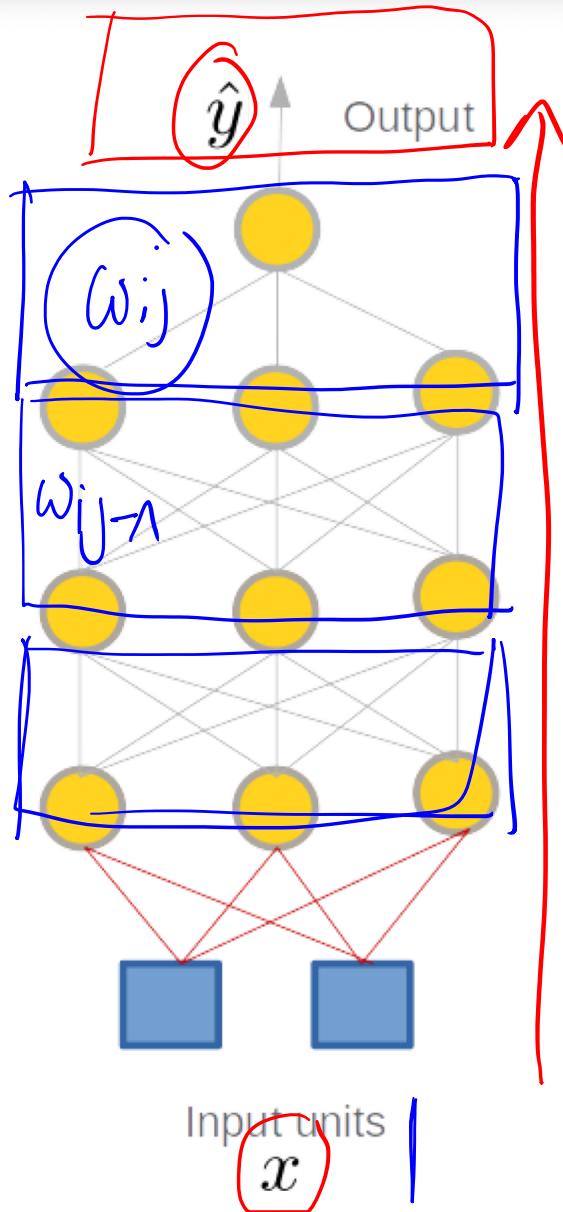

Back Propagation in Depth

Heni Ben Amor, Ph.D.
Assistant Professor
Arizona State University

Intuition of Back Propagation



1. Calculate Error at Final Layer

$$E = \frac{1}{2} (y - \hat{y})^2$$

2. Go back through layers and calculate changes to all weights

3. Repeat step 2 for all previous layers

BP = Gradient Descent

| Calculate gradient of network

| Update weights according to gradient descent

Gradient

∇E
nabla

$$\nabla E = \left(\frac{\partial E}{w_{1,1}}, \frac{\partial E}{w_{1,2}}, \dots, \frac{\partial E}{w_{l,k}} \right)$$

