

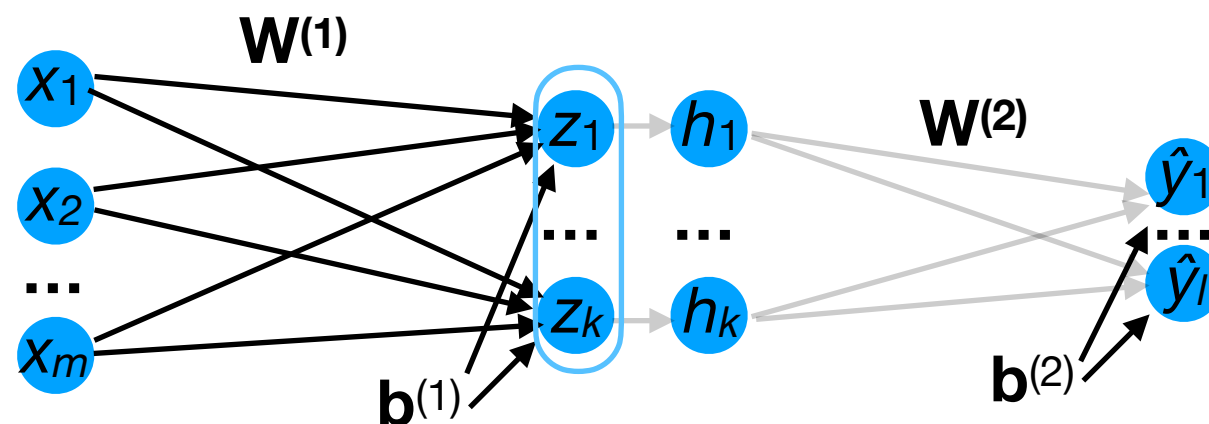
CS 4342: Class 21

Jacob Whitehill

Forwards and backwards propagation

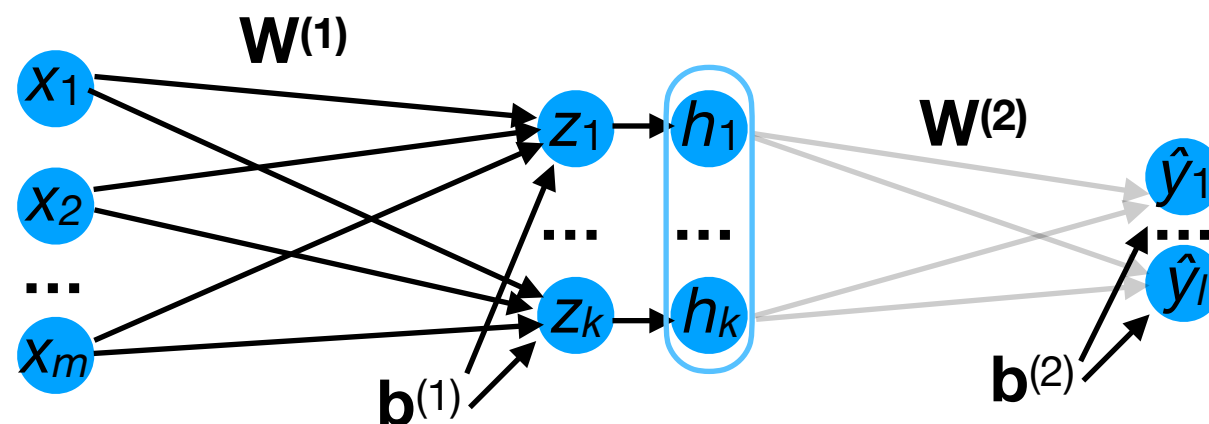
Computing the gradients

- Consider the 3-layer NN below:
 - From \mathbf{x} , $\mathbf{W}^{(1)}$, and $\mathbf{b}^{(1)}$, we can compute \mathbf{z} .



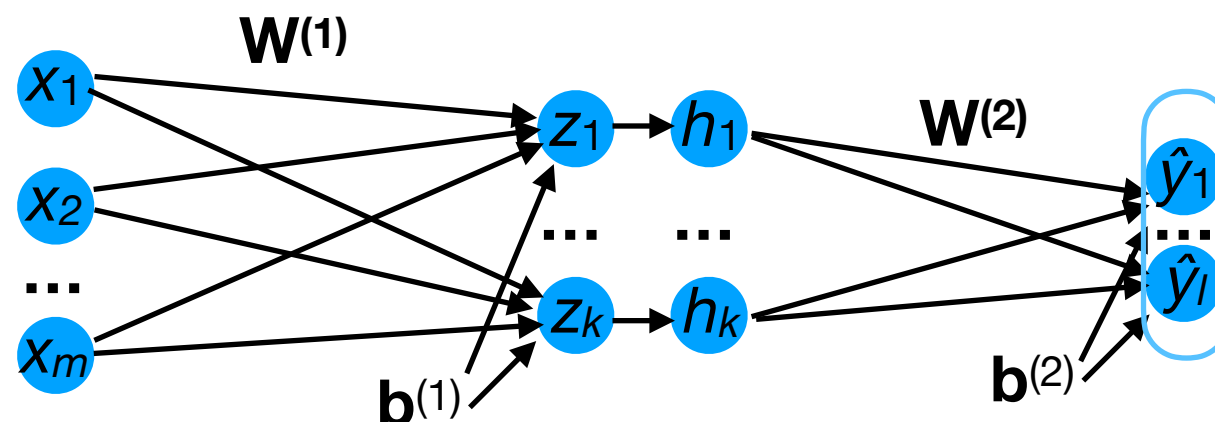
Computing the gradients

- Consider the 3-layer NN below:
 - From \mathbf{x} , $\mathbf{W}^{(1)}$, and $\mathbf{b}^{(1)}$, we can compute \mathbf{z} .
 - From \mathbf{z} and σ , we can compute $\mathbf{h} = \sigma(\mathbf{z})$.



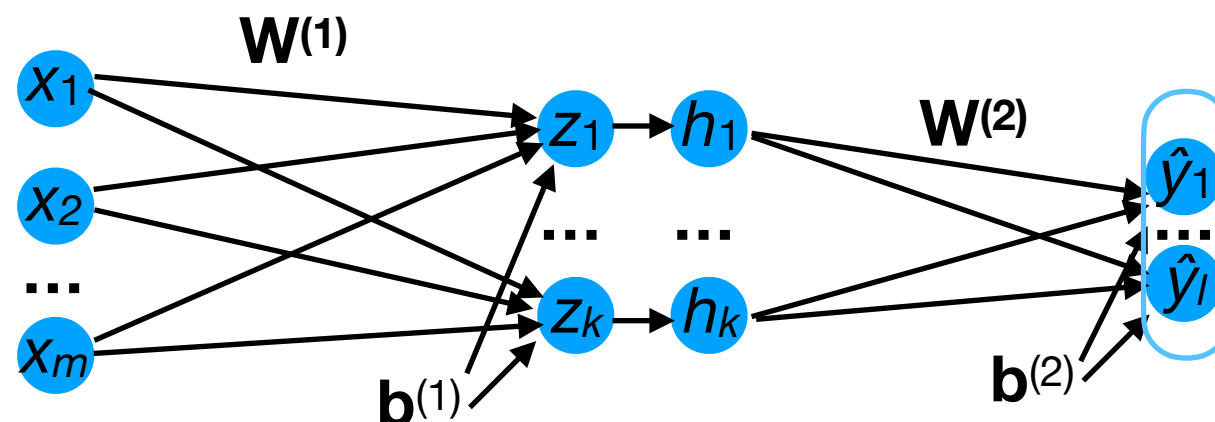
Computing the gradients

- Consider the 3-layer NN below:
 - From \mathbf{x} , $\mathbf{W}^{(1)}$, and $\mathbf{b}^{(1)}$, we can compute \mathbf{z} .
 - From \mathbf{z} and σ , we can compute $\mathbf{h} = \sigma(\mathbf{z})$.
 - From \mathbf{h} , $\mathbf{W}^{(2)}$, and $\mathbf{b}^{(2)}$, we can compute $\hat{\mathbf{y}}$.



