

# How the backpropagation algorithm works

Srikumar Ramalingam

School of Computing

University of Utah

# Reference

Most of the slides are taken from the second chapter of the online book by Michael Nielson:

- [neuralnetworksanddeeplearning.com](http://neuralnetworksanddeeplearning.com)

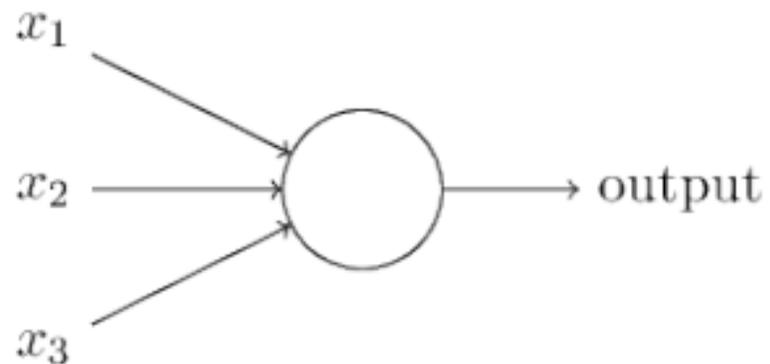
# Introduction

- First discovered in 1970.
- First influential paper in 1986:

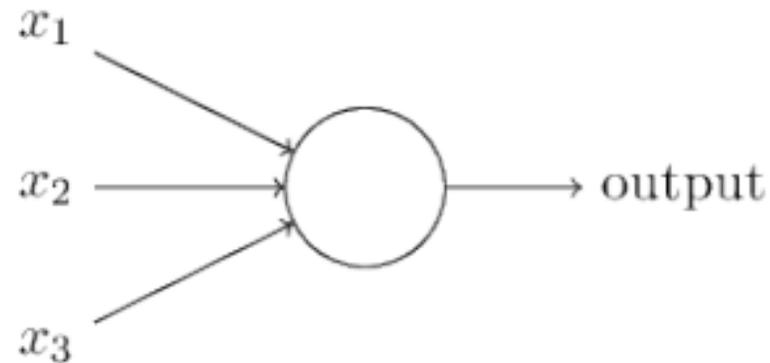
Rumelhart, Hinton and Williams, Learning representations by back-propagating errors, Nature, 1986.

# Perceptron (Reminder)

$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$



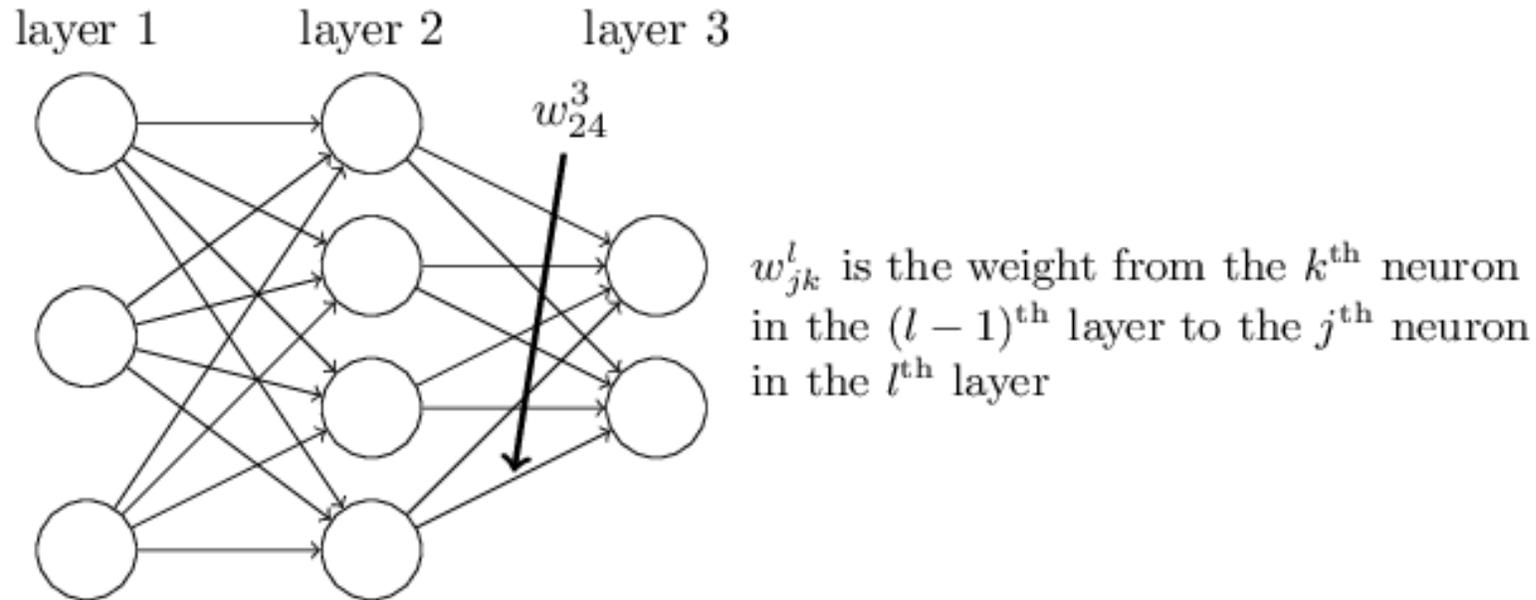
# Sigmoid neuron (Reminder)