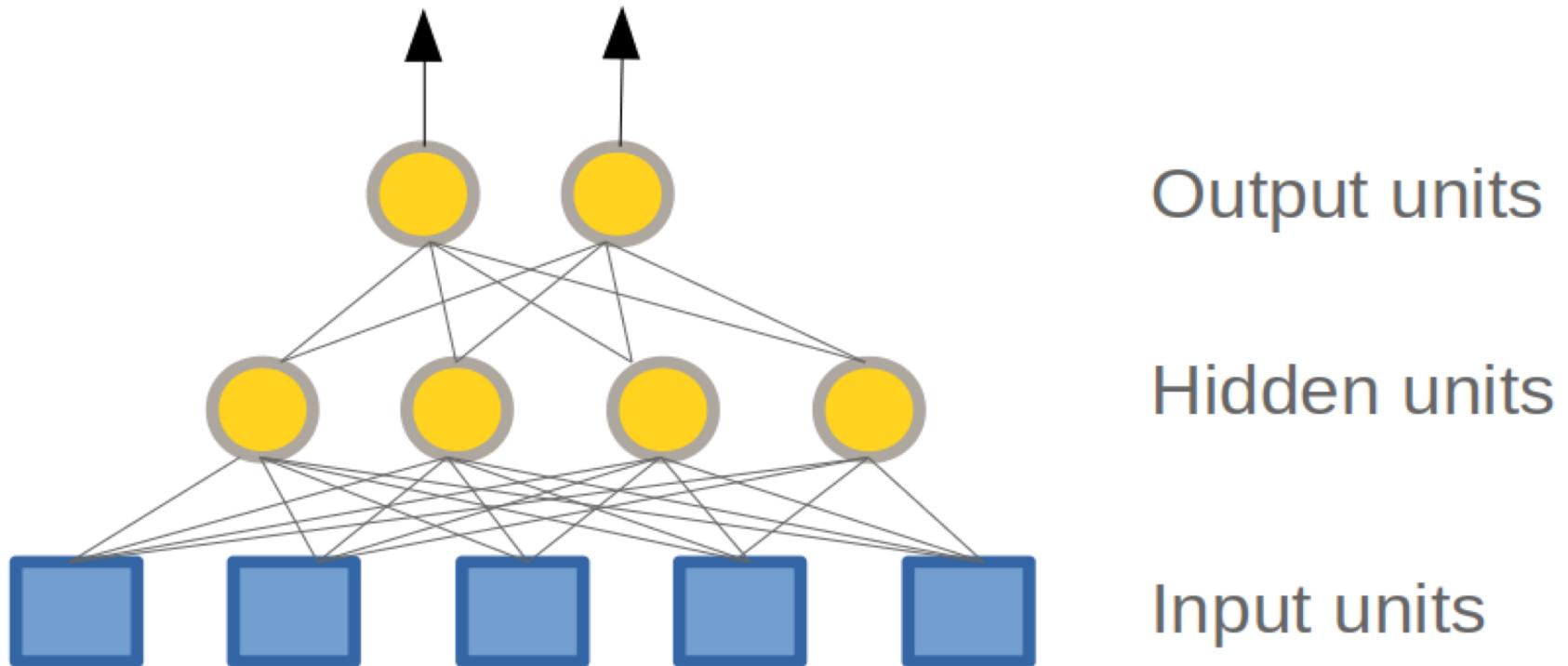
| Update equation

$$w_{i,j} \leftarrow w_{i,j} - \alpha \frac{\partial E}{w_{i,j}}$$

Learning Rate

Gradient

What is the Gradient of a Network?

