

# Artificial Neural Networks

## Lecture 5: Multi-Layer Networks

Dr Juan Neirotti

[j.p.neirotti@aston.ac.uk](mailto:j.p.neirotti@aston.ac.uk)

2023

# Multi-Layer Networks

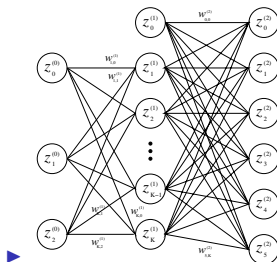


Figure: Feed-forward network with 3 layers, input, hidden and output.

- ▶ In figure 1 we have a feed-forward networks with an input layer with  $2+1$  units, a hidden layer with  $K+1$  units and an output layer with 6 units.
- ▶ The input units  $z_k^{(0)}$  can be identified with the inputs  $x_k$ , which are the entries of the input vector  $\mathbf{x} \in \mathcal{X}$ .

# Multi-Layer Networks

- ▶ The general form of the activation functions is:

$$z_k^{(\ell+1)} = \sigma \left( \left\langle \mathbf{w}_k^{(\ell)} \mid \mathbf{z}^{(\ell)} \right\rangle + w_0^{(\ell)} \right),$$

where:

$$\begin{aligned} \lim_{a \rightarrow -\infty} \sigma(a) &= 0 \quad (\text{or } -1), \\ \lim_{a \rightarrow \infty} \sigma(a) &= 1. \end{aligned}$$

Functions like these are called sigmoidal.

# Regression Problem: Cybenko's Theorem

- Theorem: Let  $\sigma(x)$  be a bounded sigmoidal function, and  $f : \mathbb{R} \rightarrow \mathbb{R}$  continuous, satisfying  $\lim_{x \rightarrow -\infty} f(x) = A$  and  $\lim_{x \rightarrow \infty} f(x) = B$ , where  $A$  and  $B$  are constants, then for any  $\varepsilon > 0$ , there exists  $N$ ,  $c_i$ ,  $y_i$ , and  $\theta_i$ , such that:

$$\left| f(x) - \sum_{i=1}^N c_i \sigma(y_i x + \theta_i) \right| < \varepsilon$$

holds for all  $x \in \mathbb{R}$ .

# Proof

- ▶ By continuity of  $f$  we have that for all  $\varepsilon > 0$  there must exist  $M \in \mathbb{N}$  such that  $|f(x) - A| < \frac{\varepsilon}{4}$  if  $x < -M$ ;  $|f(x) - B| < \frac{\varepsilon}{4}$  if  $x > M$ ; and  $|f(x') - f(x'')| < \frac{\varepsilon}{4}$  if  $|x'| \leq M$ ,  $|x''| \leq M$ , and  $|x' - x''| \leq \frac{1}{M}$ .
- ▶ Let us divide the interval  $[-M, M]$  into  $N = 2M^2$  equal segments of length  $\frac{1}{M}$ , and let define:  $x_i = -M + \frac{i}{M}$ .
- ▶ Let us also define  $t_i = \frac{x_i + x_{i+1}}{2}$ , which is the center of the interval  $[x_i, x_{i+1}]$ .
- ▶ Now let us construct:

$$g(x) = f(-M) + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \sigma(K(x - t_{i-1})).$$

- ▶ By the properties of the sigmoidal function there exists  $W > 0$  such that  $1 - \sigma(u) < \frac{1}{M^2}$  if  $u > W$  and  $\sigma(u) < \frac{1}{M^2}$  if  $u < -W$ . Let us choose  $K = 2MW$ .
- ▶ If  $x < -M$ , then

$$|f(x) - f(-M)| < |f(x) - A| + |A - f(-M)| < \frac{\varepsilon}{4} + \frac{\varepsilon}{4} = \frac{\varepsilon}{2}.$$

- ▶ Also:

$$K(x - t_{i-1}) = -K|x + M| - W(2i - 1) < -W$$

$$\begin{aligned} |g(x) - f(-M)| &= \left| \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \sigma(K(x - t_{i-1})) \right| \\ &< \sum_{i=1}^N |f(x_i) - f(x_{i-1})| \frac{1}{M^2} \\ &< 2M^2 \frac{\varepsilon}{4} \frac{1}{M^2} = \frac{\varepsilon}{2}. \end{aligned}$$

