

# Lecture 4

# Neural Network and

# Backpropagation algorithm

Lê Anh Cường

TDTU - 2020

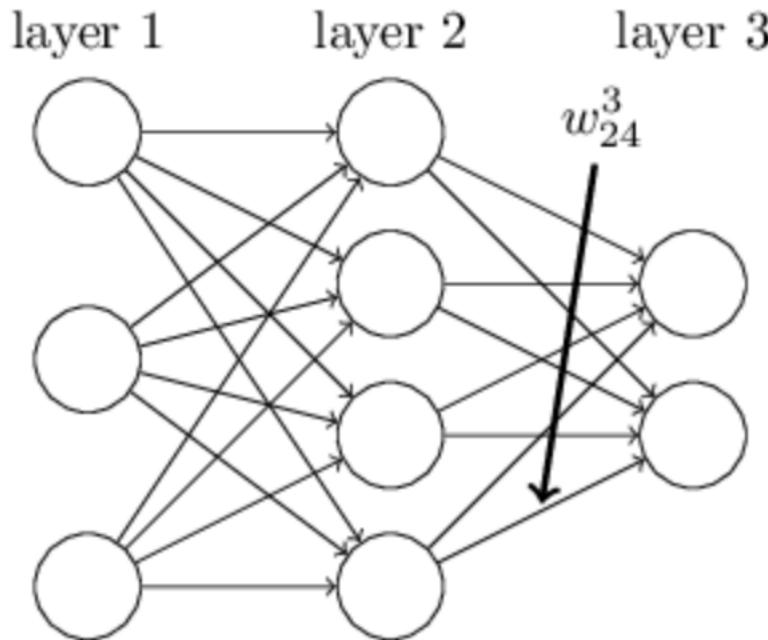
# References

- <http://neuralnetworksanddeeplearning.com/chap2.html>

# Objective

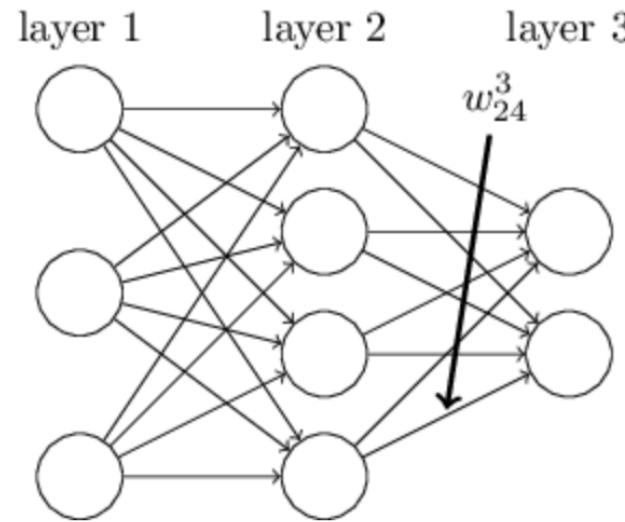
- Neural Networks architecture
- The Backpropagation algorithm

# Notations:



$w_{jk}^l$  is the weight from the  $k^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

# Formulae:



$w_{jk}^l$  is the weight from the  $k^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

