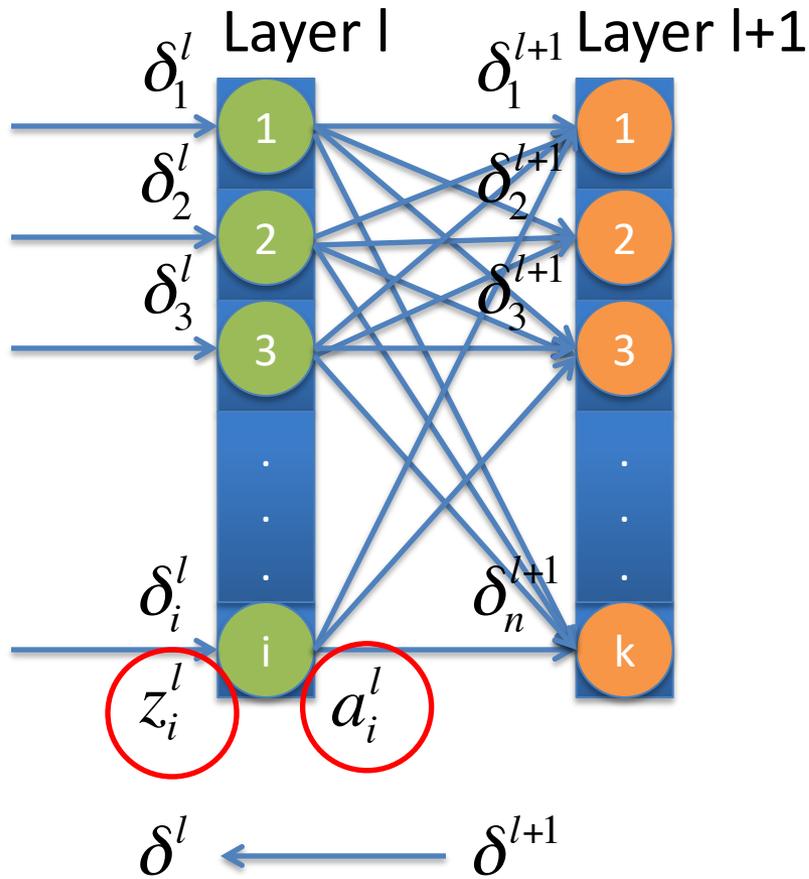# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term

Layer l   Layer l+1

$\delta_1^l$

$\delta_2^l$

$\delta_3^l$

$\delta_1^{l+1}$

$\delta_2^{l+1}$

$\delta_3^{l+1}$

$\delta_i^l$

$\delta_n^{l+1}$

$z_i^l$   $a_i^l$

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

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

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

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

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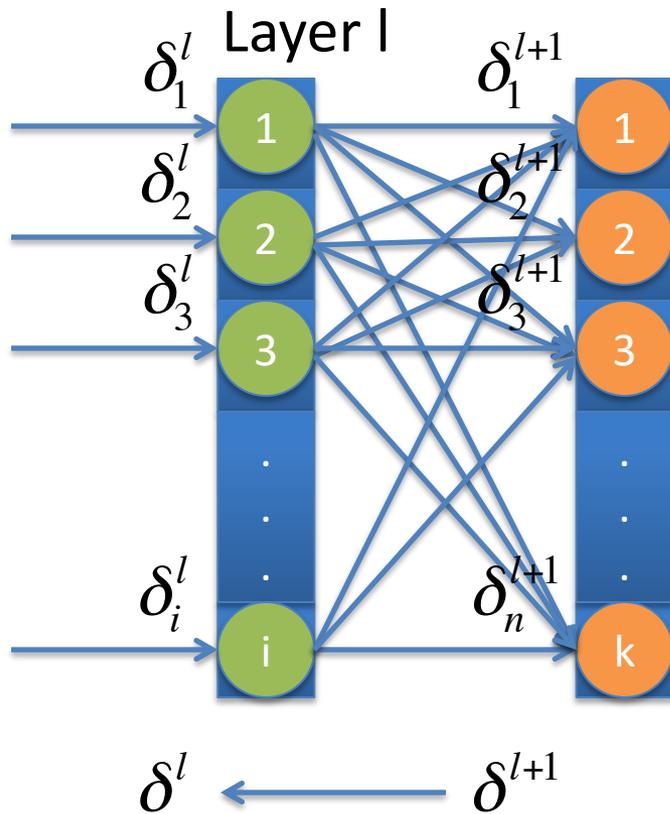
Layer l

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

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

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

$$\delta_i^l = \sigma'(z_i^l) \sum_k w_{ki}^{l+1} \delta_k^{l+1}$$
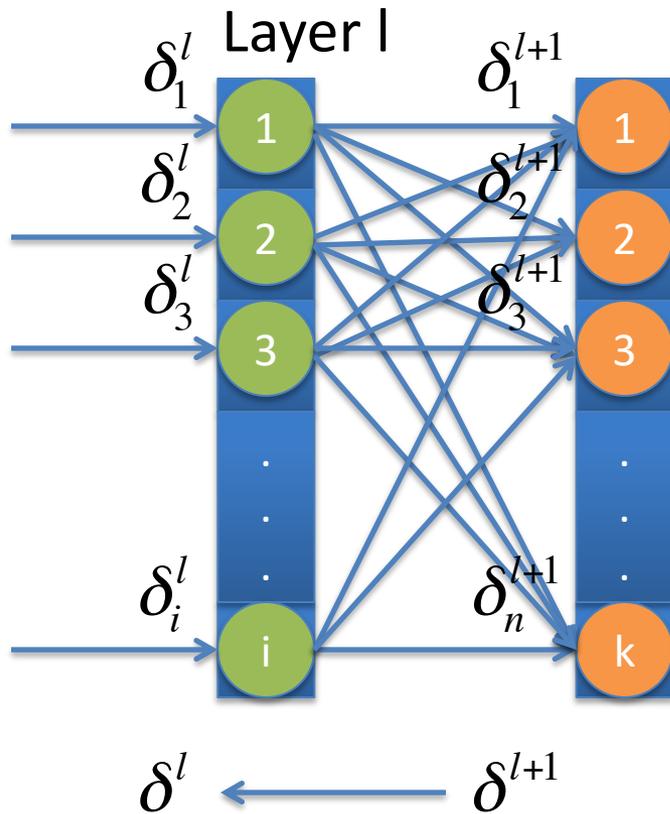
# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term

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

Layer l

$\delta_1^l$

$\delta_2^l$

$\delta_3^l$

$\delta_i^l$

1

2

3

.
.
.

i

$\delta_1^{l+1}$

$\delta_2^{l+1}$

$\delta_3^{l+1}$

$\delta_n^{l+1}$

1

2

3

.
.
.

