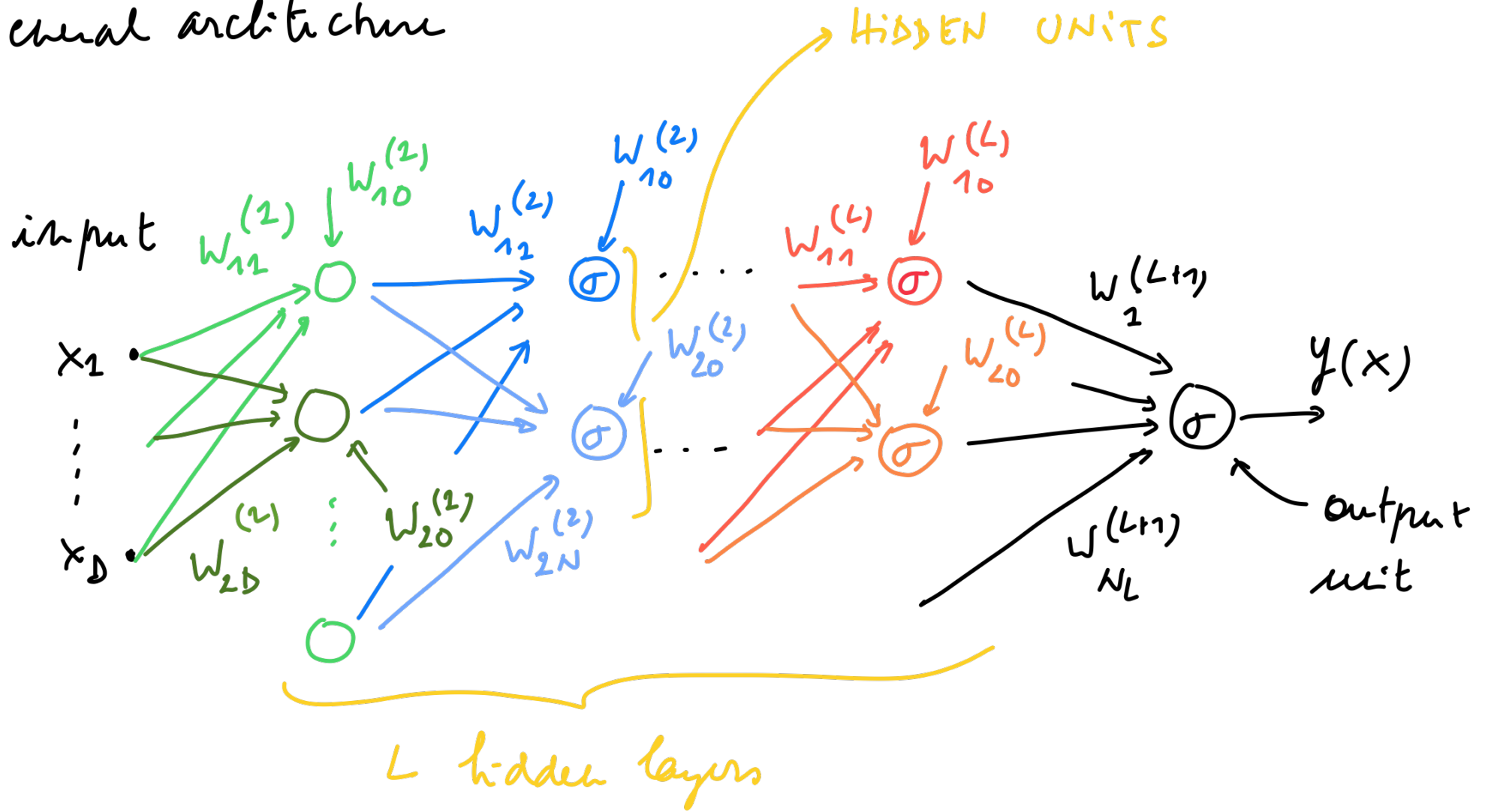


Neural Networks

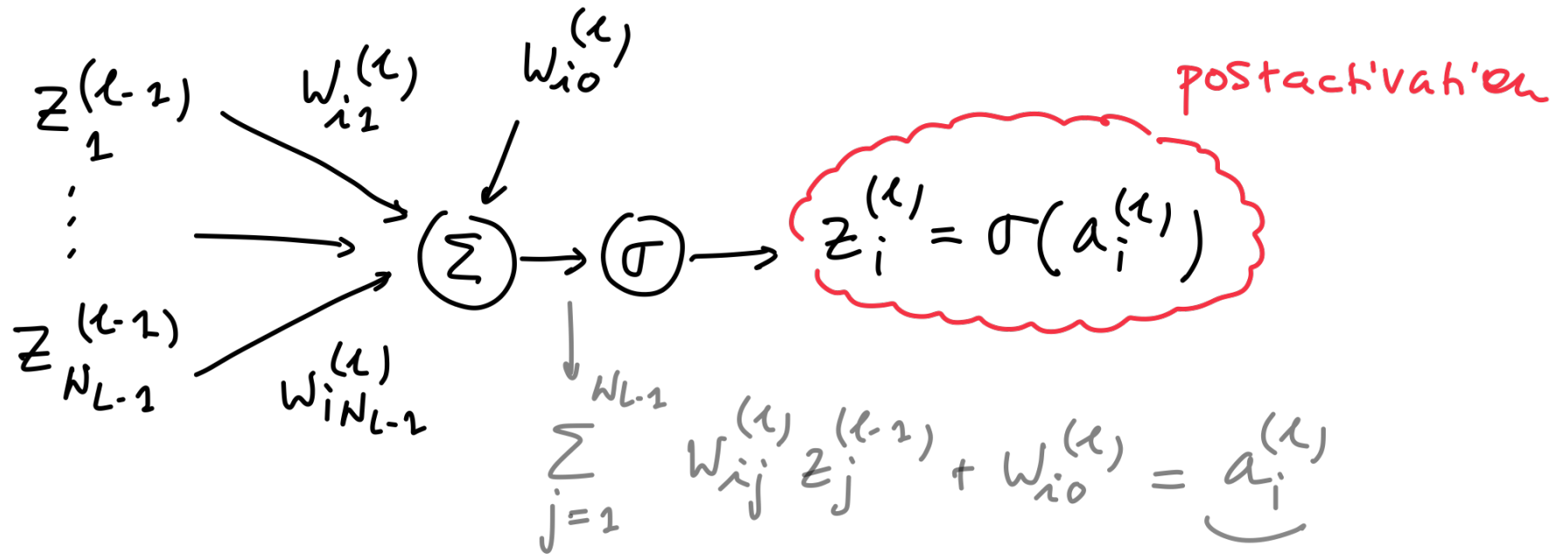
- Reminders (forward propagation, architecture)
- Derivation of gradient through backpropagation
- Recitation (implementation in numpy)

General architecture



$W_{ij}^{(l)}$ is the weight in the i -th neuron from layer l that multiplies the j -th input to this neuron

Zooming on a particular neuron i from layer l , we get



postactivation

preactivation

(whatever is passed to the activation function)

Using our network to encode class probabilities,

$$p(t^{(i)} = 1 | x^{(i)}) = y(x^{(i)}; w)$$

$$p(t^{(i)} = 0 | x^{(i)}) = 1 - y(x^{(i)}; w)$$

→ We can then train the network via a Maximum likelihood approach, corresponding to a minimization of the log loss.

$$L(w) = - \sum_{i=1}^N t^{(i)} \log(y(x^{(i)}; w)) + (1 - t^{(i)}) \log(1 - y(x^{(i)}; w))$$

→ in particular, if we use SGD, optimizing one sample at a time, the loss reduces to

