

Backpropagation lecture notes

Nicolò Navarin

The notation from the “Deep Learning” book follows (we will try to simplify this notation)

Calculus

$$\frac{dy}{dx}$$

Derivative of y with respect to x

$$\frac{\partial y}{\partial x}$$

Partial derivative of y with respect to x

$$\nabla_{\mathbf{x}} y$$

Gradient of y with respect to \mathbf{x}

$$\nabla_{\mathbf{X}} y$$

Matrix derivatives of y with respect to \mathbf{X}

$$\nabla_{\mathbf{x}} y$$

Tensor containing derivatives of y with respect to \mathbf{X}

$$\frac{\partial f}{\partial \mathbf{x}}$$

Jacobian matrix $\mathbf{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$

Notation in this lecture notes: variables and constants are in lower case, vectors are bold, matrices are uppercase bold. We will tend to use the partial derivative/Jacobian matrix notation from the book “Mathematics for Machine Learning”.

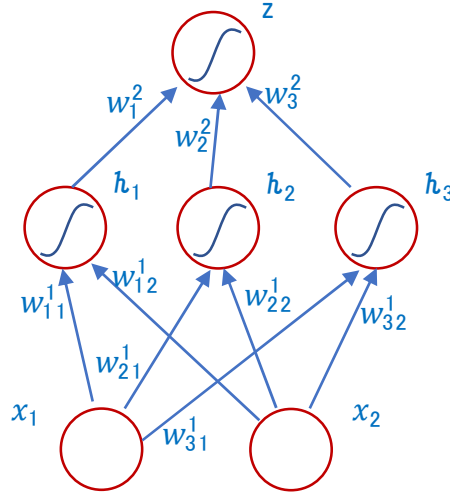
Recall:

$$\sigma(\text{net}) = \frac{1}{1 + e^{-\text{net}}} \quad \sigma'(\text{net}) = \frac{d \sigma(\text{net})}{d \text{net}} = \sigma(\text{net})(1 - \sigma(\text{net}))$$

Chain Rule:

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

We are given a training set $Tr = \{(\mathbf{x}^{(1)}, t^{(1)}), \dots, (\mathbf{x}^{(N_{Tr})}, t^{(N_{Tr})})\}$ and a network:



Forward pass

computes z and all the required variables for each example $\mathbf{x}^{(i)}$. The variables to be stored will be clear when computing the derivatives.

In the following we omit the superscript (i) indicating the specific training example for ease of notation.

$$\begin{aligned} a_1^1 &= w_{11}^1 x_1 + w_{12}^1 x_2 + b_1^1, a_2^1 = w_{21}^1 x_1 + w_{22}^1 x_2 + b_2^1, a_3^1 = w_{31}^1 x_1 + w_{32}^1 x_2 + b_3^1 \\ h_1 &= \sigma(a_1^1), h_2 = \sigma(a_2^1), h_3 = \sigma(a_3^1) \\ a^2 &= w_1^2 h_1 + w_2^2 h_2 + b^2 \\ z &= \sigma(a^2) \end{aligned}$$

Matrix notation:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \mathbf{W}^1 = \begin{bmatrix} w_{11}^1 & w_{12}^1 \\ w_{21}^1 & w_{22}^1 \\ w_{31}^1 & w_{32}^1 \end{bmatrix}, \mathbf{b}^1 = \begin{bmatrix} b_1^1 \\ b_2^1 \\ b_3^1 \end{bmatrix}, \mathbf{a}^1 = \begin{bmatrix} a_1^1 \\ a_2^1 \\ a_3^1 \end{bmatrix} = \mathbf{W}^1 \mathbf{x} + \mathbf{b}^1 = \begin{bmatrix} w_{11}^1 & w_{12}^1 \\ w_{21}^1 & w_{22}^1 \\ w_{31}^1 & w_{32}^1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^1 \\ b_2^1 \\ b_3^1 \end{bmatrix}$$

$$\mathbf{h} = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \sigma(\mathbf{a}^1), \quad \mathbf{W}^2 = \begin{bmatrix} w_1^2 \\ w_2^2 \\ w_3^2 \end{bmatrix},$$

$$z = (\mathbf{W}^2)^T \mathbf{h} + b^2 = \begin{bmatrix} w_1^2 & w_2^2 & w_3^2 \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = w_1^2 h_1 + w_2^2 h_2 + w_3^2 h_3 + b^2$$

