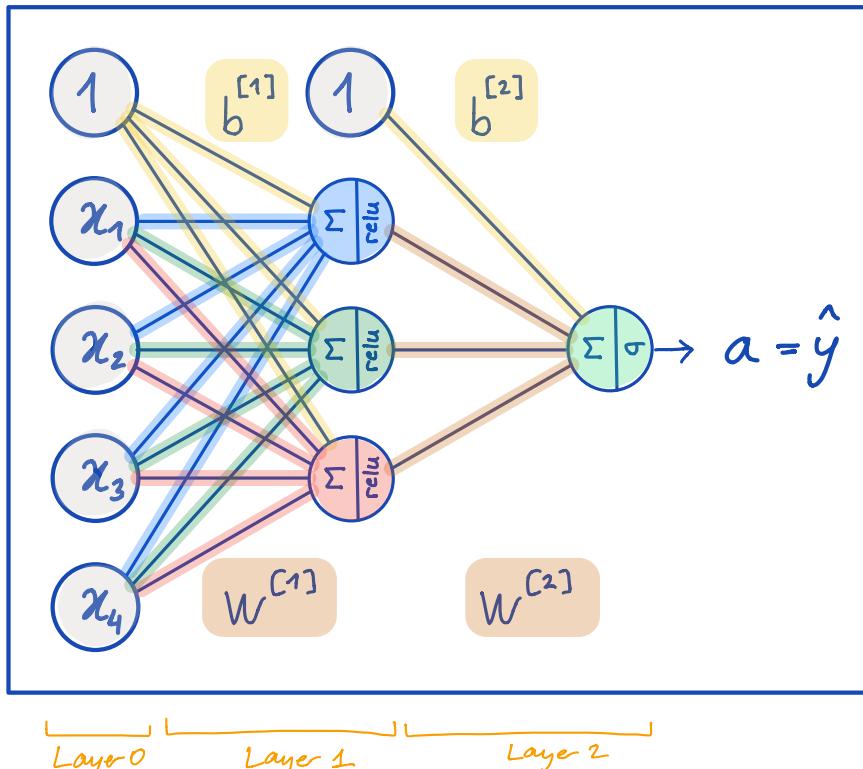


# BINÄRE KLASSEFIKASYON MIT DNN



$$\begin{bmatrix} \text{blue} \\ \text{green} \\ \text{red} \end{bmatrix} \cdot \begin{bmatrix} \text{grey} \end{bmatrix} + \begin{bmatrix} \text{yellow} \end{bmatrix} = \begin{bmatrix} \text{blue} \\ \text{green} \\ \text{red} \\ \text{yellow} \end{bmatrix} = \underline{z}^{[1]}$$

$$\text{relu}\left(W^{[1]} \cdot x + b^{[1]}\right) = a^{[1]}$$

$$\begin{bmatrix} \text{grey} \end{bmatrix} \cdot \begin{bmatrix} \text{blue} \\ \text{green} \\ \text{red} \end{bmatrix} + \begin{bmatrix} \text{yellow} \end{bmatrix} = \begin{bmatrix} \text{green} \end{bmatrix} = \underline{z}^{[2]}$$

$$\sigma\left(W^{[2]} \cdot a^{[1]} + b^{[2]}\right) = a^{[2]}$$

Layer 0      Layer 1      Layer 2

*Broadcasting*

$$\text{relu}\left(\begin{bmatrix} \text{blue} \\ \text{green} \\ \text{red} \end{bmatrix} \cdot \begin{bmatrix} \text{grey} & \text{grey} & \text{grey} & \dots & \text{grey} & \text{grey} & \text{grey} \end{bmatrix}_{(3 \times 7)} + \begin{bmatrix} \text{yellow} \end{bmatrix}_{3 \times 7}\right) = \begin{bmatrix} \text{blue} \\ \text{green} \\ \text{red} \\ \text{yellow} \end{bmatrix}_{(3 \times 7)}$$

$$\text{relu}\left(W^{[1]} \cdot x + b^{[1]}\right) = A^{[1]}$$

*Broadcasting*

$$\sigma\left(\begin{bmatrix} \text{grey} \end{bmatrix} \cdot \begin{bmatrix} \text{blue} & \text{blue} & \text{blue} & \dots & \text{blue} & \text{blue} & \text{blue} \end{bmatrix}_{(1 \times 7)} + \begin{bmatrix} \text{yellow} \end{bmatrix}_{(1 \times 7)}\right) = \begin{bmatrix} \text{green} & \text{green} & \text{green} & \dots & \text{green} & \text{green} & \text{green} \end{bmatrix}_{(1 \times 7)}$$