k

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

Universität Stuttgart

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



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

$$\delta_1^l = \sigma'(z_1^l) \sum_k w_{k1}^{l+1} \delta_k^{l+1}$$

$$\delta_2^l = \sigma'(z_2^l) \sum_k w_{k2}^{l+1} \delta_k^{l+1}$$

$$\delta_3^l = \sigma'(z_3^l) \sum_k w_{k3}^{l+1} \delta_k^{l+1}$$

Layer l

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

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

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

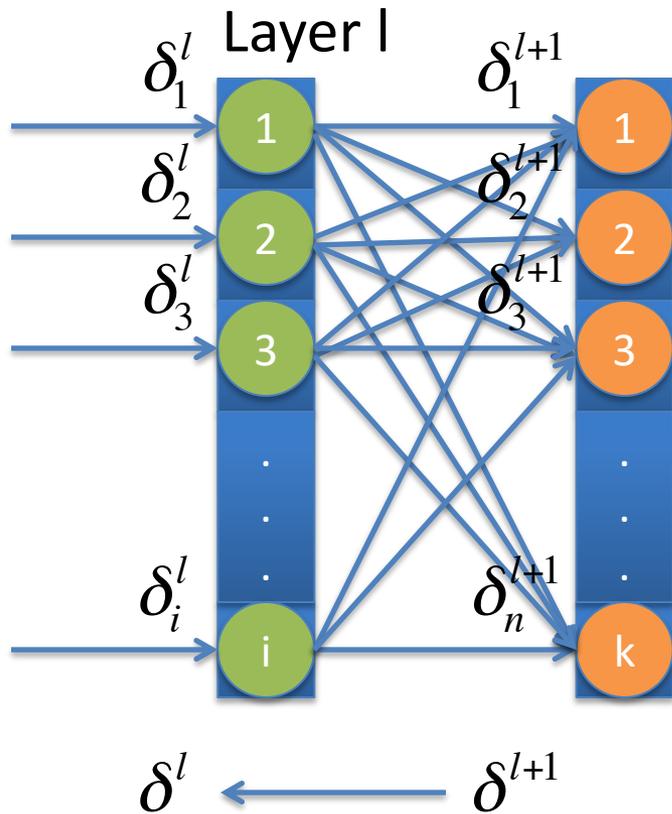# Compute $\dfrac{\partial C}{\partial w_{ij}^l}$ - Second term

Layer l

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

$$\delta_1^l = \sigma'(z_1^l) \sum_k w_{k1}^{l+1} \delta_k^{l+1}$$

$$\delta_2^l = \sigma'(z_2^l) \sum_k w_{k2}^{l+1} \delta_k^{l+1}$$

$$\delta_3^l = \sigma'(z_3^l) \sum_k w_{k3}^{l+1} \delta_k^{l+1}$$

$$W^{l+1} = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & & \end{bmatrix}$$

$N_l$

$N_{l+1}$

**Universität Stuttgart**

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

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

Layer l

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

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

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

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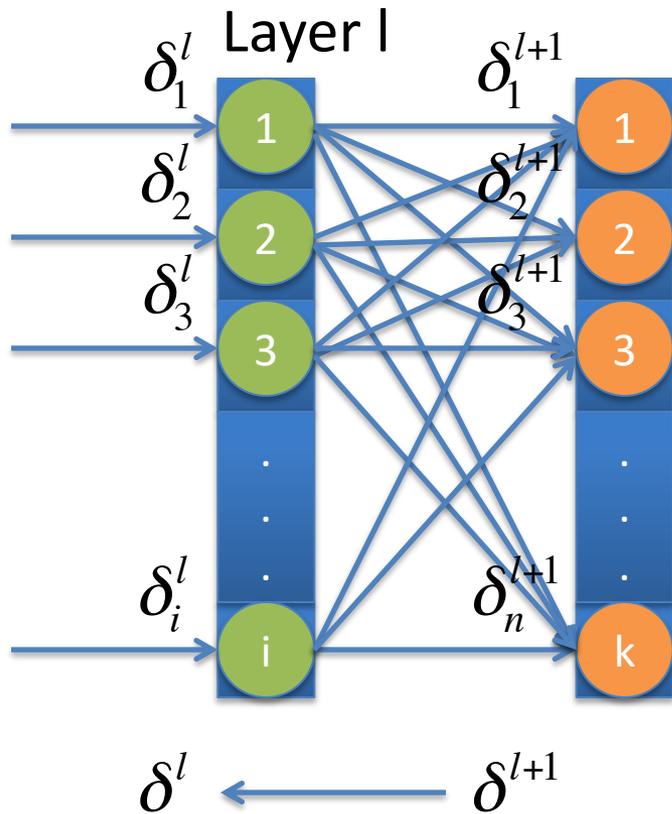
$\delta^l \longleftarrow \delta^{l+1}$

$$\delta_1^l = \sigma'(z_1^l) \sum_k w_{k1}^{l+1} \delta_k^{l+1}$$

$$\delta_2^l = \sigma'(z_2^l) \sum_k w_{k2}^{l+1} \delta_k^{l+1}$$

$$\delta_3^l = \sigma'(z_3^l) \sum_k w_{k3}^{l+1} \delta_k^{l+1}$$

$$W^{l+1} = \begin{bmatrix} w_{11}^l & w_{12}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & & \end{bmatrix}$$

$N_l$

$N_{l+1}$

Universität Stuttgart

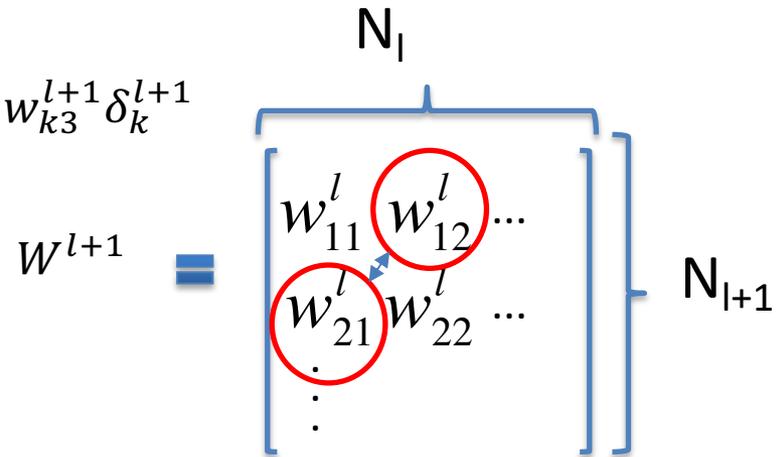# Compute $\frac{\partial C}{\partial w_{ij}^l}$ - Second term



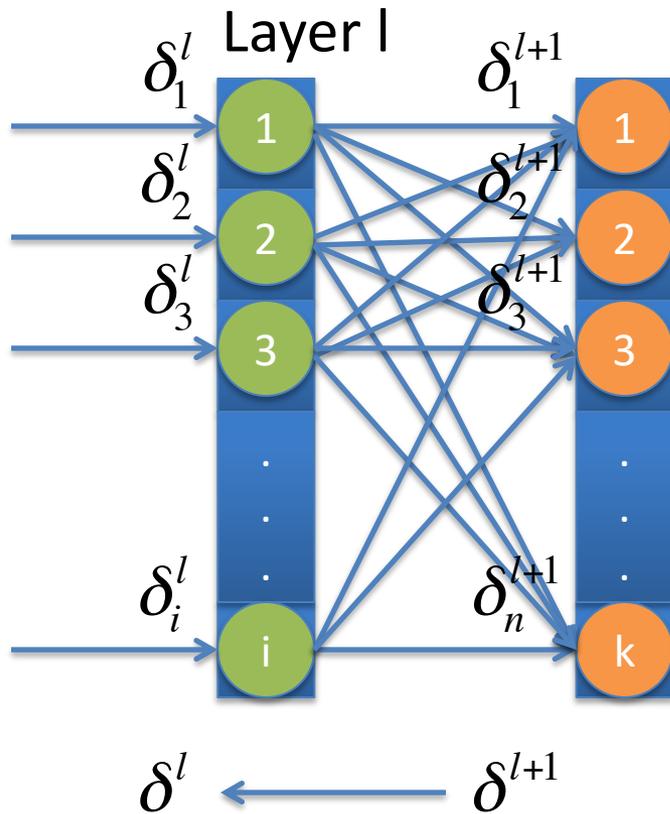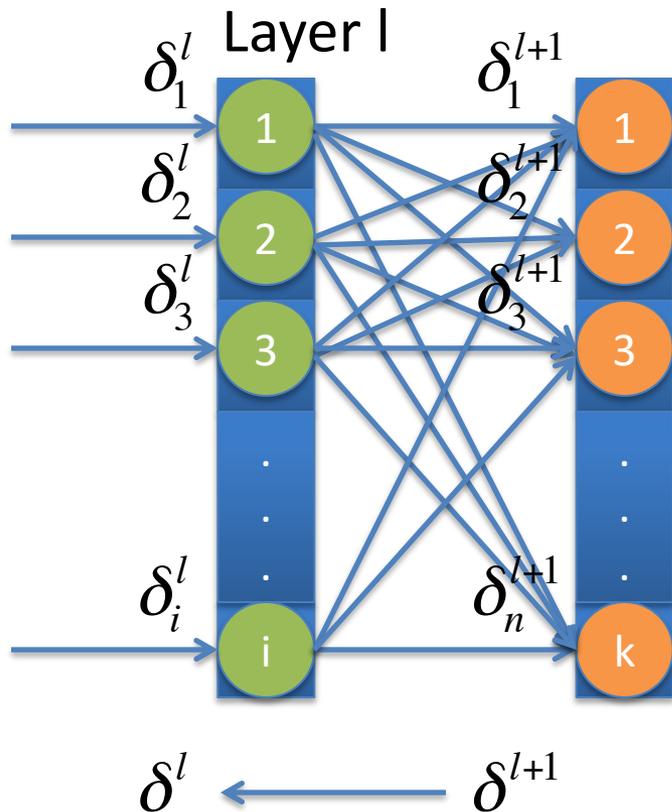$$\delta_i^l = \sigma'(z_i^l)\sum_k w_{ki}^{l+1}\delta_k^{l+1}$$