Computing the gradients

- This process is known as **forward propagation**.
 - It produces all the intermediary (\mathbf{h} , \mathbf{z}) and final ($\hat{\mathbf{y}}$) network outputs.



Computing the gradients

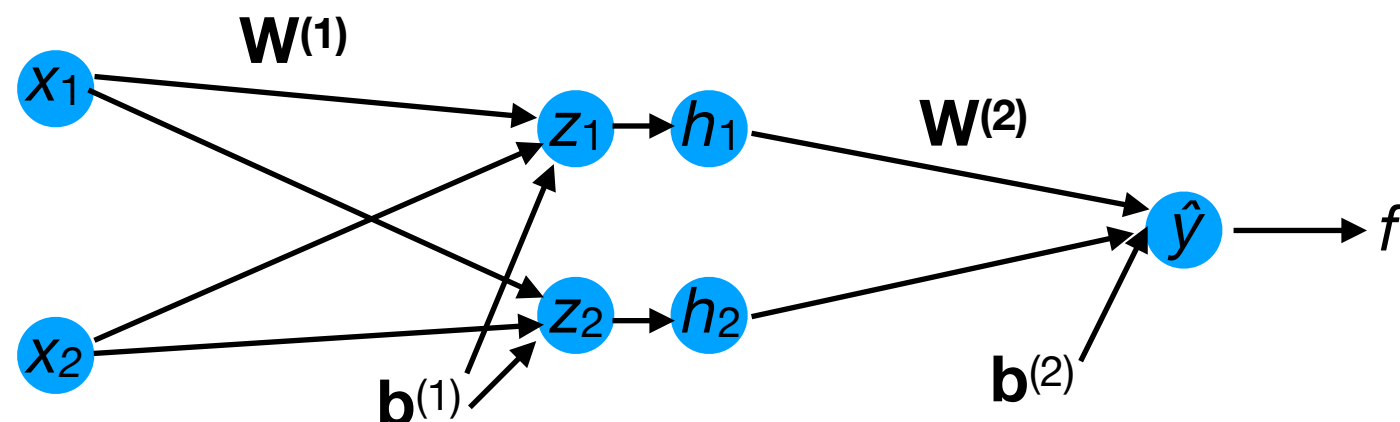
- Now, let's look at how to compute each gradient term:

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

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

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

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

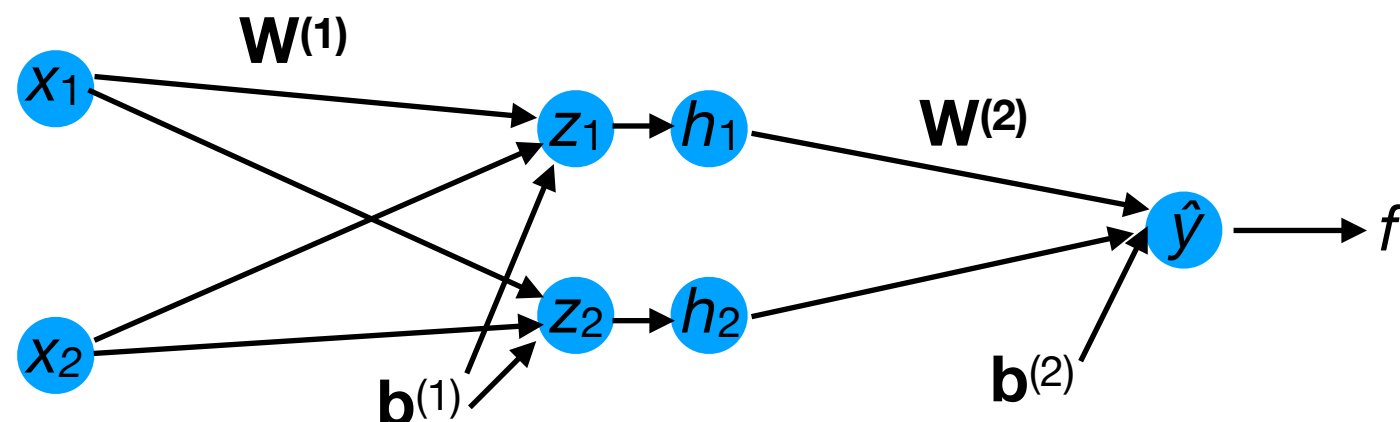


Computing the gradients

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

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

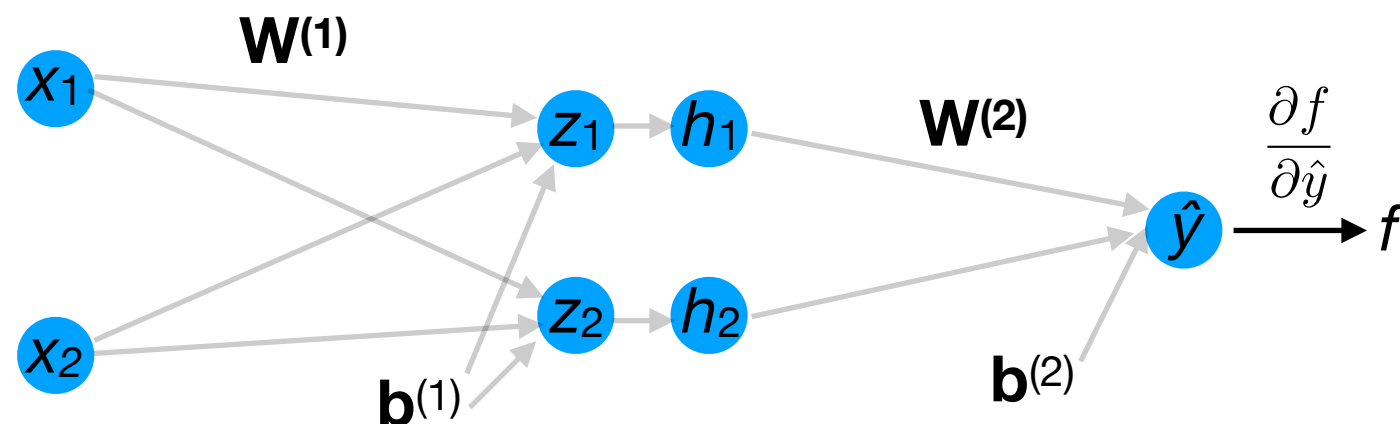
Redundant computation



Computing the gradients

- Here's how we can compute all these *efficiently*:

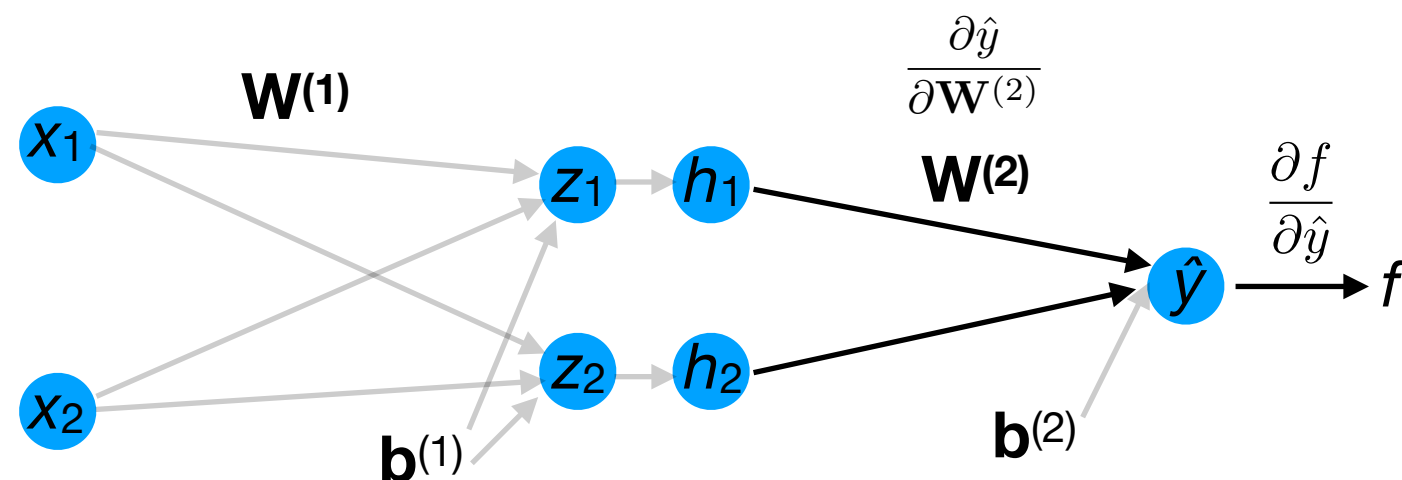
$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}}$$