$$\sigma\left(W^{[2]} \cdot A^{[1]} + b^{[2]}\right) = A^{[2]} = \hat{y}$$

# CROSS-ENTROPY LOSS

$$A^{[2]} = \hat{y} = [\text{●●●●●●●}] \quad \text{IST}$$

$$y = [\text{●●●●●●●}] \quad \text{SOLL}$$

---

$$\text{LOSS} = \text{DISTANCE}(\hat{y}, y) = [ \begin{array}{ccccccc} (1) & (2) & (3) & \dots & (6) & (7) \\ \text{●} & \text{●} & \text{●} & \text{●} & \text{●} & \text{●} & \text{●} \end{array} ]$$

- $\mathcal{L}(\text{●}, \text{●}) = \mathcal{L}(\hat{y}, y) = -y \cdot \log(\hat{y}) - (1-y) \cdot \log(1-\hat{y}) = \text{●}$

- $\mathcal{L}([\text{●●●●●●●}], [\text{●●●●●●●}])$

$$= \mathcal{L}(\hat{y}, y) = -Y * \log(\hat{y}) - (1-Y) * \log(1-\hat{y})$$

$$= [\text{●●●●●●●}]$$

elementweises Multiplizieren

- $Y * \log(\hat{y}) = [\text{●●●●●●●}] * \log([\text{●●●●●●●}])$

$$(1 \times 7) = (1 \times 7) * (1 \times 7)$$

# KOSTENFUNKTION

$$J(\theta) = \frac{1}{m} \cdot \text{np.sum}(\mathcal{L}(\hat{Y}, Y))$$

$$= \frac{1}{7} \cdot \underbrace{\text{np.sum}}_{\substack{\text{Summe über alle Datenpunkte} \\ \text{BATCH GRADIENT DESCENT}}}([● ● ● ● ● ● ●])$$

Summe über alle Datenpunkte  
BATCH GRADIENT DESCENT

# DER GRADIENT

$$\theta = [\underset{\text{Parameter}}{\text{---}}]^T$$

$$\frac{\partial J}{\partial \theta} = d\theta = \left[ \underset{\substack{\text{für aktuelles } \theta \text{ auswerten}}}{\text{---}} \right]^T$$
$$= \nabla J = \left[ \frac{\partial J}{\partial w_{11}^{[1]}} \dots \frac{\partial J}{\partial w_{34}^{[1]}} \frac{\partial J}{\partial b_1^{[1]}} \dots \frac{\partial J}{\partial b_3^{[1]}} \frac{\partial J}{\partial w_{11}^{[2]}} \dots \frac{\partial J}{\partial w_{13}^{[2]}} \frac{\partial J}{\partial b_{11}^{[2]}} \right]^T \underset{\substack{\text{for current } w}}{}$$