$$\delta_1^l = \sigma'(z_1^l)\sum_k w_{k1}^{l+1}\delta_k^{l+1}$$

$$\delta_2^l = \sigma'(z_2^l)\sum_k w_{k2}^{l+1}\delta_k^{l+1}$$

$$\delta_3^l = \sigma'(z_3^l)\sum_k w_{k3}^{l+1}\delta_k^{l+1}$$

$N_l$

$$W^{l+1^T} = \begin{bmatrix} w_{11}^l & w_{21}^l & \cdots \\ w_{21}^l & w_{22}^l & \cdots \\ \vdots & & \end{bmatrix}$$
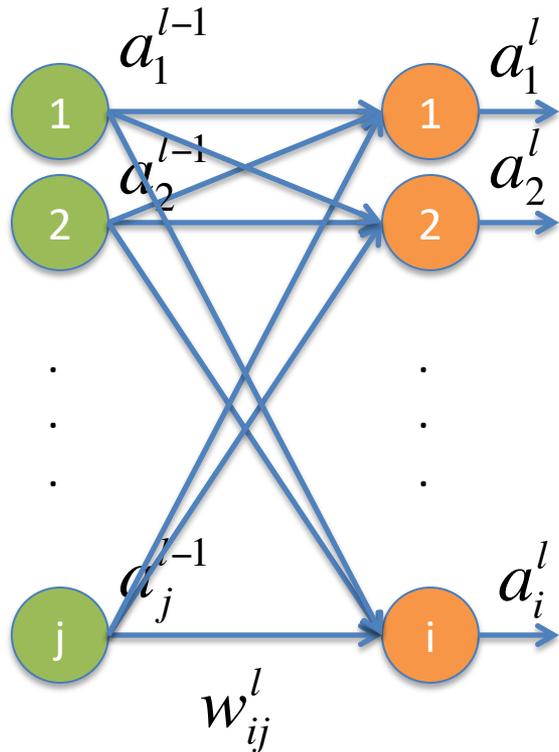
$N_{l+1}$

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

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

$$\left\{ \begin{array}{ll} a_j^{l-1} & l > 1 \\ x_j & l = 1 \end{array} \right.$$

$$\delta_i^l$$

**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

# Loss function

Ouput layer L

$z_1^L$

$\hat{y}_1$

$z_2^L$

$\hat{y}_2$

$\hat{y}_t$

$z_1^L$

$\hat{y}_N$

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

- Sofmax function:

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

- Cross Entropy (CE):

$$c(\theta) = C_x(\theta) = -\log y_t$$

Universität Stuttgart

# Cross Entropy Loss function

Ouput layer L
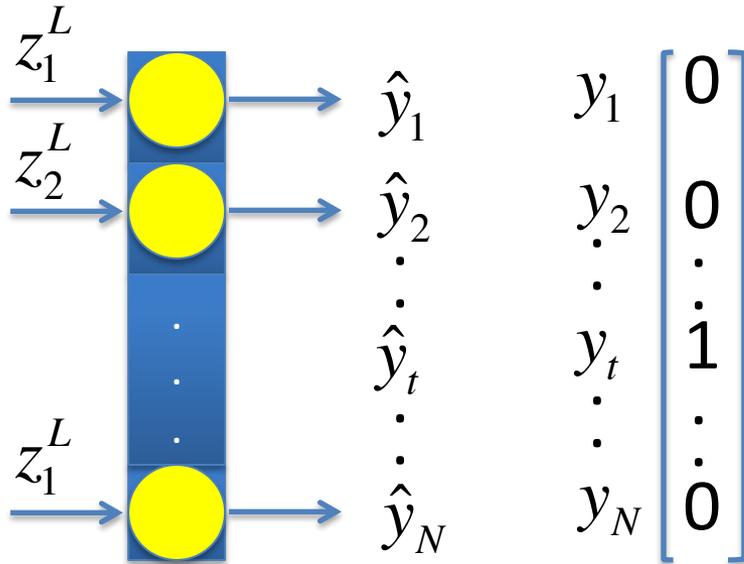
$z_1^L$ → ⬤ → $\hat{y}_1$

$z_2^L$ → ⬤ → $\hat{y}_2$

$\vdots$

→ $\hat{y}_t$

$z_1^L$ → ⬤ → $\hat{y}_N$

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

$$C_x(\theta) = -\log y_t \qquad \delta_i^L = \frac{\delta C_x}{\delta z_i^L}$$

- Case 1: $\quad i = t$

$$\delta_i^L = \frac{\delta C_x}{\delta z_t^L} = \frac{\delta C_x}{\delta y_t}\frac{\delta y_t}{\delta z_t^L}$$

$$= -\frac{1}{\hat{y}_t}\,\hat{y}_t(1-\hat{y}_t) = \hat{y}_t - 1$$

- Case 2: $\quad i \neq t$

$$\delta_i^L = \frac{\delta C_x}{\delta z_i^L} = -\frac{1}{\hat{y}_t}(-\hat{y}_t\hat{y}_i) = \hat{y}_i$$

