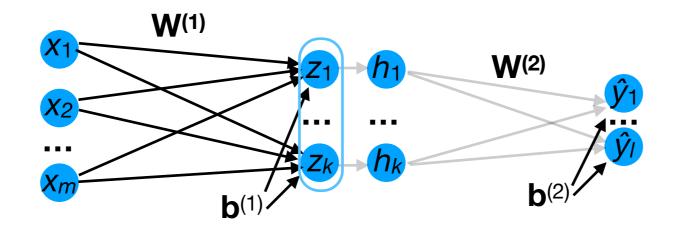
#### CS 453X: Class 21

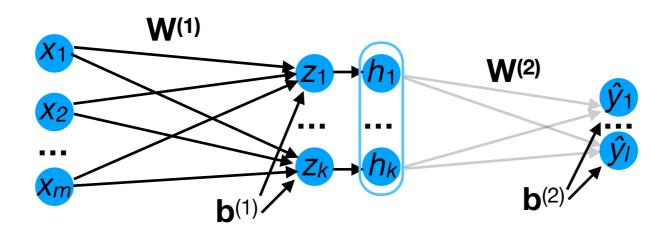
Jacob Whitehill

# More on forwards and backwards propagation

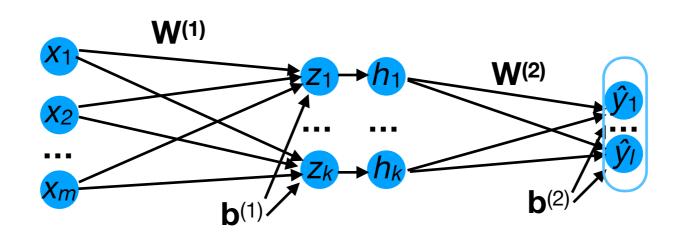
- Jacobian matrices and the chain rule provide a recipe for how to compute all the gradient terms efficiently.
- Consider the 3-layer NN below:
  - From  $\mathbf{x}$ ,  $\mathbf{W}^{(1)}$ , and  $\mathbf{b}^{(1)}$ , we can compute  $\mathbf{z}$ .



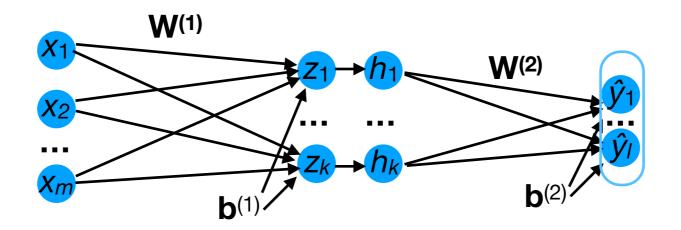
- Jacobian matrices and the chain rule provide a recipe for how to compute all the gradient terms efficiently.
- Consider the 3-layer NN below:
  - From  $\mathbf{x}$ ,  $\mathbf{W}^{(1)}$ , and  $\mathbf{b}^{(1)}$ , we can compute  $\mathbf{z}$ .
  - From **z** and  $\sigma$ , we can compute **h** =  $\sigma$ (**z**).



- Jacobian matrices and the chain rule provide a recipe for how to compute all the gradient terms efficiently.
- Consider the 3-layer NN below:
  - From  $\mathbf{x}$ ,  $\mathbf{W}^{(1)}$ , and  $\mathbf{b}^{(1)}$ , we can compute  $\mathbf{z}$ .
  - From **z** and  $\sigma$ , we can compute **h** =  $\sigma$ (**z**).
  - From h,  $W^{(2)}$ , and  $b^{(2)}$ , we can compute  $\hat{y}$ .

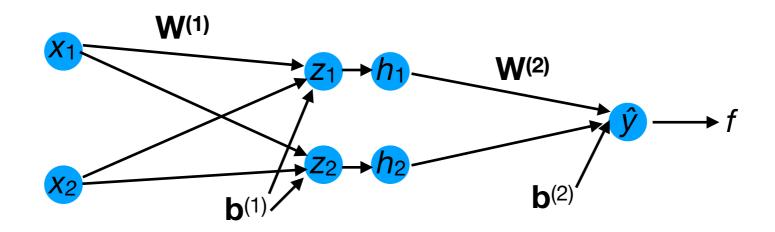


- Jacobian matrices and the chain rule provide a recipe for how to compute all the gradient terms efficiently.
- This process is known as forward propagation.
  - It produces all the intermediary (h, z) and final (ŷ) network outputs.



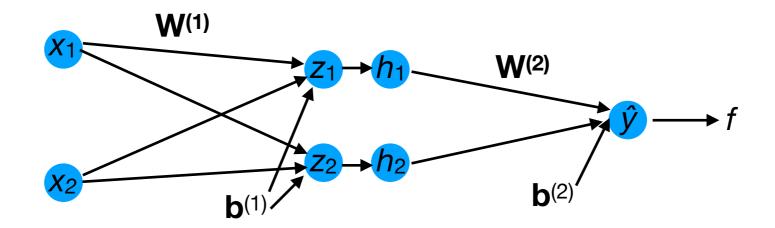
Now, let's look at how to compute each gradient term:

$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}} 
\frac{\partial f}{\partial \mathbf{b}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}} 
\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} 
\frac{\partial f}{\partial \mathbf{b}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}^{(1)}}$$

