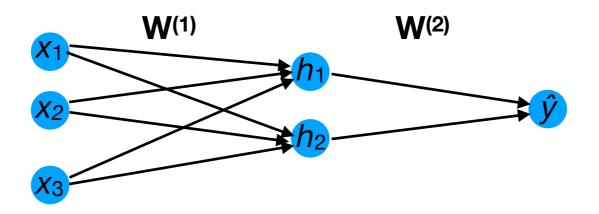
CS 453X: Class 19

Jacob Whitehill

More on neural networks

- Neural networks can have multiple neurons per layer.
- Between each adjacent pair of layers (input-hidden and hidden-output), there is a matrix of (synaptic) weights:



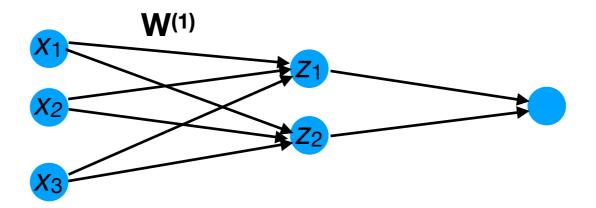
Input layer

Hidden layer

Output layer

 We can compute the pre-activation values z of the hidden layer as:

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

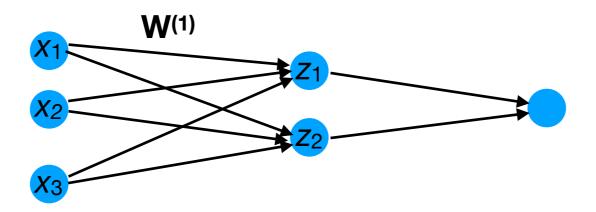


Input layer

Hidden layer Output layer

 We can compute the pre-activation values z of the hidden layer as:

$$\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x}$$
 $\mathbf{W}^{(1)}$ is 2 x 3.

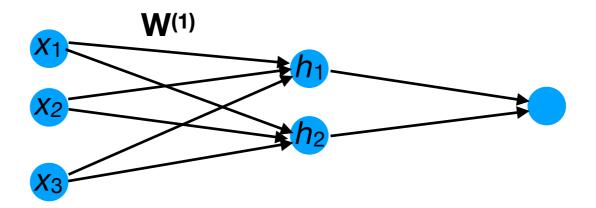


Input layer

Hidden layer Output layer

• We can then pass z to the activation function σ and compute the hidden neuron values **element-wise**, i.e.:

$$\mathbf{h}_j = \sigma(\mathbf{z}_j)$$

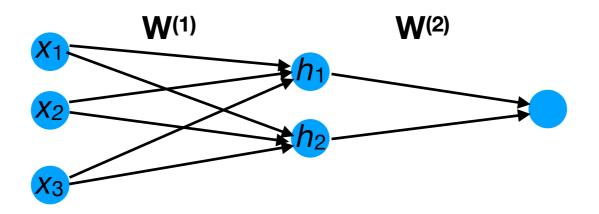


Input layer

Hidden layer

Output layer

Next, we pass h to the next layer...



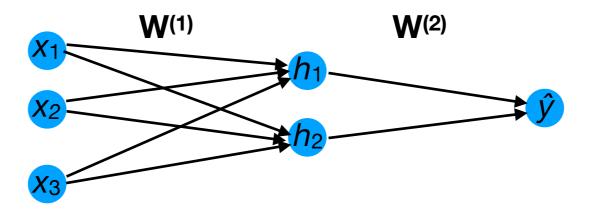
Input layer

Hidden layer

Output layer

...and compute the product:

$$\hat{\mathbf{y}} = \mathbf{W}^{(2)}\mathbf{h}$$



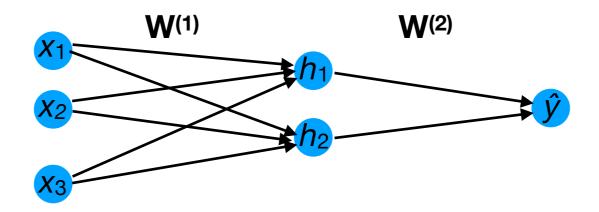
Input layer

Hidden layer Output layer

...and compute the product:

$$\hat{\mathbf{y}} = \mathbf{W}^{(2)}\mathbf{h}$$

 $W^{(2)}$ is 1 x 2.