$$a_j^l = \sigma \left( \sum_k w_{jk}^l a_k^{l-1} + b_j^l \right),$$

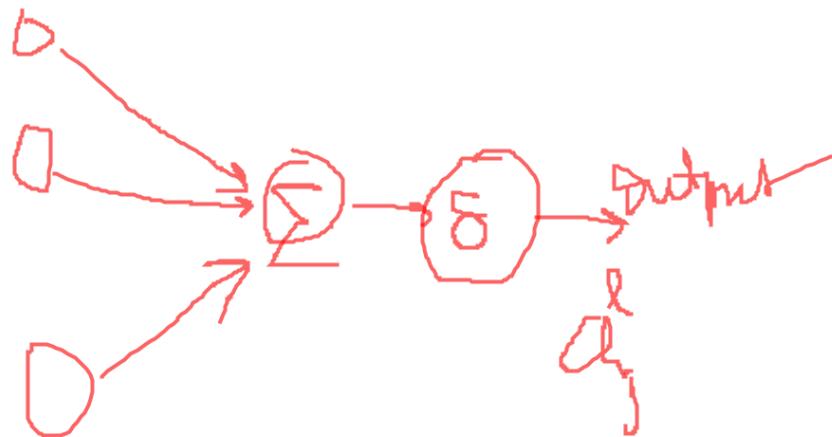
$$a^l = \sigma(w^l a^{l-1} + b^l).$$

$$z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$$

$$z^l \equiv w^l a^{l-1} + b^l$$

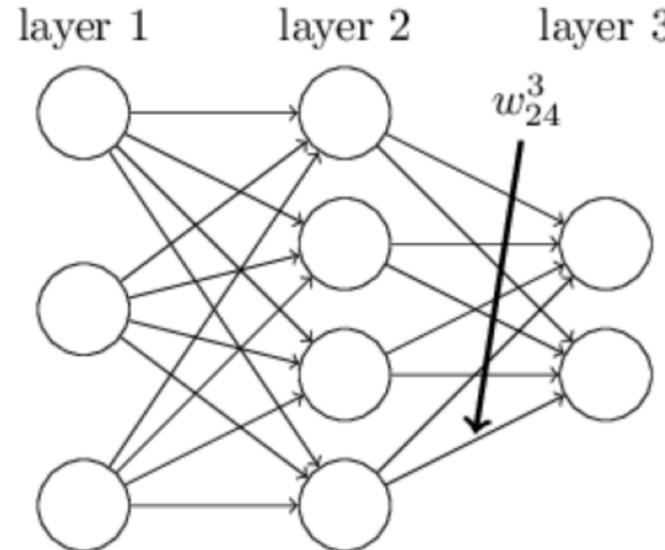
$$a^l = \sigma(z^l)$$

# Formulae:



$$a_j^l = \sigma \left( \underbrace{\sum_k w_{jk}^l a_k^{l-1}}_{\text{red underline}} + b_j^l \right),$$

$$a^l = \sigma(w^l a^{l-1} + b^l).$$



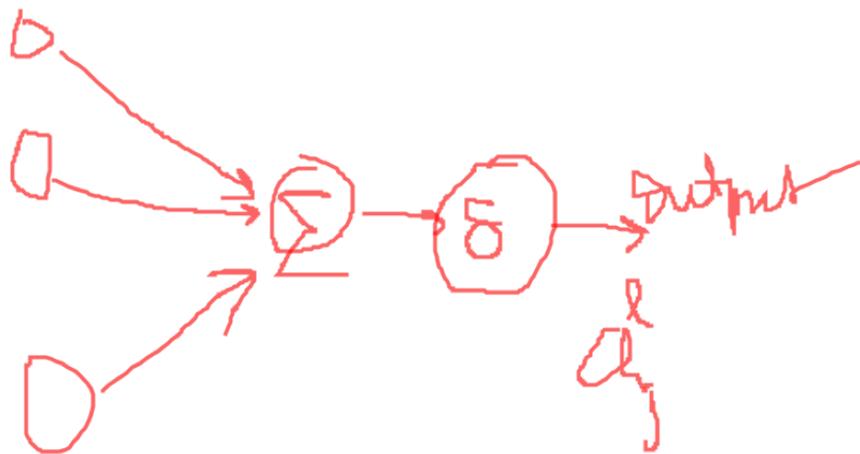
$w_{jk}^l$  is the weight from the  $k^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