Gradient Computation

Let us consider the MSE loss $J = \frac{1}{2} \frac{1}{N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)})^2$

Let's compute the gradient with respect to the output weights w_1^2 (same as linear regression)

$$\begin{aligned}
 \frac{\partial J}{\partial w_1^2} &= \frac{\partial}{\partial w_1^2} \frac{1}{2 N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)})^2 = \text{Derivative is linear: derivative of sum is the sum of derivatives} \\
 &= \frac{1}{2 N_{Tr}} \sum_{p \in Tr} \frac{\partial}{\partial w_1^2} (t^{(p)} - z^{(p)})^2 = \text{Power rule: } \frac{d f(x)^n}{dx} = n(f(x))^{n-1} f'(x) \\
 &= \frac{1}{2 N_{Tr}} \sum_{p \in Tr} 2(t^{(p)} - z^{(p)}) \frac{\partial}{\partial w_1^2} (t^{(p)} - z^{(p)}) \quad z^{(p)} = \sigma(w_1^2 h_1 + w_2^2 h_2 + w_3^2 h_3 + b^2) = \sigma((\mathbf{w}^2)^T \mathbf{h}^p) \\
 &= \frac{1}{N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)}) \frac{\partial}{\partial w_1^2} (t^{(p)} - \sigma((\mathbf{w}^2)^T \mathbf{h}^p)) \quad \text{Sum rule: } \frac{\partial t^{(p)}}{\partial w_1^2} = 0 \\
 &= -\frac{1}{N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)}) \sigma'((\mathbf{w}^2)^T \mathbf{h}^p) h_1^p = \frac{dJ}{dz} \frac{dz}{da^2} \frac{da^2}{dw_1^2} \\
 &\quad \text{Computed during backward pass} \quad \text{Computed during forward pass}
 \end{aligned}$$

For sigmoid activation function

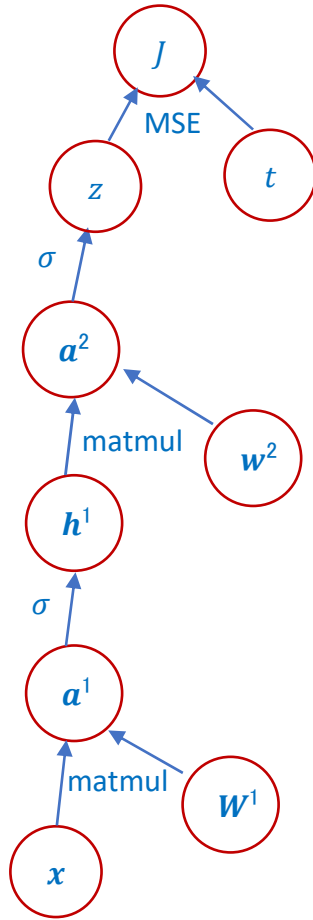
$$\sigma'((\mathbf{w}^2)^T \mathbf{h}^p) = \frac{d \sigma((\mathbf{w}^2)^T \mathbf{h}^p)}{d ((\mathbf{w}^2)^T \mathbf{h}^p)} = \sigma((\mathbf{w}^2)^T \mathbf{h}^p) (1 - \sigma((\mathbf{w}^2)^T \mathbf{h}^p))$$

$$\frac{\partial J}{\partial w_1^2} = -\frac{1}{N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)}) \sigma((\mathbf{w}^2)^T \mathbf{h}^p) (1 - \sigma((\mathbf{w}^2)^T \mathbf{h}^p)) h_1^p$$