**Universität Stuttgart**

# Live Voting

**Universität Stuttgart**

# Overview

- Review
  - Backpropagation
  - Cross-entropy loss function

- Activation function
  - Sigmoid
  - Tanh
  - ReLU

- Setting a Network

- Learning: practical usages

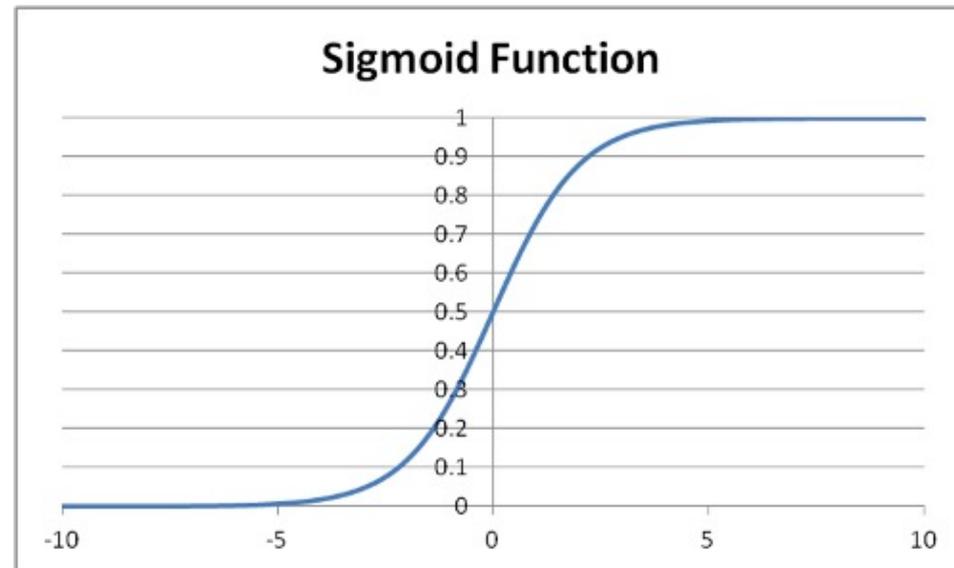Universität Stuttgart

# Activation function

- For the hidden layers:
  - Sigmoid
  - Tanh
  - ReLu and its variation

- For the output layer
  - Softmax

Universität Stuttgart

# Sigmoid function

- Frequently used activation function
- It is still used as baseline system

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

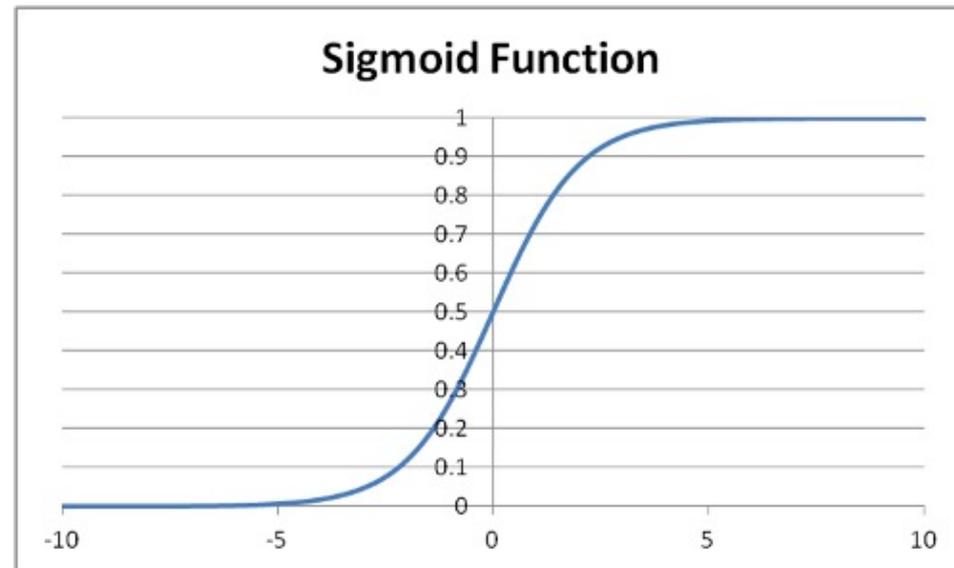$$\sigma'(z) = ?$$



**Sigmoid Function**

**Universität Stuttgart**

# Sigmoid function

- Frequently used activation function
- It is still used as baseline system

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

$$\sigma'(z) = -1.(1+e^{-z})^{-2}(-e^{-z})$$

$$= \frac{e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}}$$

$$= \frac{1}{1+e^{-z}} \cdot (1 - \frac{1}{1+e^{-z}}.) = \sigma(z)(1-\sigma(z))$$



Sigmoid Function

**Universität Stuttgart**

# Tanh function

- In NLP applications, tanh function is more often used than sigmoid function

$$\tanh(z) = \frac{\sinh(z)}{\cosh(z)} = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$
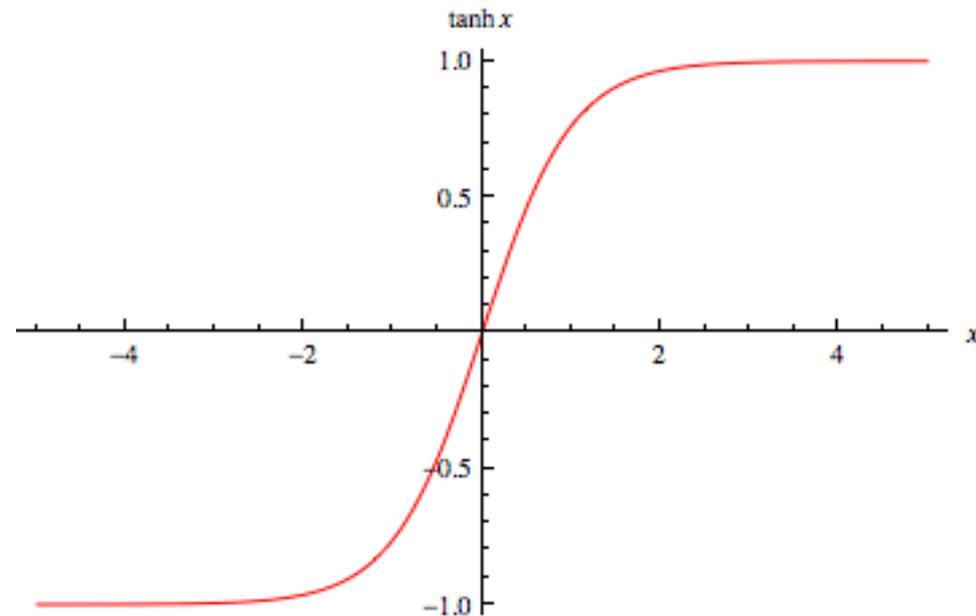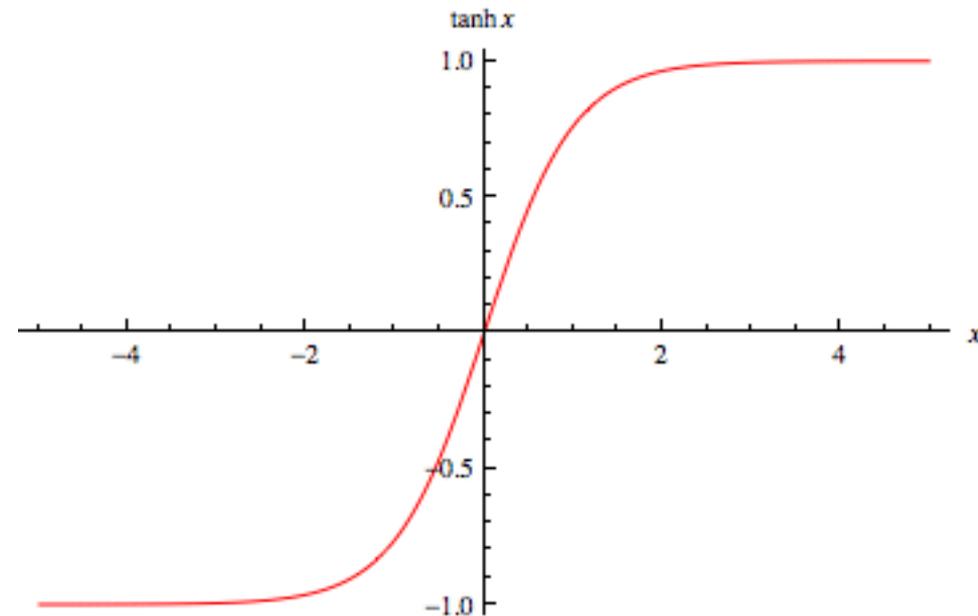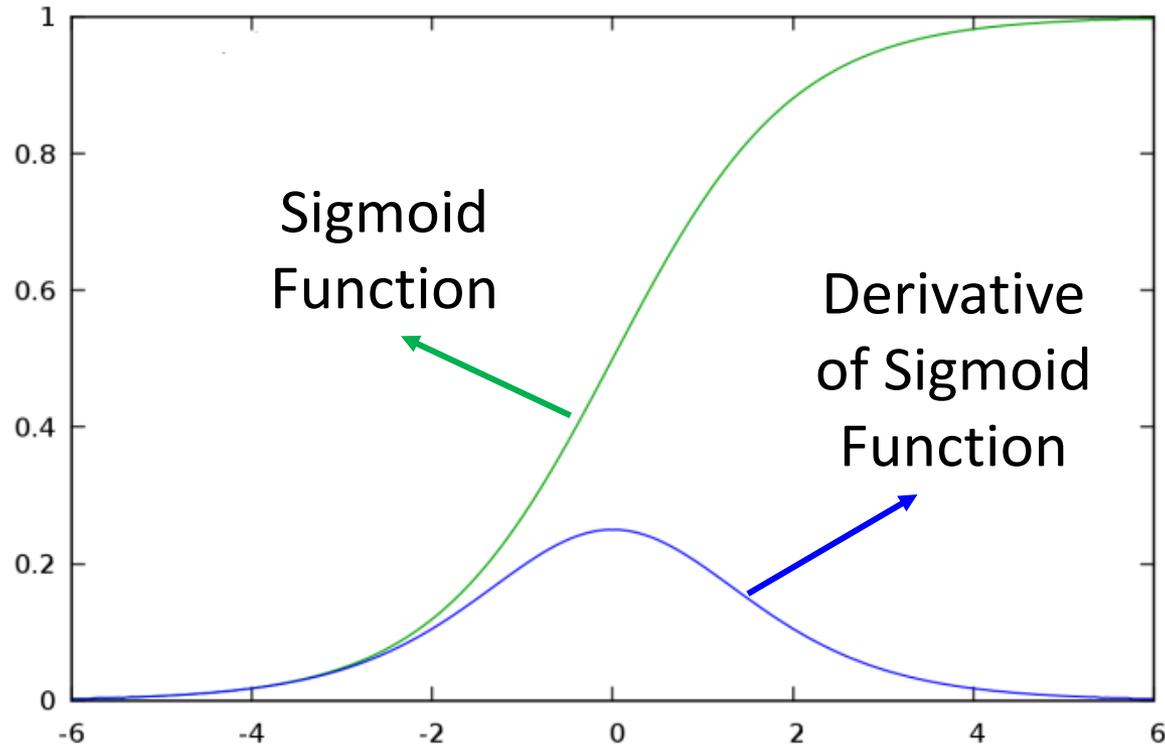
$$\tanh'(z) = ?$$