- In consequence

$$|g(x) - f(x)| < |g(x) - f(-M)| + |f(-M) - f(x)| < \varepsilon \text{ for all } x < -M.$$

- If  $x > M$ , then

$$|f(x) - f(M)| < |f(x) - B| + |B - f(M)| < \frac{\varepsilon}{4} + \frac{\varepsilon}{4} = \frac{\varepsilon}{2}.$$

- Also:

$$K(x - t_{i-1}) \geq K(x - M) + \frac{K}{2M} > W.$$

► Then:

$$\begin{aligned}g(x) &= f(-M) + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \{\sigma(K(x - t_{i-1})) - 1 + 1\} \\&= f(-M) + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \{\sigma(K(x - t_{i-1})) - 1\} \\&\quad + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \\&= f(-M) + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \{\sigma(K(x - t_{i-1})) - 1\} \\&\quad + f(x_N) - f(x_0) \\&= f(M) + \sum_{i=1}^N [f(x_i) - f(x_{i-1})] \{\sigma(K(x - t_{i-1})) - 1\} \\|g(x) - f(M)| &< \sum_{i=1}^N |f(x_i) - f(x_{i-1})| \{1 - \sigma(K(x - t_{i-1}))\} \\&< 2M^2 \frac{\varepsilon}{4} \frac{1}{M^2} = \frac{\varepsilon}{2}.\end{aligned}$$



- ▶ In consequence

$$|g(x) - f(x)| < |g(x) - f(M)| + |f(M) - f(x)| < \varepsilon \text{ for all } x > M.$$

- ▶ If  $-M < x < M$  there exists a  $0 \leq k \leq 2M^2$  such that  $x \in [x_{k-1}, x_k]$  and then  $|x - t_{i-1}| \leq \frac{1}{2M}$  if  $i = k$  and  $|x - t_{i-1}| > \frac{1}{2M}$  if  $i \neq k$ .
- ▶ Furthermore, if  $i < k$  then  $x > t_{i-1}$  and:

$$K(x - t_{i-1}) > \frac{K}{2M} = W,$$

and

$$1 - \sigma(K(x - t_{i-1})) < \frac{1}{M^2}.$$

- ▶ If  $i > k$  then  $x < t_{i-1}$  and:

$$K(x - t_{i-1}) < -\frac{K}{2M} = -W,$$

and

$$\sigma(K(x - t_{i-1})) < \frac{1}{M^2}.$$

- Consequently we have:

$$\begin{aligned}
 g(x) &= f(-M) + \sum_{i=1}^{k-1} [f(x_i) - f(x_{i-1})] \sigma(K(x - t_{i-1})) + \\
 &\quad + [f(x_k) - f(x_{k-1})] \sigma(K(x - t_{k-1})) + \\
 &\quad + \sum_{i=k+1}^N [f(x_i) - f(x_{i-1})] \sigma(K(x - t_{i-1})) \\
 &= f(-M) + \sum_{i=1}^{k-1} [f(x_i) - f(x_{i-1})] \{\sigma(K(x - t_{i-1})) - 1\} + \\
 &\quad + f(x_{k-1}) - f(x_0) + \\
 &\quad + [f(x_k) - f(x_{k-1})] \sigma(K(x - t_{k-1})) + \\
 &\quad + \sum_{i=k+1}^N [f(x_i) - f(x_{i-1})] \sigma(K(x - t_{i-1}))
 \end{aligned}$$