Similar procedure for w_2^2, w_3^2 and for the bias values

Let us now compute the gradient w.r.t. the weights of the first layer, e.g. w_1^1

$$\begin{aligned}
 \frac{\partial J}{\partial w_{11}^1} &= \frac{\partial}{\partial w_{11}^1} \frac{1}{2 N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)})^2 = \text{Sum rule} \\
 &= \frac{1}{2 N_{Tr}} \sum_{p \in Tr} \frac{\partial}{\partial w_{11}^1} (t^{(p)} - z^{(p)})^2 = \text{Power rule and } \frac{\partial t^{(p)}}{\partial w_{11}^1} = 0 \text{ and sum rule} \\
 &= \frac{1}{2 N_{Tr}} \sum_{p \in Tr} 2(t^{(p)} - z^{(p)}) \frac{\partial}{\partial w_{11}^1} (-z^{(p)}) = z^{(p)} = \sigma((\mathbf{w}^2)^T \mathbf{h}^p) \text{ and chain rule} \\
 &= -\frac{1}{2 N_{Tr}} \sum_{p \in Tr} 2(t^{(p)} - z^{(p)}) \sigma'((\mathbf{w}^2)^T \mathbf{h}^p) \frac{\partial}{\partial w_{11}^1} ((\mathbf{w}^2)^T \mathbf{h}^p) (\mathbf{w}^2)^T \mathbf{h}^p = w_1^2 h_1 + w_2^2 h_2 + w_3^2 h_3 + b^2 \\
 &= -\frac{1}{2 N_{Tr}} \sum_{p \in Tr} 2(t^{(p)} - z^{(p)}) \sigma'((\mathbf{w}^2)^T \mathbf{h}^p) \left(\frac{\partial}{\partial w_{11}^1} w_1^2 h_1^{(p)} + \frac{\partial}{\partial w_{11}^1} w_2^2 h_2^{(p)} + \frac{\partial}{\partial w_{11}^1} w_3^2 h_3^{(p)} + \frac{\partial}{\partial w_{11}^1} b^2 \right) \\
 &\quad \text{Product rule + } h_1^p = \sigma(w_1^1 x + b_1^1) \\
 &= \frac{1}{N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)}) \sigma'((\mathbf{w}^2)^T \mathbf{h}^p) w_1^2 \frac{\partial}{\partial w_{11}^1} \sigma(w_1^1 x^{(p)} + b_1^1) = \text{Chain rule} \quad [w_{11}^1 \quad w_{12}^1] \\
 &= -\frac{1}{N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)}) \sigma'((\mathbf{w}^2)^T \mathbf{h}^p) w_1^2 \sigma'(w_1^1 x^{(p)} + b_1^1) \left(\frac{\partial}{\partial w_{11}^1} w_{11}^1 x_1^{(p)} + \frac{\partial}{\partial w_{11}^1} w_{12}^1 x_2^{(p)} \right. \\
 &\quad \left. + \frac{\partial}{\partial w_{11}^1} b_1^1 \right) = 0 \\
 &= -\frac{1}{N_{Tr}} \sum_{p \in Tr} (t^{(p)} - z^{(p)}) \sigma'((\mathbf{w}^2)^T \mathbf{h}^p) w_1^2 \sigma'(w_1^1 x^{(p)} + b_1^1) x_1^{(p)}
 \end{aligned}$$



Gradient Computation in matrix notation

Let us from now on consider a single example for simplicity.