Computing the gradients

- Here's how we can compute all these *efficiently*:

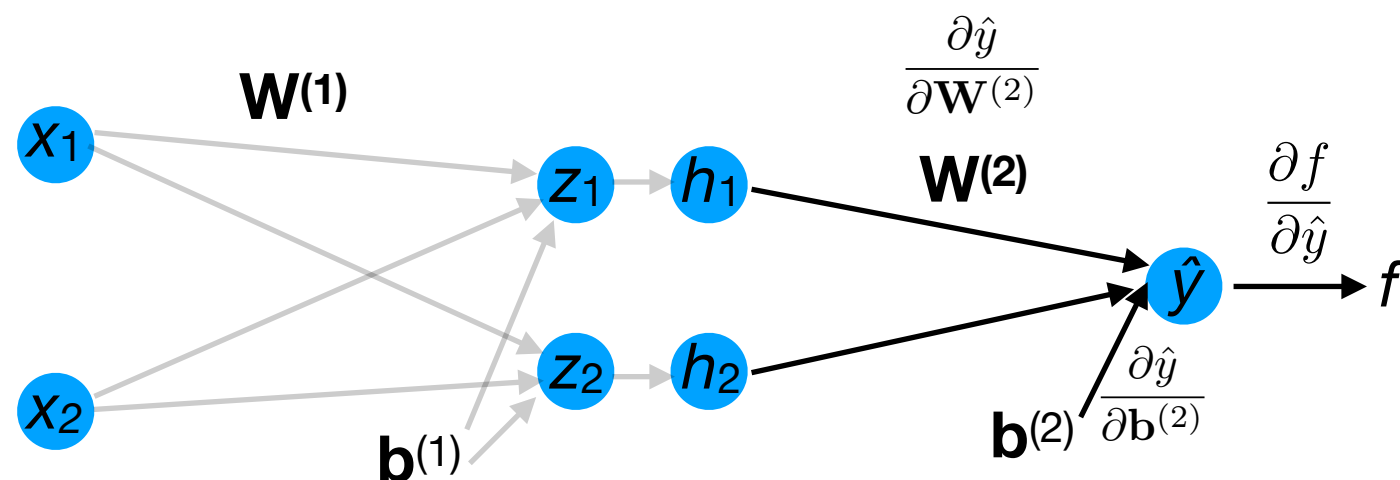
$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}}$$



Computing the gradients

- Here's how we can compute all these *efficiently*:

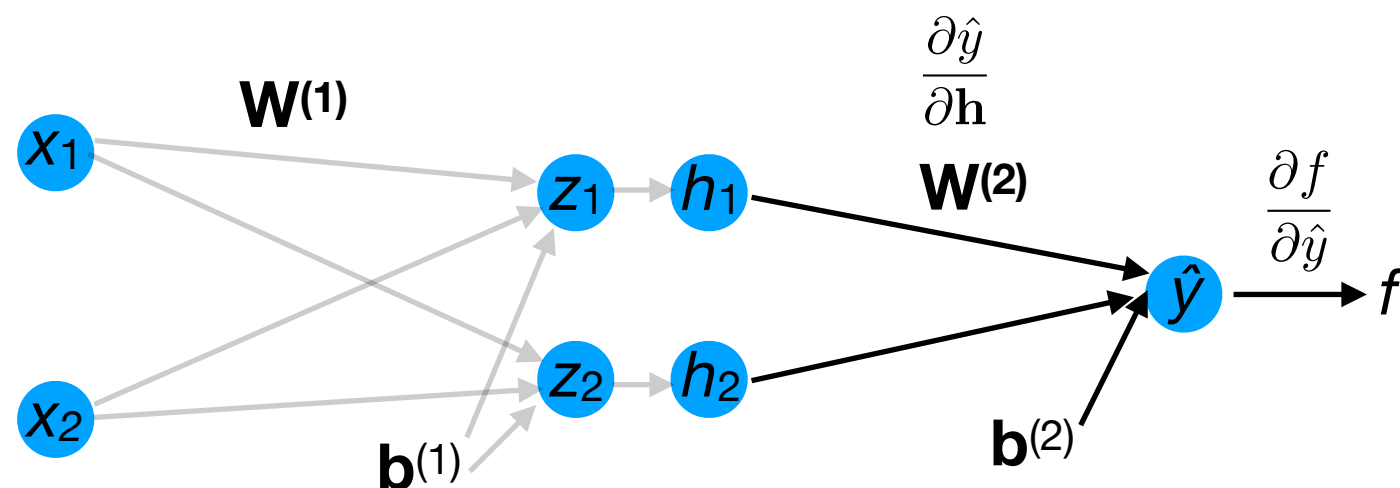
$$\frac{\partial f}{\partial \mathbf{W}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}}$$
$$\frac{\partial f}{\partial \mathbf{b}^{(2)}} = \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}}$$



Computing the gradients

- Here's how we can compute all these *efficiently*:

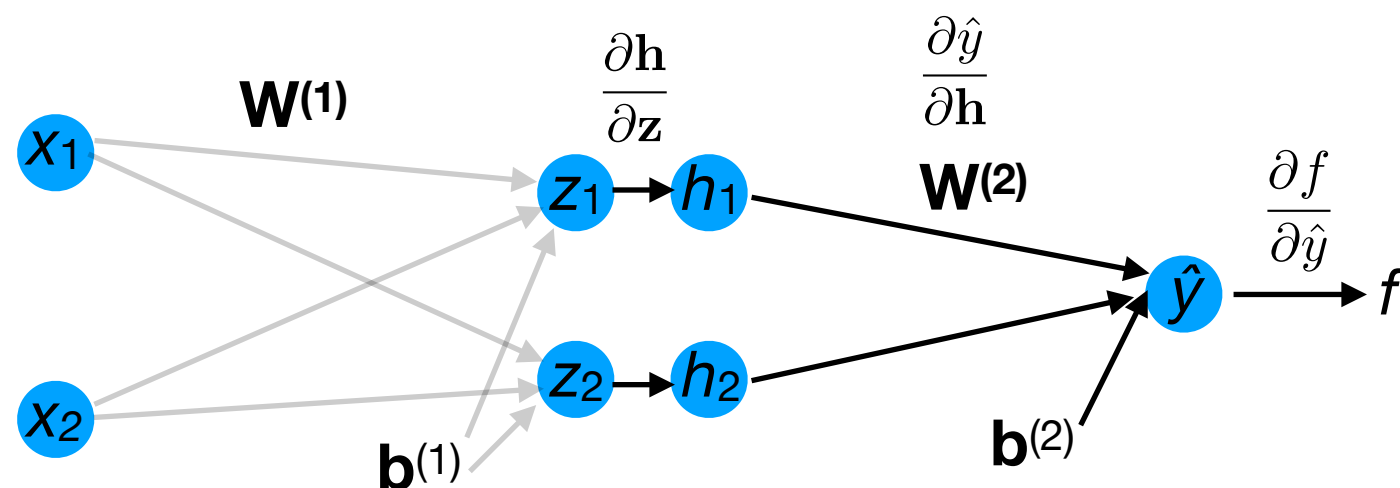
$$\begin{aligned}\frac{\partial f}{\partial \mathbf{W}^{(2)}} &= \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}} \\ \frac{\partial f}{\partial \mathbf{b}^{(2)}} &= \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}} \\ \frac{\partial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