# Tanh function

- In NLP applications, tanh function is more often used than sigmoid function

$$\tanh(z) = \frac{\sinh(z)}{\cosh(z)} = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

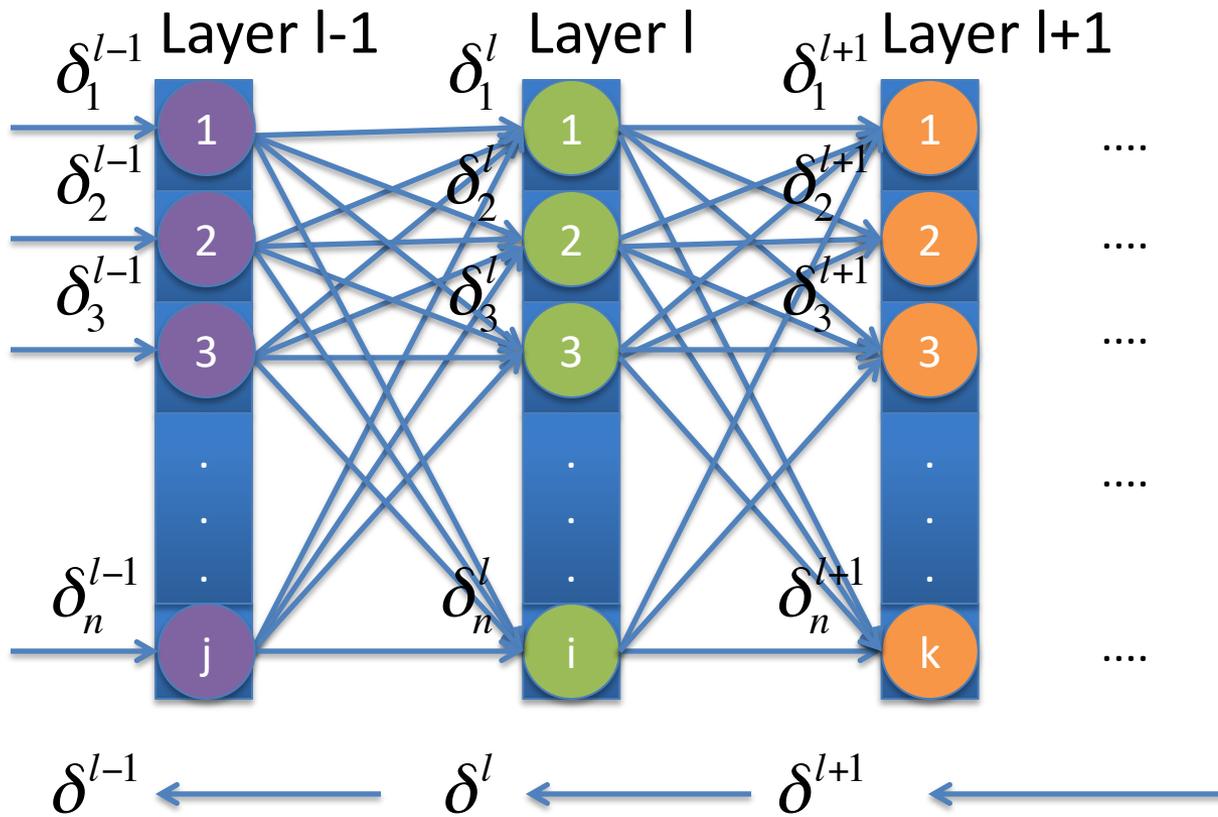$$\tanh'(z) = \ldots = 1 - \tanh^2(z)$$

**Universität Stuttgart**

# Problem of Sigmoid



Sigmoid Function

Derivative of Sigmoid Function

Universität Stuttgart

# Vanishing Gradient Problem

- Backward pass

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

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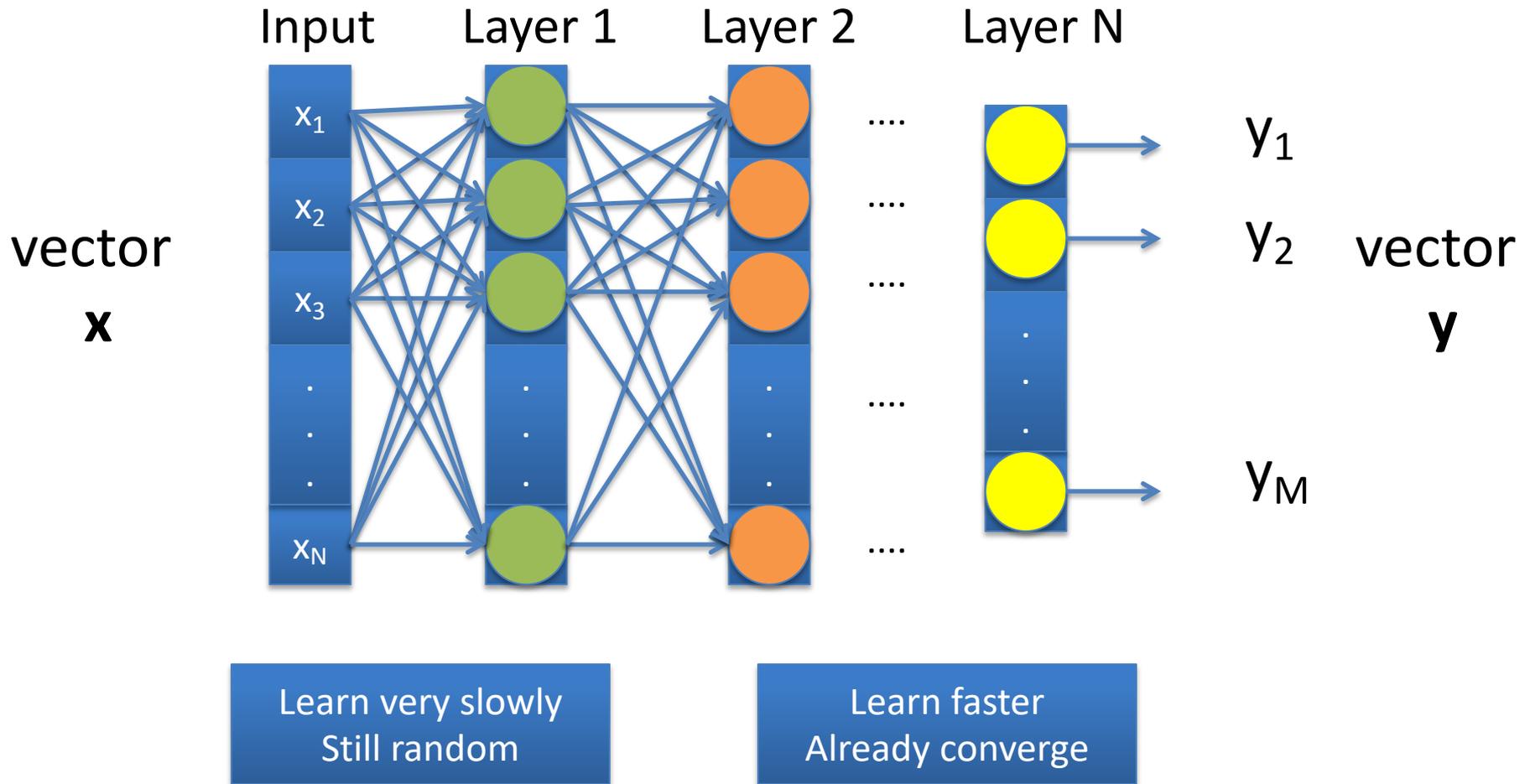
$\delta_1^{l+1}$ $\delta_2^{l+1}$ $\delta_3^{l+1}$ $\delta_n^{l+1}$

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

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

Gradient is getting smaller

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

# Vanishing Gradient Problem

vector **x**

| Input | Layer 1 | Layer 2 | | Layer N | | vector **y** |

Input: $x_1$, $x_2$, $x_3$, ., ., ., $x_N$

.... $y_1$

.... $y_2$

.... $y_M$

Learn very slowly
Still random

Learn faster
Already converge

# Recitifier Linear Unit (ReLU)

$$\sigma(z)$$

$a$

$a = z$

$a = 0$

$z$

$$\sigma'(z)$$

$a$

$1$

$0$

$z$

**Universität Stuttgart**

# ReLU

- Backward Pass

Layer l      Layer l+1

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

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

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

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

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

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

$\sigma(z)$

$a$

$a = z$

$a = 0$

$z$

$\sigma'(z)$

$a$

$1$

$0$

$z$

Universität Stuttgart