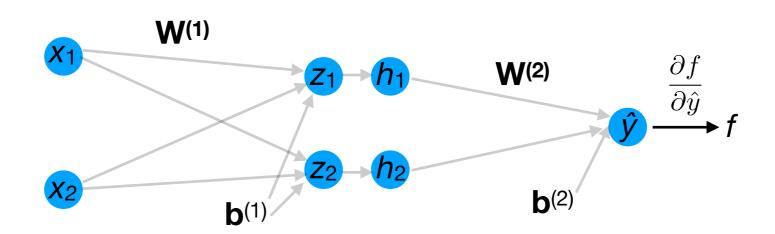


Now, let's look at how to compute each gradient term:

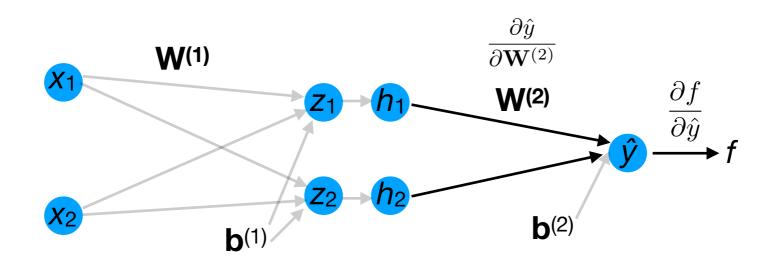
$$\begin{array}{ll} \frac{\partial f}{\partial \mathbf{W}^{(2)}} & = & \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}} \\ \frac{\partial f}{\partial \mathbf{b}^{(2)}} & = & \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}} & \mathbf{computation} \\ \frac{\partial f}{\partial \mathbf{W}^{(1)}} & = & \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} \\ \frac{\partial f}{\partial \mathbf{b}^{(1)}} & = & \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{b}^{(1)}} \end{array}$$



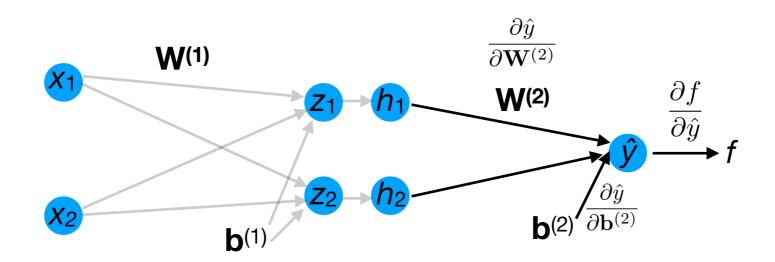
$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}}$$



$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}}$$



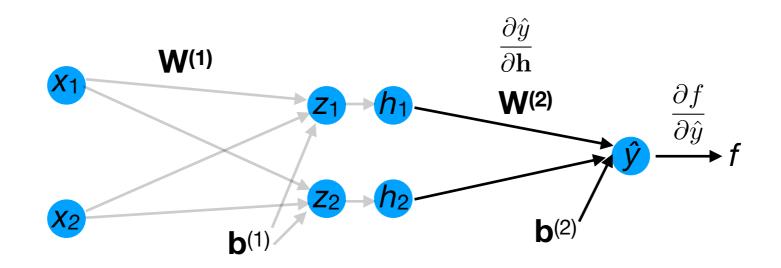
$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}}$$
$$\frac{\partial f}{\partial \mathbf{b}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}}$$



$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}}$$

$$\frac{\partial f}{\partial \mathbf{b}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}}$$

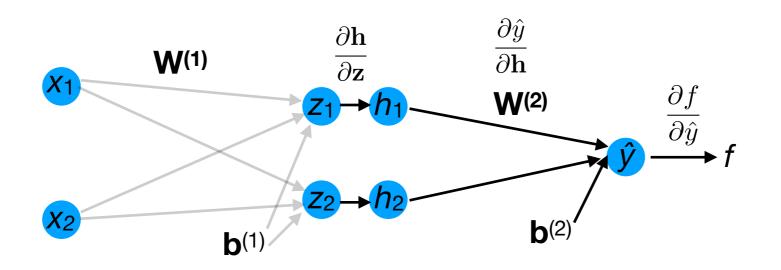
$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}}$$



$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}}$$

$$\frac{\partial f}{\partial \mathbf{b}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}}$$

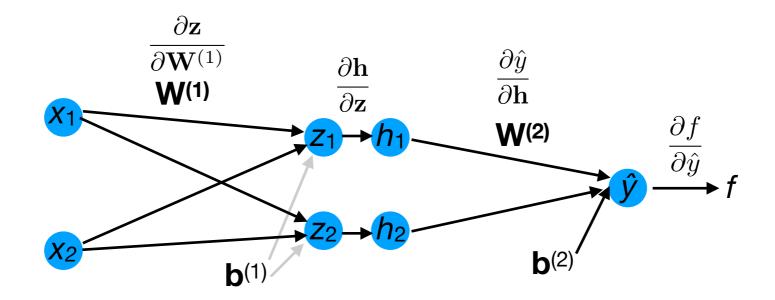
$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}}$$



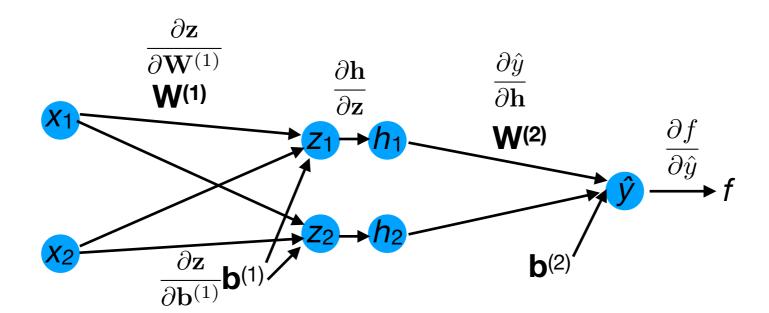
$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}}$$

$$\frac{\partial f}{\partial \mathbf{b}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}}$$

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$



- This process is known as backwards propagation ("backprop"):
  - It produces the gradient terms of all the weight matrices and bias vectors.