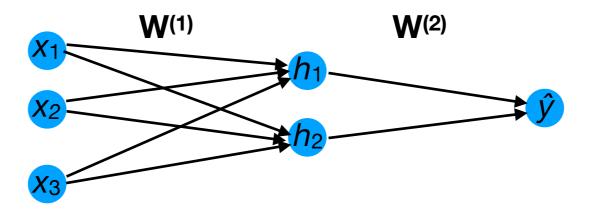


Input layer

Hidden layer Output layer

 Note that the final layer could also have an activation function (if we wanted one), e.g.:

$$\hat{\mathbf{y}} = \sigma \left(\mathbf{W}^{(2)} \mathbf{h} \right)$$



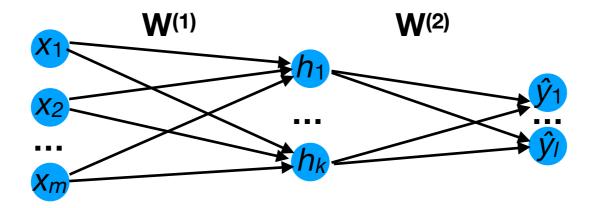
Input layer

Hidden layer

Output layer

Multiple output neurons

We can also have a NN with multiple output neurons, e.g.:



Input layer

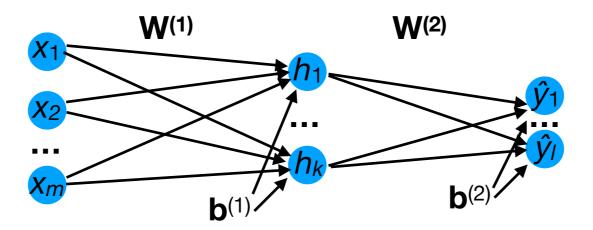
Hidden layer

Output layer

Bias terms

 We typically include a bias term for every neuron, so that the layers' values are computed as:

$$\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}$$
 and $\hat{\mathbf{y}} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$

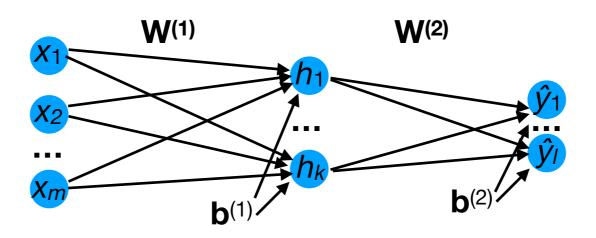


Input layer

Hidden layer Output layer

Exercise: bias terms

• What will $\hat{\mathbf{y}}$ be for $\mathbf{x} = \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} \\ -1 & 2 & 3 \end{bmatrix}$ $\mathbf{b}^{(1)} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$ $\hat{y} = g(\mathbf{x}) = \mathbf{W}^{(2)} \sigma \left(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right) + \mathbf{b}^{(2)}$ $\mathbf{W}^{(2)} = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$ $\mathbf{b}^{(2)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