$$z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$$

$$z^l \equiv w^l a^{l-1} + b^l$$

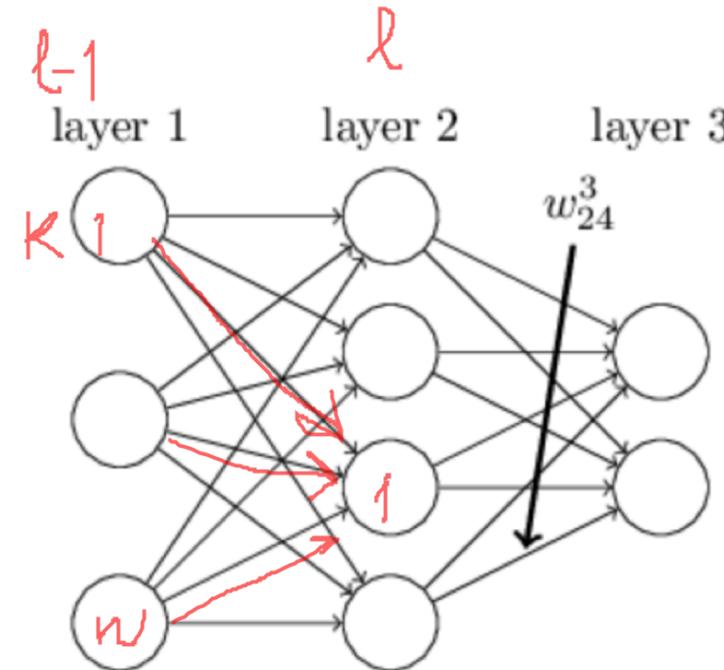
$$a^l = \sigma(z^l)$$

# Formulae:



$$a_j^l = \sigma \left( \underbrace{\sum_k w_{jk}^l a_k^{l-1}}_{\text{red underline}} + \underbrace{b_j^l}_{\text{red underline}} \right),$$

$$a^l = \sigma(w^l a^{l-1} + b^l).$$



$w_{jk}^l$  is the weight from the  $k^{\text{th}}$  neuron in the  $(l - 1)^{\text{th}}$  layer to the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer

$$z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$$

$$z^l \equiv w^l a^{l-1} + b^l$$

$$a^l = \sigma(z^l)$$

# Loss/Cost function

$$C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2,$$

$$C = \frac{1}{2} \|y - a^L\|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2,$$

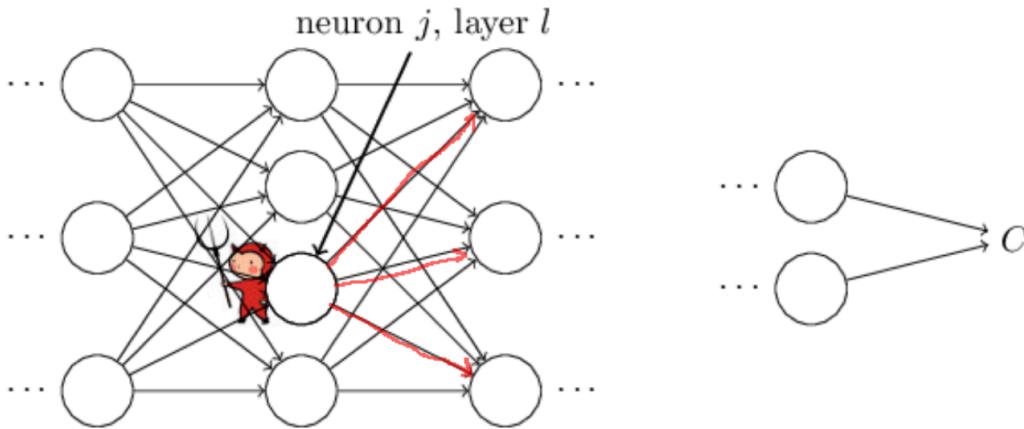
the cost for a single training example is  $C_x = \frac{1}{2} \|y - a^L\|^2$ .