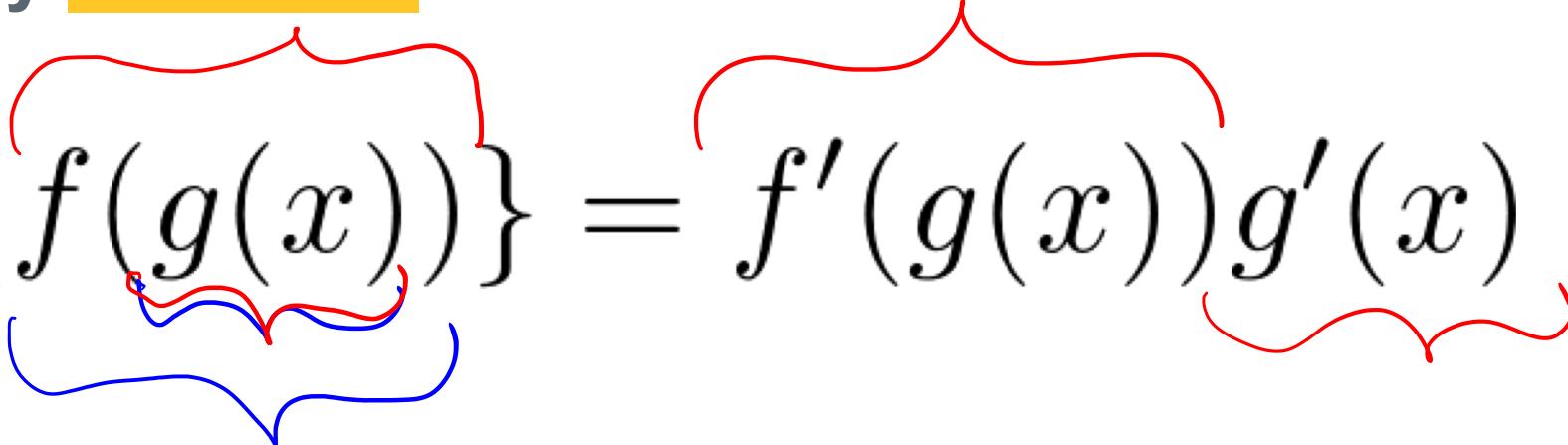


Gradient Calculation

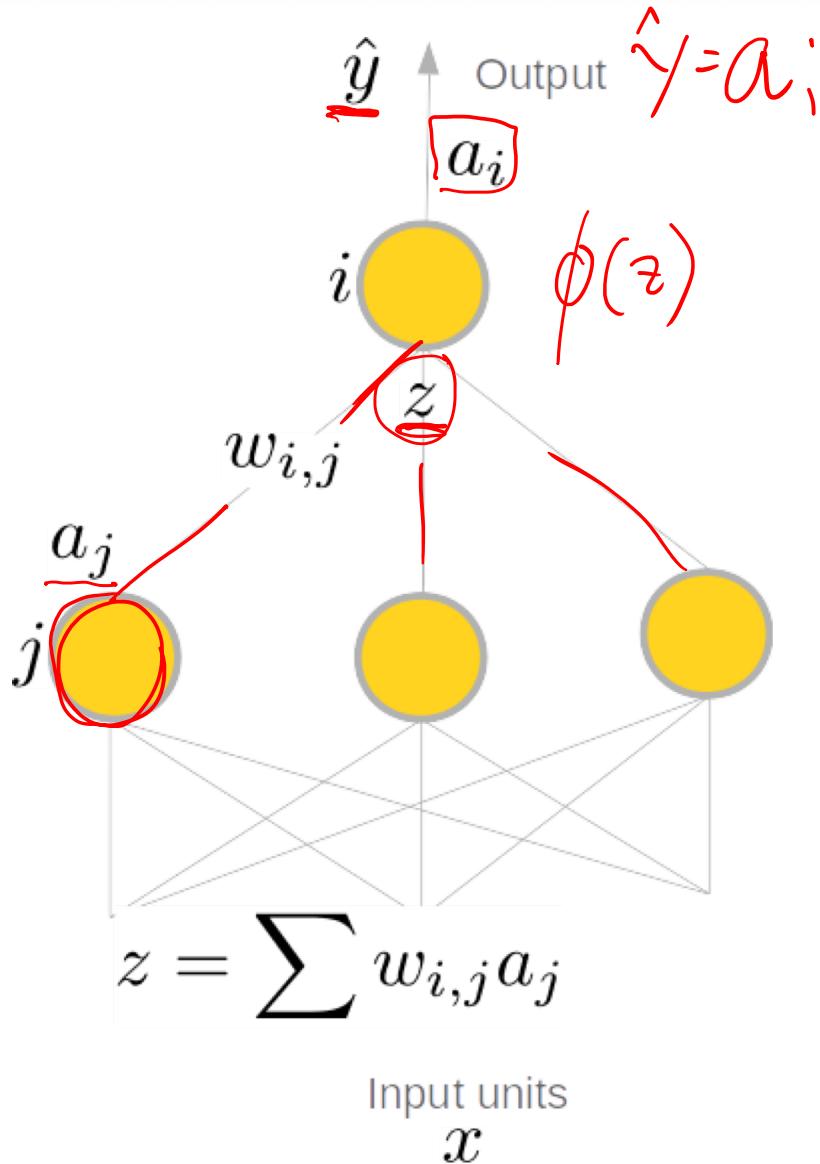
| How to get partial derivatives:

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots, \frac{\partial E}{\partial w_l} \right)$$