- A sigmoid neuron can take real numbers  $(x_1, x_2, x_3)$  within 0 to 1 and returns a number within 0 to 1. The weights  $(w_1, w_2, w_3)$  and the bias term  $b$  are real numbers.

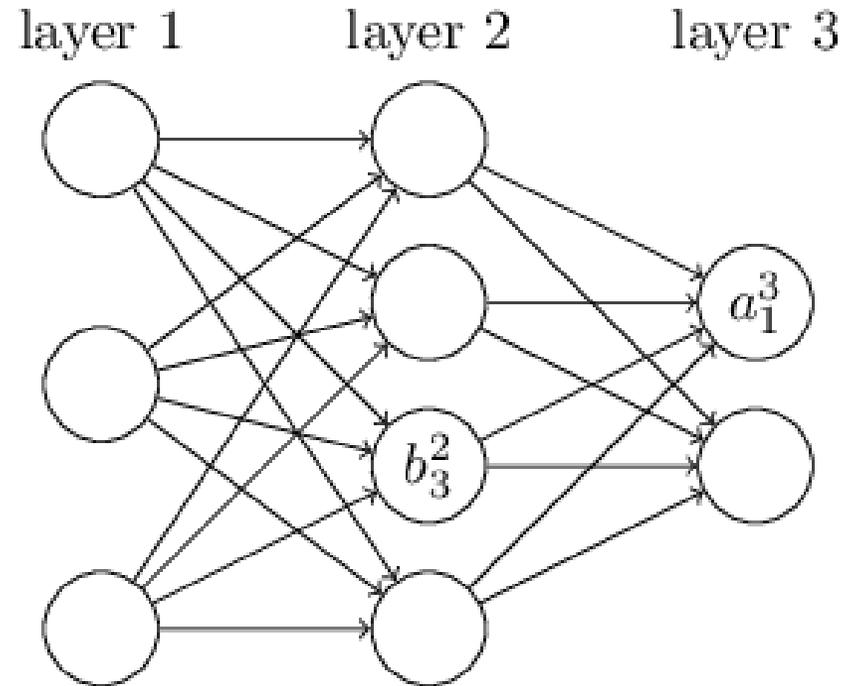
Sigmoid function  $\sigma(z) \equiv \frac{1}{1 + e^{-z}}$

# Matrix equations for neural networks



The indices “j” and “k” seem a little counter-intuitive!

# Layer to layer relationship



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

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

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

- $b_j^l$  is the bias term in the  $j$ th neuron in the  $l$ th layer.
- $a_j^l$  is the activation in the  $j$ th neuron in the  $l$ th layer.
- $z_j^l$  is the weighted input to the  $j$ th neuron in the  $l$ th layer.

# Cost function from the network

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

Groundtruth for each input

Output activation vector for a specific training sample  $x$ .