# ReLU

- Backward Pass

Layer l          Layer l+1

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

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

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

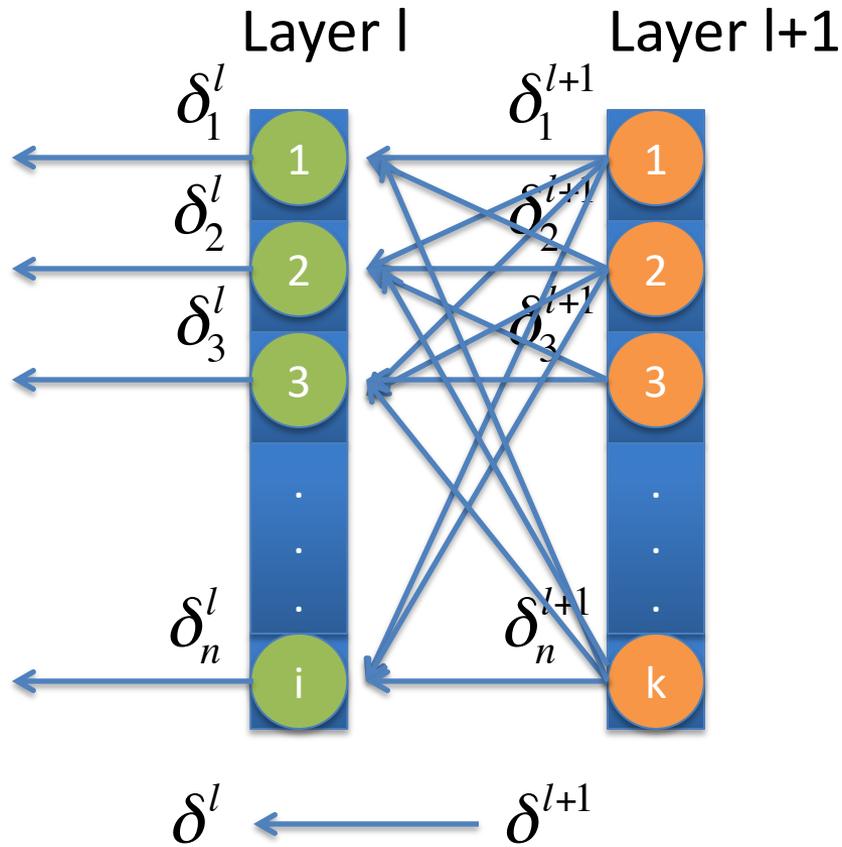$\delta_n^l$                    $\delta_n^{l+1}$

$$\delta^l = \sigma'(z^l)(W^{l+1})^T \delta^{l+1}$$
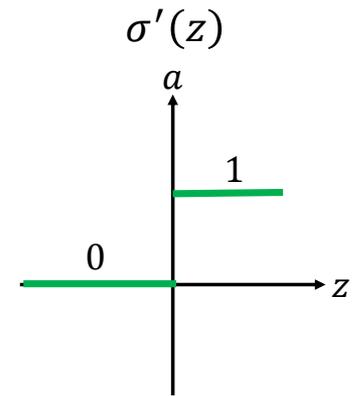
- All the weights connected to not-activated neurons have zero gradient
- Only activated neurons impact the training

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

# ReLU

- Backward Pass

Layer l       Layer l+1

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

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

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

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

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

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

$\sigma(z)$       $\sigma'(z)$

$a = z$

$a = 0$

$1$

$0$

The Dying ReLU Problem

Universität Stuttgart

# ReLU - Variants

*Leaky ReLU*



$$a = z$$
$$a = 0.01z$$

*Parametric ReLU*



$$a = z$$
$$a = \alpha z$$

α also learned by gradient descent

**Universität Stuttgart**

# ReLU - Variants

Exponential Linear Unit (ELU)

Scaled ELU (SELU)

$$a = z$$

$$a = \alpha(e^z - 1)$$

$$a = z$$

$$\times \lambda$$

$$a = \alpha(e^z - 1)$$

$$\times \lambda$$

$$\alpha = 1.6732632423543772848170429916717$$

$$\lambda = 1.0507009873554804934193349852946$$

Universität Stuttgart

# Live Voting

# Overview

- Review
  - Backpropagation
  - Cross-entropy loss function
- Activation function
  - Sigmoid
  - Tanh
  - ReLU
- Setting a Network
- Learning: practical usages

