



PERCEPTRONS AND MULTILAYER PERCEPTRONS

THE BACKPROPAGATION ALGORITHM

Vincent Barra
 LIMOS, UMR 6158 CNRS, Université Clermont Auvergne

FORWARD PASS

Train the MLP \leftrightarrow Minimize a loss function ℓ over a training set Z

$$\ell = \sum_{(\mathbf{x}, y) \in Z} \mathcal{L}(f_{(\mathbf{w}, \mathbf{b})}(\mathbf{x}), y)$$

We use a gradient optimization process: if $(\hat{\mathbf{w}}, \hat{\mathbf{b}})$ are the optimal parameters of the MLP, then

$$\sum_{(\mathbf{x}, y) \in Z} \nabla_{(\mathbf{w}, \mathbf{b})} \mathcal{L}(f_{(\hat{\mathbf{w}}, \hat{\mathbf{b}})}(\mathbf{x}), y) = \mathbf{0}$$

Let consider a single example \mathbf{x} , and let $s^i, i \in \llbracket 1 \cdots L \rrbracket$ denote the consecutive post synaptic potential computed at each layer.

$$\mathbf{x} \xrightarrow[\mathbf{w}^1, \mathbf{b}^1]{} s^1 \xrightarrow[\sigma]{} \mathbf{x}^1 \xrightarrow[\mathbf{w}^2, \mathbf{b}^2]{} s^2 \xrightarrow[\sigma]{} \cdots \xrightarrow[\mathbf{w}^L, \mathbf{b}^L]{} s^L \xrightarrow[\sigma]{} f_{(\mathbf{w}, \mathbf{b})}(\mathbf{x})$$

thus (**forward pass**): $\mathbf{x}^{(0)} = \mathbf{x}$ and for $i \in \llbracket 1 \cdots L \rrbracket$;

$$s^{(i)} = \mathbf{w}^{(i)T} \mathbf{x}^{(i-1)} + b^{(i)} \text{ and } \mathbf{x}^{(i)} = \sigma(s^{(i)})$$

BACKPROPAGATION

The core principle of the back-propagation algorithm is the chain rule in \mathbb{R} :

$$(f_1 \circ f_2)' = (f_1' \circ f_2)f_2'$$

in \mathbb{R}^d :

$$J_{f_n \circ \dots \circ f_1}(\mathbf{x}) = J_{f_n}(f_{n-1}(\dots(\mathbf{x}))) \cdots J_{f_3}(f_2(f_1(\mathbf{x}))) J_{f_2}(f_1(\mathbf{x})) J_{f_1}(\mathbf{x})$$

where $J_{f_i}(\mathbf{x})$ is the Jacobian of f at \mathbf{x} .

BACKPROPAGATION

Recall: $x^{(i-1)} \xrightarrow{\mathbf{w}^L, \mathbf{b}^L} s^L \xrightarrow{\sigma} f(\mathbf{w}, \mathbf{b})(\mathbf{x})$

Then since for all components j : $x_j^{(i)} = \sigma(s_j^{(i)})$

$$\frac{\partial \mathcal{L}}{\partial s_j^{(i)}} = \frac{\partial \mathcal{L}}{\partial x_j^{(i)}} \frac{\partial x_j^{(i)}}{\partial s_j^{(i)}} = \frac{\partial \mathcal{L}}{\partial x_j^{(i)}} \sigma'(s_j^{(i)})$$

and since $x_j^{(i-1)}$ impacts \mathcal{L} only through the $s_j^{(i)}$:

$$s_j^{(i)} = \sum_k w_{j,k}^{(i)} x_k^{(i-1)} + b_j^{(i)}$$

BACKPROPAGATION

Recall: $x^{(i-1)} \xrightarrow{\mathbf{w}^L, \mathbf{b}^L} s^L \xrightarrow{\sigma} f(\mathbf{w}, \mathbf{b})(\mathbf{x})$