# of input samples

for each input sample

The diagram shows the cost function formula  $C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2$  with four blue arrows pointing to specific parts of the formula and their corresponding text labels. An arrow points from the text '# of input samples' to the denominator  $2n$ . Another arrow points from the text 'for each input sample' to the summation variable  $x$ . A third arrow points from the text 'Groundtruth for each input' to the term  $y(x)$ . A fourth arrow points from the text 'Output activation vector for a specific training sample  $x$ .' to the term  $a^L(x)$ .

# Backpropagation and stochastic gradient descent

- The goal of the backpropagation algorithm is to compute the gradients  $\frac{\partial C}{\partial w}$  and  $\frac{\partial C}{\partial b}$  of the cost function  $C$  with respect to each and every weight and bias parameters. Note that backpropagation is only used to compute the gradients.

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

- Stochastic gradient descent is the training algorithm.

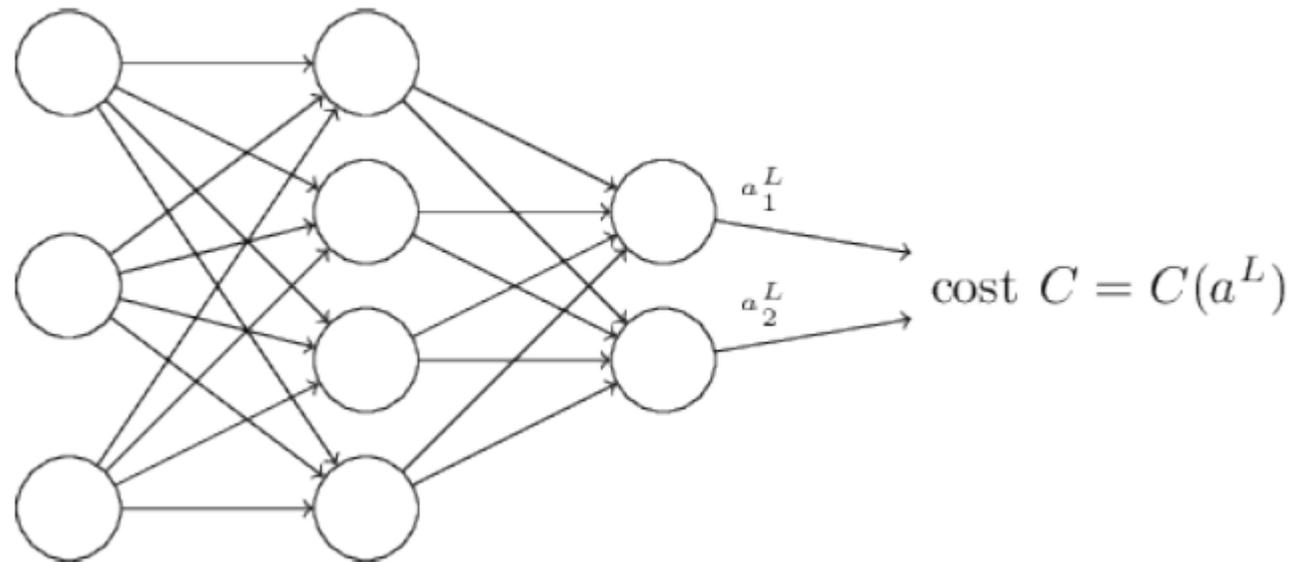
# Assumptions on the cost function

1. We assume that the cost function can be written as the average over the cost functions from individual training samples:  $C = \frac{1}{n} \sum_x C_x$ . The cost function for the individual training sample is given by  $C_x = \frac{1}{2} |y(x) - a^L(x)|^2$ .

- why do we need this assumption? Backpropagation will only allow us to compute the gradients with respect to a single training sample as given by  $\frac{\partial C_x}{\partial w}$  and  $\frac{\partial C_x}{\partial b}$ . We then recover  $\frac{\partial C}{\partial w}$  and  $\frac{\partial C}{\partial b}$  by averaging the gradients from the different training samples.

# Assumptions on the cost function (continued)

2. We assume that the cost function can be written as a function of the output from the neural network. We assume that the input  $x$  and its associated correct labeling  $y(x)$  are fixed and treated as constants.