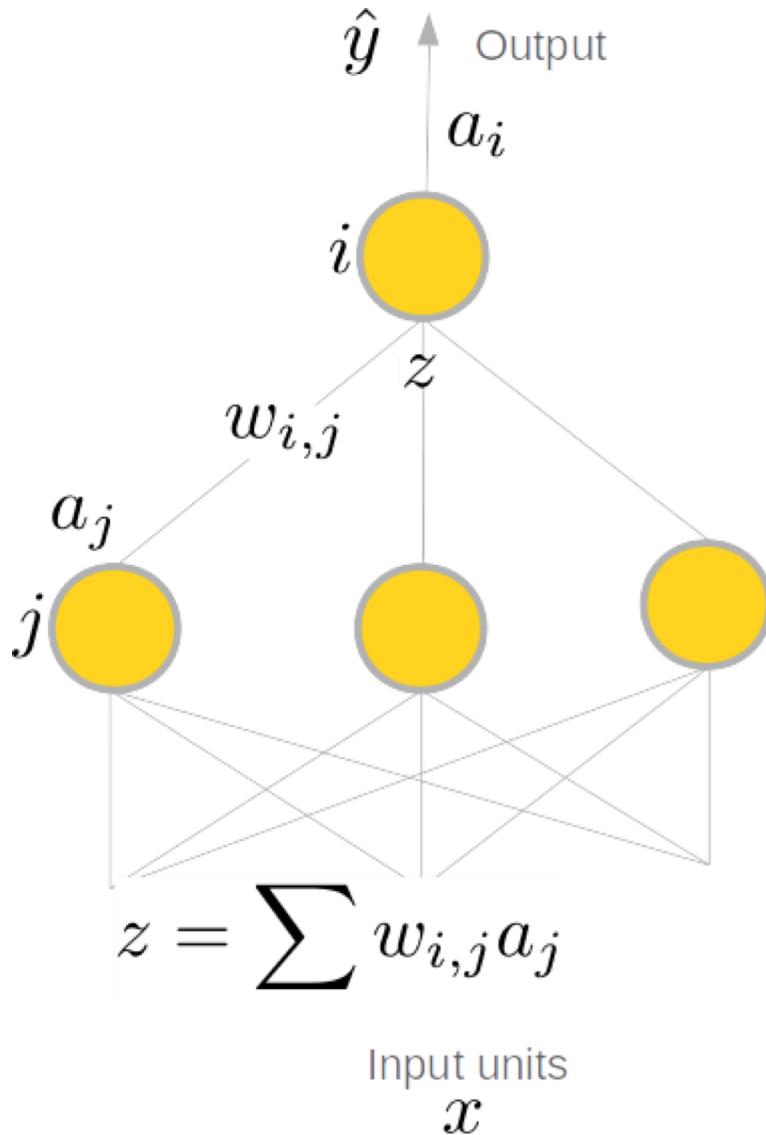
| Apply **chain rule**:

$$D\{f(g(x))\} = f'(g(x))g'(x)$$


In Depth Look at Back Propagation



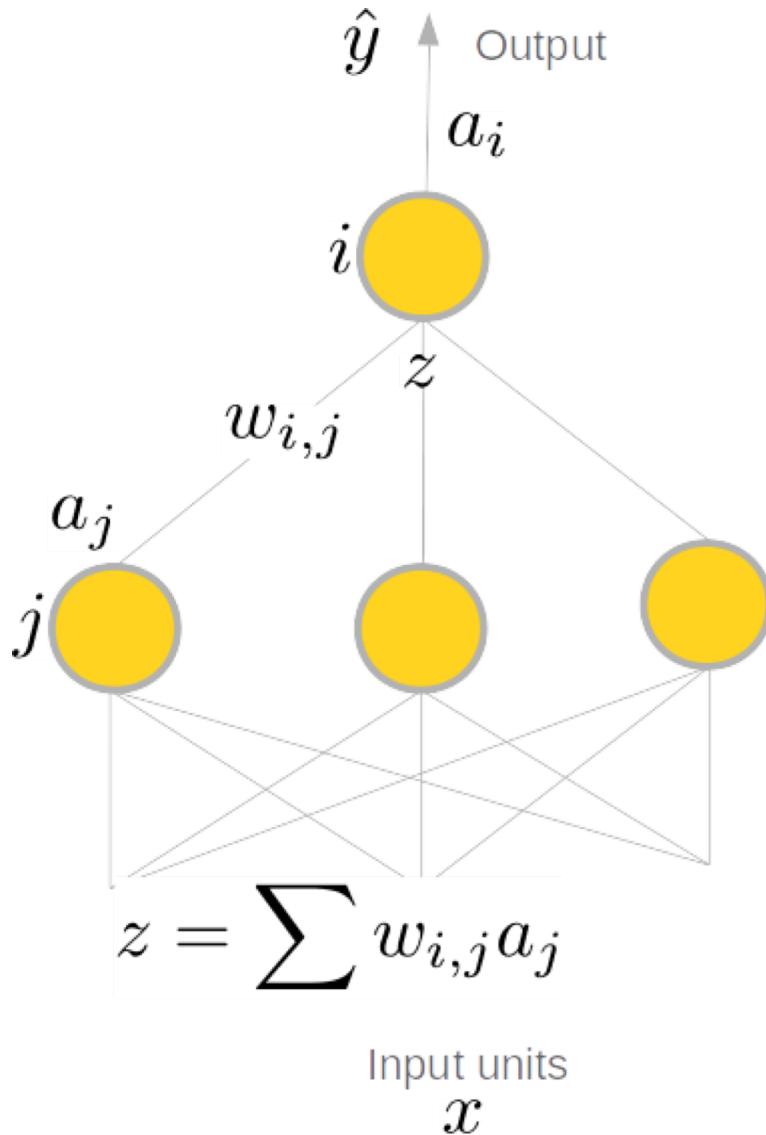
In Depth Look at Back Propagation



$$E = \frac{1}{2} (\underline{y - \hat{y}})^2$$

discrepancy

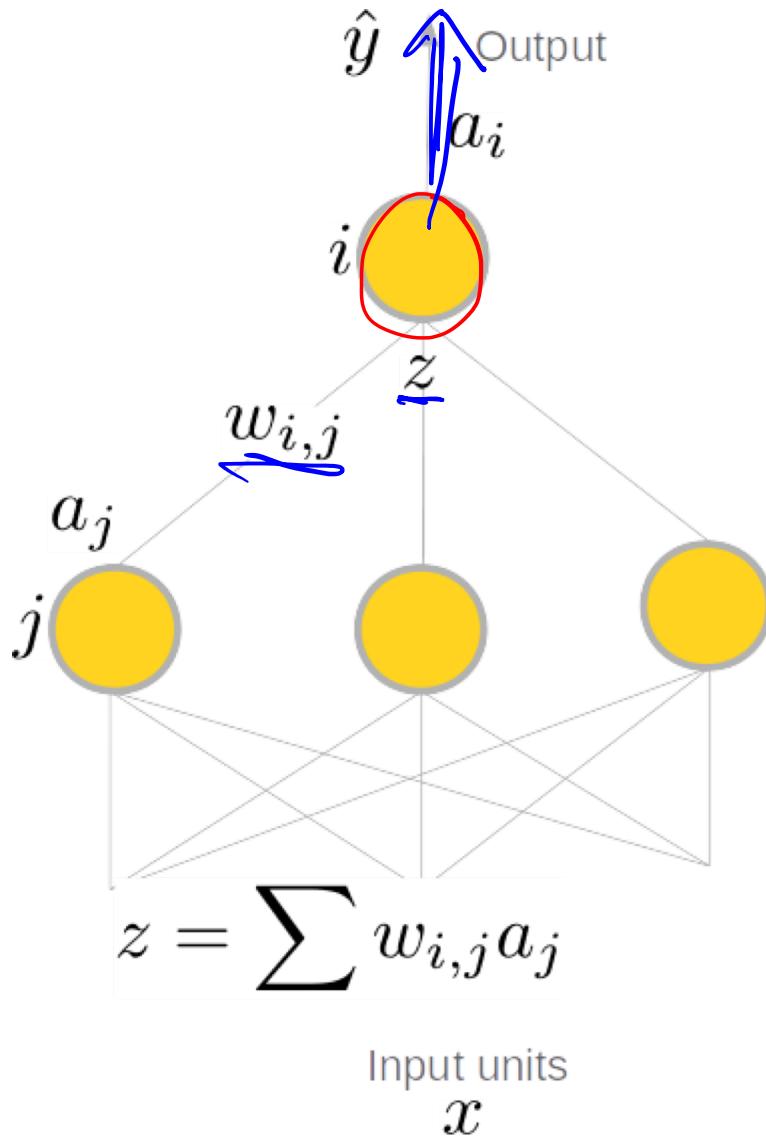
In Depth Look at Back Propagation



$$E = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{\partial E}{\partial \underline{w_{i,j}}}$$