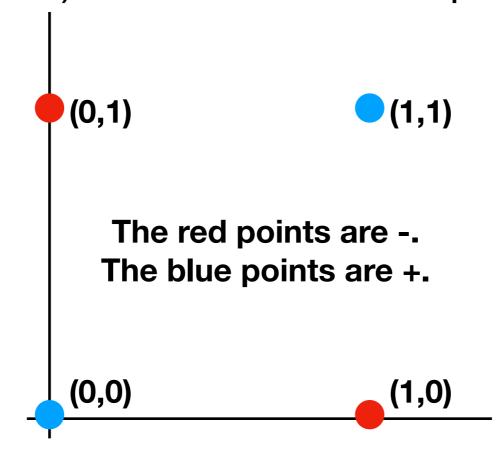


Input layer

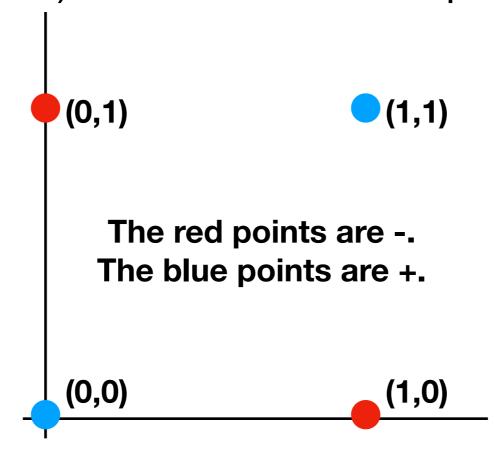
Hidden layer

Output layer

 Recall that no linear decision boundary (e.g., linear SVM, linear regression) can solve the XOR problem.

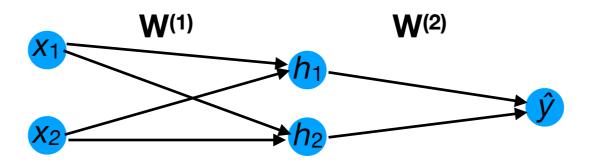


 Recall that no linear decision boundary (e.g., linear SVM, linear regression) can solve the XOR problem.

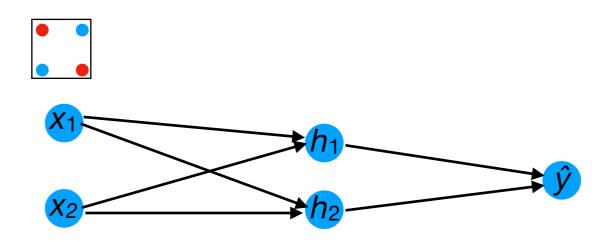


Let's see how using a hidden layer can help us solve it...

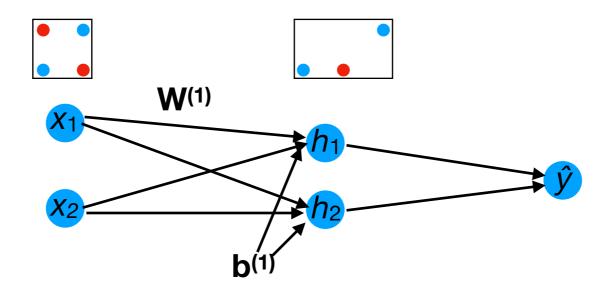
- We want to use a NN to define a function f such that:
 - f(0,0) = 0f(0,1) = 1f(1,0) = 1f(1,1) = 0



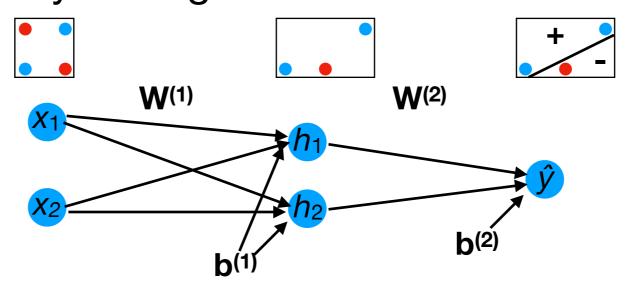
 Here's how a 3-layer NN with a non-linear (ReLU) activation function can solve it:



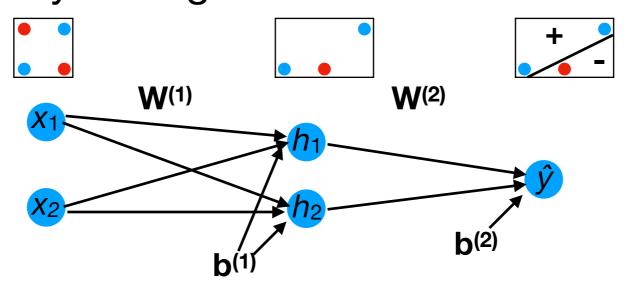
- Here's how a 3-layer NN with a non-linear (ReLU) activation function can solve it:
 - **W**⁽¹⁾, **b**⁽¹⁾ will "collapse" the two data points onto one point in the "hidden" 2-D space.