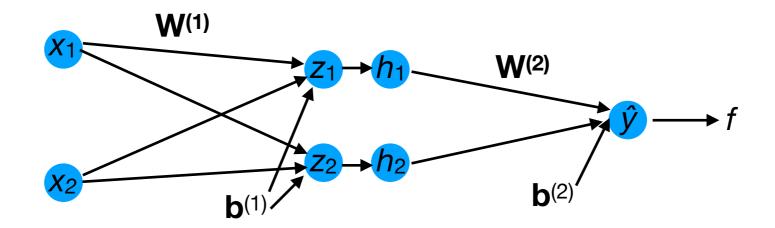


Where do these come from?

$$abla_{\mathbf{W}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)}^{\top}$$
 $abla_{\mathbf{b}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$ 
 $abla_{\mathbf{b}^{(1)}} f_{\mathrm{CE}} = \mathbf{g} \mathbf{x}^{\top}$ 
 $abla_{\mathbf{b}^{(1)}} f_{\mathrm{CE}} = \mathbf{g}$ 

where

$$\mathbf{g}^{\top} = \left( (\hat{\mathbf{y}}^{\top} - \mathbf{y}) \mathbf{W}^{(2)} \right) \odot \text{relu}' (\mathbf{z}^{(1)}^{\top})$$



• Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

• How does  $\hat{y}$  depend on **h**?

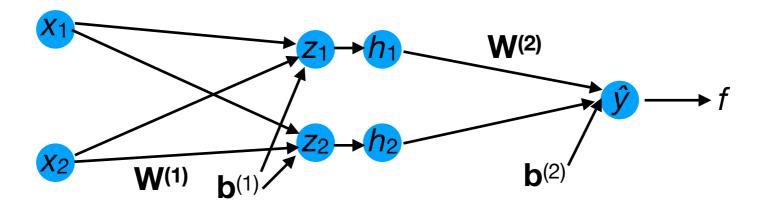
$$\hat{y} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)} \\
= \mathbf{W}^{(2)}_{1}\mathbf{h}_{1} + \mathbf{W}^{(2)}_{2}\mathbf{h}_{2} + \mathbf{b}^{(2)} \\
\Rightarrow \frac{\partial \hat{y}}{\partial \mathbf{h}} = \begin{bmatrix} \frac{\partial \hat{y}}{\partial \mathbf{h}_{1}} & \frac{\partial \hat{y}}{\partial \mathbf{h}_{2}} \end{bmatrix} \\
= \begin{bmatrix} \mathbf{W}^{(2)}_{1} & \mathbf{W}^{(2)}_{2} \end{bmatrix} \\
= \mathbf{W}^{(2)} \\
\mathbf{W}^{(2)}_{1} & \mathbf{h}^{(1)}_{2} & \mathbf{h}^{(2)}_{2}$$

Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

How does h depend on z?

$$\mathbf{h} = \begin{bmatrix} \operatorname{relu}(\mathbf{z}_1) \\ \operatorname{relu}(\mathbf{z}_2) \end{bmatrix}$$



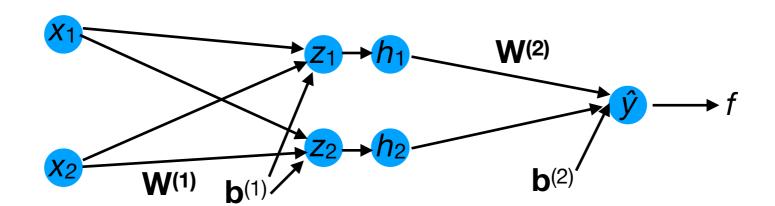
• Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

How does h depend on z?

$$\mathbf{h} = \begin{bmatrix} \operatorname{relu}(\mathbf{z}_1) \\ \operatorname{relu}(\mathbf{z}_2) \end{bmatrix}$$

$$\Longrightarrow \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial \mathbf{h}_1}{\partial \mathbf{z}_1} & \frac{\partial \mathbf{h}_1}{\partial \mathbf{z}_2} \\ \frac{\partial \mathbf{h}_2}{\partial \mathbf{z}_1} & \frac{\partial \mathbf{h}_2}{\partial \mathbf{z}_2} \end{bmatrix}$$



• Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

How does h depend on z?

$$\mathbf{h} = \begin{bmatrix} \operatorname{relu}(\mathbf{z}_1) \\ \operatorname{relu}(\mathbf{z}_2) \end{bmatrix}$$

$$\Rightarrow \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial \mathbf{h}_1}{\partial \mathbf{z}_1} & \frac{\partial \mathbf{h}_1}{\partial \mathbf{z}_2} \\ \frac{\partial \mathbf{h}_2}{\partial \mathbf{z}_1} & \frac{\partial \mathbf{h}_2}{\partial \mathbf{z}_2} \end{bmatrix}$$

$$= \begin{bmatrix} \operatorname{relu}'(\mathbf{z}_1) & 0 \\ 0 & \operatorname{relu}'(\mathbf{z}_2) \end{bmatrix}$$

$$\mathbf{x}_1$$

$$\mathbf{z}_1$$

$$\mathbf{h}_1$$

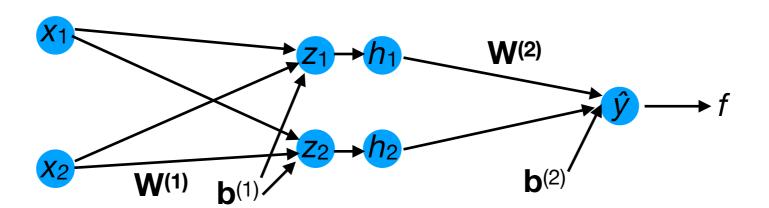
$$\mathbf{w}_2$$

• Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

How does z depend on W<sup>(1)</sup>?