► Therefore

$$\begin{aligned} A &= |g(x) - f(x_{k-1}) - [f(x_k) - f(x_{k-1})] \sigma(K(x - t_{k-1}))| \\ &< \sum_{i=1}^{k-1} |f(x_i) - f(x_{i-1})| \sigma(K(x - t_{i-1})) + \\ &\quad + \sum_{i=k+1}^N |f(x_i) - f(x_{i-1})| \{1 - \sigma(K(x - t_{i-1}))\} \\ &< \sum_{i=1}^{k-1} \frac{\varepsilon}{4} \frac{1}{M^2} + \sum_{i=k+1}^N \frac{\varepsilon}{4} \frac{1}{M^2} = \frac{2M^2 - 1}{2M^2} \frac{\varepsilon}{2} < \frac{\varepsilon}{2}. \end{aligned}$$

► Finally:

$$\begin{aligned} |g(x) - f(x)| &< |f(x_{k-1}) - f(x)| + \frac{\varepsilon}{4} \sigma(K(x - t_{k-1})) + \frac{\varepsilon}{2} \\ &< \frac{\varepsilon}{4} + \frac{\varepsilon}{4} + \frac{\varepsilon}{2} = \varepsilon. \clubsuit \end{aligned}$$

# Classification Problem: Grand Mother Cells

- ▶ Consider  $\mathbf{x} \in \mathcal{X} = \{\pm 1\}^d \subset \mathbb{R}^d$ .
- ▶ Let us verify that any dichotomy  $\{\mathcal{X}^+, \mathcal{X}^-\}$  of  $\mathcal{X}$  can be realized by a neural network with one hidden layer.
- ▶ Let us suppose that  $|\mathcal{X}^+| \leq |\mathcal{X}^-|$ . Let us consider the vector  $\mathbf{x}_k^+ \in \mathcal{X}^+$ , for any  $1 \leq k \leq |\mathcal{X}^+|$ .

Observe that

$$\langle \mathbf{x} | \mathbf{x}_k^+ \rangle \begin{cases} < d - 1 & \mathcal{X} \ni \mathbf{x} \neq \mathbf{x}_k^+ \\ = d & \mathcal{X} \ni \mathbf{x} = \mathbf{x}_k^+. \end{cases}$$

- Therefore

$$\text{sgn} \left( \langle \mathbf{x} | \mathbf{x}_k^+ \rangle - d + \frac{1}{2} \right) = \begin{cases} 1 & \mathbf{x} = \mathbf{x}_k^+ \\ -1 & \mathbf{x} \neq \mathbf{x}_k^+. \end{cases}$$

- The network with  $|\mathcal{X}^+| = K \leq 2^{d-1}$  units in its hidden layer, each one of them implementing the activation function

$$z_k^{(1)} = \text{sgn} \left( \langle \mathbf{x} | \mathbf{x}_k^+ \rangle - d + \frac{1}{2} \right), \text{ i.e. } \mathbf{w}_k^{(1)} = \mathbf{x}_k^+ \text{ and } w_{0,k} = -d + \frac{1}{2}, \text{ and an output layer that implements}$$

$$z^{(2)} = \text{sgn} \left( \sum_{k=1}^K z_k^{(1)} + K - \frac{1}{2} \right), \text{ i.e. } \langle \mathbf{w}^{(2)} | = \overbrace{(1, \dots, 1)}^K, \text{ and } w_0^{(2)} = K - \frac{1}{2}, \text{ produces the correct classification of the dichotomy } \{\mathcal{X}^+, \mathcal{X}^-\}.$$

- The problem with this approach is that the total number of units needed may be  $2^{d-1}$ , which could be very large indeed.

# General Feed-Forward Networks (with a hidden layer)

- ▶ We consider the case of a network with one hidden layer and one linear output.
- ▶ Linear outputs are usually used for regression problems (binary outputs are used for classification)
- ▶ Let us consider the data set  $\mathcal{D} = \{(\mathbf{y}_\ell, \mathbf{t}_\ell)\}_{\ell=1}^L$  with  $\mathbf{y}_\ell \in \mathcal{Y}$  and  $\mathbf{t}_\ell \in \mathbb{R}^o$  for all  $1 \leq \ell \leq L$ , and let us suppose there is a continuous function  $\mathbf{h} : \mathcal{Y} \rightarrow \mathbb{R}^o$  such that  $\mathbf{h}(\mathbf{y}_\ell) = \mathbf{t}_\ell$  for all  $1 \leq \ell \leq L$ .
- ▶ Observe that  $\langle \mathbf{h} | = (h_1, \dots, h_o)$  where  $h_n : \mathcal{Y} \rightarrow \mathbb{R}$  for all  $1 \leq n \leq o$ .