Let's compute the gradient with respect to the output weights **vector** \mathbf{w}^2

$$\begin{aligned} \frac{\partial J}{\partial \mathbf{w}^2} &= \nabla_{\mathbf{w}^2} J = \frac{\partial J}{\partial z} \frac{dz}{da^2} \frac{\partial a^2}{\partial \mathbf{w}^2} = \nabla_z J \sigma'(a^2) \nabla_{\mathbf{w}^2} a^2 = [-(t - z)] [\sigma(a^2)(1 - \sigma(a^2))] \nabla_{\mathbf{w}^2} a^2 = \\ &= [-(t - z)] [\sigma(a^2)(1 - \sigma(a^2))] [h_1 \quad h_2 \quad h_3] = [-(t - z)] [\sigma(a^2)(1 - \sigma(a^2))] \mathbf{h}^T \end{aligned}$$

Notice the shape of $\frac{\partial J}{\partial \mathbf{w}^2}$. Since we defined the gradients as row vectors, we expect $\frac{\partial J}{\partial \mathbf{w}^2}$ to have a similar shape to \mathbf{w}^2 . Actually, it should have the same shape as $(\mathbf{w}^2)^T$, so live in $\mathbb{R}^{1 \times 3}$. This is obtained since a^2 is a number and \mathbf{w}^2 is a vector in \mathbb{R}^3 , so for our definition of gradient $\nabla_{\mathbf{w}^2} a^2$ have to be a row vector of size $\mathbb{R}^{1 \times 3}$.

Let's now compute the gradient with respect to the hidden weights **matrix** \mathbf{W}^1 .

We will require the computation of the gradient of vectors with respect to matrices, please see

$$\begin{aligned} &\mathbb{R} \quad \mathbb{R} \quad \mathbb{R}^{1 \times 3} \quad \mathbb{R}^{3 \times 3} \quad \mathbb{R}^{3 \times (3 \times 2)} \\ \frac{\partial J}{\partial \mathbf{W}^1} &= \frac{\partial J}{\partial z} \frac{dz}{da^2} \frac{\partial a^2}{\partial \mathbf{h}} \frac{d\mathbf{h}}{da^1} \frac{\partial a^1}{\partial \mathbf{W}^1} \end{aligned}$$

Let us study the dimensions of each term.

$\frac{\partial J}{\partial z} \in \mathbb{R}$ since both J and z are numbers.

$\frac{dz}{da^2} \in \mathbb{R}$ since both z and a^2 are numbers. This is the derivative of the output after applying the output activation function w.r.t. the output pre-activation.

$\frac{\partial a^2}{\partial \mathbf{h}} \in \mathbb{R}^{1 \times 3}$ is a row vector since a^2 is a number and \mathbf{h} is a vector. In particular, it will be defined as:

$$\frac{\partial a^2}{\partial \mathbf{h}} = \begin{bmatrix} \frac{\partial a^2}{\partial h_1} & \frac{\partial a^2}{\partial h_2} & \frac{\partial a^2}{\partial h_3} \end{bmatrix} = [w_1^2 \quad w_2^2 \quad w_3^2]$$

$\frac{d\mathbf{h}}{da^1} \in \mathbb{R}^{3 \times 3}$ is a (Jacobian) matrix since both \mathbf{h} and \mathbf{a}^1 are vectors. This matrix however has a particularity: it is diagonal since h_i only depends on a_i^1 (pointwise non-linearity). Thus:

$$\begin{aligned} \frac{d\mathbf{h}}{da^1} &= \begin{bmatrix} \nabla_{\mathbf{a}^1} h_1 \\ \nabla_{\mathbf{a}^1} h_2 \\ \nabla_{\mathbf{a}^1} h_3 \end{bmatrix} = \begin{bmatrix} \frac{\partial h_1}{\partial a_1^1} & \frac{\partial h_1}{\partial a_2^1} & \frac{\partial h_1}{\partial a_3^1} \\ \frac{\partial h_2}{\partial a_1^1} & \frac{\partial h_2}{\partial a_2^1} & \frac{\partial h_2}{\partial a_3^1} \\ \frac{\partial h_3}{\partial a_1^1} & \frac{\partial h_3}{\partial a_2^1} & \frac{\partial h_3}{\partial a_3^1} \end{bmatrix} = \begin{bmatrix} \frac{\partial h_1}{\partial a_1^1} & 0 & 0 \\ 0 & \frac{\partial h_2}{\partial a_2^1} & 0 \\ 0 & 0 & \frac{\partial h_3}{\partial a_3^1} \end{bmatrix} \\ &= \begin{bmatrix} \sigma'(a_1^1) & 0 & 0 \\ 0 & \sigma'(a_2^1) & 0 \\ 0 & 0 & \sigma'(a_3^1) \end{bmatrix} \end{aligned}$$