Computing the gradients

- Here's how we can compute all these *efficiently*:

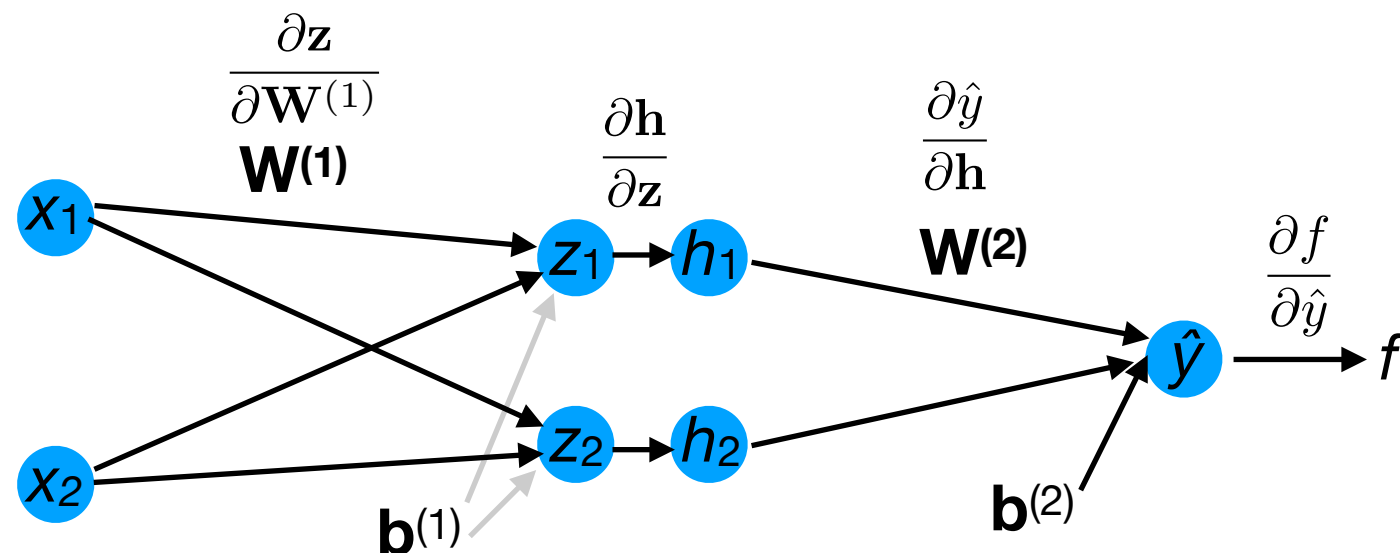
$$\begin{aligned}\frac{\partial f}{\partial \mathbf{W}^{(2)}} &= \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}} \\ \frac{\partial f}{\partial \mathbf{b}^{(2)}} &= \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}} \\ \frac{\partial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



Computing the gradients

- Here's how we can compute all these *efficiently*:

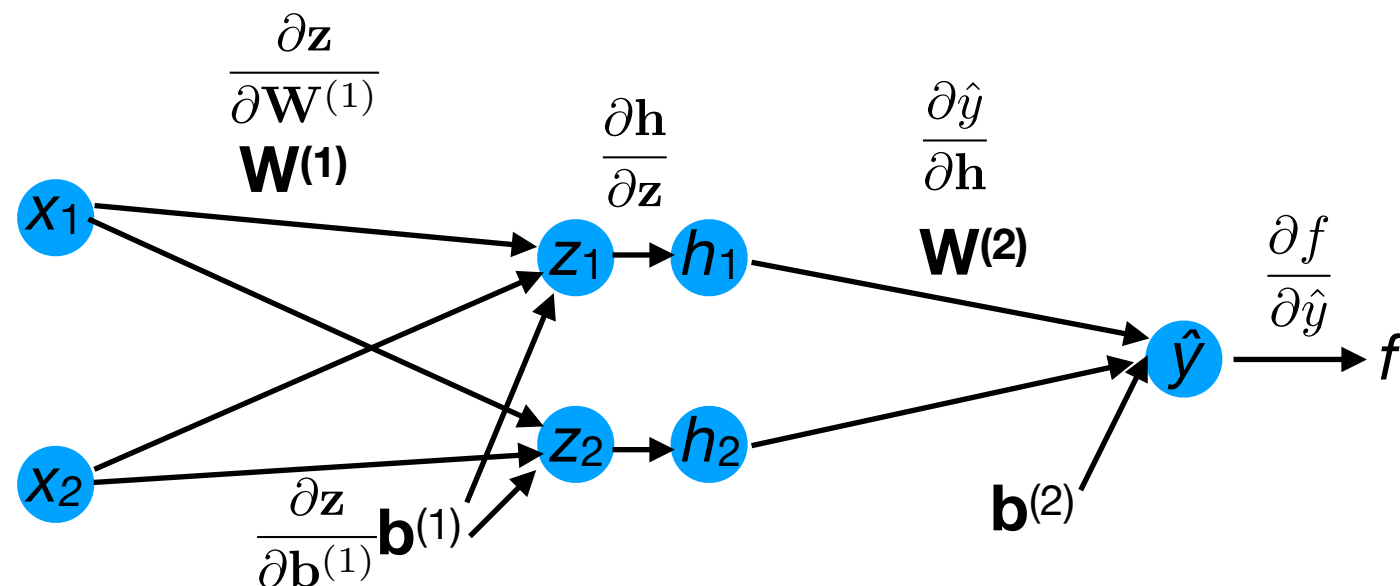
$$\begin{aligned}\frac{\partial f}{\partial \mathbf{W}^{(2)}} &= \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{W}^{(2)}} \\ \frac{\partial f}{\partial \mathbf{b}^{(2)}} &= \frac{\partial f}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}} \\ \frac{\partial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



Computing the gradients

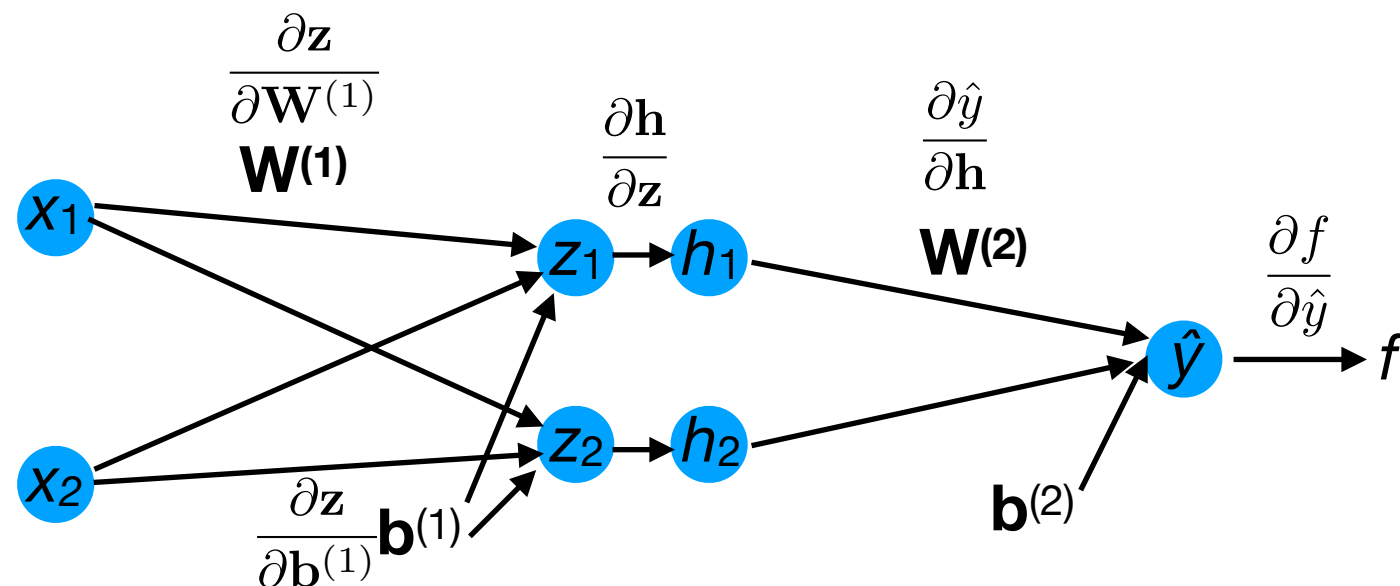
- Here's how we can compute all these *efficiently*:

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



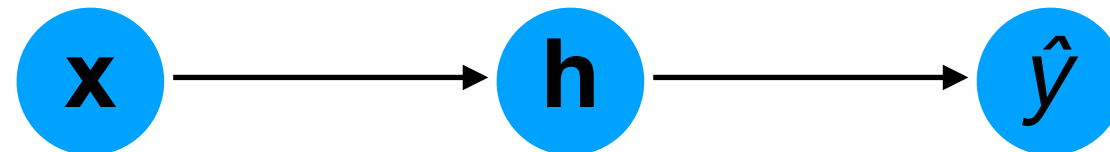
Computing the gradients

- This process is known as **backwards propagation** (“**backprop**”):
 - It produces the gradient terms of all the weight matrices and bias vectors.
 - It requires first conducting forward propagation.

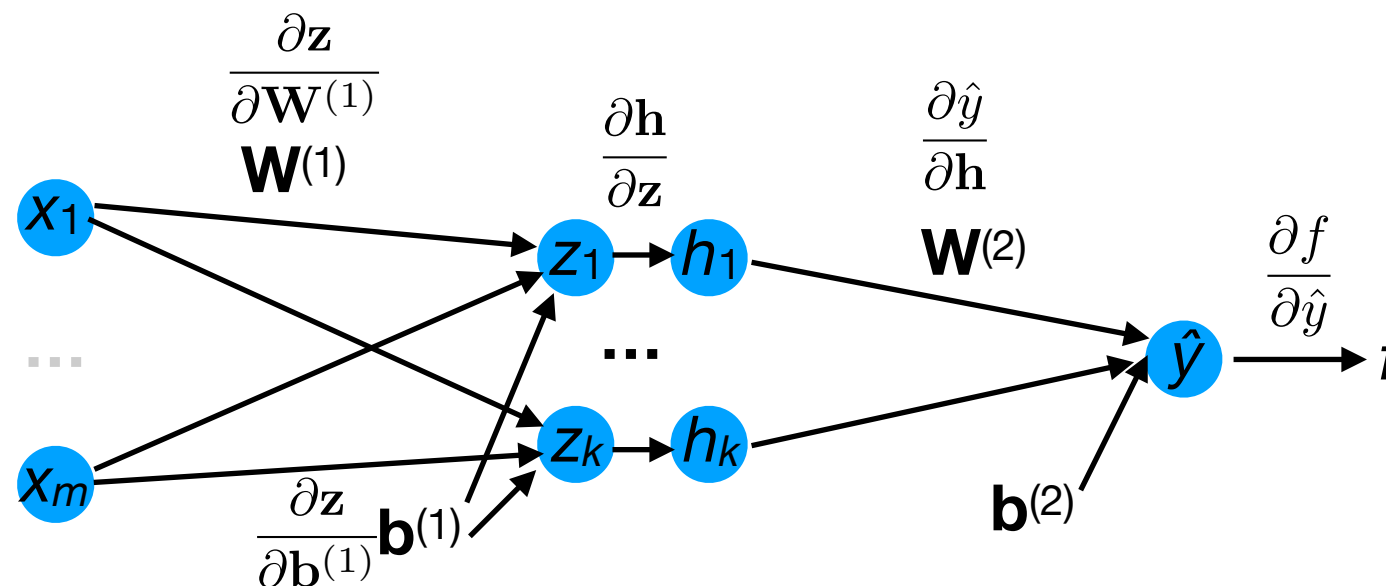
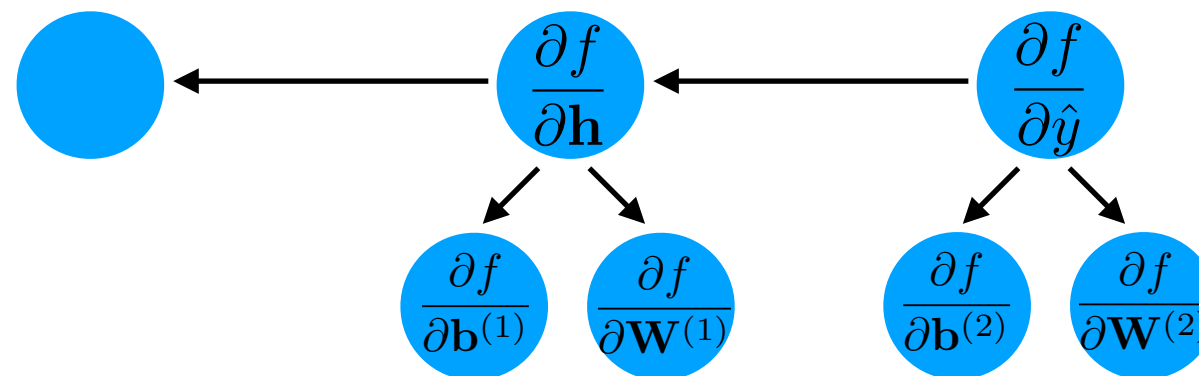


Computing the gradients

Forward propagation



Backward propagation



Computing the gradients

- Where do these come from?

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

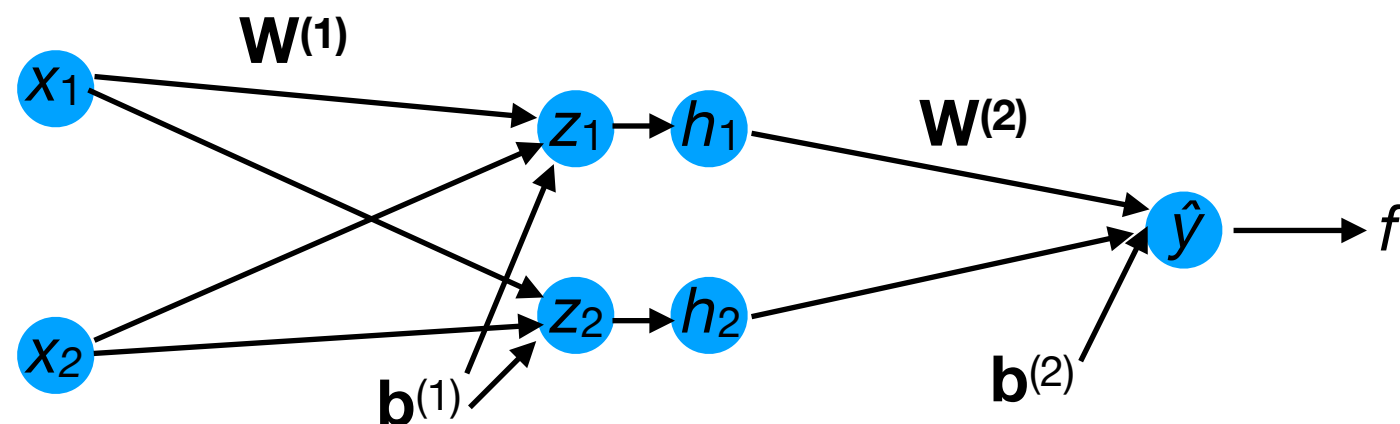
$$\nabla_{\mathbf{b}^{(2)}} f_{\text{CE}} = (\hat{\mathbf{y}} - \mathbf{y})$$