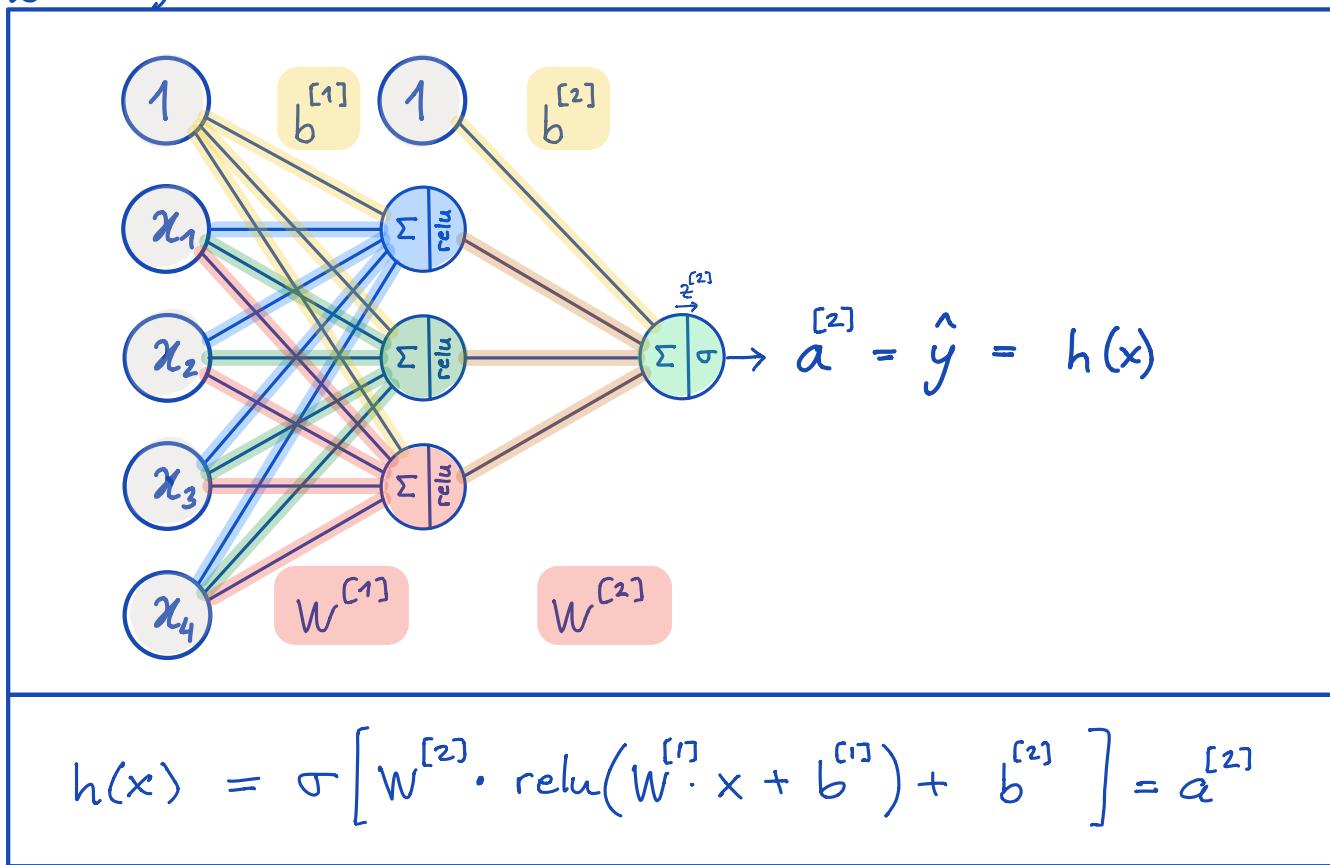
Richtung des steilsten Aufstiegs ! (Gehe in entgegengesetzte Richtung) !

$$\theta := \theta - \alpha \cdot \nabla J$$

Gewichtsaktualisierungen

# • Backpropagation - Multiple Layers

2-Layer Network (Graph Representation)



(Ausgabe als Funktion von  $x$ )

$$J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(a^{(i)}, y^{(i)})$$

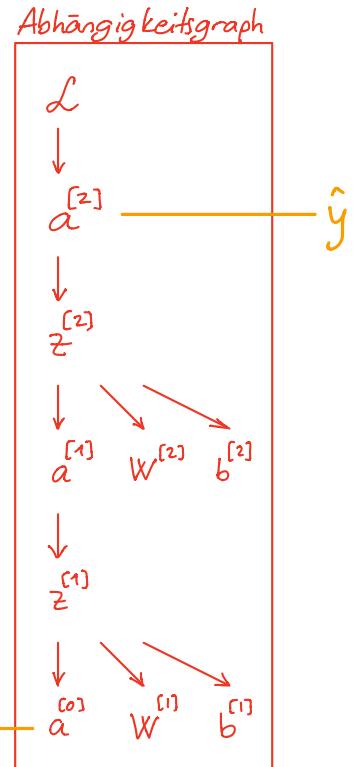
$$\mathcal{L} = - \left( y \cdot \log(a) + (1-y) \cdot \log(1-a) \right)$$