# Hadamard product

- Let  $s$  and  $t$  are two vectors. The Hadamard product is given by:

$$s \odot t$$

$$(s \odot t)_j = s_j t_j$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 * 3 \\ 2 * 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 8 \end{bmatrix}$$

Such elementwise multiplication is also referred to as schur product.

# Backpropagation

- Our goal is to compute the partial derivatives  $\frac{\partial C}{\partial w_{jk}^l}$  and  $\frac{\partial C}{\partial b_j^l}$ .
- We compute some intermediate quantities while doing so:

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}$$

# Four equations of the BP (backpropagation)

**Summary: the equations of backpropagation**

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

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

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

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

# Chain Rule in differentiation

- In order to differentiate a function  $z = f(g(x))$  w.r.t  $x$ , we can do the following:

$$\text{Let } y = g(x), \quad z = f(y), \quad \frac{dz}{dx} = \frac{dz}{dy} \times \frac{dy}{dx}$$

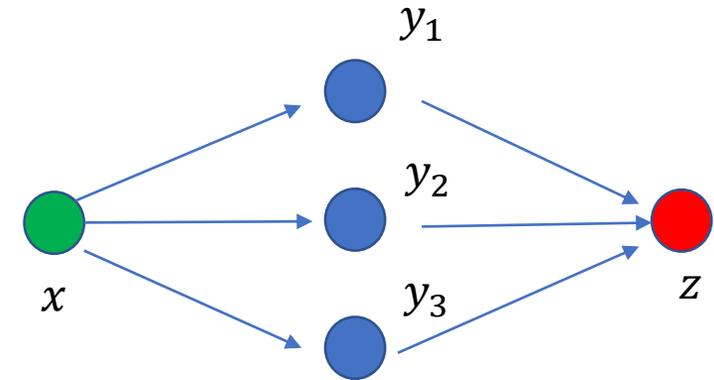
# Chain Rule in differentiation (vector case)

Let  $x \in \mathbb{R}^m$ ,  $y \in \mathbb{R}^n$ ,  $g$  maps from  $\mathbb{R}^m$  to  $\mathbb{R}^n$ , and  $f$  maps from  $\mathbb{R}^n$  to  $\mathbb{R}$ . If  $y = g(x)$  and  $z = f(y)$ , then

$$\frac{\partial z}{\partial x_i} = \sum_k \frac{\partial z}{\partial y_k} \frac{\partial y_k}{\partial x_i}$$

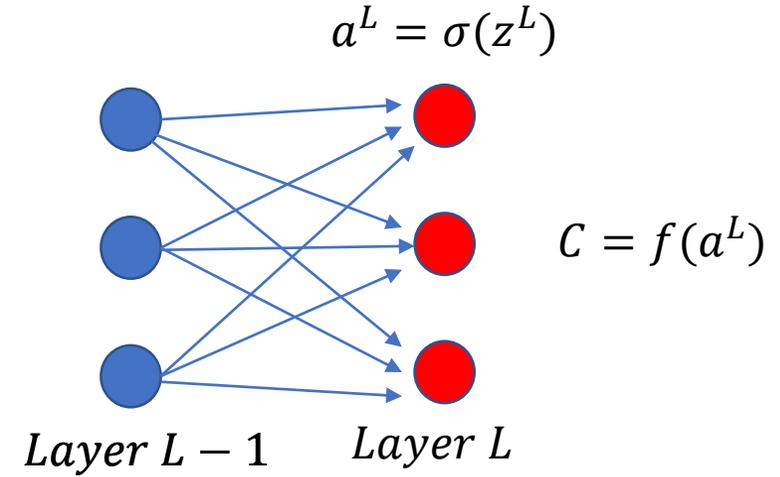
# Chain Rule in differentiation (computation graph)

$$\frac{\partial z}{\partial x} = \sum_{\substack{j: x \in \text{Parent}(y_j), \\ y_j \in \text{Ancestor}(z)}} \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x}$$



# BP1

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



Here L is the last layer.

$$\delta^L = \frac{\partial C}{\partial z^L},$$

$$\sigma'(z^L) = \frac{\partial(\sigma(z^L))}{\partial z^L},$$

$$\nabla_a C = \frac{\partial C}{\partial a^L} = \left( \frac{\partial C}{\partial a_1^L}, \frac{\partial C}{\partial a_2^L}, \dots, \frac{\partial C}{\partial a_n^L} \right)^T$$