- Here's how a 3-layer NN with a non-linear (ReLU) activation function can solve it:
 - **W**⁽¹⁾, **b**⁽¹⁾ will "collapse" the two data points onto one point in the "hidden" 2-D space.
 - Since the + and data are now linearly separated, W⁽²⁾,
 b⁽²⁾ can easily distinguish the two classes.

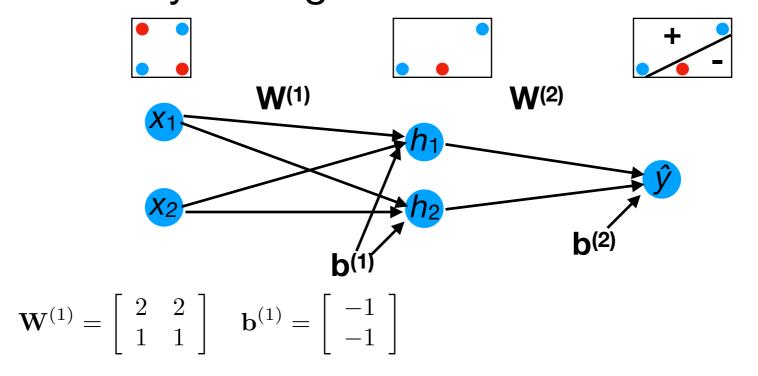


- Here's how a 3-layer NN with a non-linear (ReLU) activation function can solve it:
 - **W**⁽¹⁾, **b**⁽¹⁾ will "collapse" the two data points onto one point in the "hidden" 2-D space.
 - Since the + and data are now linearly separated, W⁽²⁾,
 b⁽²⁾ can easily distinguish the two classes.



What values for W⁽¹⁾, b⁽¹⁾ will make the 4 data points linearly separable? (Hint: set b = $[-1 -1]^T$; W contains only 1s and 2s.)

- Here's how a 3-layer NN with a non-linear (ReLU) activation function can solve it:
 - **W**⁽¹⁾, **b**⁽¹⁾ will "collapse" the two data points onto one point in the "hidden" 2-D space.
 - Since the + and data are now linearly separated, W⁽²⁾,
 b⁽²⁾ can easily distinguish the two classes.

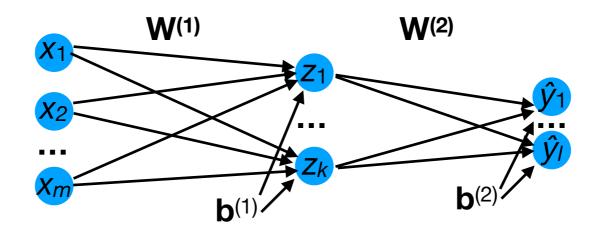


(There are other solutions as well.)

- Note that the ability of the 3-layer NN to solve the XOR problem relies crucially on the non-linear ReLU activation function.
 - Note that other non-linear functions would also work.
- Without non-linearity, a multi-layer NN is no more powerful than a 2-layer network!

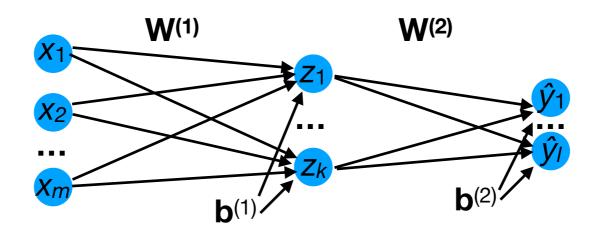
Suppose we define a 3-layer NN without non-linearity:

$$g(\mathbf{x}) = \mathbf{W}^{(2)} \left(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right) + \mathbf{b}^{(2)}$$