$$\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$



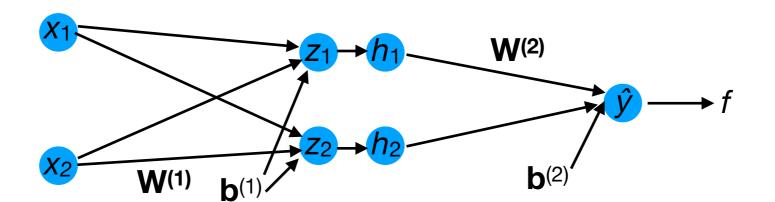
Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

How does z depend on W<sup>(1)</sup>?

$$\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$

$$\begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{W}_1^{(1)} & \mathbf{W}_2^{(1)} \\ \mathbf{W}_3^{(1)} & \mathbf{W}_4^{(1)} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{b}_1^{(1)} \\ \mathbf{b}_2^{(1)} \end{bmatrix}$$

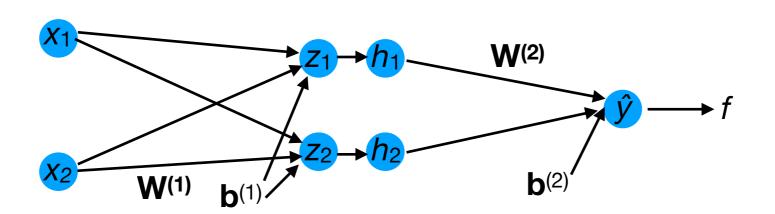


Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

How does z depend on W<sup>(1)</sup>?

were a vector.



Let's derive each gradient term in turn:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

How does z depend on W<sup>(1)</sup>?

$$\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$

$$\begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{W}_1^{(1)} & \mathbf{W}_2^{(1)} \\ \mathbf{W}_3^{(1)} & \mathbf{W}_4^{(1)} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{b}_1^{(1)} \\ \mathbf{b}_2^{(1)} \end{bmatrix}$$

$$\Rightarrow \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} = \begin{bmatrix} \frac{\partial \mathbf{z}_1}{\partial \mathbf{W}_1^{(1)}} & \frac{\partial \mathbf{z}_1}{\partial \mathbf{W}_2^{(1)}} & \frac{\partial \mathbf{z}_1}{\partial \mathbf{W}_3^{(1)}} & \frac{\partial \mathbf{z}_1}{\partial \mathbf{W}_4^{(1)}} \\ \frac{\partial \mathbf{z}_2}{\partial \mathbf{W}_1^{(1)}} & \frac{\partial \mathbf{z}_2}{\partial \mathbf{W}_2^{(1)}} & \frac{\partial \mathbf{z}_2}{\partial \mathbf{W}_3^{(1)}} & \frac{\partial \mathbf{z}_2}{\partial \mathbf{W}_4^{(1)}} \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & 0 & 0 \\ 0 & 0 & \mathbf{x}_1 & \mathbf{x}_2 \end{bmatrix}$$

$$\mathbf{X}_1 \qquad \mathbf{X}_2 \qquad \mathbf{X}_1 \qquad \mathbf{X}_2$$

$$\mathbf{X}_1 \qquad \mathbf{X}_2 \qquad \mathbf{X}_1 \qquad \mathbf{X}_2 \qquad \mathbf{X}_2 \qquad \mathbf{X}_1 \qquad \mathbf{X}_2 \qquad \mathbf{X}_2 \qquad \mathbf{X}_2 \qquad \mathbf{X}_2 \qquad \mathbf{X}_2 \qquad \mathbf{X}_2 \qquad \mathbf{X}_3 \qquad \mathbf{X}_4 \qquad \mathbf{X}_4 \qquad \mathbf{X}_5 \qquad \mathbf{X}_5 \qquad \mathbf{X}_5 \qquad \mathbf{X}_5 \qquad \mathbf{X}_5 \qquad \mathbf{X}_5 \qquad \mathbf{X}_6 \qquad \mathbf{X}_7 \qquad \mathbf{X}_7$$

We can now finally derive the gradient update for W<sup>(1)</sup>:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}}$$

We can now finally derive the gradient update for W<sup>(1)</sup>:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} 
= (\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)} \begin{bmatrix} \operatorname{relu}'(\mathbf{z}_1) & 0 \\ 0 & \operatorname{relu}'(\mathbf{z}_2) \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & 0 & 0 \\ 0 & 0 & \mathbf{x}_1 & \mathbf{x}_2 \end{bmatrix}$$

We can now finally derive the gradient update for W<sup>(1)</sup>:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} 
= (\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)} \begin{bmatrix} \operatorname{relu}'(\mathbf{z}_{1}) & 0 \\ 0 & \operatorname{relu}'(\mathbf{z}_{2}) \end{bmatrix} \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} 
= (((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot [\operatorname{relu}'(\mathbf{z}_{1}) & \operatorname{relu}'(\mathbf{z}_{2}) ]) \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix}$$

since multiplying by a diagonal matrix is equivalent to element-wise (Hadamard) product.