**Universität Stuttgart**

# Computation of the final output



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

**Universität Stuttgart**

# The Input



- For each of the dimension compute the mean $m_i$ $\sigma_i$

$$x_i^r \leftarrow \frac{x_i^r - m_i}{\sigma_i}$$

- Mean = 0
- Variance = 1

Universität Stuttgart

# Setting up the architecture

You need to answer the following questions:

- How many layers?
  - Shallow vs. deep

- How many nodes for each hidden layer?
  - For the output layer, #nodes = #target classes

- Which activation function?
  - Sigmoid vs. Tanh vs. …

# Fat + short vs. Thin + tall ??



Shallow

Deep

Universität Stuttgart

# Fat + short vs. Thin + tall: ASR

- Word error rate (WER)

| L x N | WER (%) |
|-------|---------|
| 1 x 2k | 24.2 |
| 2 x 2k | 20.4 |
| 3 x 2k | 18.4 |
| 4 x 2k | 17.8 |
| 5 x 2k | 17.2 |
| 7x2k | 17.1 |

| 1 x N | WER (%) |
|-------|---------|
| 1 x 3772 | 22.5 |
| 1 x 4,634 | 22.6 |
| 1 x 16k | 22.1 |

Frank Seide, Gang Li and Dong Yu. Conversational Speech Transcriptions Using Context-Dependent Deep Neural Networks. In Proc. of Interspeech, 2011
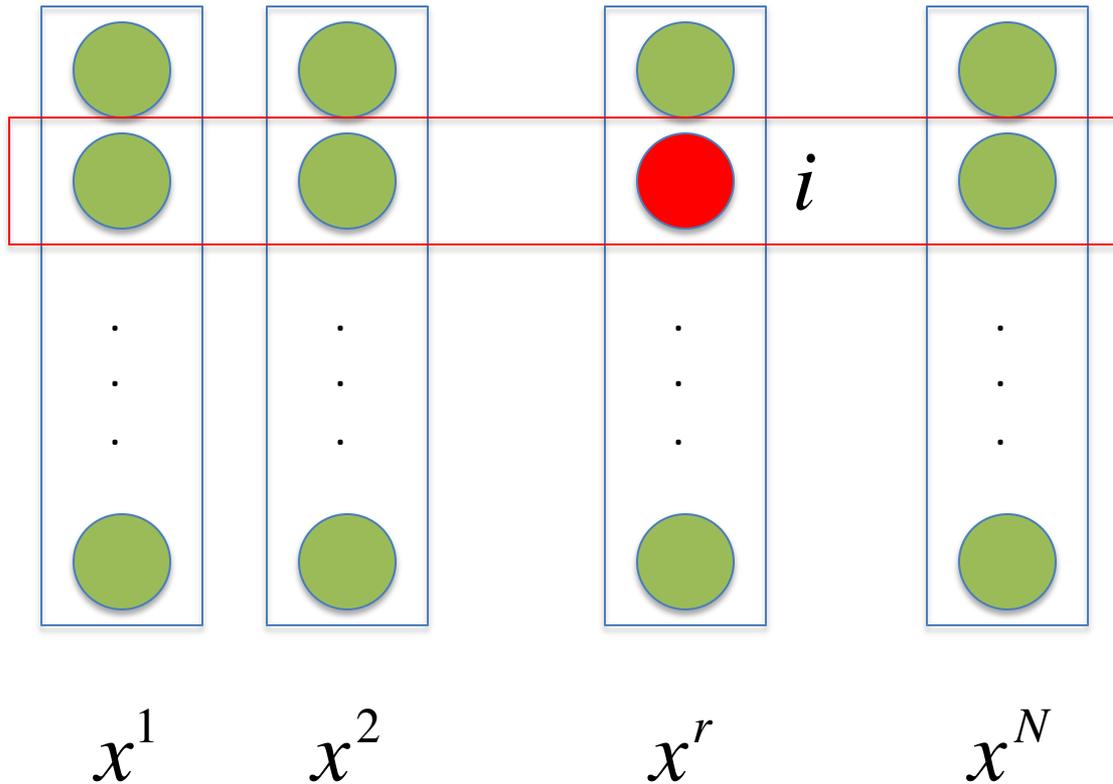
# Overview

- Review
  - Backpropagation
  - Cross-entropy loss function

- Activation function
  - Sigmoid
  - Tanh
  - ReLU

- Setting a Network

- Learning: practical usages

**Universität Stuttgart**

# Learning: Practical Usages

- Parameter Initialization

- Learning Rate

- Stochastic Gradient Descent and Mini-batch

- Learning Recipe

**Universität Stuttgart**

# Parameter Initialization

- For gradient descent, we need to intialize the parameters $\theta_0$
- It has strong impact on the model quality
  - Different parameter initialization $\longrightarrow$ different results

- Some suggestions
  - Do not set the parameters equally
  - Set the parameters randomly, e.g. uniform (-r,r) distribution
$$r = \sqrt{6 \big/ (fan\_in + fan\_out)}$$
  - Train several models with different parameter initializations and combine them

Universität Stuttgart

# Stochastic Gradient Descent

- Gradient Descent:

$$\theta_i \leftarrow \theta_{i-1} - \eta \nabla C(\theta_{i-1})$$

  – In which

$$\nabla C(\theta_{i-1}) = \frac{1}{R} \sum_r \nabla C^r(\theta)$$

- Stochastic Gradient Descent:

  – Pick an example $x^r$

$$\theta_i \leftarrow \theta_{i-1} - \eta \nabla C^r(\theta_{i-1})$$

Universität Stuttgart

# Stochastic Gradient Descent

- *What is an epoch?*

- Training data: $(x_1, y_1), (x_2, y_2), ..., (x_R, y_R)$

- When using the stochastic gradient descent:

  - Starting with $\theta_0$

  - Pick $(x_1, y_1)$ $\qquad \theta_1 \leftarrow \theta_0 - \eta \nabla C^1(\theta_0)$

    $(x_2, y_2)$ $\qquad \theta_2 \leftarrow \theta_1 - \eta \nabla C^2(\theta_1)$ $\qquad$ Seen all the training data

    $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ One epoch

    $(x_R, y_R)$ $\qquad \theta_R \leftarrow \theta_{R-1} - \eta \nabla C^R(\theta_{R-1})$

$(x_1, y_1)$ $\qquad \theta_{R+1} \leftarrow \theta_R - \eta \nabla C^{R+1}(\theta_R)$