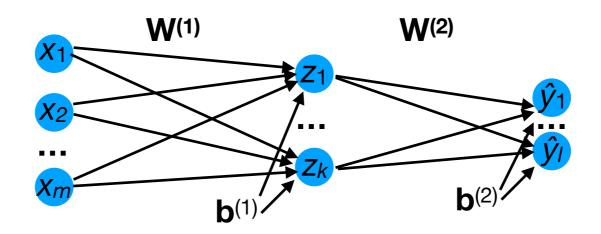
Then we can simplify g to be:

$$g(\mathbf{x}) = \mathbf{W}^{(2)} \left(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right) + \mathbf{b}^{(2)}$$



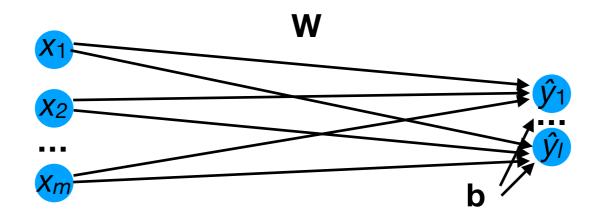
Then we can simplify g to be:

$$g(\mathbf{x}) = \mathbf{W}^{(2)} \left(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right) + \mathbf{b}^{(2)}$$
$$= \left(\mathbf{W}^{(2)} \mathbf{W}^{(1)} \right) \mathbf{x} + \mathbf{W}^{(2)} \mathbf{b}^{(1)} + \mathbf{b}^{(2)}$$



• Then we can simplify g to be:

$$g(\mathbf{x}) = \mathbf{W}^{(2)} \left(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)} \right) + \mathbf{b}^{(2)}$$
$$= \left(\mathbf{W}^{(2)} \mathbf{W}^{(1)} \right) \mathbf{x} + \mathbf{W}^{(2)} \mathbf{b}^{(1)} + \mathbf{b}^{(2)}$$
$$= \mathbf{W} \mathbf{x} + \mathbf{b}$$



Training neural networks

Training neural networks

- While training neural networks by hand is (arguably) fun, it is completely impractical except for toy examples.
- How can we find good values for the weights and bias terms automatically?

Training neural networks

- While training neural networks by hand is (arguably) fun, it is completely impractical except for toy examples.
- How can we find good values for the weights and bias terms automatically?
 - Gradient descent.

- Here is how we can conduct gradient descent for the XOR problem...
- Let's first define:

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}, \mathbf{W}^{(2)} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix}, \mathbf{b}^{(1)} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \mathbf{b}^{(2)} = \begin{bmatrix} b_3 \end{bmatrix}$$

Then we can define g so that:

$$\hat{y} = g\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \mathbf{W}^{(2)}\sigma\left(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}\right) + \mathbf{b}^{(2)}$$

- Here is how we can conduct gradient descent for the XOR problem...
- Let's first define:

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}, \mathbf{W}^{(2)} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix}, \mathbf{b}^{(1)} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \mathbf{b}^{(2)} = \begin{bmatrix} b_3 \end{bmatrix}$$

Then we can define g so that:

$$\hat{y} = g\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \mathbf{W}^{(2)}\sigma\left(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}\right) + \mathbf{b}^{(2)}$$

$$= \begin{bmatrix} w_5 \\ w_6 \end{bmatrix}^\mathsf{T}\sigma\left(\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}\right) + b_3$$

- Here is how we can conduct gradient descent for the XOR problem...
- Let's first define:

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}, \mathbf{W}^{(2)} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix}, \mathbf{b}^{(1)} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}, \mathbf{b}^{(2)} = \begin{bmatrix} b_3 \end{bmatrix}$$