We can now finally derive the gradient update for W<sup>(1)</sup>:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} 
= (\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)} \begin{bmatrix} \operatorname{relu}'(\mathbf{z}_{1}) & 0 \\ 0 & \operatorname{relu}'(\mathbf{z}_{2}) \end{bmatrix} \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} 
= (((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot [\operatorname{relu}'(\mathbf{z}_{1}) & \operatorname{relu}'(\mathbf{z}_{2}) ]) \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} 
= [\mathbf{g}_{1} \quad \mathbf{g}_{2}] \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix}$$

To simplify notation, let's define a new vector that equals the first few terms.

We can now finally derive the gradient update for W<sup>(1)</sup>:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} 
= (\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)} \begin{bmatrix} \operatorname{relu}'(\mathbf{z}_{1}) & 0 \\ 0 & \operatorname{relu}'(\mathbf{z}_{2}) \end{bmatrix} \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} 
= (((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot [ \operatorname{relu}'(\mathbf{z}_{1}) & \operatorname{relu}'(\mathbf{z}_{2}) ]) \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} 
= [ \mathbf{g}_{1} & \mathbf{g}_{2} ] \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} 
= [ \mathbf{g}_{1}\mathbf{x}_{1} & \mathbf{g}_{1}\mathbf{x}_{2} & \mathbf{g}_{2}\mathbf{x}_{1} & \mathbf{g}_{2}\mathbf{x}_{2} ]$$

We can now finally derive the gradient update for W<sup>(1)</sup>:

$$\frac{\partial f}{\partial \mathbf{W}^{(1)}} = \frac{\partial f}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} \\
= (\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)} \begin{bmatrix} \operatorname{relu}'(\mathbf{z}_{1}) & 0 \\ 0 & \operatorname{relu}'(\mathbf{z}_{2}) \end{bmatrix} \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} \\
= (((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot [ \operatorname{relu}'(\mathbf{z}_{1}) & \operatorname{relu}'(\mathbf{z}_{2}) ]) \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} \\
= [ \mathbf{g}_{1} & \mathbf{g}_{2} ] \begin{bmatrix} \mathbf{x}_{1} & \mathbf{x}_{2} & 0 & 0 \\ 0 & 0 & \mathbf{x}_{1} & \mathbf{x}_{2} \end{bmatrix} \\
= [ \mathbf{g}_{1}\mathbf{x}_{1} & \mathbf{g}_{1}\mathbf{x}_{2} & \mathbf{g}_{2}\mathbf{x}_{1} & \mathbf{g}_{2}\mathbf{x}_{2} \end{bmatrix} \\
\Rightarrow \nabla_{\mathbf{W}^{(1)}} f = \mathbf{g}\mathbf{x}^{\top}$$

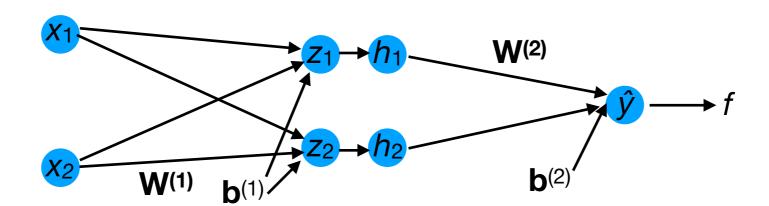
**Outer product** 

### Weight initialization

#### NNs and convexity

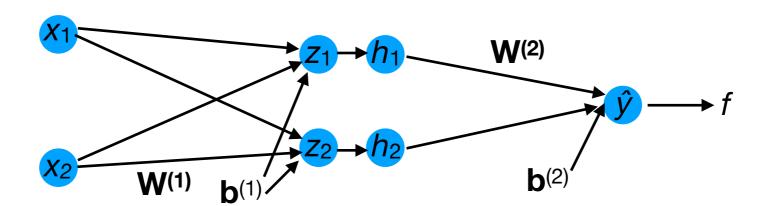
- Neural networks are a non-convex ML model.
- Hence, the values you use to initialize the weights and bias terms can make a big difference on the accuracy of the network.
  - The network might end up in a worse local minimum of the cost function.

- Suppose we initialize all the weights and bias terms of a 3-layer NN to be 0.
- What will happen during SGD?



- Suppose we initialize all the weights and bias terms of a 3-layer NN to be 0.
- What will happen during SGD?

During forwards propagation, z and h will be 0. Hence, ŷ will also be 0.



- Suppose we initialize all the weights and bias terms of a 3-layer NN to be 0.
- What will happen during SGD?

#### **During backwards propagation, we have:**

$$\nabla_{\mathbf{W}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)}^{\top}$$

$$\nabla_{\mathbf{b}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$$

$$\nabla_{\mathbf{W}^{(1)}} f_{\mathrm{CE}} = \mathbf{g} \mathbf{x}^{\top}$$

$$\nabla_{\mathbf{b}^{(1)}} f_{\mathrm{CE}} = \mathbf{g}$$

$$\mathbf{g}^{\top} = ((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot \mathrm{relu}'(\mathbf{z}^{(1)}^{\top})$$

$$\mathbf{z}_{1} \rightarrow h_{1} \qquad \mathbf{w}_{2}$$

$$\mathbf{z}_{2} \rightarrow h_{2}$$

$$\mathbf{z}_{1} \rightarrow h_{2}$$

- Suppose we initialize all the weights and bias terms of a 3-layer NN to be 0.
- What will happen during SGD?