$$L(w) = - t^{(i)} \log(y(x^{(i)}; w)) - (1 - t^{(i)}) \log(1 - y(x^{(i)}; w))$$

To implement the training of the network through SGD, we need the derivatives $\frac{\partial \mathcal{L}}{\partial w_{ij}^{(l)}}$

→ use the fact that $a_i^{(l)} = \sum_{j=2}^{N_L} w_{ij}^{(l)} z_j^{(l-1)} + w_{i0}^{(l)}$

→ From this $\frac{\partial \mathcal{L}}{\partial w_{ij}^{(l)}} = \frac{\partial \mathcal{L}}{\partial a_i^{(l)}} \cdot \frac{\partial a_i^{(l)}}{\partial w_{ij}^{(l)}} = \frac{\partial \mathcal{L}}{\partial a_i^{(l)}} \cdot z_j^{(l-1)}$

We will now use $\delta_i^{(l)} = \frac{\partial \mathcal{L}}{\partial a_i^{(l)}}$

From this we get $\frac{\partial \mathcal{L}}{\partial w_{ij}^{(l)}} = \delta_i^{(l)} z_j^{(l-1)}$

$z_j \rightarrow$ can be derived through forward propagation

$\delta_i^{(k)} \rightarrow$ can be derived through back propagation

step 1 $\delta_{out} = \frac{\partial \mathcal{L}}{\partial a_{out}} \rightarrow$ can be obtained easily

