

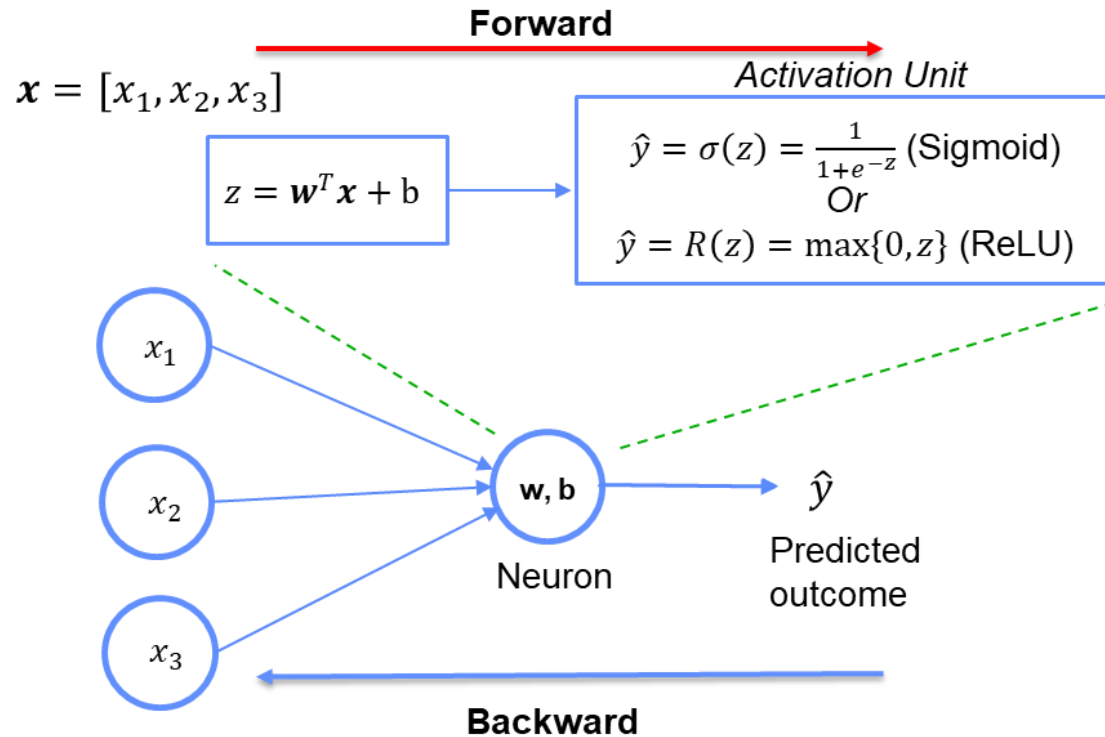
Discussion Neural Network

ECE 498 Data Science & Engr Spring 2020

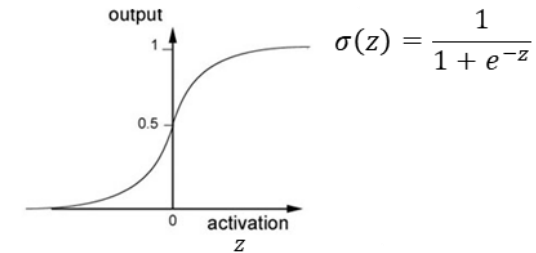
04/24/2020

Perceptron

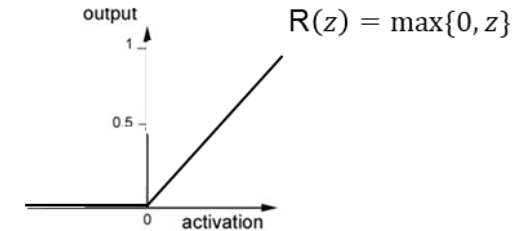
- The core of the neural network is perceptron model



Sigmoid Function



ReLU Function



Update Rule (Backward):

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla J(\mathbf{w}_t)$$

η : Learning rate

Loss

$$J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N L(\mathbf{w} \cdot \mathbf{x}^{(i)}, y^{(i)})$$