$$C = \frac{1}{n} \sum_x C_x$$

$$y = y(x)$$

$$a^L = a^L(x)$$

# Error propagation

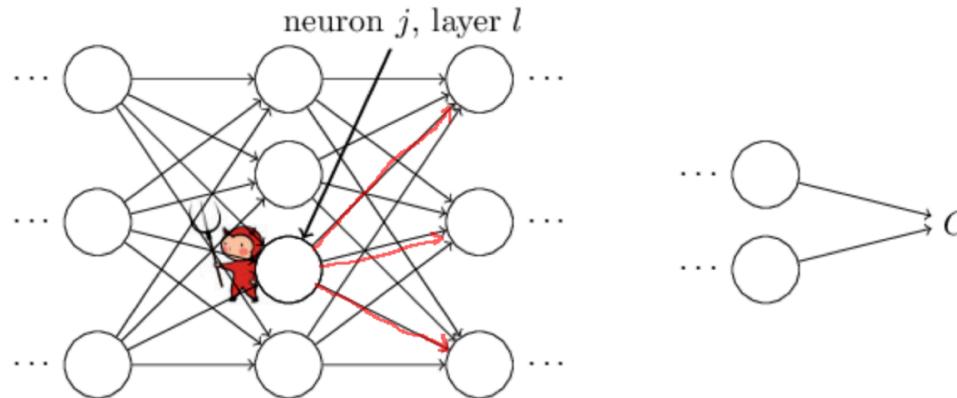


The demon sits at the  $j^{\text{th}}$  neuron in layer  $l$ . As the input to the neuron comes in, the demon messes with the neuron's operation. It adds a little change  $\Delta z_j^l$  to the neuron's weighted input, so that instead of outputting  $\sigma(z_j^l)$ , the neuron instead outputs  $\sigma(z_j^l + \Delta z_j^l)$ . This change propagates through later layers in the network, finally causing the overall cost to change by an amount  $\frac{\partial C}{\partial z_j^l} \Delta z_j^l$ .

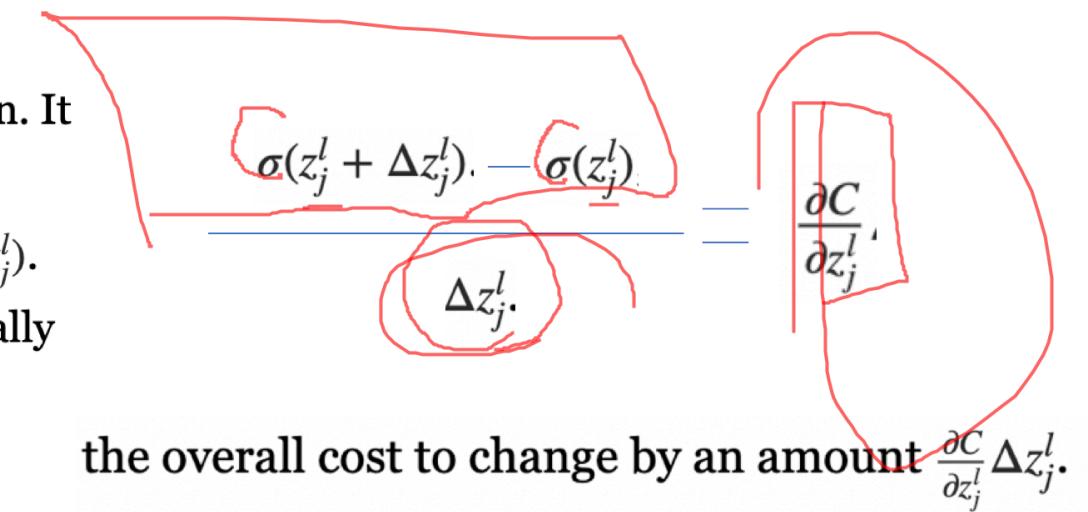
$$\frac{\sigma(z_j^l + \Delta z_j^l) - \sigma(z_j^l)}{\Delta z_j^l} = \frac{\partial C}{\partial z_j^l},$$

the overall cost to change by an amount  $\frac{\partial C}{\partial z_j^l} \Delta z_j^l$ .

# Error propagation



The demon sits at the  $j^{\text{th}}$  neuron in layer  $l$ . As the input to the neuron comes in, the demon messes with the neuron's operation. It adds a little change  $\Delta z_j^l$  to the neuron's weighted input, so that instead of outputting  $\sigma(z_j^l)$ , the neuron instead outputs  $\sigma(z_j^l + \Delta z_j^l)$ . This change propagates through later layers in the network, finally causing the overall cost to change by an amount  $\frac{\partial C}{\partial z_j^l} \Delta z_j^l$ .



partial derivatives  $\partial C/\partial w_{jk}^l$  and  $\partial C/\partial b_j^l$ .

we define the error  $\delta_j^l$  of neuron  $j$  in layer  $l$

$$\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}.$$

partial derivatives  $\frac{\partial C}{\partial w_{jk}^l}$  and  $\frac{\partial C}{\partial b_j^l}$

we define the error  $\delta_j^l$  of neuron  $j$  in layer  $l$

$$\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}.$$

$$w_{jk}^l = w_{jk}^l - \mu^* - \frac{\partial C}{\partial w_{jk}^l}$$

partial derivatives  $\partial C/\partial w_{jk}^l$  and  $\partial C/\partial b_j^l$ .