Let us now focus on the term $\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1}$. $\mathbf{a}^1 \in \mathbb{R}^3$ is a vector of size 3. $\mathbf{W}^1 \in \mathbb{R}^{3 \times 2}$ is the matrix of weights of the first layer. $\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1}$ is then a tensor in $\mathbb{R}^{3 \times (3 \times 2)}$.

In fact, by definition, the gradient is the collection of partial derivatives:

$$\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1} = \begin{bmatrix} \frac{\partial a_1^1}{\partial \mathbf{W}^1} \\ \frac{\partial a_2^1}{\partial \mathbf{W}^1} \\ \frac{\partial a_3^1}{\partial \mathbf{W}^1} \end{bmatrix}, \text{ where } \frac{\partial a_i^1}{\partial \mathbf{W}^1} \in \mathbb{R}^{1 \times (3 \times 2)}$$

Let us write down explicitly $a_i^1 = w_{i1}^1 x_1 + w_{i2}^1 x_2 + b_i^1$

We know that $\frac{\partial a_i^1}{\partial w_{ij}^1} = x_j$, that is $\frac{\partial a_i^1}{\partial w_{i1}^1} = x_1$ and $\frac{\partial a_i^1}{\partial w_{i2}^1} = x_2$

We can also compute the partial derivative of a_i^1 with respect to a row of \mathbf{W}^1 , that is:

$$\frac{\partial a_i^1}{\partial \mathbf{w}_{i,:}^1} = \mathbf{x}^T \in \mathbb{R}^{1 \times 2}. \text{ Moreover, the tensor is sparse since } \frac{\partial a_i^1}{\partial w_{k \neq i,:}^1} = \mathbf{0}^T \in \mathbb{R}^{1 \times 2}$$

Let us now look more closely at $\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1}$

$$\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1} = \begin{bmatrix} \frac{\partial a_1^1}{\partial w^1} \\ \frac{\partial a_2^1}{\partial w^1} \\ \frac{\partial a_3^1}{\partial w^1} \end{bmatrix}, \text{ and } \frac{\partial a_i^1}{\partial \mathbf{W}^1} = \begin{bmatrix} \frac{\partial a_i^1}{\partial w_{11}^1} & \frac{\partial a_i^1}{\partial w_{12}^1} \\ \frac{\partial a_i^1}{\partial w_{21}^1} & \frac{\partial a_i^1}{\partial w_{22}^1} \\ \frac{\partial a_i^1}{\partial w_{31}^1} & \frac{\partial a_i^1}{\partial w_{32}^1} \end{bmatrix} = \begin{bmatrix} \frac{\partial a_i^1}{\partial w_{1,:}^1} \\ \frac{\partial a_i^1}{\partial w_{2,:}^1} \\ \frac{\partial a_i^1}{\partial w_{3,:}^1} \end{bmatrix}$$

Then we have that

$$\frac{\partial a_1^1}{\partial \mathbf{W}^1} = \begin{bmatrix} x_1 & x_2 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}, \frac{\partial a_2^1}{\partial \mathbf{W}^1} = \begin{bmatrix} 0 & 0 \\ x_1 & x_2 \\ 0 & 0 \end{bmatrix} \text{ and } \frac{\partial a_3^1}{\partial \mathbf{W}^1} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ x_1 & x_2 \end{bmatrix}$$