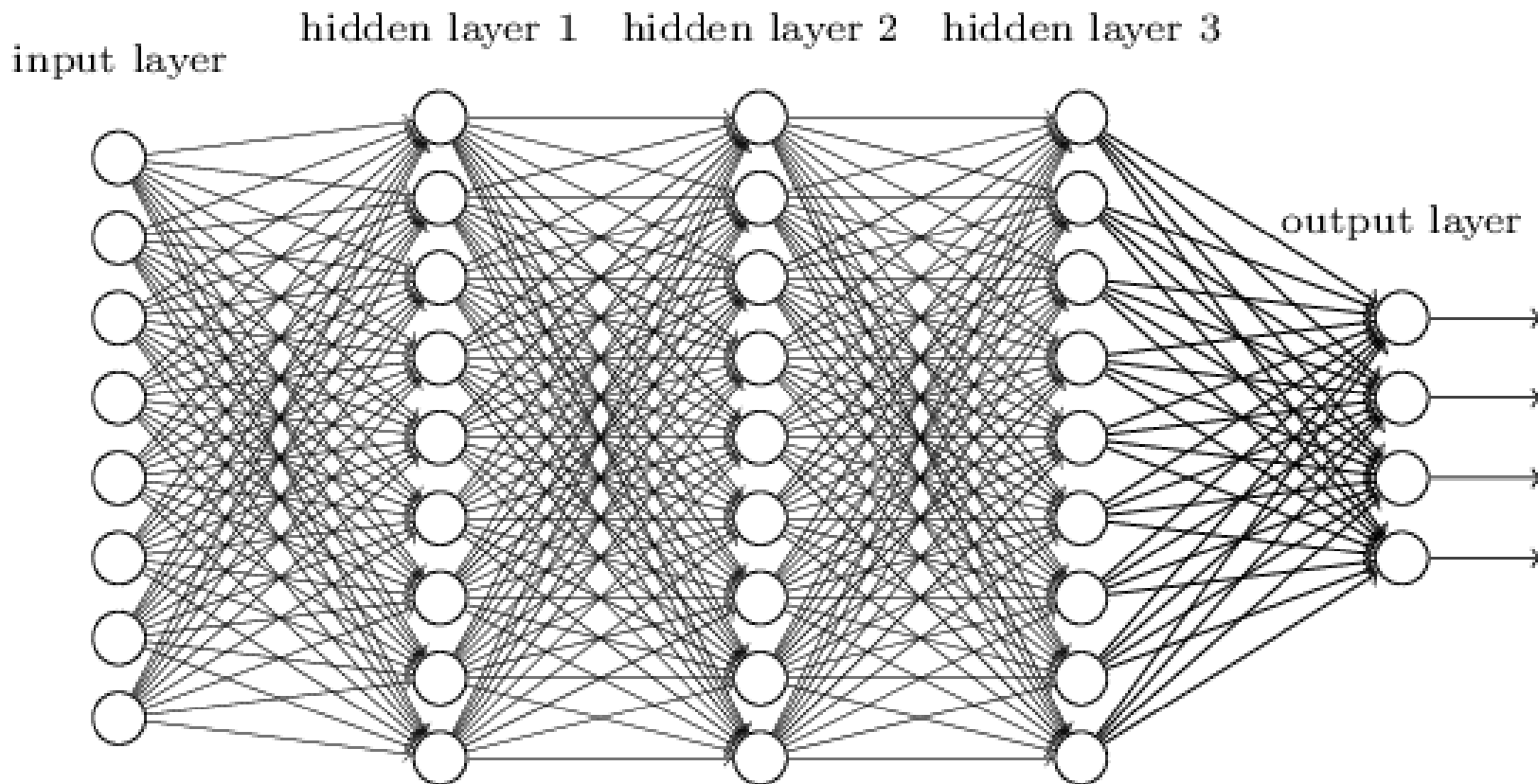
N : number of samples $\mathbf{x}^{(i)}$: feature of i^{th} sample

Computing Gradient

$$\nabla J(\mathbf{w}_0) = \left(\frac{\partial J(\mathbf{w})}{\partial w_0}, \frac{\partial J(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial J(\mathbf{w})}{\partial w_n} \right)_{\mathbf{w}_0}$$

Neural Network

- When we stack perceptron together (universal function approximator)
- we can have many hidden layers, and each layer of arbitrary size (# of neurons)



Why Activation Function?

- Recall that at each layer:

$$z = \mathbf{w}^T \mathbf{x} + b$$

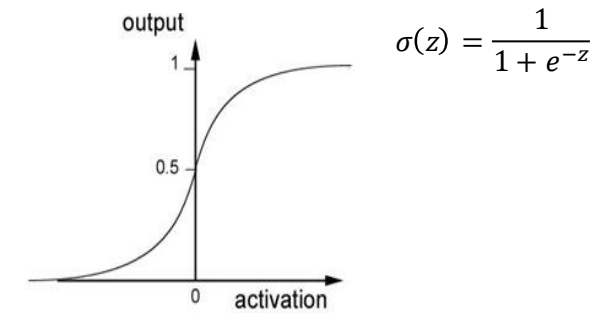
 \longrightarrow

$$\hat{y} = \sigma(z) = \frac{1}{1+e^{-z}} \text{ (Sigmoid)}$$

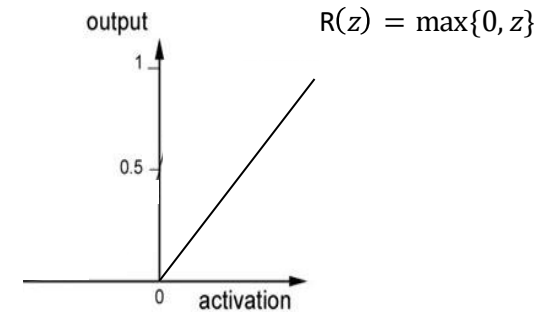
Or

$$\hat{y} = R(z) = \max\{0, z\} \text{ (ReLU)}$$

Sigmoid Function



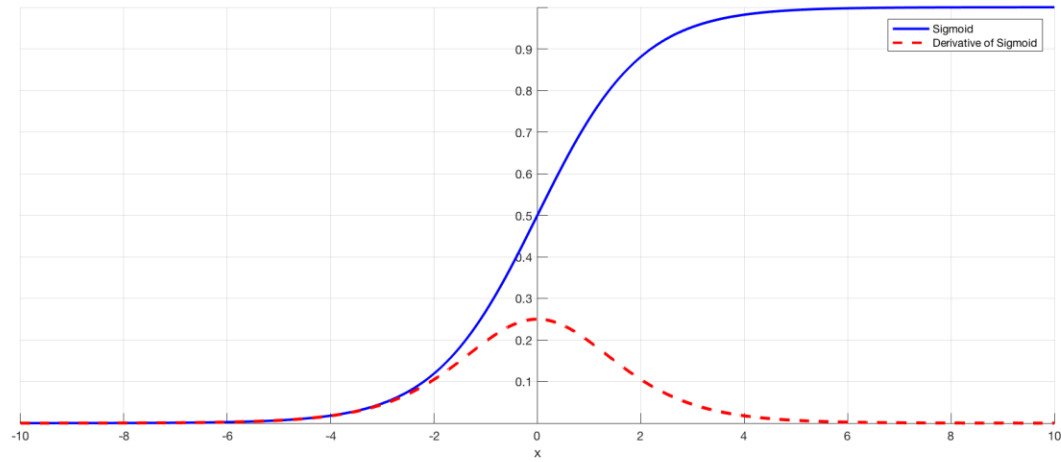
ReLU Function



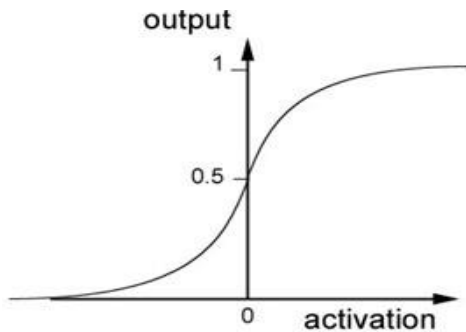
- Activation function adds non-linearity into the network.

ReLU?

- ReLU vs Sigmoid

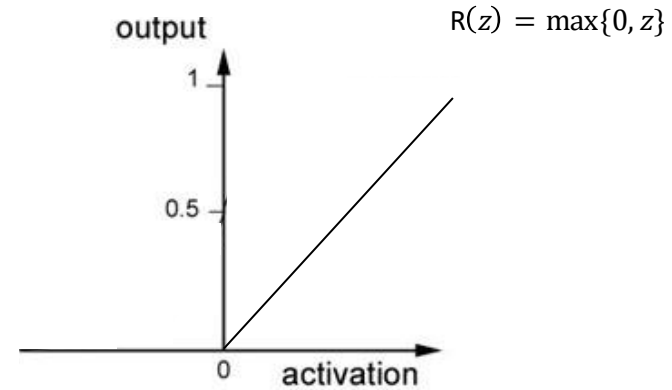


Sigmoid Function



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

ReLU Function

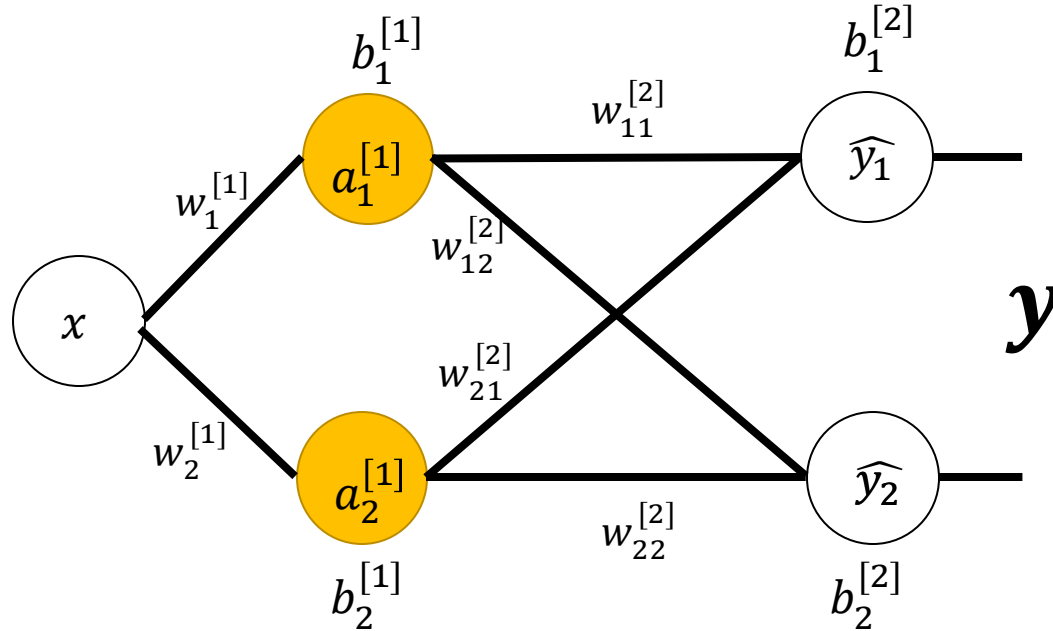


- Vanishing gradient problem in deep networks.
- But could run into numerical issue