$$a^{[2]} = \sigma(z^{[2]})$$

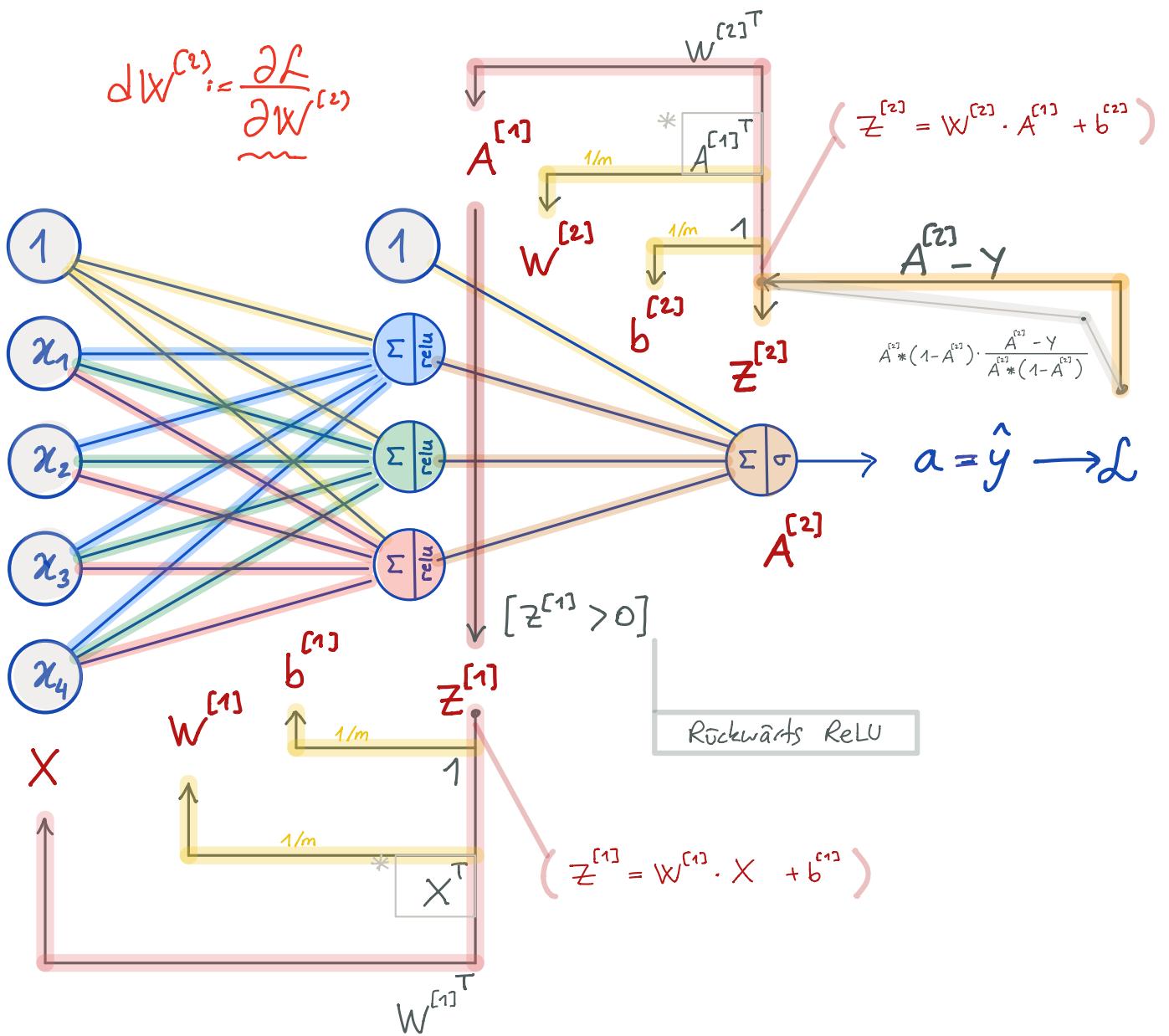
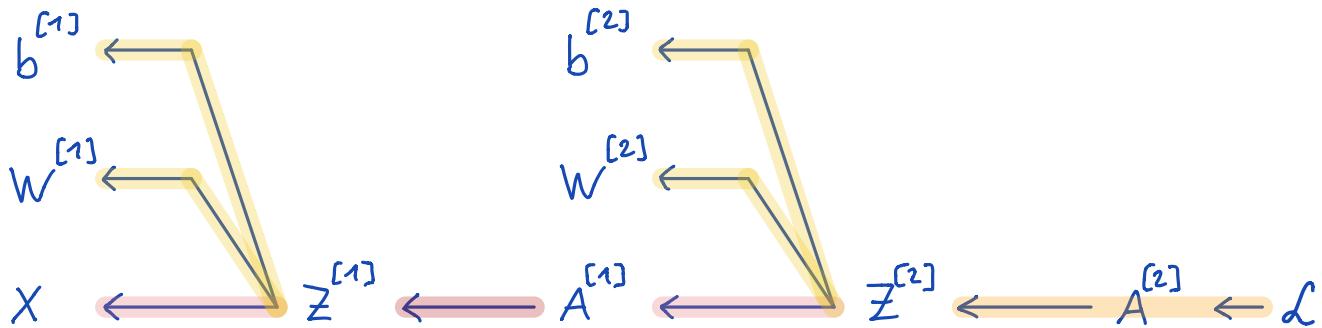
$$z^{[2]} = W^{[2]} \cdot a^{[1]} + b^{[2]}$$

$$a^{[1]} = \text{ReLU}(z^{[1]})$$

$$z^{[1]} = W^{[1]} \cdot a^{[0]} + b^{[1]}$$



# BACKWARD PASS



\* : The right term to be multiplied with  $dZ^{[2]}$ , but not  $\partial Z / \partial W$ ! See notes below!