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

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

where

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



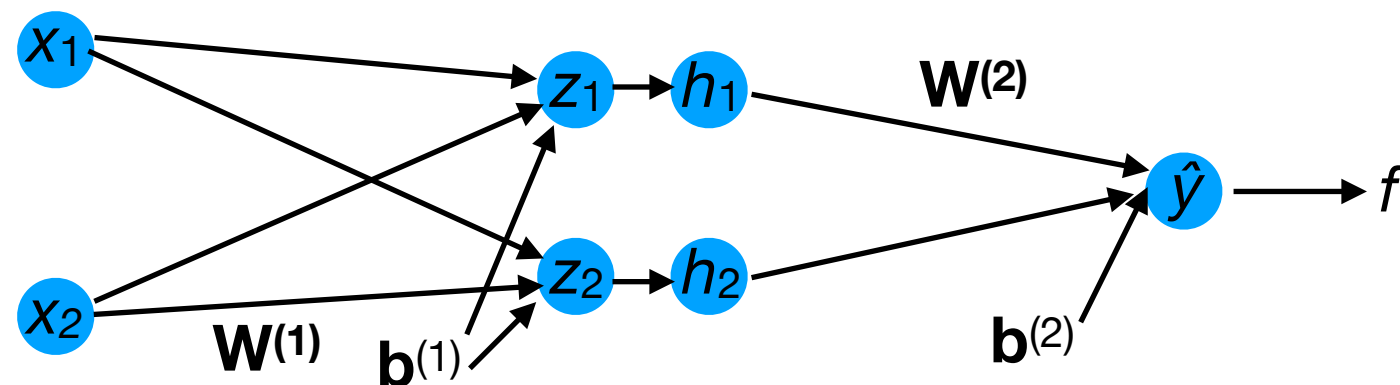
Computing the gradients

- Let's derive each gradient term in turn:

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

- How does f depend on $\hat{\mathbf{y}}$?

$$\frac{\partial f}{\partial \hat{\mathbf{y}}} = (\hat{\mathbf{y}} - \mathbf{y})^\top$$



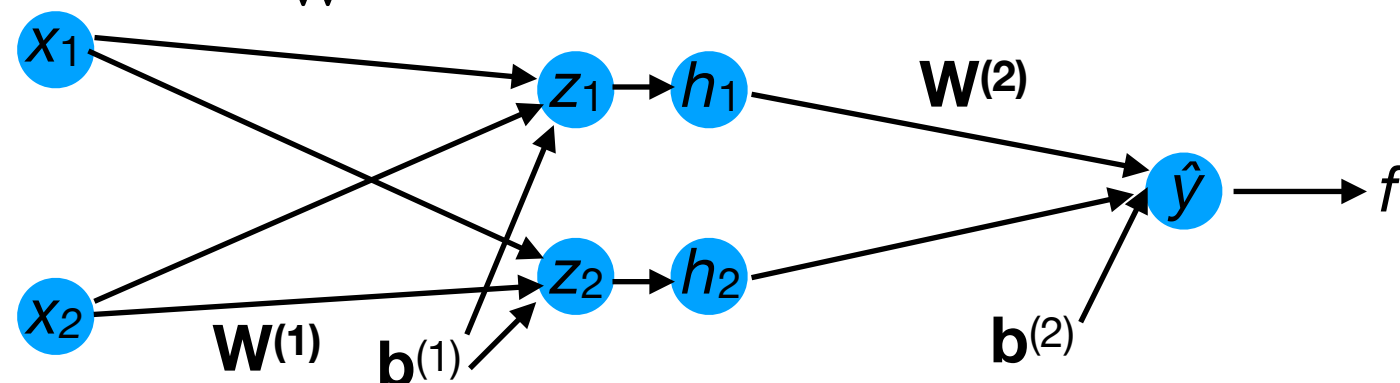
Computing the gradients

- 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 \mathbf{h} ?

$$\begin{aligned}\hat{y} &= \mathbf{W}^{(2)} \mathbf{h} + \mathbf{b}^{(2)} \\ &= \mathbf{W}_1^{(2)} h_1 + \mathbf{W}_2^{(2)} h_2 + \mathbf{b}^{(2)} \\ \Rightarrow \frac{\partial \hat{y}}{\partial \mathbf{h}} &= \begin{bmatrix} \frac{\partial \hat{y}}{\partial h_1} & \frac{\partial \hat{y}}{\partial h_2} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{W}_1^{(2)} & \mathbf{W}_2^{(2)} \end{bmatrix} \\ &= \mathbf{W}^{(2)}\end{aligned}$$



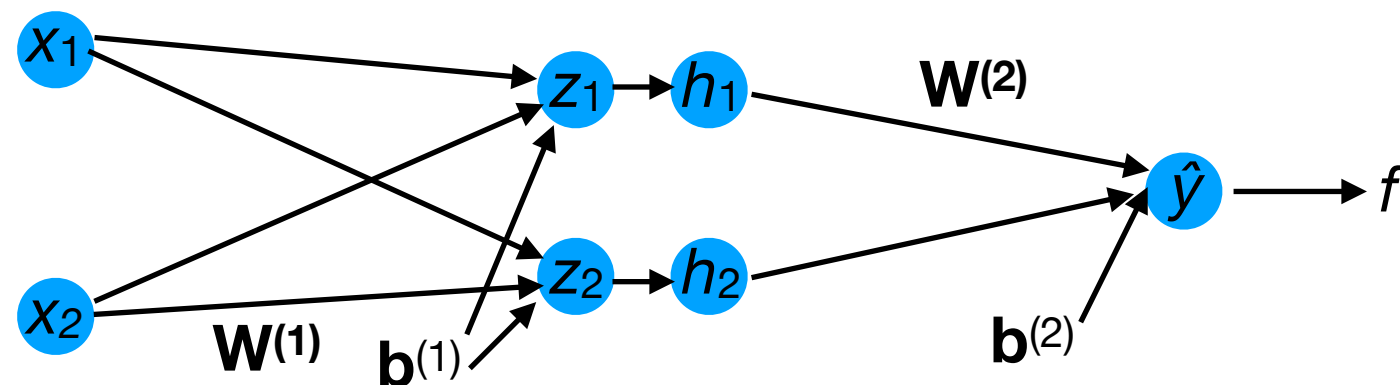
Computing the gradients

- 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 \mathbf{h} depend on \mathbf{z} ?

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



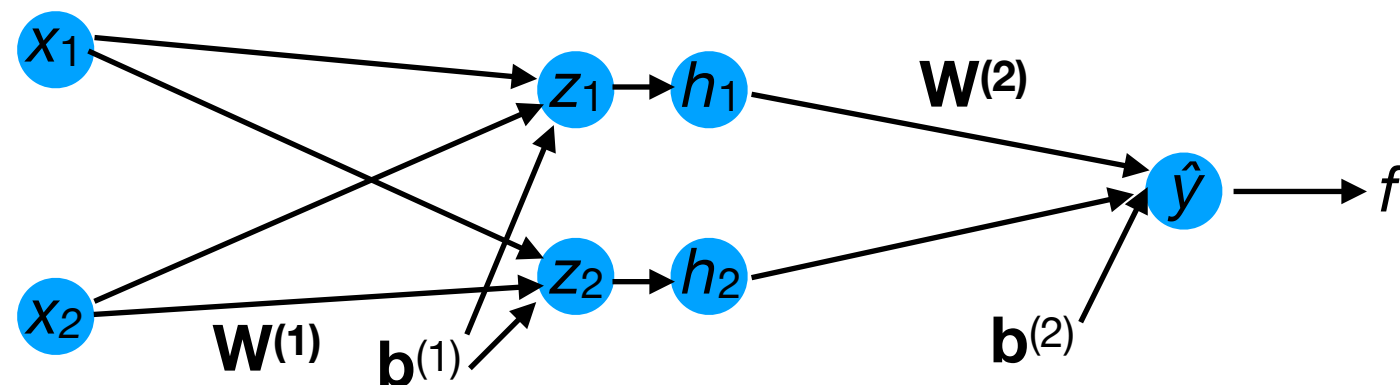
Computing the gradients

- 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 \mathbf{h} depend on \mathbf{z} ?