Validation vs Test set

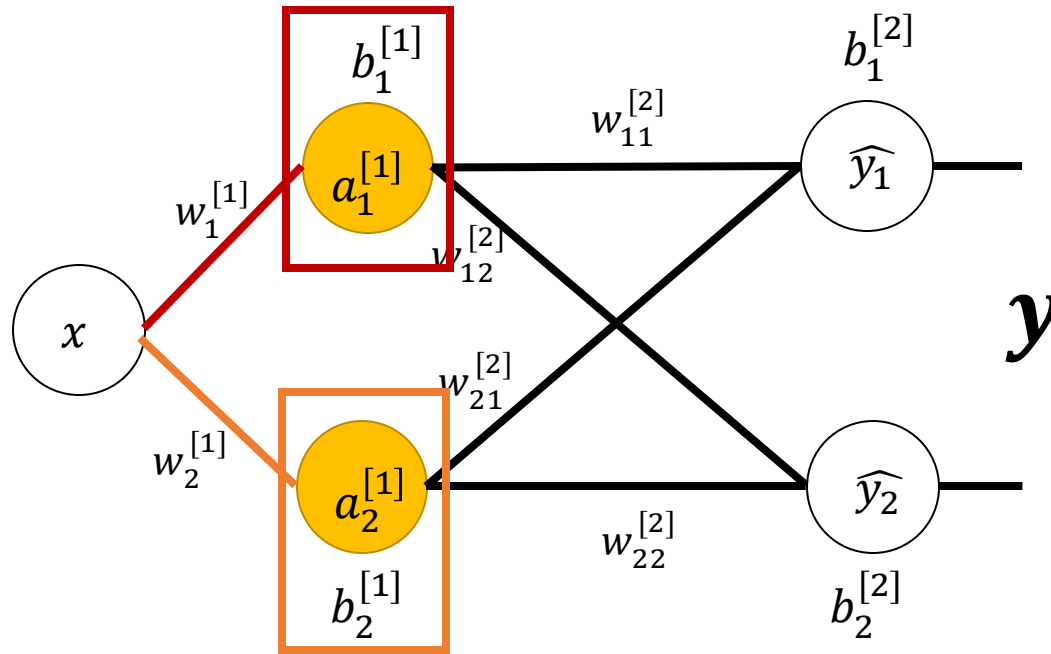
- Validation == Test Set? Not really strictly speaking
- Validation set: Dataset you validate your model on during training for model selection
- Test set: Dataset you test your model on after everything is fixed for unbiased-generalization error, **test set should not affect training procedure, including model selection**

Example



- Yellow --- with activation function
- Activation --- Sigmoid --- $\sigma(z) = a$
- Loss --- L_2 squared loss, also known as the squared L-2 norm
- Squared $L_2(\hat{\mathbf{y}}, \mathbf{y}) = ||(\hat{\mathbf{y}} - \mathbf{y})||_2^2$

Forward Pass 1

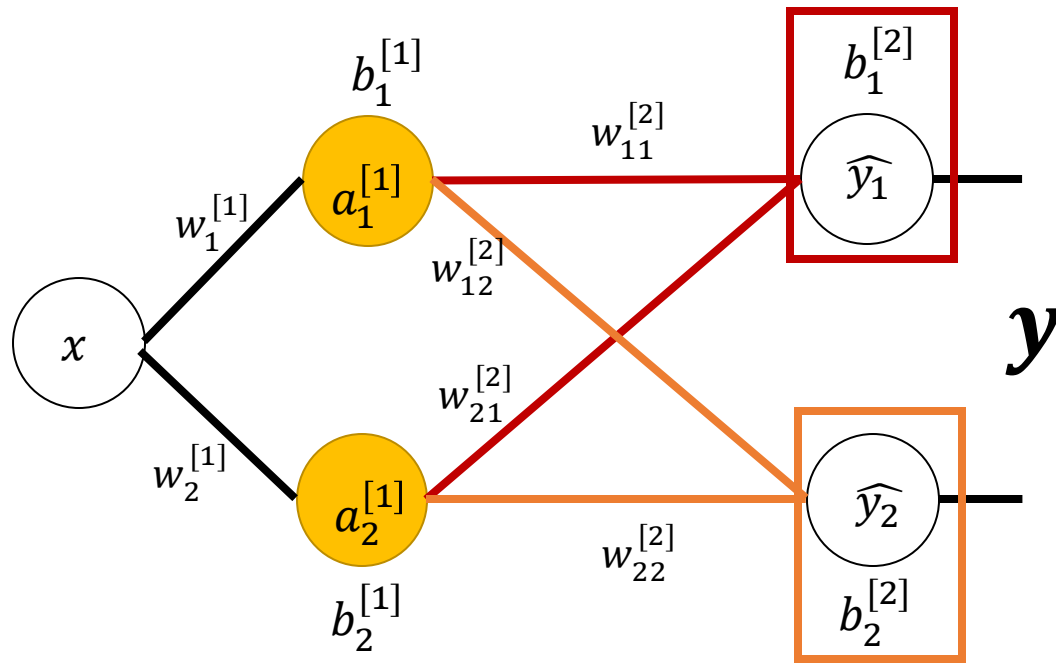


- Yellow --- with activation function
- Activation --- Sigmoid
- $\sigma(z_i) = (1 + e^{-z_i})^{-1} = a_i$
- Loss --- L_2 squared loss, also known as the squared L-2 norm
- Squared $L_2(\hat{y}, y) = ||(\hat{y} - y)||_2^2$

$$z_1^{[1]} = x w_1^{[1]} + b_1^{[1]}$$
$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = x w_2^{[1]} + b_2^{[1]}$$
$$a_2^{[1]} = \sigma(z_2^{[1]})$$