# • DIE ABLEITUNGEN

①  $\frac{\partial \mathcal{L}}{\partial a}$  für

$$\mathcal{L}(a, y) = -y \cdot \log_e(a) - (1-y) \cdot \log(1-a)$$

•  $\frac{\partial \mathcal{L}}{\partial a} = \frac{-y}{a} - \frac{(1-y)}{(1-a)} \cdot (-1) = \frac{-y}{a} + \frac{(1-y)}{(1-a)}$

$$\frac{-y(1-a) + a(1-y)}{a(1-a)} = \frac{-y + ya + a - ay}{a(1-a)} = \frac{a - y}{a(1-a)}$$

②  $\frac{da}{dz}$  für  $a = \sigma(z) = \frac{1}{1+e^{-z}}$   $\Rightarrow \frac{1}{a} = 1+e^{-z}$

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

$$= -1 \cdot a^2 \cdot \left(1 - \frac{1}{a}\right) \quad A * (1-A)$$

$$= -a^2 \cdot \left(\frac{a-1}{a}\right) = a(1-a)$$

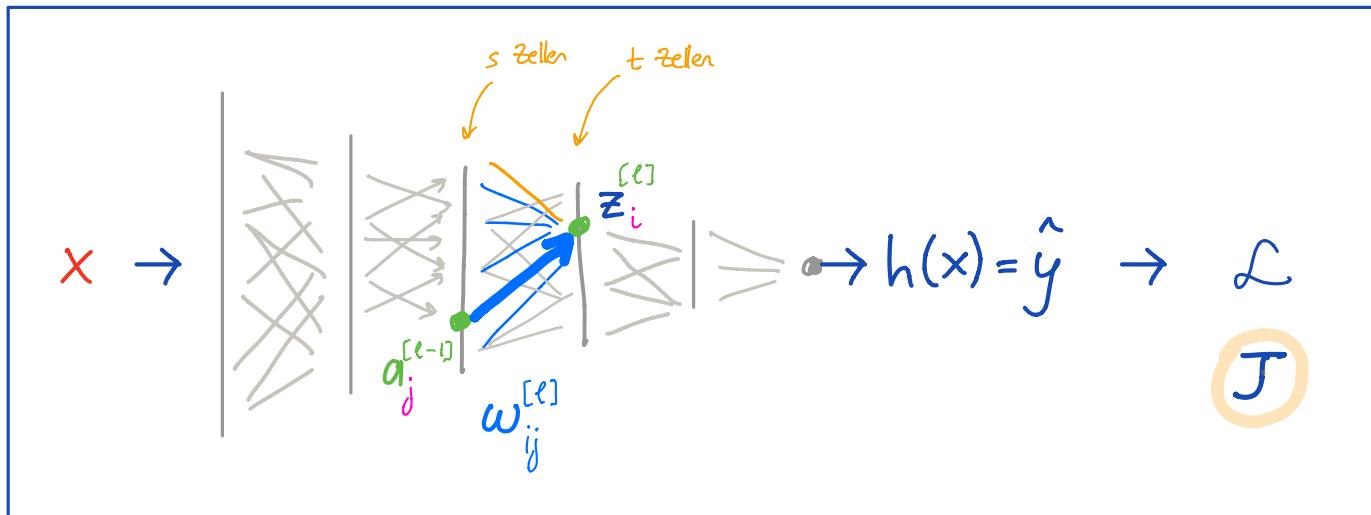
•  $\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z} = \frac{a-y}{a(1-a)} \cdot \frac{a-1}{a} = a-y \quad A^{[2]} - Y$