- By Cybenko's approach we should be able to construct a network of sigmoidal units such that the  $K$ th approximation:

$$z_n^{(2)}(\mathbf{y}) = w_{0,n}^{(2)} + \sum_{k=1}^K w_{k,n}^{(2)} \sigma \left( w_{k,0}^{(1)} + \sum_{i=1}^d w_{k,i}^{(1)} y_i \right),$$

is as close as  $h_n(\mathbf{y})$  as required.

- Inputs  $\mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^d$  are augmented to  $\mathbf{x} \in \{1\} \times \mathcal{Y} \subset \mathbb{R}^{d+1}$ .
- Let us identify the variables  $\mathbf{z}^{(0)} = \mathbf{x}$ . Observe that  $z_0^{(0)} = 1$  always.
- The  $k$ th *activation* of the hidden layer is given by

$$a_k^{(1)} = \left\langle \mathbf{x} \mid \mathbf{w}_k^{(1)} \right\rangle = \left\langle \mathbf{z}^{(0)} \mid \mathbf{w}_k^{(1)} \right\rangle$$

for a suitable  $\mathbf{w}_k^{(1)} \in \Omega^{(1)} \subset \mathbb{R}^{d+1}$ .

- ▶ The  $k$ th *activation function* ( $1 \leq k \leq K$ ) is implemented by a sigmoidal function  $\sigma$  :

$$z_k^{(1)} = \sigma(a_k^{(1)}),$$

where

$$\lim_{x \rightarrow \infty} \sigma(x) = 1$$

$$\lim_{x \rightarrow -\infty} \sigma(x) = -1.$$

- ▶ We define  $z_0^{(1)} = 1$ .
- ▶ The linear output is  $z_n^{(2)} = a_n^{(2)} = \left\langle z^{(1)} \mid \mathbf{w}_n^{(2)} \right\rangle$  for a suitable  $\mathbf{w}_n^{(2)} \in \Omega^{(2)} \subset \mathbb{R}^{K+1}$ .



# Error Back Propagation (EBP)

- ▶ An adequate loss function for the problem is defined as:

$$\mathcal{L} \left( \{\mathbf{w}_n^{(2)}\}_{n=1}^O, \{\mathbf{w}_k^{(1)}\}_{k=1}^K \right) = \frac{1}{2} \sum_{\ell=1}^L \left\langle \mathbf{t}_\ell - \mathbf{z}^{(2)}(\mathbf{x}_\ell) \mid \mathbf{t}_\ell - \mathbf{z}^{(2)}(\mathbf{x}_\ell) \right\rangle = \sum_{\ell=1}^L \mathcal{L}_\ell.$$

- ▶ Let us compute the derivatives of the Loss function with respect to the vectors  $\mathbf{w}^{(2)}$  and  $\mathbf{w}_k^{(1)}$  :

$$\begin{aligned} \left| \nabla_{\mathbf{w}_n^{(2)}} \mathcal{L}_\ell \right\rangle &= \frac{\partial \mathcal{L}_\ell}{\partial a_n^{(2)}} \nabla_{\mathbf{w}_n^{(2)}} a_n^{(2)} \\ \delta_n^{(2)} &:= \frac{\partial \mathcal{L}_\ell}{\partial a_n^{(2)}} = z_n^{(2)}(\mathbf{x}_\ell) - t_{n,\ell} \\ \left| \nabla_{\mathbf{w}_n^{(2)}} a_n^{(2)} \right\rangle &= \left| \nabla_{\mathbf{w}_n^{(2)}} \left\langle \mathbf{z}^{(1)} \mid \mathbf{w}_n^{(2)} \right\rangle \right\rangle = \left| \mathbf{z}^{(1)} \right\rangle. \end{aligned}$$