Forward Pass 2

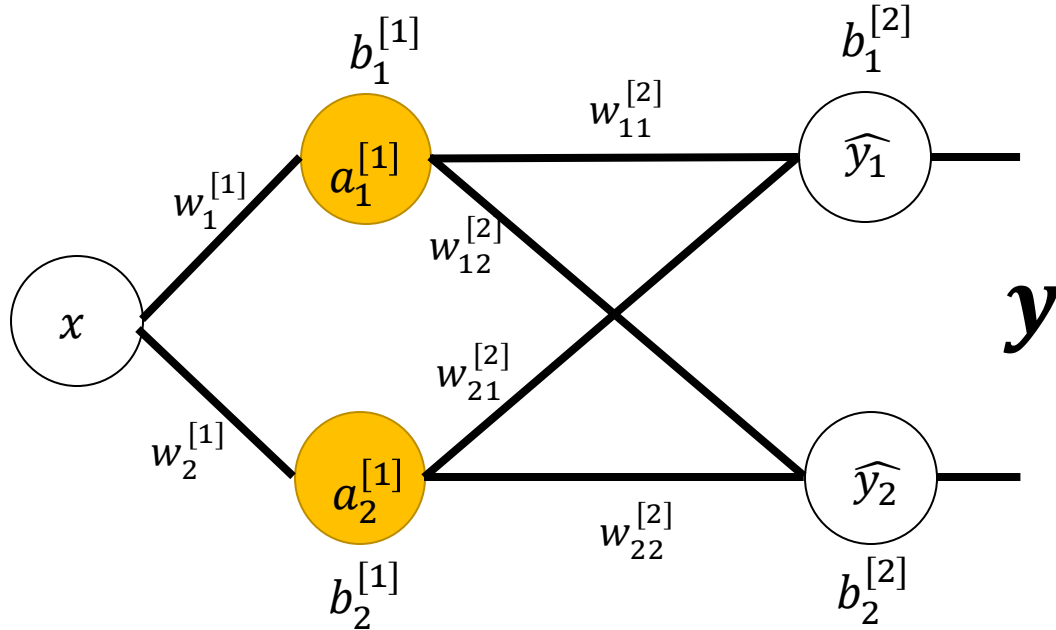


- Yellow --- with activation function
- Activation --- Sigmoid
- $\sigma(z_i) = (1 + e^{-z_i})^{-1} = a_i$
- Loss --- L_2 squared loss, also known as the squared L-2 norm
- Squared $L_2(\hat{\mathbf{y}}, \mathbf{y}) = ||(\hat{\mathbf{y}} - \mathbf{y})||_2^2$

$$\hat{y}_1 = w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + b_1^{[2]}$$

$$\hat{y}_2 = w_{12}^{[2]} a_1^{[1]} + w_{22}^{[2]} a_2^{[1]} + b_2^{[2]}$$

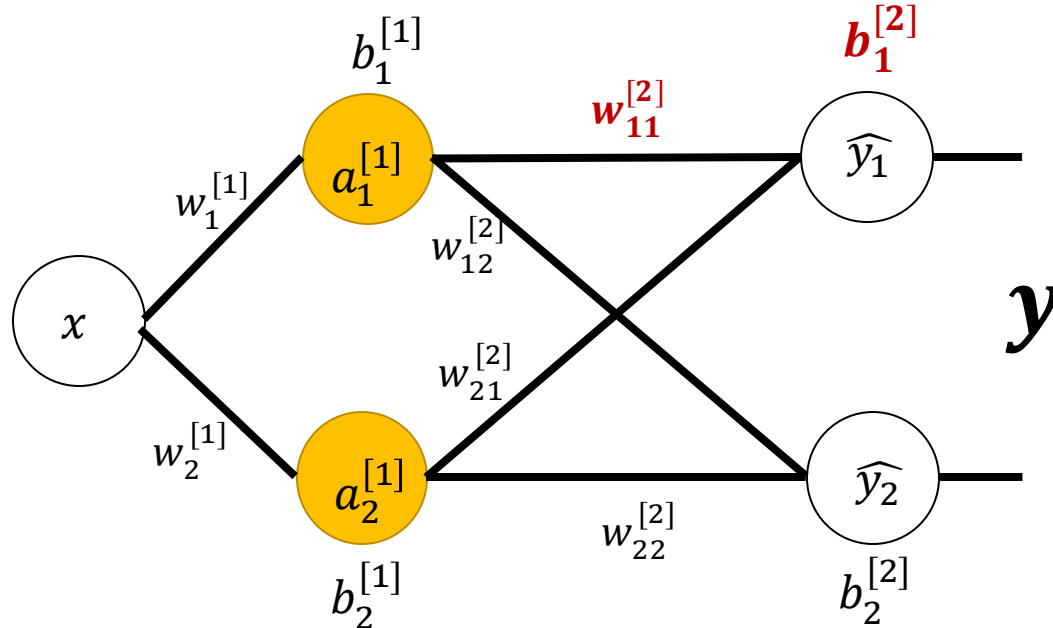
Calculate Loss



- Yellow --- with activation function
- Activation --- Sigmoid
- $\sigma(z_i) = (1 + e^{-z_i})^{-1} = a_i$
- Loss --- L_2 squared loss, also known as the squared L-2 norm
- Squared $L_2(\hat{\mathbf{y}}, \mathbf{y}) = ||(\hat{\mathbf{y}} - \mathbf{y})||_2^2$

$$L = ||(\hat{\mathbf{y}} - \mathbf{y})||_2^2 = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2$$

Back Propagation 1



- Yellow --- with activation function
- Activation --- Sigmoid
- $\sigma(z_i) = (1 + e^{-z_i})^{-1} = a_i$
- Loss --- L_2 squared loss, also known as the squared L-2 norm
- Squared $L_2(\hat{y}, y) = ||(\hat{y} - y)||_2^2$

$$L = ||(\hat{y} - y)||_2^2 = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2$$