Then we can define g so that:

$$\hat{y} = g\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \mathbf{W}^{(2)}\sigma\left(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}\right) + \mathbf{b}^{(2)}$$

$$= \begin{bmatrix} w_5 \\ w_6 \end{bmatrix}^\mathsf{T}\sigma\left(\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}\right) + b_3$$

$$= w_5\sigma(w_1x_1 + w_2x_2 + b_1) + w_6\sigma(w_3x_1 + w_4x_2 + b_2) + b_3$$

• From \hat{y} , we can compute the cost (f_{MSE}):

$$f_{\text{MSE}}(\hat{y}; w_1, w_2, w_3, w_4, w_5, w_6, b_1, b_2, b_3) = \frac{1}{2} \sum_{i=1}^{n} (\hat{y}^{(i)} - y^{(i)})^2$$

• From \hat{y} , we can compute the cost (f_{MSE}):

$$f_{\text{MSE}}(\hat{y}; w_1, w_2, w_3, w_4, w_5, w_6, b_1, b_2, b_3) = \frac{1}{2} \sum_{i=1}^{n} (\hat{y}^{(i)} - y^{(i)})^2$$

 We then then calculate the derivative of f_{MSE} w.r.t. each parameter p using the chain rule as:

$$\frac{\partial f_{\text{MSE}}}{\partial p} = \frac{1}{2} \sum_{i=1}^{n} \frac{\partial f_{\text{MSE}}}{\partial \hat{y}^{(i)}} \frac{\hat{y}^{(i)}}{\partial p}$$

where:

$$\frac{\partial f_{\text{MSE}}}{\partial \hat{y}^{(i)}} = (\hat{y}^{(i)} - y^{(i)})$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial}{\partial w_1} \left[w_5 \sigma(w_1 x_1 + w_2 x_2 + b_1) + w_6 \sigma(w_3 x_1 + w_4 x_2 + b_2) + b_3 \right]$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial}{\partial w_1} \left[w_5 \sigma(w_1 x_1 + w_2 x_2 + b_1) + w_6 \sigma(w_3 x_1 + w_4 x_2 + b_2) + b_3 \right] \\
= w_5$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial}{\partial w_1} [w_5 \sigma(w_1 x_1 + w_2 x_2 + b_1) + w_6 \sigma(w_3 x_1 + w_4 x_2 + b_2) + b_3]
= w_5 \sigma'(w_1 x_1 + w_2 x_2 + b_1)$$

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial}{\partial w_1} [w_5 \sigma(w_1 x_1 + w_2 x_2 + b_1) + w_6 \sigma(w_3 x_1 + w_4 x_2 + b_2) + b_3]
= w_5 \sigma'(w_1 x_1 + w_2 x_2 + b_1) x_1$$

• Now we just have to differentiate $\hat{y} = g(\mathbf{x})$ w.r.t each parameter p:, e.g.:

$$\frac{\partial \hat{y}}{\partial w_1} = \frac{\partial}{\partial w_1} [w_5 \sigma(w_1 x_1 + w_2 x_2 + b_1) + w_6 \sigma(w_3 x_1 + w_4 x_2 + b_2) + b_3]
= w_5 \sigma'(w_1 x_1 + w_2 x_2 + b_1) x_1$$

where:

$$\sigma'(z) = \left\{ \begin{array}{ll} 0 & \text{if} & z \leq 0 \\ 1 & \text{if} & z > 0 \end{array} \right. \quad \text{for ReLU}$$

• Hence:

$$\frac{\partial \hat{y}}{\partial w_1} = \begin{cases} 0 & \text{if } w_1 x_1 + w_2 x_2 + b_1 \le 0\\ w_5 x_1 & \text{if } w_1 x_1 + w_2 x_2 + b_1 > 0 \end{cases}$$

since:

$$\sigma'(z) = \left\{ \begin{array}{ll} 0 & \text{if} & z \leq 0 \\ 1 & \text{if} & z > 0 \end{array} \right. \quad \text{for ReLU}$$

We also need to derive the other partial derivatives:

$$\frac{\partial \hat{y}}{\partial w_1}, \dots \frac{\partial \hat{y}}{\partial w_6}, \frac{\partial \hat{y}}{\partial b_1}, \dots \frac{\partial \hat{y}}{\partial w_3}$$