$$\mathbf{h} = \begin{bmatrix} \text{relu}(\mathbf{z}_1) \\ \text{relu}(\mathbf{z}_2) \end{bmatrix}$$
$$\Rightarrow \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial h_1}{\partial z_1} & \frac{\partial h_1}{\partial z_2} \\ \frac{\partial h_2}{\partial z_1} & \frac{\partial h_2}{\partial z_2} \end{bmatrix}$$



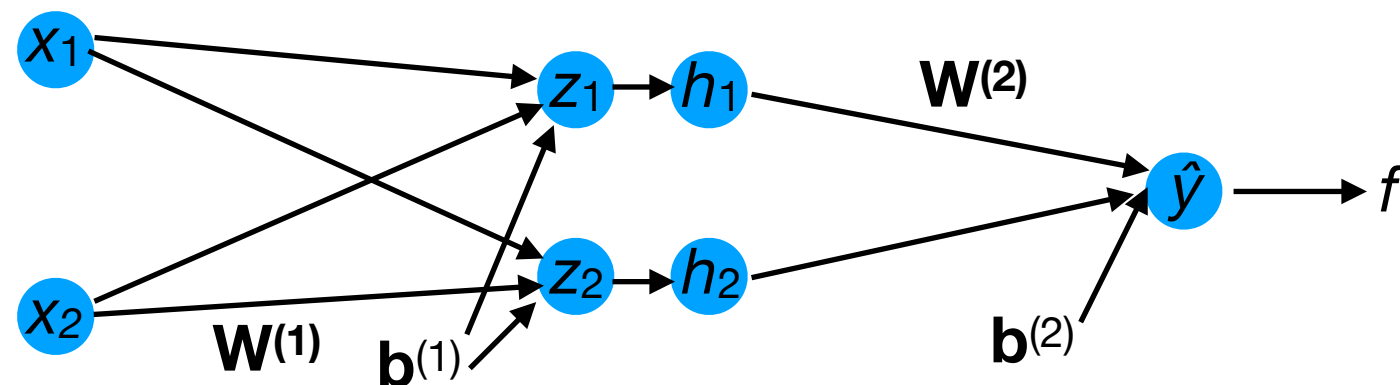
Computing the gradients

- 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 \mathbf{h} depend on \mathbf{z} ?

$$\begin{aligned} \mathbf{h} &= \begin{bmatrix} \text{relu}(\mathbf{z}_1) \\ \text{relu}(\mathbf{z}_2) \end{bmatrix} \\ \Rightarrow \frac{\partial \mathbf{h}}{\partial \mathbf{z}} &= \begin{bmatrix} \frac{\partial h_1}{\partial \mathbf{z}_1} & \frac{\partial h_1}{\partial \mathbf{z}_2} \\ \frac{\partial h_2}{\partial \mathbf{z}_1} & \frac{\partial h_2}{\partial \mathbf{z}_2} \end{bmatrix} \\ &= \begin{bmatrix} \text{relu}'(\mathbf{z}_1) & 0 \\ 0 & \text{relu}'(\mathbf{z}_2) \end{bmatrix} \end{aligned}$$



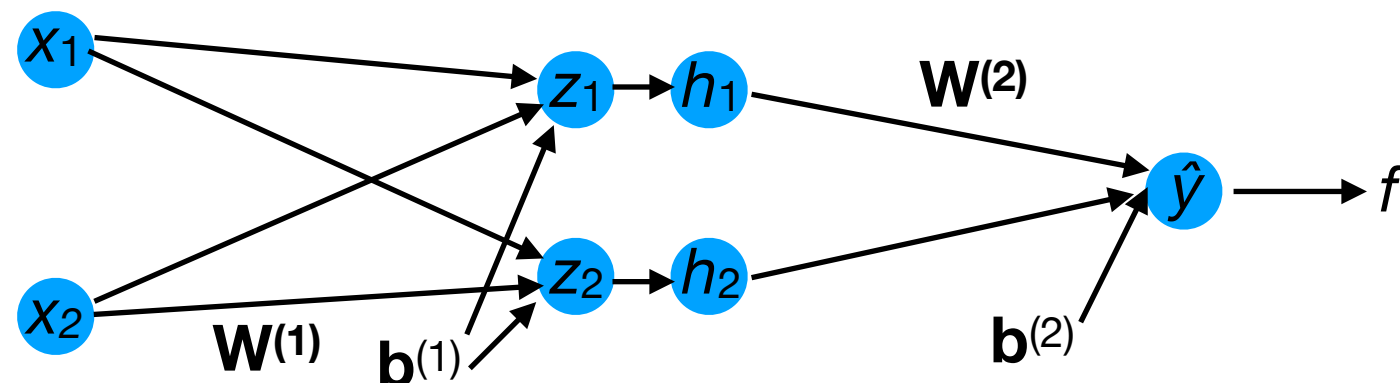
Computing the gradients

- 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 \mathbf{z} depend on $\mathbf{W}^{(1)}$?

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



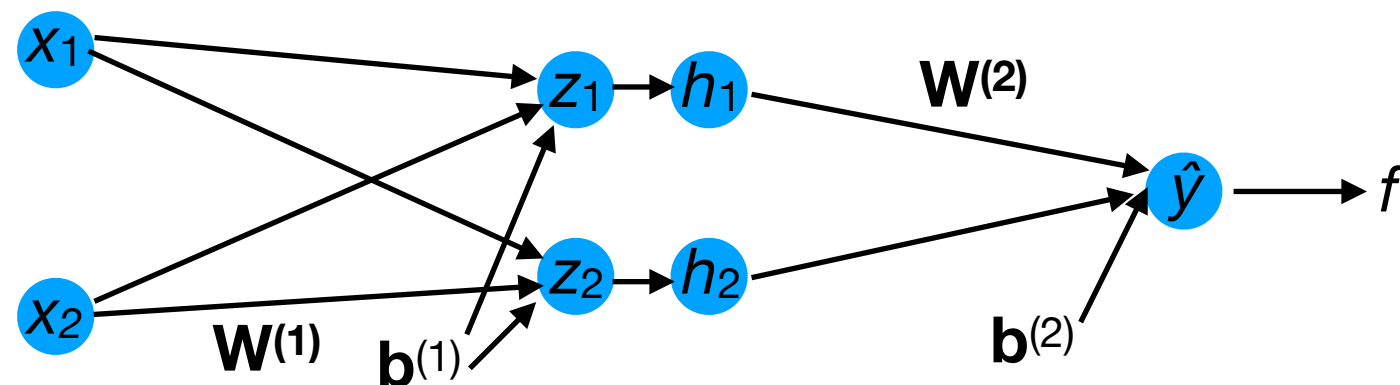
Computing the gradients

- 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 \mathbf{z} depend on $\mathbf{W}^{(1)}$?

$$\mathbf{z} = \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}$$
$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} \mathbf{W}_1^{(1)} & \mathbf{W}_2^{(1)} \\ \mathbf{W}_3^{(1)} & \mathbf{W}_4^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \end{bmatrix}$$