Lets update $w_{11}^{[2]}$ and $b_1^{[2]}$ using **gradient descent**

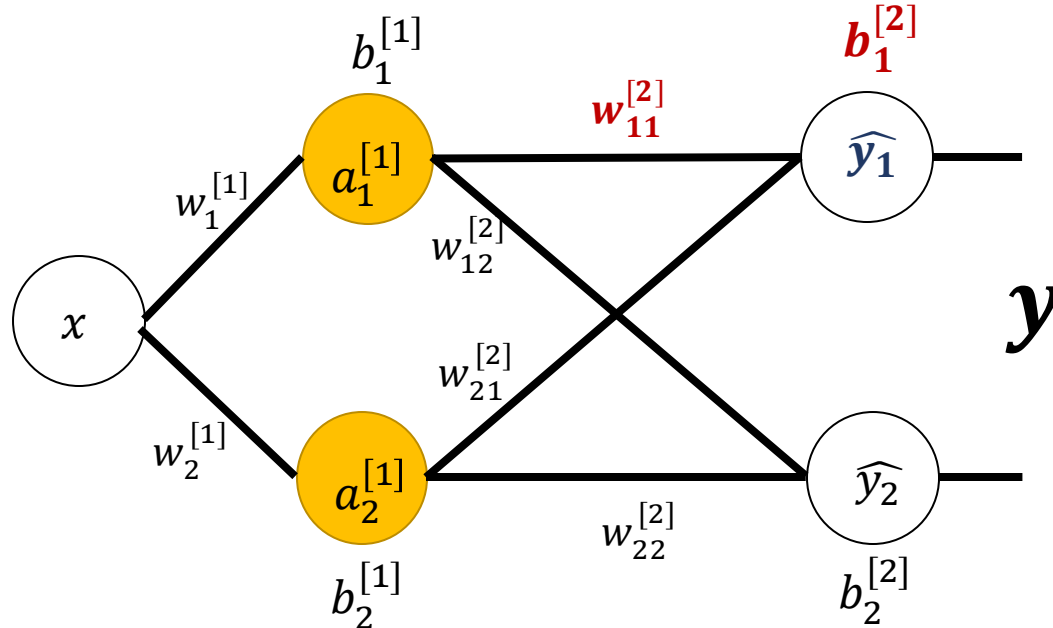
Compute the gradients:

$$\frac{\partial L}{\partial w_{11}^{[2]}}$$

and

$$\frac{\partial L}{\partial b_1^{[2]}}$$

Back Propagation 1



Compute the gradients with chain rule:

$$\frac{\partial L}{\partial w_{11}^{[2]}} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial w_{11}^{[2]}}$$

Similarly

$$\frac{\partial L}{\partial b_1^{[2]}} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial b_1^{[2]}}$$

- Yellow --- with activation function
- Activation --- Sigmoid
- $\sigma(z_i) = (1 + e^{-z_i})^{-1} = a_i$
- Loss --- L_2 squared loss, also known as the squared L-2 norm
- Squared $L_2(\hat{y}, y) = ||(\hat{y} - y)||_2^2$

$$L = ||(\hat{y} - y)||_2^2 = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2$$

$$\hat{y}_1 = w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + b_1^{[2]}$$

$$\hat{y}_2 = w_{12}^{[2]} a_1^{[1]} + w_{22}^{[2]} a_2^{[1]} + b_2^{[2]}$$

$$\frac{\partial L}{\partial \hat{y}_1} = 2(\hat{y}_1 - y_1)$$

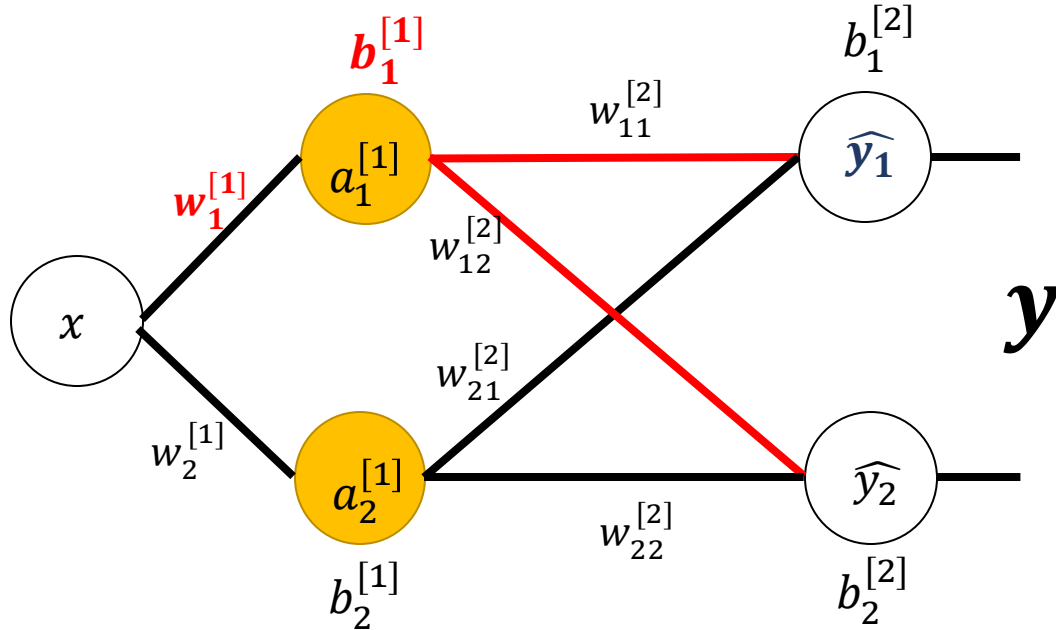
$$\frac{\partial \hat{y}_1}{\partial w_{11}^{[2]}} = a_1^{[1]}$$

(We get $a_1^{[1]}$ from the forward pass)

$$\frac{\partial L}{\partial w_{11}^{[2]}} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial w_{11}^{[2]}} = 2(\hat{y}_1 - y_1) a_1^{[1]}$$

$$\frac{\partial L}{\partial b_1^{[2]}} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial b_1^{[2]}} = 2(\hat{y}_1 - y_1)$$

Back Propagation 2



Compute the gradients with chain rule:

$$\frac{\partial L}{\partial w_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial w_1^{[1]}}$$