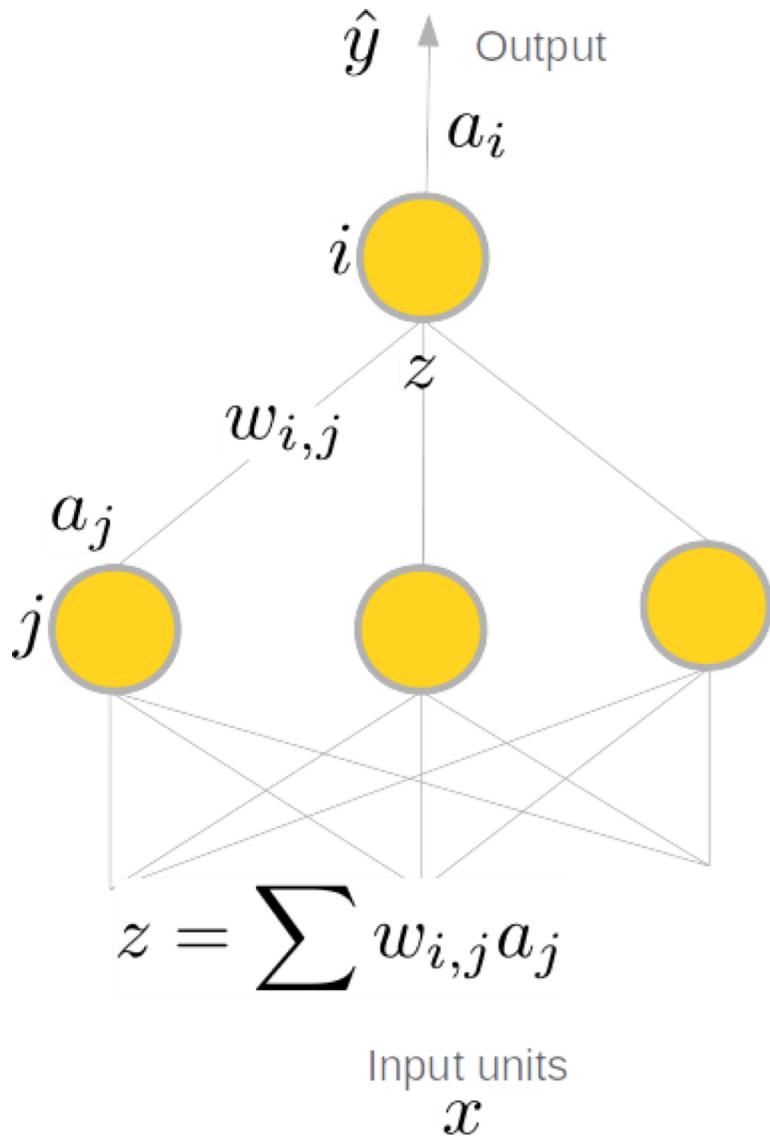
In Depth Look at Back Propagation



$$E = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{\partial E}{\partial w_{i,j}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