- ▶ Thus

$$\left| \nabla_{\mathbf{w}_n^{(2)}} \mathcal{L}_\ell \right\rangle = \delta_n^{(2)} \left| \mathbf{z}^{(1)} \right\rangle.$$

- The derivatives with respect the weights linking hidden units to inputs are:

$$\begin{aligned}
 \left| \nabla_{\mathbf{w}_k^{(1)}} \mathcal{L}_\ell \right\rangle &= \frac{\partial \mathcal{L}_\ell}{\partial a_k^{(1)}} \left| \nabla_{\mathbf{w}_k^{(1)}} a_k^{(1)} \right\rangle \\
 \left| \nabla_{\mathbf{w}_k^{(1)}} a_k^{(1)} \right\rangle &= \left| \nabla_{\mathbf{w}_k^{(1)}} \left\langle \mathbf{x} \left| \mathbf{w}_k^{(1)} \right\rangle \right\rangle = |\mathbf{x}\rangle \\
 \frac{\partial a_n^{(2)}}{\partial a_k^{(1)}} &= \frac{\partial}{\partial a_k^{(1)}} \left\langle \mathbf{z}^{(1)} \left| \mathbf{w}_n^{(2)} \right\rangle = w_{n,k}^{(2)} \sigma' \left( a_k^{(1)} \right) \\
 \delta_k^{(1)} &:= \frac{\partial \mathcal{L}_\ell}{\partial a_k^{(1)}} = \sum_{n=1}^o \frac{\partial \mathcal{L}_\ell}{\partial a_n^{(2)}} \frac{\partial a_n^{(2)}}{\partial a_k^{(1)}} = \sigma' \left( a_k^{(1)} \right) \sum_{n=1}^o \delta_n^{(2)} w_{n,k}^{(2)},
 \end{aligned}$$

- Thus

$$\left| \nabla_{\mathbf{w}_k^{(1)}} \mathcal{L}_\ell \right\rangle = \sigma' \left( a_k^{(1)} \right) \sum_{n=1}^o \delta_n^{(2)} w_{n,k}^{(2)} |\mathbf{x}\rangle.$$

Observe that:

$$\delta_j^{(\text{out})} = \frac{\partial \mathcal{L}_\ell}{\partial z_j^{(\text{out})}} \frac{dz_j^{(\text{out})}}{da_j^{(\text{out})}}$$

$$\delta_j^{(m)} = \frac{dz_j^{(m)}}{da_j^{(m)}} \sum_{i=1}^{K^{(m+1)}} \delta_i^{(m+1)} w_{i,j}^{(m+1)}$$

$$\left| \nabla_{\mathbf{w}_k^{(m)}} a_k^{(m)} \right\rangle = \left| \mathbf{z}^{(m-1)} \right\rangle$$

$$\left| \nabla_{\mathbf{w}_k^{(m)}} \mathcal{L}_\ell \right\rangle = \delta_k^{(m)} \left| \mathbf{z}^{(m-1)} \right\rangle.$$

- ▶ Information is propagated from the input nodes towards the output (all activations  $z^{(m)}$  are computed). Forward step.
- ▶ Deltas ( $\delta_k^{(m)}$ ) are computed from output to input to complete the derivatives. Backward step.
- ▶ Observe that the derivatives of the sigmoidal functions can be expressed in terms of the sigmoids themselves:

$$\frac{d \tanh(x)}{dx} = 1 - \tanh^2(x),$$

$$\frac{d}{dx} \left( \frac{2}{1 + e^{-x}} - 1 \right) = \frac{2e^{-x}}{(1 + e^{-x})^2} = \frac{1}{2} \left[ 1 - \left( \frac{2}{1 + e^{-x}} - 1 \right)^2 \right].$$

# Gradient descent

- ▶ The EBP algorithm is completed by applying the gradients to the weight-update algorithm:

$$\left| \mathbf{w}_k^{(K)} \right\rangle_{\ell+1} = \left| \mathbf{w}_k^{(K)} \right\rangle_{\ell} - \eta_{\ell} \left| \nabla_{\mathbf{w}_k^{(K)}} \mathcal{L}_{\ell} \right\rangle$$