Note that $\frac{\partial \mathcal{L}}{\partial a_{out}} = -\frac{\partial}{\partial a_{out}} (t^{(i)} \log(\sigma(a_{out})) + (1-t^{(i)}) \log(1-\sigma(a_{out})))$

$$\rightarrow \frac{\partial \mathcal{L}}{\partial a_{out}} = -t^{(i)} \frac{\sigma'(a_{out})}{\sigma(a_{out})} + -(1-t^{(i)}) \frac{(-\sigma'(a_{out}))}{1-\sigma(a_{out})} \quad (*)$$

Recall that $\sigma'(x) = \sigma(x)(1-\sigma(x))$

Substituting this in (*)

$$\delta^{\text{out}} = \frac{\partial \mathcal{L}}{\partial a_{\text{out}}} = -t^{(i)}(1 - \sigma(a_{\text{out}})) + (1 - t^{(i)})\sigma(a_{\text{out}})$$

$$= \sigma(a_{\text{out}}) - t^{(i)}$$

$$\delta^{\text{out}} = f(\lambda^{(i)}; w) - t^{(i)}$$

→ Question? Given δ^{out} , how can we derive $\delta_i^{(l)}$ for $l < L$?

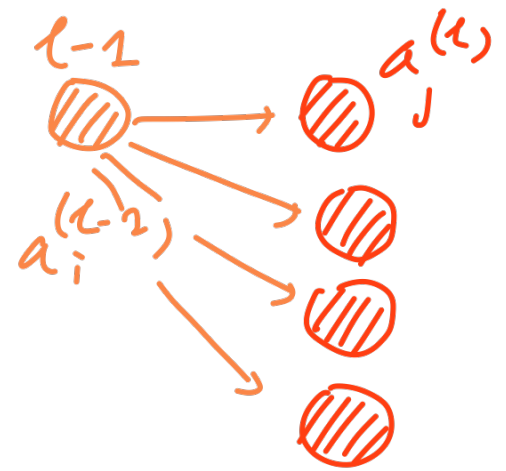
Using the chain rule and the fact that the preactivation $a_i^{(l-1)}$ is "fed" to all the neurons at layer l , we can

write

$$\delta_i^{(l-1)} \frac{\partial \mathcal{L}}{\partial a_i^{(l-1)}} = \sum_{j=1}^{N_L} \frac{\partial \mathcal{L}}{\partial a_j^{(l)}} \frac{\partial a_j^{(l)}}{\partial a_i^{(l-1)}} \delta_j^{(l)}$$

$$\delta_i^{(l-2)} = \sum_{j=2}^{N_L} \delta_j^{(l)}$$

$$\frac{\partial a_j^{(l)}}{\partial a_i^{(l-2)}} \delta_j^{(l)}$$



$$a_j^{(l)} = \sum_{i=1}^{N_{L-2}} w_{ji}^{(l)} z_i^{(l-2)} + w_{j0}^{(l)}$$

$$a_j^{(l)} = \sum_{i=2}^{N_{L-2}} w_{ji}^{(l)} \sigma(a_i^{(l-2)}) + w_{j0}^{(l)}$$

$$\delta_i^{(l-2)} = \sum_{j=2}^{N_L} \delta_j^{(l)} w_{ji}^{(l)} \sigma'(a_i^{(l-2)})$$

$$\delta_i^{(l-2)} = \sigma'(a_i^{(l-2)}) \sum_{j=2}^{N_L} \delta_j^{(l)} w_{ji}^{(l)}$$

BACK PROPAGATION
OF THE $\delta_i^{(l)}$

→ Derivatives w.r.t intercepts $w_{i0}^{(l)}$ can be obtained similarly by augmenting the post activation vector z with a 1 (i.e. generating \tilde{z} as in regression and classification)

→ GENERAL BACKPROPAGATION ALGORITHM

Step 1 Forward propagate input $x^{(l)}$ through the network and derive all pre and post activations $a_i^{(l)}, z_i^{(l)}$

Step 2 Compute δ^{out}

Step 3 Backpropagate the δ^{out} to derive all $\delta_i^{(l)}$

Step 4 Get the gradient as $\frac{\partial \mathcal{L}}{\partial w_{ij}^{(l)}} = \delta_i^{(l)} z_j^{(l-1)}$

Step 5 apply the update

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta \frac{\partial \mathcal{L}}{\partial w_{ij}^{(l)}}$$