In Depth Look at Back Propagation

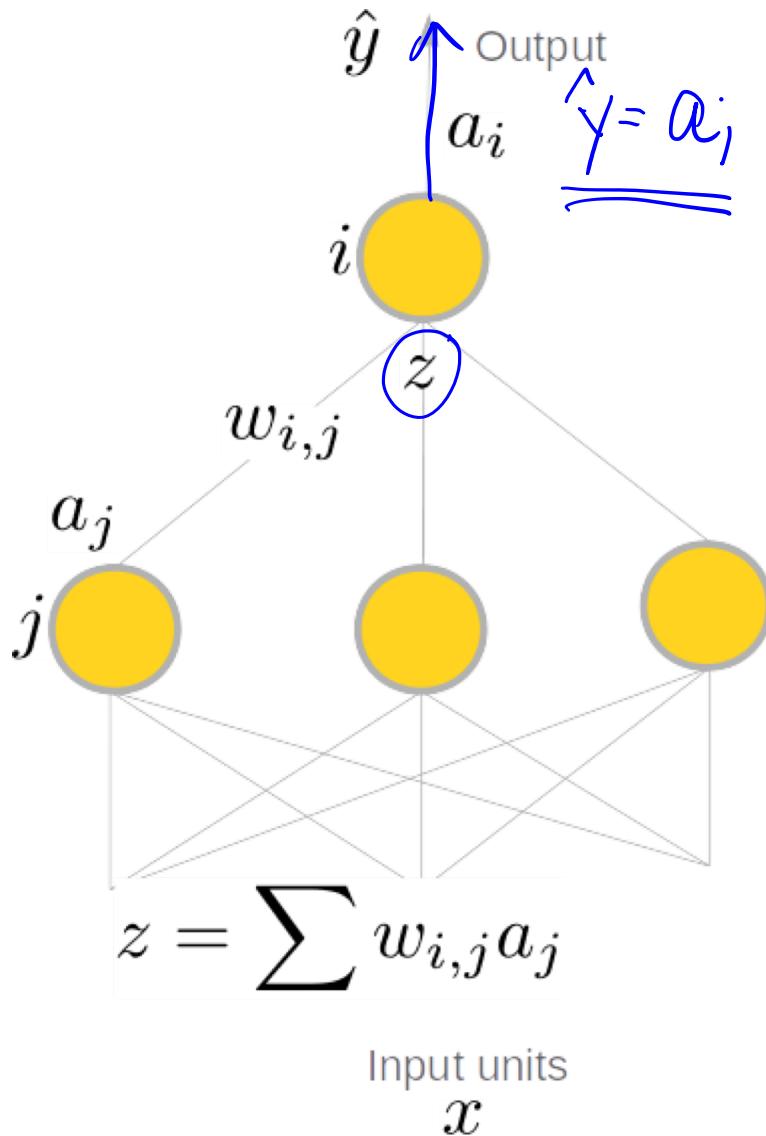


$$E = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{\partial E}{\partial w_{i,j}} = \boxed{\frac{\partial E}{\partial \hat{y}}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

$$\frac{\partial E}{\partial w_{i,j}} = -(y - \hat{y}) \boxed{\frac{\partial \hat{y}}{\partial z}} \frac{\partial z}{\partial w_{i,j}}$$

In Depth Look at Back Propagation



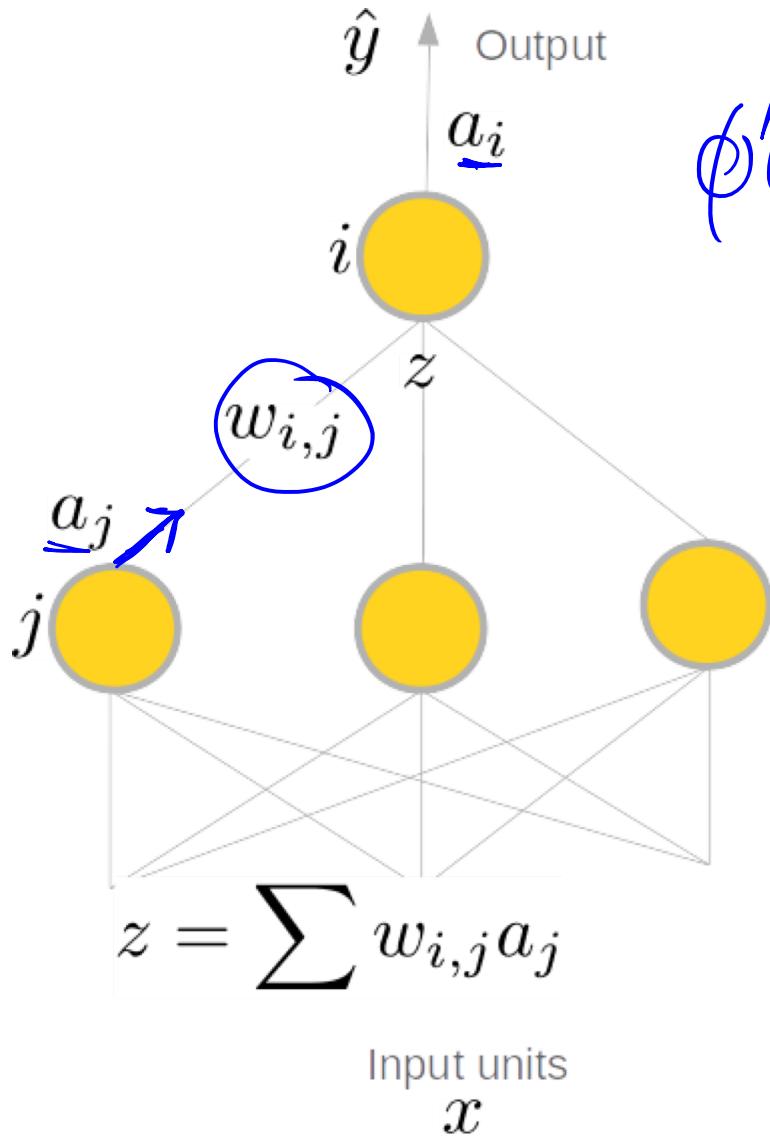
$$E = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{\partial E}{\partial w_{i,j}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

$$\frac{\partial E}{\partial w_{i,j}} = -(y - \hat{y}) \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

$$\frac{\partial E}{\partial w_{i,j}} = -(y - a_i) \frac{\partial \phi(z)}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