Computing the gradients

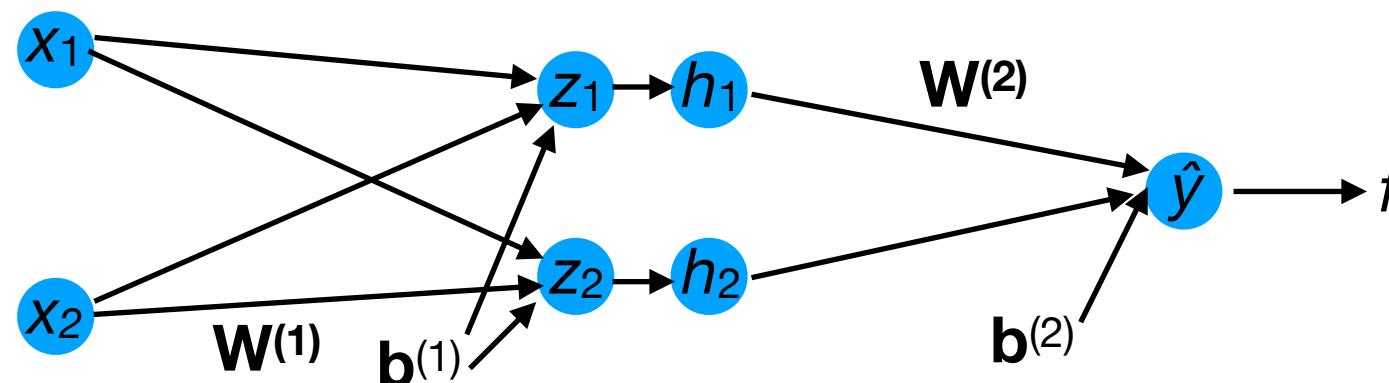
- Let's derive each gradient term in turn:

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

- How does \mathbf{z} depend on $\mathbf{W}^{(1)}$?

$$\begin{aligned} \mathbf{z} &= \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \\ \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{W}_1^{(1)} & \mathbf{W}_2^{(1)} \\ \mathbf{W}_3^{(1)} & \mathbf{W}_4^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \end{bmatrix} \\ \Rightarrow \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} &= \begin{bmatrix} \frac{\partial z_1}{\partial \mathbf{W}_1^{(1)}} & \frac{\partial z_1}{\partial \mathbf{W}_2^{(1)}} & \frac{\partial z_1}{\partial \mathbf{W}_3^{(1)}} & \frac{\partial z_1}{\partial \mathbf{W}_4^{(1)}} \\ \frac{\partial z_2}{\partial \mathbf{W}_1^{(1)}} & \frac{\partial z_2}{\partial \mathbf{W}_2^{(1)}} & \frac{\partial z_2}{\partial \mathbf{W}_3^{(1)}} & \frac{\partial z_2}{\partial \mathbf{W}_4^{(1)}} \end{bmatrix} \end{aligned}$$

For Jacobian matrix, we have to treat $\mathbf{W}^{(1)}$ as if it were a vector.



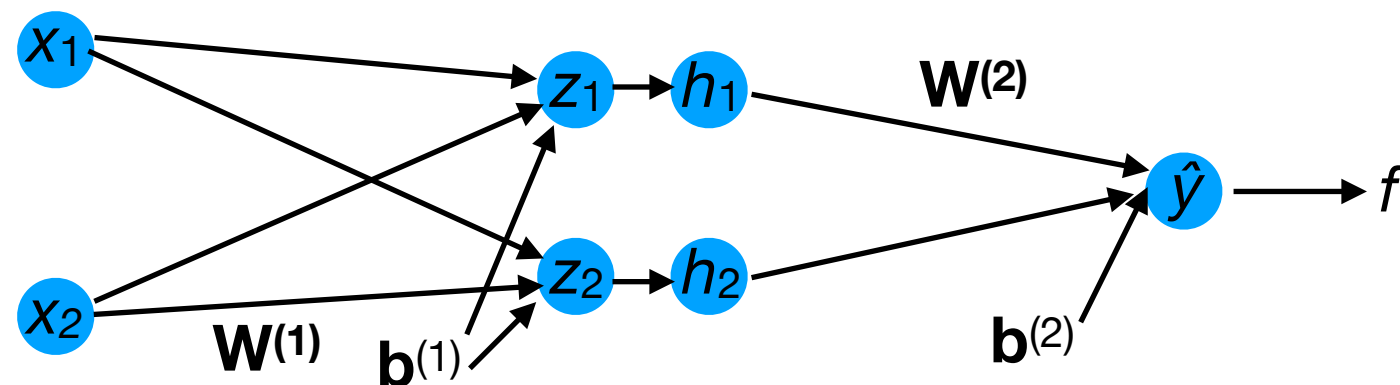
Computing the gradients

- 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 \mathbf{z} depend on $\mathbf{W}^{(1)}$?

$$\begin{aligned} \mathbf{z} &= \mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \\ \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} &= \begin{bmatrix} \mathbf{W}_1^{(1)} & \mathbf{W}_2^{(1)} \\ \mathbf{W}_3^{(1)} & \mathbf{W}_4^{(1)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \end{bmatrix} \\ \Rightarrow \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} &= \begin{bmatrix} \frac{\partial z_1}{\partial \mathbf{W}_1^{(1)}} & \frac{\partial z_1}{\partial \mathbf{W}_2^{(1)}} & \frac{\partial z_1}{\partial \mathbf{W}_3^{(1)}} & \frac{\partial z_1}{\partial \mathbf{W}_4^{(1)}} \\ \frac{\partial z_2}{\partial \mathbf{W}_1^{(1)}} & \frac{\partial z_2}{\partial \mathbf{W}_2^{(1)}} & \frac{\partial z_2}{\partial \mathbf{W}_3^{(1)}} & \frac{\partial z_2}{\partial \mathbf{W}_4^{(1)}} \end{bmatrix} \\ &= \begin{bmatrix} x_1 & x_2 & 0 & 0 \\ 0 & 0 & x_1 & x_2 \end{bmatrix} \end{aligned}$$



Analytical simplification

- We can now finally derive the gradient update for $\mathbf{W}^{(1)}$:

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

Analytical simplification

- We can now finally derive the gradient update for $\mathbf{W}^{(1)}$:

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

Analytical simplification

- We can now finally derive the gradient update for $\mathbf{W}^{(1)}$:

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

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

Analytical simplification

- We can now finally derive the gradient update for $\mathbf{W}^{(1)}$:

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

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

Analytical simplification

- We can now finally derive the gradient update for $\mathbf{W}^{(1)}$:

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

Analytical simplification

- We can now finally derive the gradient update for $\mathbf{W}^{(1)}$:

$$\begin{aligned}\frac{\partial f}{\partial \mathbf{W}^{(1)}} &= \frac{\partial f}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}^{(1)}} \\&= (\hat{\mathbf{y}} - \mathbf{y})^\top \mathbf{W}^{(2)} \begin{bmatrix} \text{relu}'(\mathbf{z}_1) & 0 \\ 0 & \text{relu}'(\mathbf{z}_2) \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & 0 & 0 \\ 0 & 0 & \mathbf{x}_1 & \mathbf{x}_2 \end{bmatrix} \\&= \left(\left((\hat{\mathbf{y}} - \mathbf{y})^\top \mathbf{W}^{(2)} \right) \odot \begin{bmatrix} \text{relu}'(\mathbf{z}_1) & \text{relu}'(\mathbf{z}_2) \end{bmatrix} \right) \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & 0 & 0 \\ 0 & 0 & \mathbf{x}_1 & \mathbf{x}_2 \end{bmatrix} \\&= \begin{bmatrix} \mathbf{g}_1 & \mathbf{g}_2 \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 & \mathbf{x}_2 & 0 & 0 \\ 0 & 0 & \mathbf{x}_1 & \mathbf{x}_2 \end{bmatrix} \\&= \begin{bmatrix} \mathbf{g}_1 \mathbf{x}_1 & \mathbf{g}_1 \mathbf{x}_2 & \mathbf{g}_2 \mathbf{x}_1 & \mathbf{g}_2 \mathbf{x}_2 \end{bmatrix} \\ \Rightarrow \nabla_{\mathbf{W}^{(1)}} f &= \mathbf{g} \mathbf{x}^\top\end{aligned}$$

Outer product