Then since for all components j : $x_j^{(i)} = \sigma(s_j^{(i)})$

$$\frac{\partial \mathcal{L}}{\partial s_j^{(i)}} = \frac{\partial \mathcal{L}}{\partial x_j^{(i)}} \frac{\partial x_j^{(i)}}{\partial s_j^{(i)}} = \frac{\partial \mathcal{L}}{\partial x_j^{(i)}} \sigma'(s_j^{(i)})$$

and since $s_j^{(i)} = \sum_k w_{j,k}^{(i)} x_k^{(i-1)} + b_j^{(i)}$, then

$$\frac{\partial \mathcal{L}}{\partial x_k^{(i-1)}} = \sum_j \frac{\partial \mathcal{L}}{\partial s_j^{(i)}} \frac{\partial s_j^{(i)}}{\partial x_k^{(i-1)}} = \sum_j \frac{\partial \mathcal{L}}{\partial s_j^{(i)}} w_{j,k}^{(i)}$$

$$\text{and } \frac{\partial \mathcal{L}}{\partial w_{j,k}^{(i)}} = \frac{\partial \mathcal{L}}{\partial s_j^{(i)}} \frac{\partial s_j^{(i)}}{\partial w_{j,k}^{(i)}} = \frac{\partial \mathcal{L}}{\partial s_j^{(i)}} x_k^{(i-1)} \text{ and } \frac{\partial \mathcal{L}}{\partial b_j^{(i)}} = \frac{\partial \mathcal{L}}{\partial s_j^{(i)}}$$

BACKPROPAGATION - SUMMARY

So $\frac{\partial \mathcal{L}}{\partial x_k^{(L)}}$ can be recursively computed (**backward propagated**):

$$\frac{\partial \mathcal{L}}{\partial s_j^{(l)}} = \frac{\partial \mathcal{L}}{\partial x_j^{(i)}} \sigma'(s_j^{(i)})$$

and

$$\frac{\partial \mathcal{L}}{\partial x_k^{(i-1)}} = \sum_j \frac{\partial \mathcal{L}}{\partial s_j^{(i)}} w_{j,k}^{(i)}$$

and then compute the derivatives w.r.t to w and b

$$\frac{\partial \mathcal{L}}{\partial w_{j,k}^{(i)}} = \frac{\partial \mathcal{L}}{\partial s_j^{(i)}} x_k^{(i-1)} \text{ and } \frac{\partial \mathcal{L}}{\partial b_j^{(i)}} = \frac{\partial \mathcal{L}}{\partial s_j^{(i)}}$$

BACKPROPAGATION - MATRIX FORMULATION

Notations:

► for $f : \mathbb{R}^d \rightarrow \mathbb{R}^c$, $f = (f_1 \cdots f_c)^T$, $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$

$$\left[\frac{\partial f}{\partial \mathbf{x}} \right] = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_c}{\partial x_1} & \cdots & \frac{\partial f_c}{\partial x_d} \end{pmatrix}$$

► for $f : \mathbb{R}^{d \times c} \rightarrow \mathbb{R}$,

$$\left[\frac{\partial f}{\partial \mathbf{w}} \right] = \begin{pmatrix} \frac{\partial f}{\partial w_{1,1}} & \cdots & \frac{\partial f}{\partial w_{1,c}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial w_{d,1}} & \cdots & \frac{\partial f}{\partial w_{d,c}} \end{pmatrix}$$

BACKPROPAGATION - MATRIX FORMULATION

- 1 Forward pass: $\mathbf{x}^{(0)} = \mathbf{x}$ and for $i \in \llbracket 1 \cdots L \rrbracket$;

$$\mathbf{s}^{(i)} = \mathbf{w}^{(i)T} \mathbf{x}^{(i-1)} + \mathbf{b}^{(i)} \text{ and } \mathbf{x}^{(i)} = \sigma(\mathbf{s}^{(i)})$$

- 2 Backward pass:

$$\begin{cases} \text{Compute output layer from } \mathcal{L} \\ \text{Hidden layers } i < L \end{cases} \quad \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(L)}} \\ \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(i)}} \end{bmatrix} = \mathbf{w}^{(i+1)T} \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{s}^{(i+1)}} \end{bmatrix}$$

- 3 Compute $\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{s}^{(i)}} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{(i)}} \end{bmatrix} \odot \sigma'(\mathbf{s}^{(i)})$

- 4 Compute gradient: $\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(i)}} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{s}^{(i)}} \end{bmatrix} \mathbf{x}^{(i-1)T}$ and $\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(i)}} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{s}^{(i)}} \end{bmatrix}$

- 5 Gradient descent: $\mathbf{w}^{(i)} = \mathbf{w}^{(i)} - \eta \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{w}^{(i)}} \end{bmatrix}$ and $\mathbf{b}^{(i)} = \mathbf{b}^{(i)} - \eta \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{b}^{(i)}} \end{bmatrix}$