Similarly

$$\frac{\partial L}{\partial b_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial b_1^{[1]}}$$

$$z_1^{[1]} = x \mathbf{w}_1^{[1]} + \mathbf{b}_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = x w_2^{[1]} + b_2^{[1]}$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

$$\hat{y}_1 = w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + b_1^{[2]}$$

$$\hat{y}_2 = w_{12}^{[2]} a_1^{[1]} + w_{22}^{[2]} a_2^{[1]} + b_2^{[2]}$$

$$L = ||(\hat{\mathbf{y}} - \mathbf{y})||_2^2 = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2$$

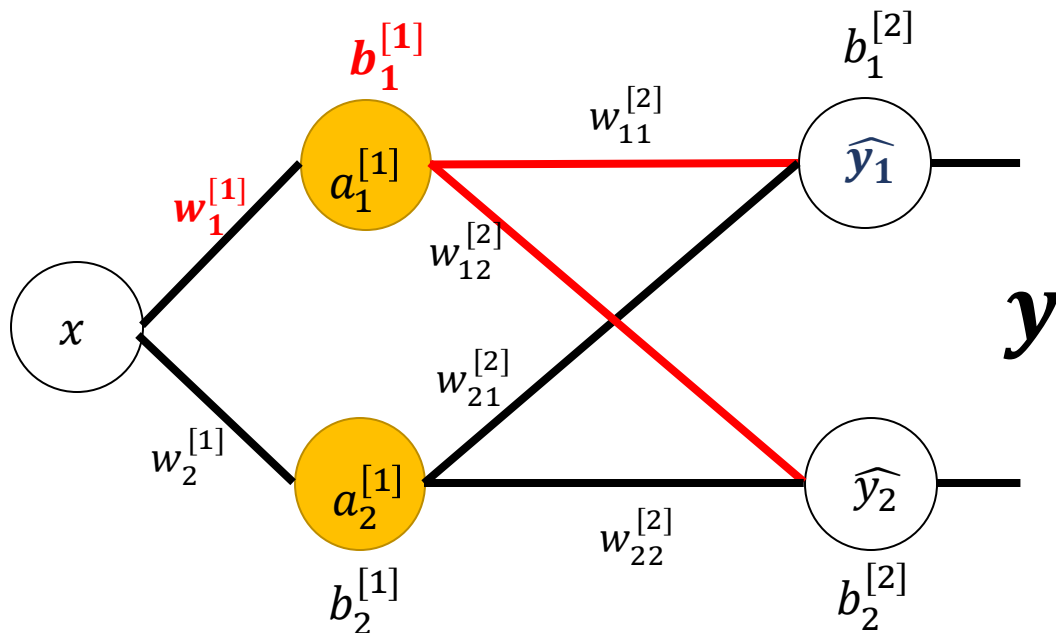
$$\frac{\partial L}{\partial \hat{\mathbf{y}}} = \left[\frac{\partial L}{\partial \hat{y}_1} = 2(\hat{y}_1 - y_1) \quad \frac{\partial L}{\partial \hat{y}_2} = 2(\hat{y}_2 - y_2) \right]$$

$$\frac{\partial \hat{\mathbf{y}}}{\partial a_1} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial a_1} = w_{11}^{[2]} \\ \frac{\partial \hat{y}_2}{\partial a_1} = w_{12}^{[2]} \end{bmatrix}$$

$$\frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} = \sigma(z_1^{[1]}) (1 - \sigma(z_1^{[1]}))$$

$$\frac{\partial z_1^{[1]}}{\partial w_1^{[1]}} = x \quad \text{and} \quad \frac{\partial z_1^{[1]}}{\partial b_1^{[1]}} = 1$$

Back Propagation 2



Compute the gradients with chain rule:

$$\frac{\partial L}{\partial w_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial w_1^{[1]}}$$

Similarly

$$\frac{\partial L}{\partial b_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial b_1^{[1]}}$$

$$z_1^{[1]} = x \mathbf{w}_1^{[1]} + \mathbf{b}_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = x w_2^{[1]} + b_2^{[1]}$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

$$\hat{y}_1 = w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + b_1^{[2]}$$

$$\hat{y}_2 = w_{12}^{[2]} a_1^{[1]} + w_{22}^{[2]} a_2^{[1]} + b_2^{[2]}$$

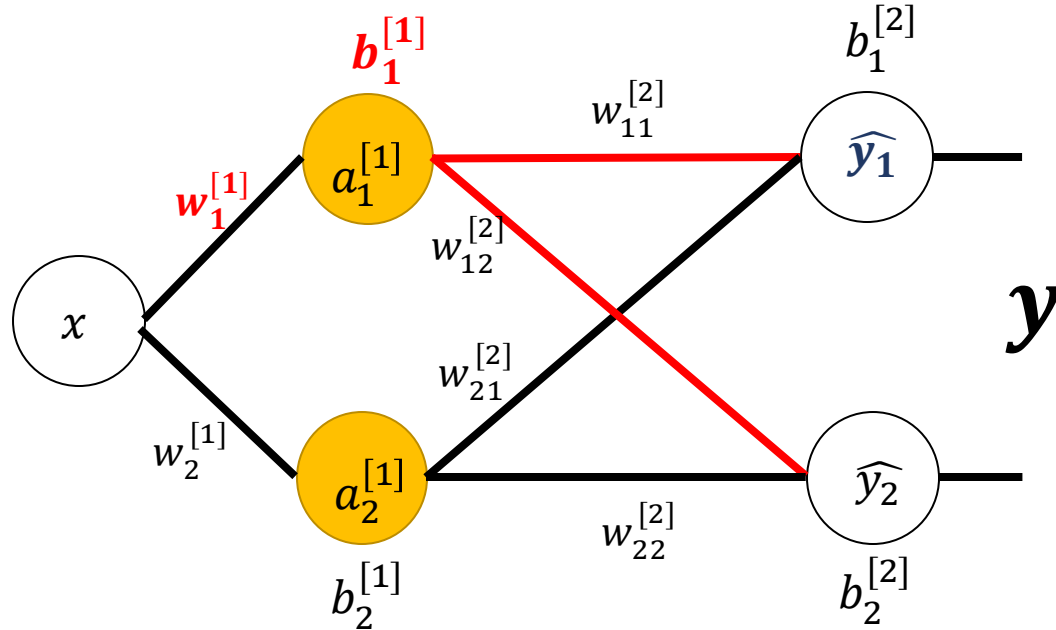
$$L = ||(\hat{\mathbf{y}} - \mathbf{y})||_2^2 = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2$$