Universität Stuttgart

# Stochastic Gradient Descent

- Mini-batch Gradient Descent:
  - Pick B examples as a batch b
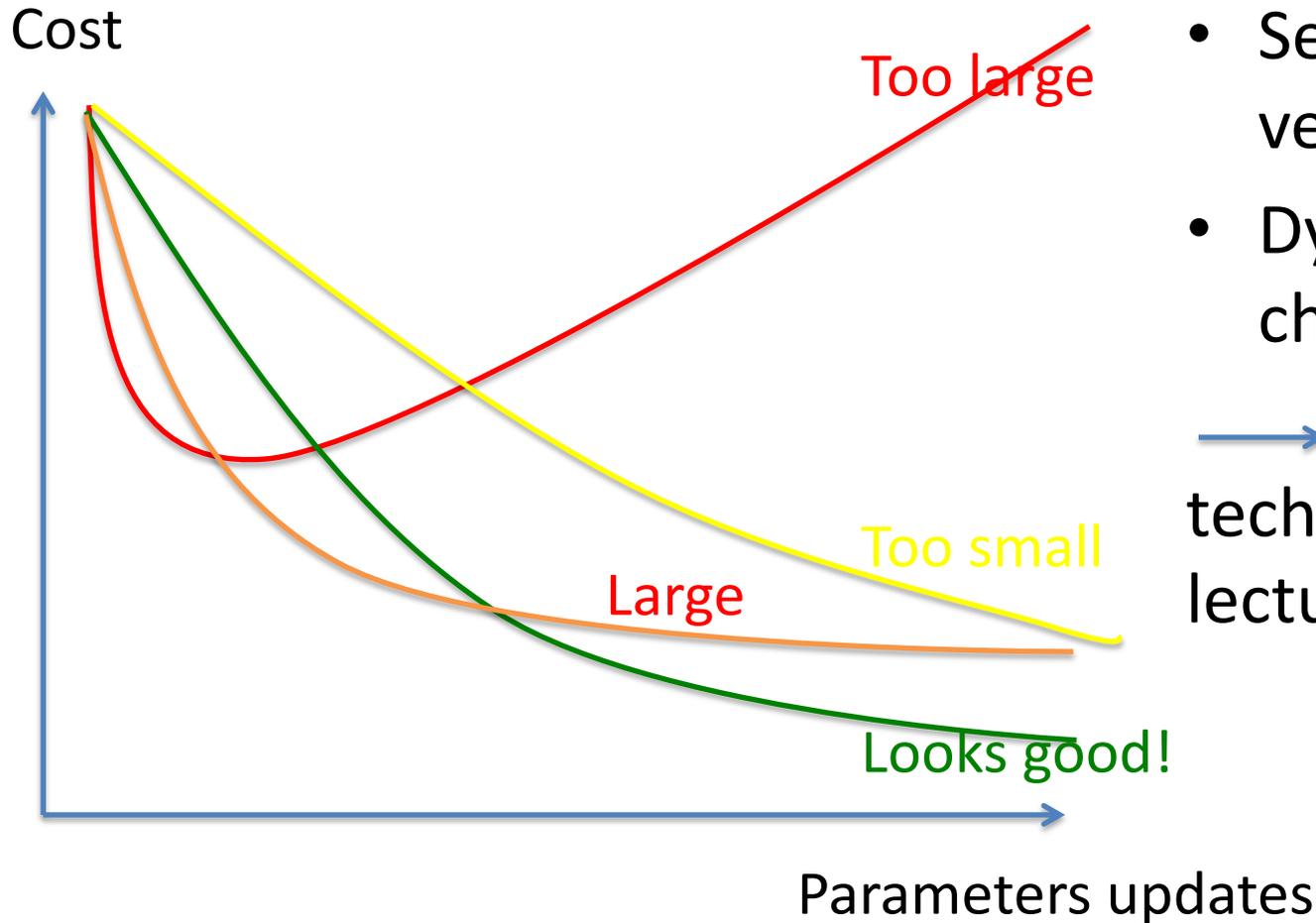  - B is the batch size

$$\theta_i \leftarrow \theta_{i-1} - \eta \frac{1}{B} \sum_{x_r \in b} \nabla C^r(\theta_{i-1})$$

- Mini-batch Gradient Descent is faster than Stochastic Gradient Descent
  - Less updates
  - Better parallelization
- <u>Important:</u> Shuffle the data after each epoch

**Universität Stuttgart**

# Learning rate

Cost

Too large

Too small

Large

Looks good!

Parameters updates

- Set the learning rate very carefully

- Dynamically changed

  → More advanced techniques in the next lectures

# Learning Rate

- Popular & Simple idea: Reduce the learning rate by some factor every few epochs
  - At the beginning, larger learning rate
  - After several epochs, reduce the learning rate

$$\eta^t = \eta / (t + 1)$$

When to reduce the learning rate?

How much should we reduce the learning rate?

Universität Stuttgart

# Learning Recipe

- Split the data in three parts

| Training Data | Cross Validation Data | Evaluation Data |
|---|---|---|

Used for training and monitoring the performance

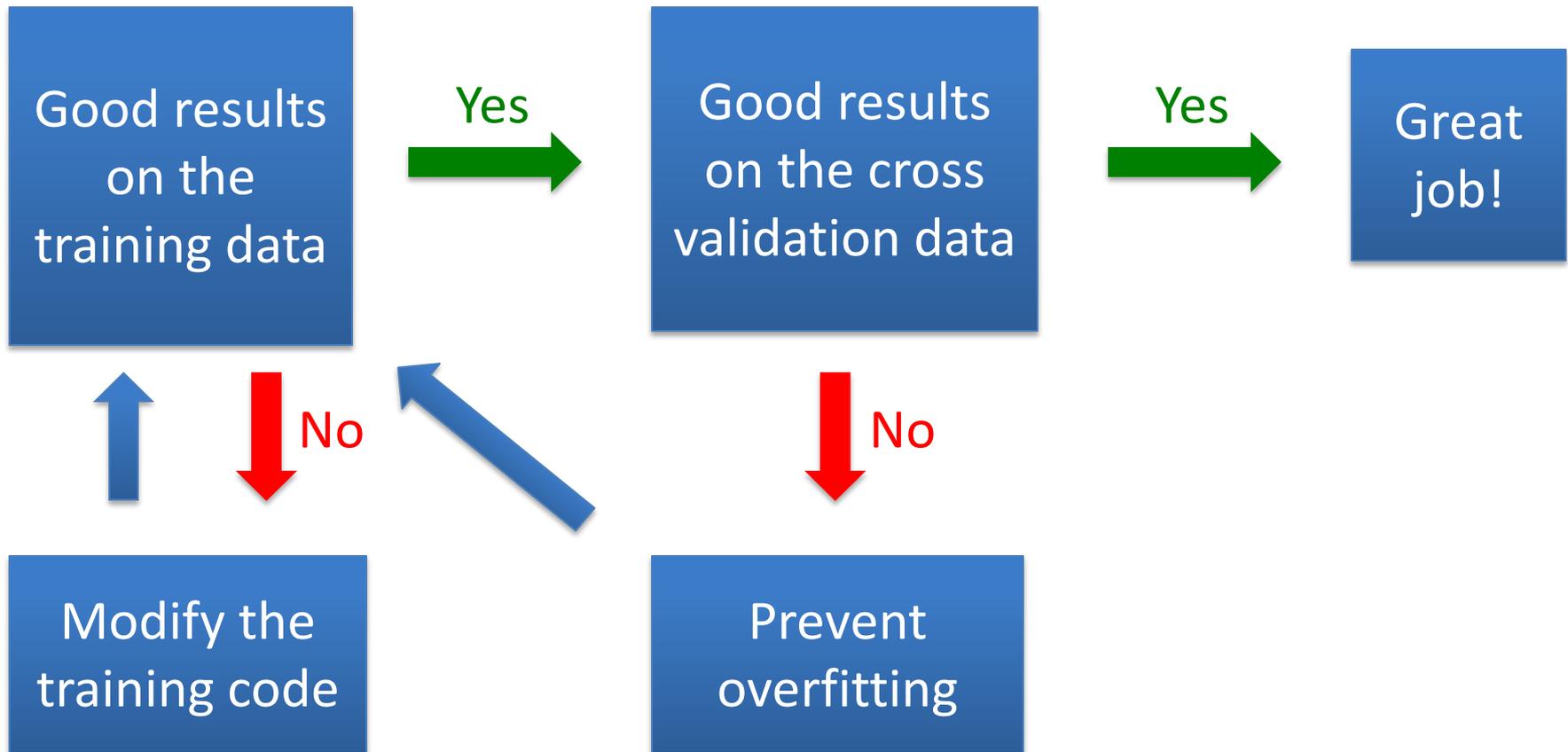Used for monitoring the performance and tuning parameters

Don't touch it till the deadline

Universität Stuttgart

# Learning Recipe

- Monitor the cost function after each update

- Evaluate your model not only on the cross validation but also on the training data
  - The model should work well at first on the training data

- The learning rate can be adjusted based on the performance of the cross validation set
  - Decrease the learning rate if no improvement can be observed on the cross validation set

Universität Stuttgart

# Learning Recipe

Universität Stuttgart

# Thanks for listening!

**Universität Stuttgart**