Proof:

$$\delta_j^L = \frac{\partial C}{\partial z_j^L} = \sum_k \frac{\partial C}{\partial a_k^L} \frac{\partial a_k^L}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L} \quad \text{when } j \neq k, \text{ the term } \frac{\partial a_k^L}{\partial z_j^L} \text{ vanishes.}$$

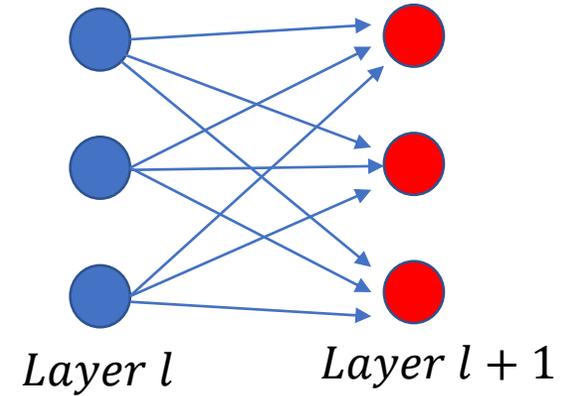
$$\delta_j^L = \frac{\partial C}{\partial a_j^L} \sigma'(z_j^L)$$

Thus we have

$$\delta^L = \frac{\partial C}{\partial a^L} \odot \sigma'(z^L)$$

# BP2

$$z^{l+1} = w^{l+1} a^l + b^{l+1}$$



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

Proof:

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} = \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1}$$

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

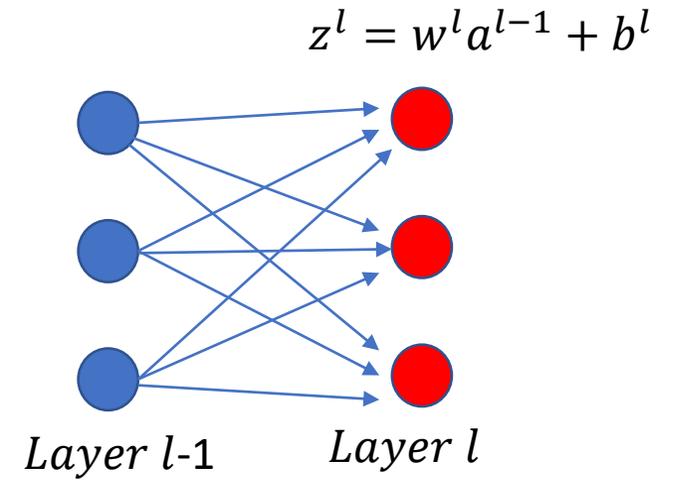
By differentiating we have:

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

$$\delta_j^l = \sum_k w_{kj}^{l+1} \delta_k^{l+1} \sigma'(z_j^l)$$

# BP3

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$



Proof:

$$\begin{aligned} \frac{\partial C}{\partial b_j^l} &= \sum_k \left( \frac{\partial C}{\partial z_k^l} \frac{\partial z_k^l}{\partial b_j^l} \right) = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial b_j^l} \\ &= \delta_j^l \frac{\partial (\sum_k (w_{jk} a_k^{l-1} + b_j^l))}{\partial b_j} \\ &= \delta_j^l \end{aligned}$$

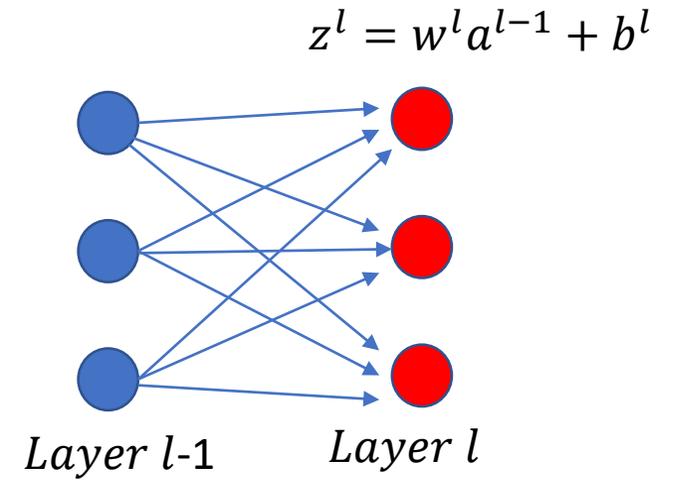
# BP4

Proof:

$$\frac{\partial C}{\partial w_{jk}^l} = \sum_m \frac{\partial C}{\partial z_m^l} \frac{\partial z_m^l}{\partial w_{jk}^l}$$

$$\begin{aligned} &= \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} \\ &= \delta_j^l \frac{\partial (\sum_k w_{jk}^l a_k^{l-1} + b_j^l)}{\partial w_{jk}^l} \\ &= \delta_j^l a_k^{l-1} \end{aligned}$$

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



# The backpropagation algorithm

1. **Input  $x$ :** Set the corresponding activation  $a^1$  for the input layer.
2. **Feedforward:** For each  $l = 2, 3, \dots, L$  compute  $z^l = w^l a^{l-1} + b^l$  and  $a^l = \sigma(z^l)$ .
3. **Output error  $\delta^L$ :** Compute the vector  $\delta^L = \nabla_a C \odot \sigma'(z^L)$ .
4. **Backpropagate the error:** For each  $l = L - 1, L - 2, \dots, 2$  compute  $\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$ .
5. **Output:** The gradient of the cost function is given by  $\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$  and  $\frac{\partial C}{\partial b_j^l} = \delta_j^l$ .

The word “backpropagation” comes from the fact that we compute the error vectors  $\delta_j^l$  in the backward direction.

# Stochastic gradient descent with BP

1. **Input a set of training examples**
2. **For each training example  $x$ :** Set the corresponding input activation  $a^{x,1}$ , and perform the following steps:
  - **Feedforward:** For each  $l = 2, 3, \dots, L$  compute  $z^{x,l} = w^l a^{x,l-1} + b^l$  and  $a^{x,l} = \sigma(z^{x,l})$ .
  - **Output error  $\delta^{x,L}$ :** Compute the vector  $\delta^{x,L} = \nabla_a C_x \odot \sigma'(z^{x,L})$ .
  - **Backpropagate the error:** For each  $l = L - 1, L - 2, \dots, 2$  compute  $\delta^{x,l} = ((w^{l+1})^T \delta^{x,l+1}) \odot \sigma'(z^{x,l})$ .
3. **Gradient descent:** For each  $l = L, L - 1, \dots, 2$  update the weights according to the rule  $w^l \rightarrow w^l - \frac{\eta}{m} \sum_x \delta^{x,l} (a^{x,l-1})^T$ , and the biases according to the rule  $b^l \rightarrow b^l - \frac{\eta}{m} \sum_x \delta^{x,l}$ .

# Gradients using finite differences

$$\frac{\partial C}{\partial w_j} \approx \frac{C(w + \epsilon e_j) - C(w)}{\epsilon}$$

Here  $\epsilon$  is a small positive number and  $e_j$  is the unit vector in the  $j$ th direction.

Conceptually very easy to implement.

In order to compute this derivative w.r.t one parameter, we need to do one forward pass – for millions of variables we will have to do millions of forward passes.

- Backpropagation can get all the gradients in just one forward and backward pass – forward and backward passes are roughly equivalent in computations.

The derivatives using finite differences would be a million times slower!!

# Backpropagation – the big picture

$$\Delta C \approx \sum_{mnp\dots q} \frac{\partial C}{\partial a_m^L} \frac{\partial a_m^L}{\partial a_n^{L-1}} \frac{\partial a_n^{L-1}}{\partial a_p^{L-2}} \cdots \frac{\partial a_q^{l+1}}{\partial a_j^l} \frac{\partial a_j^l}{\partial w_{jk}^l} \Delta w_{jk}^l$$