$$\frac{\partial L}{\partial w_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial w_1^{[1]}}$$

$$= [2(\hat{y}_1 - y_1) \quad 2(\hat{y}_2 - y_2)] \begin{bmatrix} w_{11}^{[2]} \\ w_{12}^{[2]} \end{bmatrix} \sigma(z_1^{[1]}) (1 - \sigma(z_1^{[1]})) x$$

$$= (2(\hat{y}_1 - y_1) w_{11}^{[2]} + 2(\hat{y}_2 - y_2) w_{12}^{[2]}) \sigma(z_1^{[1]}) (1 - \sigma(z_1^{[1]})) x$$

Back Propagation 2



Compute the gradients with chain rule:

$$\frac{\partial L}{\partial w_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial w_1^{[1]}}$$

Similarly

$$\frac{\partial L}{\partial b_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial b_1^{[1]}}$$

$$z_1^{[1]} = x \mathbf{w}_1^{[1]} + \mathbf{b}_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = x w_2^{[1]} + b_2^{[1]}$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

$$\hat{y}_1 = w_{11}^{[2]} a_1^{[1]} + w_{21}^{[2]} a_2^{[1]} + b_1^{[2]}$$

$$\hat{y}_2 = w_{12}^{[2]} a_1^{[1]} + w_{22}^{[2]} a_2^{[1]} + b_2^{[2]}$$

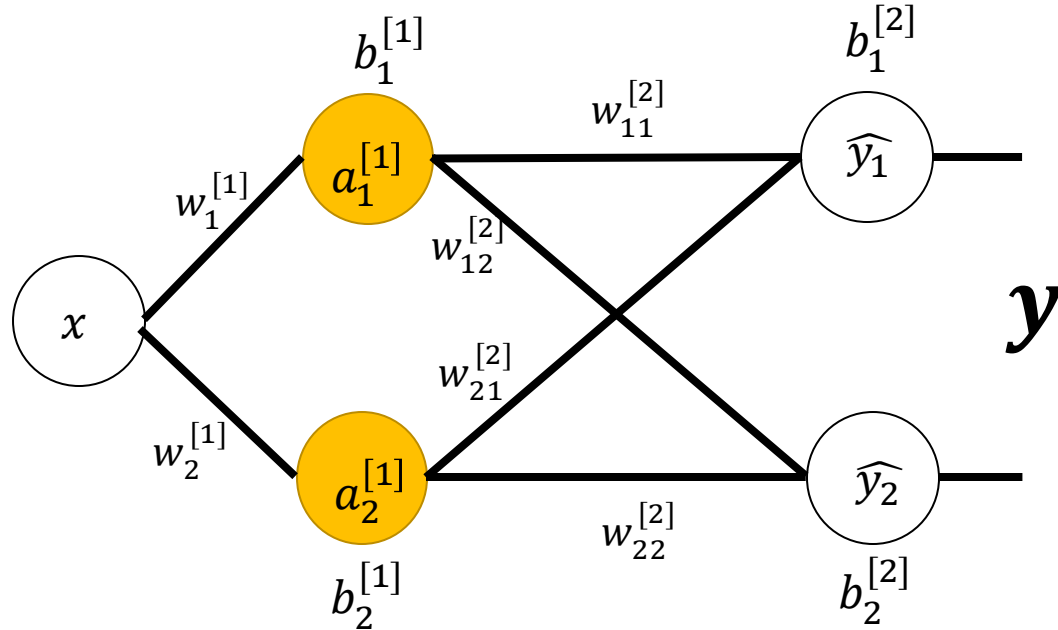
$$L = ||(\hat{\mathbf{y}} - \mathbf{y})||_2^2 = (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2$$

$$\frac{\partial L}{\partial b_1^{[1]}} = \frac{\partial L}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial a_1^{[1]}} \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial b_1^{[1]}}$$

$$= [2(\hat{y}_1 - y_1) \quad 2(\hat{y}_2 - y_2)] \begin{bmatrix} w_{11}^{[2]} \\ w_{12}^{[2]} \end{bmatrix} \sigma(z_1^{[1]}) (1 - \sigma(z_1^{[1]})) 1$$

$$= (2(\hat{y}_1 - y_1) w_{11}^{[2]} + 2(\hat{y}_2 - y_2) w_{12}^{[2]}) \sigma(z_1^{[1]}) (1 - \sigma(z_1^{[1]}))$$

Update Rule Gradient Descent



$$w_{k+1} = w_k - \eta \frac{\partial L}{\partial w}$$

Walking in the opposite direction of the gradient