•  $\frac{\partial \mathcal{L}}{\partial z^{[2]}} = A^{[2]} - Y$

This holds, as long as  
 $A^{[2]} = \sigma(z^{[2]})$

$$= [ \text{green dots} ] - [ \text{green dots} ] = [ \text{yellow dots} ]$$

$$③ \frac{\partial z}{\partial a} \text{ für } z = Wa + b$$



$$i \rightarrow \begin{bmatrix} w_{in} & \dots & w_{ij} & \dots & w_{is} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{in} & \dots & w_{ij} & \dots & w_{is} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ w_{ti} & \dots & w_{tj} & \dots & w_{ts} \end{bmatrix} \cdot \begin{bmatrix} a_1^{[l-1]} \\ \vdots \\ a_j^{[l-1]} \\ \vdots \\ a_s^{[l-1]} \end{bmatrix} + \begin{bmatrix} b_1 \\ \vdots \\ b_i \\ \vdots \\ b_t \end{bmatrix} = \begin{bmatrix} z_1^{[l]} \\ \vdots \\ z_i^{[l]} \\ \vdots \\ z_t^{[l]} \end{bmatrix}$$

$$W^{[l]} \cdot a^{[l-1]} + b^{[l]} = z^{[l]}$$

$$\frac{\partial z_i^{[l]}}{\partial a_j^{[l-1]}} = w_{ij}^{[l]}$$

$$\begin{aligned} z_1^{[l]} &= w_{in}^{[l]} \cdot a_1^{[l-1]} + \dots + w_{1j}^{[l]} \cdot a_j^{[l-1]} + \dots + w_{is}^{[l]} \cdot a_s^{[l-1]} + b_1^{[l]} \\ &\vdots \\ z_i^{[l]} &= w_{in}^{[l]} \cdot a_1^{[l-1]} + \dots + w_{ij}^{[l]} \cdot a_j^{[l-1]} + \dots + w_{is}^{[l]} \cdot a_s^{[l-1]} + b_i^{[l]} \\ &\vdots \\ z_t^{[l]} &= w_{in}^{[l]} \cdot a_1^{[l-1]} + \dots + w_{tj}^{[l]} \cdot a_j^{[l-1]} + \dots + w_{ts}^{[l]} \cdot a_s^{[l-1]} + b_t^{[l]} \end{aligned}$$

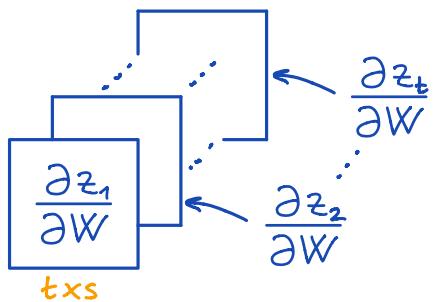
- Zwei Layout Optionen für  $\frac{\partial z}{\partial a}$ :

- Numerator (Zähler) Layout
- Denominator (Nenner) Layout

$$\circ \frac{\partial z}{\partial a} = \begin{bmatrix} \frac{\partial z_1}{\partial a_1} & \frac{\partial z_2}{\partial a_1} & \frac{\partial z_t}{\partial a_1} \\ \frac{\partial z_1}{\partial a_2} & \frac{\partial z_2}{\partial a_2} & \frac{\partial z_t}{\partial a_2} \\ \vdots & \vdots & \vdots \\ \frac{\partial z_1}{\partial a_s} & \frac{\partial z_2}{\partial a_s} & \frac{\partial z_t}{\partial a_s} \end{bmatrix}_{(s \times t)} = W^{(c)^T}$$

Hessian Formulation  
oder  
"Denominator Layout"

④  $\frac{\partial z}{\partial w}$  für  $z = Wa + b$



Vektor-nach-Matrix Ableitung:  
3D-Tensor ( $t \times s \times t$ )

$$\begin{array}{c} \text{3D-Tensor: } \dots \\ \text{Slices: } \dots = \frac{\partial z}{\partial w} = (-0-) \\ \text{Row vector: } \dots = (-a^T-) \end{array}$$

◦  $dW := \frac{\partial L}{\partial W}$  ist was wir eigentlich brauchen

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial W} \quad \left[ \begin{array}{l} \text{für Gewichtsaktualisierung} \\ W := W - \alpha \cdot dW \end{array} \right]$$

↑

Mit einem Trick können wir diese Ableitungen einfacher berechnen. (ohne  $\frac{\partial z}{\partial w}$  explizit auszurechnen).

- IDEE : Berechne den Gradienten für  $W$  Zeile-für-Zeile.

Bemerkung: Erste Zeile von  $W$  beeinflusst nur  $z_1$ !

Frage: Was wäre, wenn Schicht  $l$  nur die eine Zelle  $z_1$  hätte?