Backpropagation is about understanding how changing the weights and biases in a network changes the cost function. Ultimately, this means computing the partial derivatives  $\partial C/\partial w_{jk}^l$  and  $\partial C/\partial b_j^l$ . But to compute those, we first introduce an intermediate quantity,  $\delta_j^l$ , which we call the *error* in the  $j^{\text{th}}$  neuron in the  $l^{\text{th}}$  layer.

Backpropagation will give us a procedure to compute the error  $\delta_j^l$ , and then will relate  $\delta_j^l$  to  $\partial C/\partial w_{jk}^l$  and  $\partial C/\partial b_j^l$ .

partial derivatives  $\partial C/\partial w_{jk}^l$  and  $\partial C/\partial b_j^l$ .

we define the error  $\delta_j^l$  of neuron  $j$  in layer  $l$

$$\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}.$$

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L).$$

$$a^l = \sigma(z^l)$$

$$a_j^l = \sigma(z_j^l)$$

$$C = \frac{1}{2} \sum_j (y_j - a_j^L)^2, \text{ and so } \partial C / \partial a_j^L = (a_j^L - y_j).$$

$$\delta^L = (a^L - y) \odot \sigma'(z^L).$$

partial derivatives  $\partial C / \partial w_{jk}^l$  and  $\partial C / \partial b_j^l$ .

we define the error  $\delta_j^l$  of neuron  $j$  in layer  $l$

$$\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}.$$

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L).$$

$$a^l = \sigma(z^l)$$

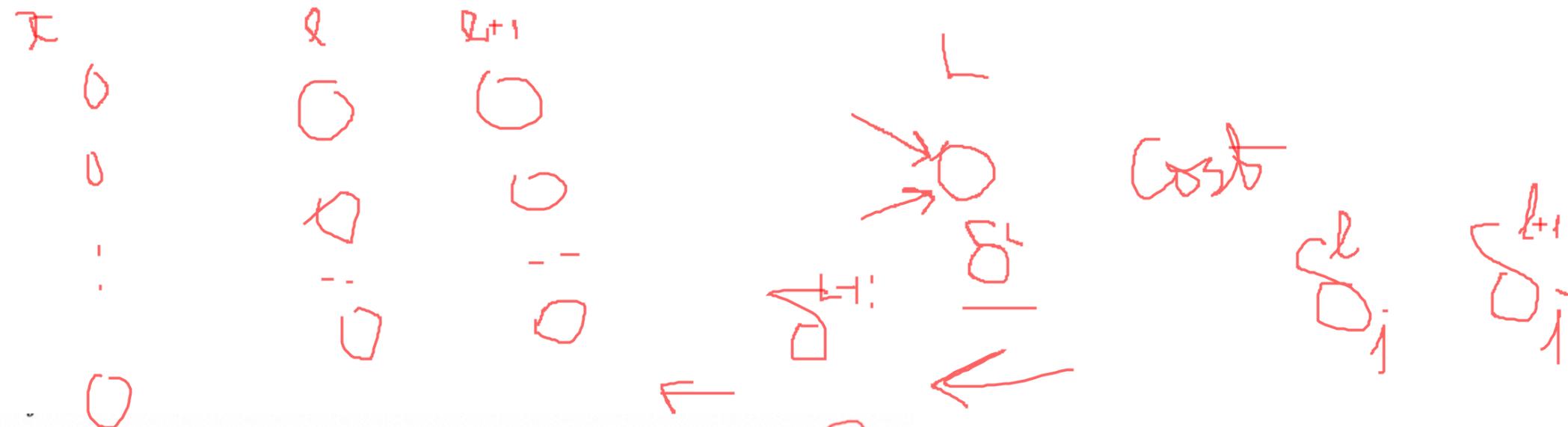
$$a^l = \sigma(z^l)$$

$$C = \frac{1}{2} \sum_j (y_j - a_j^L)^2, \text{ and so } \frac{\partial C}{\partial a_j^L} = (a_j^L - y_j).$$