**W**(1)

#### **During backwards propagation, we have:**

$$\nabla_{\mathbf{W}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)^{\top}} \mathbf{0}$$

$$\nabla_{\mathbf{b}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$$

$$\nabla_{\mathbf{W}^{(1)}} f_{\text{CE}} = \mathbf{g} \mathbf{x}^{\top} \mathbf{0}$$

$$\nabla_{\mathbf{b}^{(1)}} f_{\text{CE}} = \mathbf{g} \mathbf{0}$$

$$\mathbf{g}^{\top} = ((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot \text{relu}'(\mathbf{z}^{(1)^{\top}}) \mathbf{0}$$

# Weight initialization: example

- Because the gradients w.r.t. W<sup>(1)</sup>, W<sup>(2)</sup>, and b<sup>(1)</sup> are all 0, they will never change.
- Only b<sup>(2)</sup> will change (to the mean of the target values y).

#### **During backwards propagation, we have:**

$$\nabla_{\mathbf{W}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)^{\top}} \mathbf{0}$$

$$\nabla_{\mathbf{b}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$$

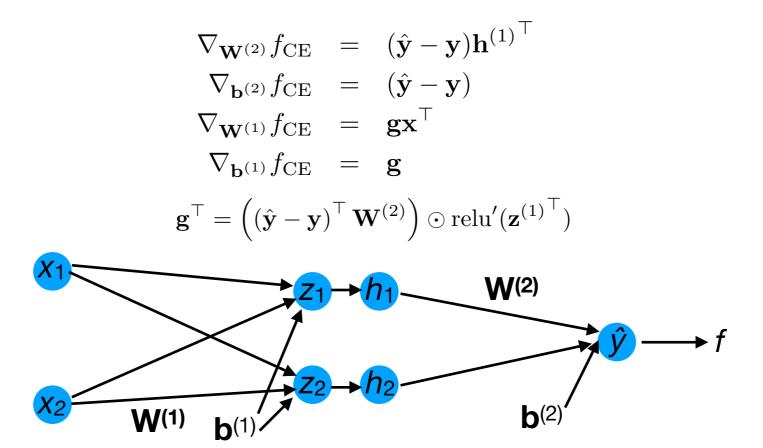
$$\nabla_{\mathbf{W}^{(1)}} f_{\text{CE}} = \mathbf{g} \mathbf{x}^{\top} \mathbf{0}$$

$$\nabla_{\mathbf{b}^{(1)}} f_{\text{CE}} = \mathbf{g} \mathbf{0}$$

$$\mathbf{g}^{\top} = ((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot \text{relu}'(\mathbf{z}^{(1)^{\top}}) \mathbf{0}$$

**W**(1)

- Suppose we initialize  $\mathbf{W}^{(1)}=\mathbf{b}^{(1)}=0$ , but  $\mathbf{W}^{(2)}$ ,  $\mathbf{b}^{(2)}$  are non-zero.
- Assume relu'(0)=0.
- What will happen during SGD?



• Since  $\mathbf{W}^{(1)} = \mathbf{b}^{(1)} = 0$ , then  $\mathbf{z}^{(1)} = \mathbf{h}^{(1)} = 0$ . Hence,  $\nabla_{\mathbf{W}^{(2)}} f_{CE} = 0$ .

$$\nabla_{\mathbf{W}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)^{\top}} \mathbf{0}$$

$$\nabla_{\mathbf{b}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$$

$$\nabla_{\mathbf{W}^{(1)}} f_{\mathrm{CE}} = \mathbf{g} \mathbf{x}^{\top}$$

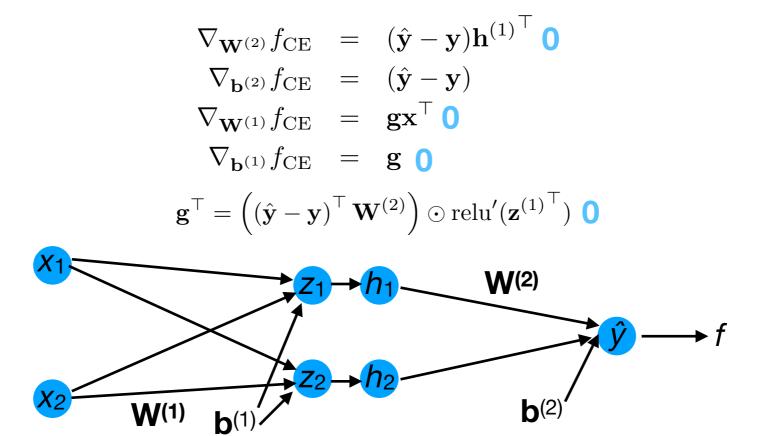
$$\nabla_{\mathbf{b}^{(1)}} f_{\mathrm{CE}} = \mathbf{g}$$

$$\mathbf{g}^{\top} = ((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot \mathrm{relu}'(\mathbf{z}^{(1)^{\top}})$$