- To compute the total change in C we need to consider all possible paths from the weight to the cost

$$\frac{\partial C}{\partial w_{jk}^l} = \sum_{mnp\dots q} \frac{\partial C}{\partial a_m^L} \frac{\partial a_m^L}{\partial a_n^{L-1}} \frac{\partial a_n^{L-1}}{\partial a_p^{L-2}} \cdots \frac{\partial a_q^{l+1}}{\partial a_j^l} \frac{\partial a_j^l}{\partial w_{jk}^l}$$

- We are computing the rate of change of C w.r.t a weight w.
- Every edge between two neurons in the network is associated with a rate factor that is just the ratio of partial derivatives of one neurons activation with respect to another neurons activation.
- The rate factor for a path is just the product of the rate factors of the edges in the path.
- The total change is the sum of the rate factors of all the paths from the weight to the cost.

Thank You

# Chain Rule in differentiation (vector case)

Let  $x \in \mathbb{R}^m$ ,  $y \in \mathbb{R}^n$ ,  $g$  maps from  $\mathbb{R}^m$  to  $\mathbb{R}^n$ , and  $f$  maps from  $\mathbb{R}^n$  to  $\mathbb{R}$ . If  $y = g(x)$  and  $z = f(y)$ , then

$$\frac{\partial z}{\partial x_i} = \sum_k \frac{\partial z}{\partial y_k} \frac{\partial y_k}{\partial x_i}$$

$$\nabla_x z = \left( \frac{\partial y}{\partial x} \right)^T \nabla_y z$$

Here  $\left( \frac{\partial y}{\partial x} \right)$  is the  $n \times m$  Jacobian matrix of  $g$ .

## DERIVATIVE RULES

$$\frac{d}{dx}(x^n) = nx^{n-1}$$

$$\frac{d}{dx}(\sin x) = \cos x$$

$$\frac{d}{dx}(\cos x) = -\sin x$$

$$\frac{d}{dx}(a^x) = \ln a \cdot a^x$$

$$\frac{d}{dx}(\tan x) = \sec^2 x$$

$$\frac{d}{dx}(\cot x) = -\operatorname{csc}^2 x$$

$$\frac{d}{dx}(f(x) \cdot g(x)) = f(x) \cdot g'(x) + g(x) \cdot f'(x)$$

$$\frac{d}{dx}(\sec x) = \sec x \tan x$$

$$\frac{d}{dx}(\operatorname{csc} x) = -\operatorname{csc} x \cot x$$

$$\frac{d}{dx}\left(\frac{f(x)}{g(x)}\right) = \frac{g(x) \cdot f'(x) - f(x) \cdot g'(x)}{(g(x))^2}$$

$$\frac{d}{dx}(\arcsin x) = \frac{1}{\sqrt{1-x^2}}$$

$$\frac{d}{dx}(\arctan x) = \frac{1}{1+x^2}$$

$$\frac{d}{dx}(f(g(x))) = f'(g(x)) \cdot g'(x)$$

$$\frac{d}{dx}(\operatorname{arc} \sec x) = \frac{1}{x\sqrt{x^2-1}}$$

$$\frac{d}{dx}(\ln x) = \frac{1}{x}$$

$$\frac{d}{dx}(\sinh x) = \cosh x$$

$$\frac{d}{dx}(\cosh x) = \sinh x$$

## INTEGRAL RULES

$$\int x^n dx = \frac{1}{n+1} x^{n+1} + c, \quad n \neq -1$$

$$\int a^x dx = \frac{1}{\ln a} a^x + c$$

$$\int \frac{1}{x} dx = \ln|x| + c$$

$$\int \frac{dx}{\sqrt{1-x^2}} = \arcsin x + c$$

$$\int \frac{dx}{1+x^2} = \arctan x + c$$

$$\int \frac{dx}{x\sqrt{x^2-1}} = \operatorname{arcsec} x + c$$

$$\int \sin x dx = -\cos x + c$$

$$\int \cos x dx = \sin x + c$$

$$\int \sec^2 x dx = \tan x + c$$

$$\int \sinh x dx = \cosh x + c$$

$$\int \csc^2 x dx = -\cot x + c$$

$$\int \sec x \tan x dx = \sec x + c$$

$$\int \csc x \cot x dx = -\csc x + c$$

$$\int \cosh x dx = \sinh x + c$$