$$\delta^L = (a^L - y) \odot \sigma'(z^L).$$

$$\frac{\partial A}{\partial C} \sum_B \frac{\partial A}{\partial B} \cdot \frac{\partial B}{\partial C}$$

$$\sigma = \underline{\sigma}(z)$$



$$\frac{\partial C}{\partial a_j^L} = \sum_B \frac{\partial C}{\partial a_B^L} \cdot \frac{\partial a_B^L}{\partial a_j^L}$$

$$\delta_j^L = \frac{\partial C}{\partial a_j^L} = \frac{\text{Sigmoid}}{1 - \text{Sigmoid}}$$

$$C = \frac{1}{2} \sum_j (y_j - a_j^L)^2, \text{ and so } \frac{\partial C}{\partial a_j^L} = (a_j^L - y_j)$$

$$\delta_j^L = (a_j^L - y_j) \odot \sigma'(z_j^L)$$

$$a = \sigma(z)$$

## An equation for the error $\delta^l$ in terms of the error in the next layer, $\delta^{l+1}$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

Suppose we know the error  $\delta^{l+1}$  at the  $l + 1^{\text{th}}$  layer. When we apply the transpose weight matrix,  $(w^{l+1})^T$ , we can think intuitively of this as moving the error *backward* through the network, giving us some sort of measure of the error at the output of the  $l^{\text{th}}$  layer. We then take the Hadamard product  $\odot \sigma'(z^l)$ . This moves the error backward through the activation function in layer  $l$ , giving us the error  $\delta^l$  in the weighted input to layer  $l$ .

$$z_k^{l+1} = \sum_j w_{kj}^{l+1} a_j^l + b_k^{l+1} = \sum_j w_{kj}^{l+1} \sigma(z_j^l) + b_k^{l+1}. \quad (43)$$

Differentiating, we obtain

$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \sigma'(z_j^l). \quad (44)$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

$$\begin{aligned} \delta_j^l &= \frac{\partial C}{\partial z_j^l} \\ &= \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} \\ &= \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1}, \end{aligned}$$

compute the error  $\delta^l$

$$\delta^L = (a^L - y) \odot \sigma'(z^L). \quad (\text{BP1})$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l), \quad (\text{BP2})$$

By combining (BP2) with (BP1) we can compute the error  $\delta^l$  for any layer in the network. We start by using (BP1) to compute  $\delta^L$ , then apply Equation (BP2) to compute  $\delta^{L-1}$ , then Equation (BP2) again to compute  $\delta^{L-2}$ , and so on, all the way back through the network.

**An equation for the rate of change of the cost with respect to any weight in the network:**

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

An equation for the rate of change of the cost with respect to any weight in the network:

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

$$w_{jk}^l = w_{jk}^l - \mu \cdot \frac{\partial C}{\partial w_{jk}^l}$$

## Summary: the equations of backpropagation

$$\delta^L = \nabla_a C \odot \sigma'(z^L) \quad (\text{BP1})$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (\text{BP2})$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad (\text{BP3})$$

$$\frac{\partial C}{\partial w_{j,k}^l} = a_k^{l-1} \delta_j^l \quad (\text{BP4})$$

$$\delta^L = \nabla_a C \odot \sigma'(z^L).$$

$$\delta^L = (a^L - y) \odot \sigma'(z^L)$$



$$a_j^l = \sigma\left(\sum_k w_{jk}^l a_k^{l-1} + b_j^l\right)$$

$$a^l=\sigma(w^la^{l-1}+b^l)$$

$$a^l=\sigma(z^l)$$

$$C = \frac{1}{2n}\sum_x \|y(x) - a^L(x)\|^2.$$

$$C = \tfrac{1}{n}\sum_x C_x$$

$$C_x = \tfrac{1}{2}\|y-a^L\|^2$$

$$a^L=a^L(x)$$

$$y=y(x)$$

$$z_j^l = \textstyle{\sum_k} w_{jk}^l a_k^{l-1} + b_j^l$$

$$z^l \equiv w^l a^{l-1} + b^l$$

$$C = \frac{1}{2}\|y-a^L\|^2 = \frac{1}{2}\sum_j(y_j-a_j^L)^2$$