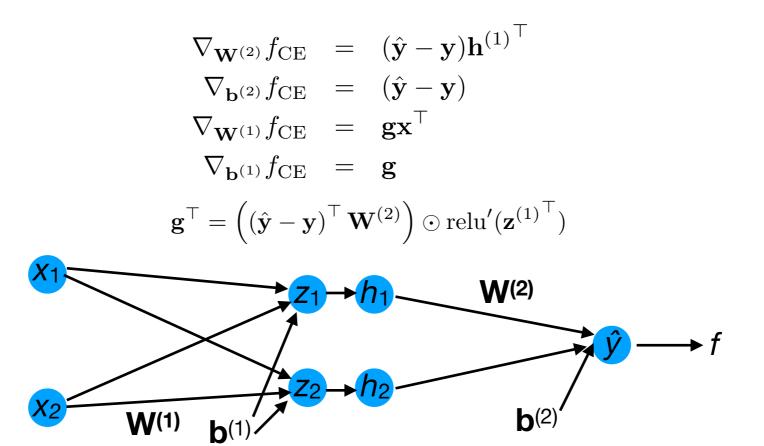
$$\mathbf{X}_{1} \qquad \mathbf{X}_{2} \qquad \mathbf{h}_{2}$$

$$\mathbf{X}_{2} \qquad \mathbf{h}_{2} \qquad \mathbf{h}_{2}$$

- Since  $\mathbf{W}^{(1)} = \mathbf{b}^{(1)} = 0$ , then  $\mathbf{z}^{(1)} = \mathbf{h}^{(1)} = 0$ . Hence,  $\nabla_{\mathbf{W}^{(2)}} f_{CE} = 0$ .
- Since relu'(0) = 0, then  $\mathbf{g}$ =0. Hence, gradients w.r.t.  $\mathbf{W}^{(1)}$  and  $\mathbf{b}^{(1)}$  are 0.
- Only  $\mathbf{b}^{(2)}$  can change (so that  $\hat{y}$  approaches mean of y).



- Suppose we initialize  $\mathbf{W}^{(2)}=\mathbf{b}^{(2)}=0$ , but  $\mathbf{W}^{(1)}$ ,  $\mathbf{b}^{(1)}$  are non-zero.
- What will happen during SGD?



• Since  $\mathbf{W}^{(2)}=0$ , then  $\mathbf{g}=0$ . Hence,  $\nabla_{\mathbf{W}^{(1)}}f_{\mathrm{CE}}$ ,  $\nabla_{\mathbf{b}^{(1)}}f_{\mathrm{CE}}=0$ .

$$\nabla_{\mathbf{W}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)^{\top}}$$

$$\nabla_{\mathbf{b}^{(2)}} f_{\mathrm{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$$

$$\nabla_{\mathbf{W}^{(1)}} f_{\mathrm{CE}} = \mathbf{g} \mathbf{x}^{\top}$$

$$\nabla_{\mathbf{b}^{(1)}} f_{\mathrm{CE}} = \mathbf{g}$$

$$\mathbf{g}^{\top} = ((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot \mathrm{relu}'(\mathbf{z}^{(1)^{\top}})$$

$$\mathbf{X}_{1} \qquad \mathbf{X}_{2} \qquad \mathbf{h}_{2}$$

$$\mathbf{b}_{(1)} \qquad \mathbf{h}_{(1)} \qquad \mathbf{b}_{(2)}$$

- Since  $\mathbf{W}^{(2)}=0$ , then  $\mathbf{g}=0$ . Hence,  $\nabla_{\mathbf{W}^{(1)}} f_{\mathrm{CE}}$ ,  $\nabla_{\mathbf{b}^{(1)}} f_{\mathrm{CE}}=0$ .
- However, **h** is non-zero. Hence,  $\nabla_{\mathbf{W}^{(2)}} f_{\mathrm{CE}}$  is nonzero =>  $\mathbf{W}^{(2)}$  will change.

$$\nabla_{\mathbf{W}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)}^{\top}$$

$$\nabla_{\mathbf{b}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$$

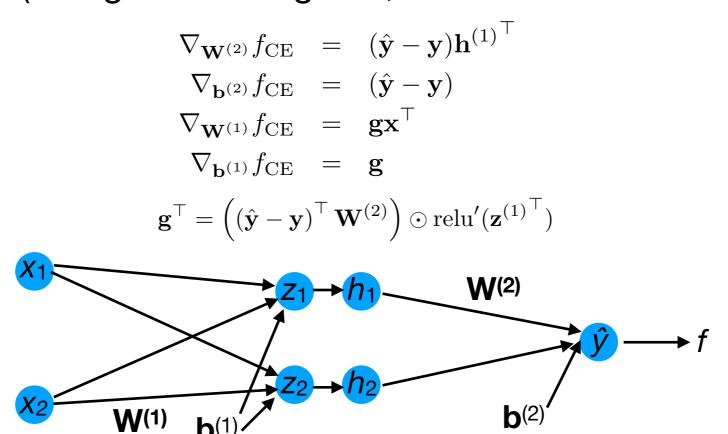
$$\nabla_{\mathbf{W}^{(1)}} f_{\text{CE}} = \mathbf{g} \mathbf{x}^{\top}$$