In Depth Look at Back Propagation



Output error
 $\phi'(z)$

$$E = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{\partial E}{\partial w_{i,j}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

$$\frac{\partial E}{\partial w_{i,j}} = -(y - \hat{y}) \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

$$\frac{\partial E}{\partial w_{i,j}} = -(y - \hat{y}) \frac{\partial \phi(z)}{\partial z} \frac{\partial z}{\partial w_{i,j}}$$

$$\frac{\partial E}{\partial w_{i,j}} = -(y - \hat{y}) \phi'(z) a_j$$

Update Equation for Output Layer


$$\frac{\partial E}{\partial w_{i,j}} = -(y - a_i) \phi'(z) a_j$$

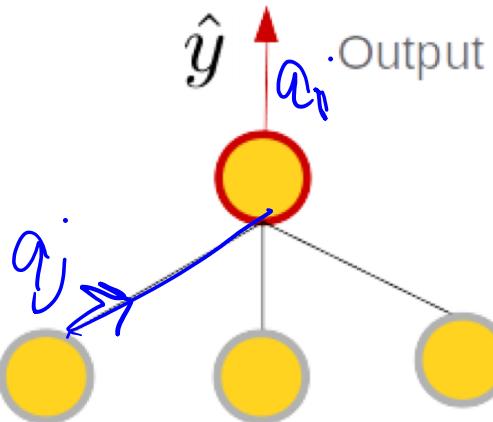
The diagram illustrates the update equation for the output layer. It shows a connection between two neurons, *i* and *j*, with a weight $w_{i,j}$. The error term $(y - a_i)$ is highlighted in blue, the derivative of the activation function $\phi'(z)$ is highlighted in orange, and the activation a_j from the previous layer is highlighted in green. Arrows point from each term to its corresponding label below the equation.

Difference to target

Derivative of Activation Function

Activation
previous layer

Output and Hidden Layer Updates in Back Propagation with Sigmoids

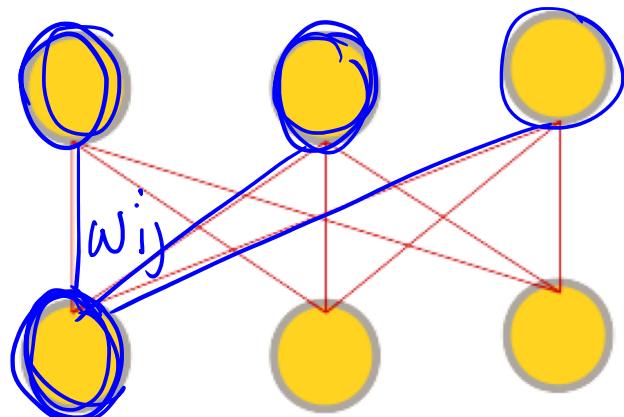


$$\frac{\partial E}{\partial w_{i,j}} = \underline{\delta_i} \underline{a_j}$$

w_{ij} leading into
output lay

$$\delta_i = \underline{a_i} (1 - \underline{a_i}) (\underline{a_i} - \underline{y})$$

Hidden Layer



$$\frac{\partial E}{\partial w_{i,j}} = \underline{\delta_i} \underline{a_j}$$

$$\delta_i = \underline{a_i} (1 - \underline{a_i})$$

$$\sum_{h \in \text{succ}(i)} \delta_h w_{i,h}$$

Back Propagation Algorithm

| Initialize all weights randomly

| Until **convergence** do: $E(NN) < \text{threshold}$

- Input example and calculate network output
- For each output unit do:

$$\delta_i = a_i(1 - a_i)(a_i - y)$$

- For each hidden unit do:

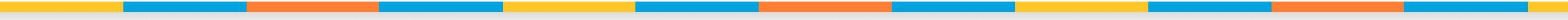
$$\delta_i = a_i(1 - a_i) \sum_{h \in \text{succ}(i)} \delta_h w_{i,h}$$

Update weights of connection between i and j:

$$w_{i,j} \leftarrow w_{i,j} - \alpha \delta_i a_j$$

← Previous activation
← learning rate α

Summary



- | Back propagation algorithm
- | Generates the weights of the network
- | Objective is to minimize the training error
- | Perform gradient descent
- | Calculate derivative of network
- | Apply **chain rule** to get partial derivatives
- | Derivatives are propagated backwards