Sei  $w_{1.}$  die erste Zeile von  $W$ . ( $w_{i.} : i\text{-te Zeile von } W$ )

$$\frac{\partial \mathcal{L}}{\partial w_{1.}} = dW_{1.} = \frac{\partial \mathcal{L}}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_{1.}} = \frac{\partial \mathcal{L}}{\partial z_1} \cdot a^T \quad \frac{\partial z_2}{\partial w_{1.}} = 0$$

$$\frac{\partial \mathcal{L}}{\partial w_{2.}} = dW_{2.} = \frac{\partial \mathcal{L}}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_{2.}} = \frac{\partial \mathcal{L}}{\partial z_2} \cdot a^T$$

:

$$\frac{\partial \mathcal{L}}{\partial w_{t.}} = dW_{t.} = \frac{\partial \mathcal{L}}{\partial z_t} \cdot \frac{\partial z_t}{\partial w_{t.}} = \frac{\partial \mathcal{L}}{\partial z_t} \cdot a^T$$

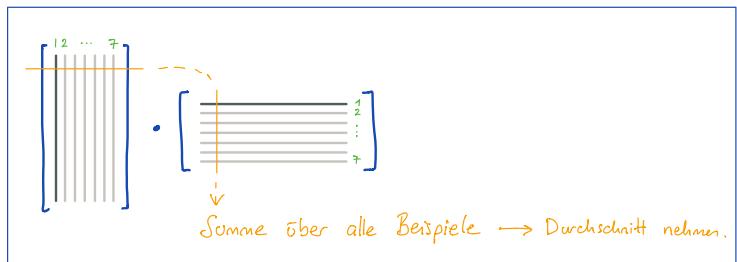
$$\Rightarrow dW = \begin{bmatrix} dz_1 \cdot a^T \\ dz_2 \cdot a^T \\ \vdots \\ dz_t \cdot a^T \end{bmatrix} = \begin{bmatrix} dz_1 \\ dz_2 \\ \vdots \\ dz_t \end{bmatrix} \cdot [ \quad a^T \quad ]$$

$$dW = dz \cdot a^T$$

- Frage: Was passiert, wenn wir mehrere Beispiele haben? ( $m > 1$ )

$$dW = \frac{1}{m} dZ \cdot A^T$$

[ $\ell$ ]    [ $\ell$ ]    [ $\ell-1$ ]



$$\textcircled{4} \quad \frac{dA}{dZ} \quad \text{für} \quad A = \text{relu}(Z)$$

$$\circ \quad \frac{dA}{dZ} = [ z > 0 ] = [ \text{True if } (z > 0) \text{ else False} ] \\ = [ 1 \text{ if } (z > 0) \text{ else } 0 ]$$

BSP.

$$\begin{array}{c}
 Z \downarrow \\
 \boxed{Z = \begin{bmatrix} 5 & -7 & 6 \\ 2 & 1 & -4 \end{bmatrix}} \\
 \downarrow (\text{relu}) \\
 A = \begin{bmatrix} 5 & 0 & 6 \\ 2 & 1 & 0 \end{bmatrix}
 \end{array}
 \quad
 \begin{array}{c}
 dZ \uparrow \\
 \boxed{dZ = \begin{bmatrix} 0.5 & 0 & 0.7 \\ 0.1 & -0.6 & 0 \end{bmatrix}} \\
 * [ z > 0 ] \uparrow (\text{rückwärts relu}) \\
 dA = \begin{bmatrix} 0.5 & 0.05 & 0.7 \\ 0.1 & -0.6 & 0.8 \end{bmatrix}
 \end{array}$$

$$\downarrow A$$

$$dA = \frac{\partial L}{\partial A}$$

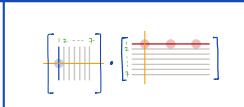
$$\Rightarrow \frac{dA}{dZ} = [ z > 0 ] = \boxed{\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}}$$

# THE BACKPROP MAP

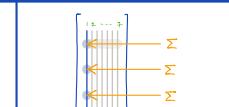
$L$

$A^{[2]}$ $(1 \times 7)$	$dA^{[2]} := \frac{\partial L}{\partial A^{[2]}} = \frac{A^{[2]} - Y}{A^{[2]} * (1 - A^{[2]})} = -Y/A^{[2]} + (1-Y)/(1-A^{[2]})$	$\uparrow$ element-wise division
-----------------------------	--	-------------------------------------

$A^{[2]}$ $A * (1 - A^{[2]})$	$dZ^{[2]} := \frac{\partial L}{\partial Z^{[2]}} = dA^{[2]} * A * (1 - A^{[2]}) = A^{[2]} - Y$
----------------------------------	--