$$\nabla_{\mathbf{b}^{(1)}} f_{\text{CE}} = \mathbf{g}$$

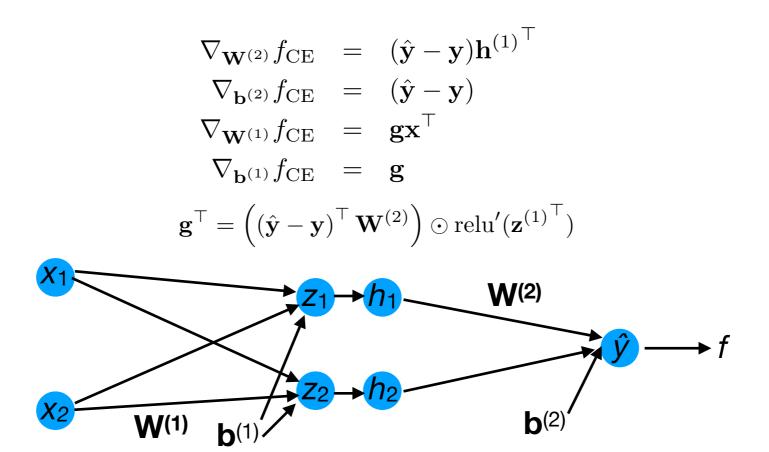
$$\mathbf{g}^{\top} = ((\hat{\mathbf{y}} - \mathbf{y})^{\top} \mathbf{W}^{(2)}) \odot \text{relu}'(\mathbf{z}^{(1)}^{\top})$$

$$\mathbf{X}_{1} \qquad \mathbf{Y}_{2} \qquad \mathbf{Y}_{3} \qquad \mathbf{Y}_{4} \qquad \mathbf{Y}$$

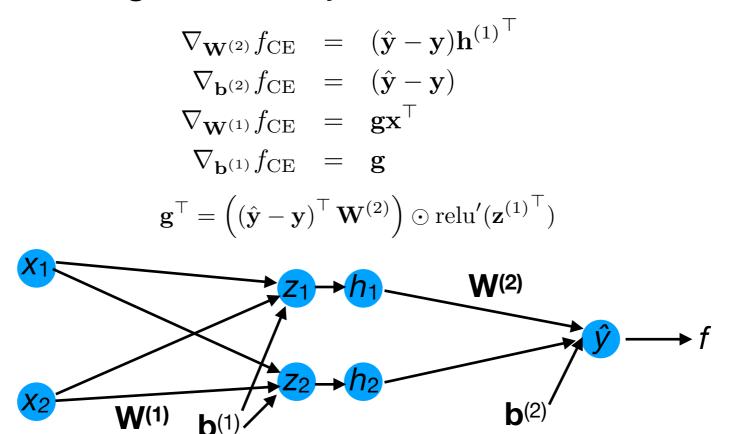
- Since  $\mathbf{W}^{(2)}=0$ , then  $\mathbf{g}=0$ . Hence,  $\nabla_{\mathbf{W}^{(1)}} f_{\mathrm{CE}}$ ,  $\nabla_{\mathbf{b}^{(1)}} f_{\mathrm{CE}}=0$ .
- However, **h** is non-zero. Hence,  $\nabla_{\mathbf{W}^{(2)}} f_{\mathrm{CE}}$  is nonzero =>  $\mathbf{W}^{(2)}$  will change.
- During the next gradient update,  $\mathbf{g}$  is non-zero =>  $\mathbf{W}^{(1)}$ ,  $\mathbf{b}^{(1)}$  will change.
- In summary: this initialization does not severely inhibit the network's performance (though initializing **W**<sup>(2)</sup>, **b**<sup>(2)</sup> to 0 is still not recommended).



- Suppose that each weight matrix & bias vector consists of the same row repeated many times.
- What will happen during SGD?



- In this case, every node of **h** (and **z**) has the same value.
- Therefore, the gradient update to each row of **W**<sup>(1)</sup> and **b**<sup>(1)</sup> has the same value.
- The NN is performing redundant computation although it has 2 hidden units, it might as well just have 1!



### Weight initialization methods

- There are various different methods of initializing the weights of a neural network.
- One common approach:
  - For weight matrix W<sup>(j)</sup>, sample each component from a 0-mean Gaussian with deviation 1/√cols(W<sup>(j)</sup>).
    - Within certain NNs, helps to ensure that the gradients are usually non-zero.

### Weight initialization methods

- There are various different methods of initializing the weights of a neural network.
- One common approach:
  - For weight matrix W<sup>(j)</sup>, sample each component from a 0-mean Gaussian with deviation 1/√cols(W<sup>(j)</sup>).
    - Within certain NNs, helps to ensure that the gradients are usually non-zero.
  - Optional: orthogonalize the rows of W<sup>(j)</sup> to reduce correlation between different units of the pre-activation layer z<sup>(j)</sup>.

#### Regularization

#### L<sub>2</sub> regularization in NNs

• To prevent the weight matrices from growing too big, we can apply an  $L_2$  regularization term to each matrix by augmenting the cross-entropy loss:

$$f_{\text{CE}}(\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}) = -\frac{1}{n} \sum_{i=1}^{n} \sum_{k=1}^{10} \mathbf{y}_{k}^{(i)} \log \hat{\mathbf{y}}_{k}^{(i)} + \frac{1}{2} \|\mathbf{W}^{(1)}\|_{\text{Fr}}^{2} + \frac{1}{2} \|\mathbf{W}^{(2)}\|_{\text{Fr}}^{2}$$

- Here, |W|<sub>Fr²</sub> means the squared Frobenius norm of W.
  - It's just the sum of squares of all the elements of W.

#### L<sub>2</sub> regularization in NNs

 In practice, this just means that the gradients have an additional term, e.g.:

$$\nabla_{\mathbf{W}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y}) \mathbf{h}^{(1)}^{\top} + \mathbf{W}^{(2)}$$
$$\nabla_{\mathbf{W}^{(1)}} f_{\text{CE}} = \mathbf{g} \mathbf{x}^{\top} + \mathbf{W}^{(1)}$$