$A^{[1]}$ $(3 \times 7) = (3 \times 1) \cdot (1 \times 7)$	$dA^{[1]} = W^{[2]T} \cdot dZ^{[2]}$	$W^{[2]}$ $A^{[1]T} \cdot *$	$1$	$Z^{[2]} = W^{[2]} \cdot A^{[1]} + b^{[2]}$
	$dW^{[2]} = \frac{1}{m} \cdot dZ^{[2]} \cdot A^{[1]T}$ $(1 \times 3) = (1 \times 7) \circ (7 \times 3)$			$db^{[2]} = \frac{1}{m} \cdot \text{np.sum}(dZ^{[2]})$ $(1 \times 1) =$ $\downarrow$ $\text{axis}=1 \quad (1 \times 7)$
				

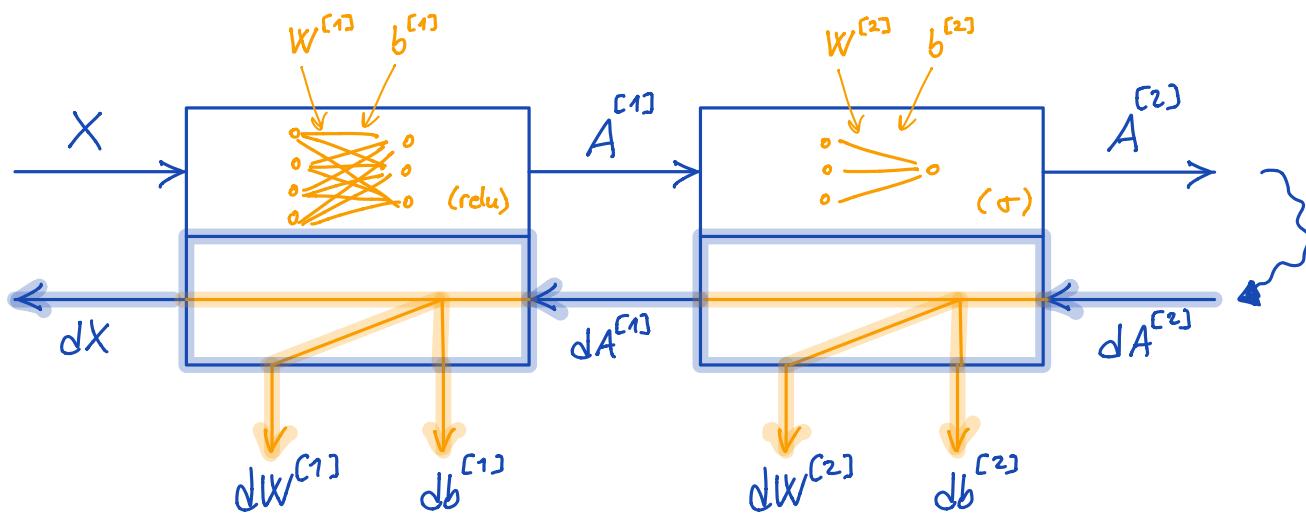
$Z^{[1]}$ $(3 \times 7)$	$dZ^{[1]} := \frac{\partial L}{\partial Z^{[1]}} = dA^{[1]} * [Z^{[1]} > 0] \quad (\text{element-wise})$
	$[Z^{[1]} > 0]$

$X$ $(4 \times 7) = (4 \times 3) \cdot (3 \times 7)$	$dX = W^{[1]T} \cdot dZ^{[1]}$	$W^{[1]}$ $X^T \cdot *$	$1$
	$dW^{[1]} = \frac{1}{m} \cdot dZ^{[1]} \cdot X^T$ $(3 \times 4) = (3 \times 7) \circ (7 \times 4)$		
			
			

\* Nicht gleich  $\partial Z / \partial W$ , aber der richtige Term in der "Kette". Siehe Notizen unten.

## • PROZESSABLAUF IM ÜBERBLICK

### FORWÄRTSLAUF



### RÜCKWÄRTSLAUF

### GEWICHTSAKTUALISIERUNGEN

$$\begin{bmatrix} W^{[1]} \\ b^{[1]} \\ W^{[2]} \\ b^{[2]} \end{bmatrix}_{\text{neu}} := \begin{bmatrix} W^{[1]} \\ b^{[1]} \\ W^{[2]} \\ b^{[2]} \end{bmatrix}_{\text{alt}} - \alpha \cdot \begin{bmatrix} dW^{[1]} \\ db^{[1]} \\ dW^{[2]} \\ db^{[2]} \end{bmatrix}$$