We can now try start to put things together and compute our gradients:

$$\frac{dJ}{d\mathbf{W}^1} = \frac{dJ}{dz} \frac{dz}{da^2} \frac{da^2}{dh} \frac{dh}{da^1} \frac{da^1}{d\mathbf{W}^1}$$

$$\begin{aligned} \frac{dJ}{d\mathbf{W}^1} &= \nabla_z J \sigma'(a^2) \nabla_h a^2 \frac{dh}{da^1} \frac{da^1}{d\mathbf{W}^1} = [-(t-z)][\sigma'(a^2)] \\ &= (w^2)^T \begin{bmatrix} \sigma'(a_1^1) & 0 & 0 \\ 0 & \sigma'(a_2^1) & 0 \\ 0 & 0 & \sigma'(a_3^1) \end{bmatrix} \frac{da^1}{d\mathbf{W}^1} \end{aligned}$$

We can simplify this multiplication since the product of a vector with a diagonal matrix is equivalent to the element-wise multiplication (indicated with \odot) with the diagonal vector.

Thus:

$$\begin{aligned} \frac{dJ}{d\mathbf{W}^1} &= [-(t-z)][\sigma'(a^2)] [w_1^2 \quad w_2^2 \quad w_3^2] \odot [\sigma'(a_1^1) \quad \sigma'(a_2^1) \quad \sigma'(a_3^1)] \frac{da^1}{d\mathbf{W}^1} \\ &= [-(t-z)][\sigma'(a^2)] [w_1^2 \sigma'(a_1^1) \quad w_2^2 \sigma'(a_2^1) \quad w_3^2 \sigma'(a_3^1)] \frac{da^1}{d\mathbf{W}^1} \end{aligned}$$

$\mathbb{R}^{1 \times 3}$

Finally, we have only the last multiplication with the tensor $\frac{da^1}{d\mathbf{W}^1}$ to take care of.

Even though the mathematics of tensors might be hard, we just need to compute a product that resembles matrix multiplication.

Luckily, there is a very simple “trick” to compute our product, and it requires knowing only how to do matrix multiplications, so we have hope!

Example of Vector/Tensor multiplication

$$\mathbf{a} = [a_1 \quad a_2 \quad a_3], \mathbf{B} = \begin{bmatrix} \mathbf{B}^1 \\ \mathbf{B}^2 \end{bmatrix}, B_i = \begin{bmatrix} b_{11}^i & b_{12}^i \\ b_{21}^i & b_{22}^i \\ b_{31}^i & b_{32}^i \end{bmatrix}$$

From the reasoning above, we know that for this example the shape we need is $\mathbf{aB} \in \mathbb{R}^{2 \times 2}$

Approach 1: compute the multiplications of \mathbf{a} with the various \mathbf{B}^i and then “collate” the results. In fact $\mathbf{aB}^i \in \mathbb{R}^{1 \times 2}$, so the shapes match.

$$\text{In other words, } \mathbf{aB} = \begin{bmatrix} \mathbf{aB}^1 \\ \mathbf{aB}^2 \end{bmatrix}$$

Approach 2: we can obtain the same result first flattening the tensor in a matrix, and then re-shaping it to account for the “lost” dimension.

We define $\bar{\mathbf{B}} = \begin{bmatrix} b_{11}^1 & b_{12}^1 & b_{11}^2 & b_{12}^2 \\ b_{21}^1 & b_{22}^1 & b_{21}^2 & b_{22}^2 \\ b_{31}^1 & b_{31}^1 & b_{31}^2 & b_{32}^2 \end{bmatrix}$. We can now easily compute $\mathbf{a}\bar{\mathbf{B}} \in \mathbb{R}^{1 \times 4}$. We just have to reshape this vector back into our desired shape:

$$\mathbf{a}\bar{\mathbf{B}} = [(ab)_1 \quad (ab)_2 \quad (ab)_3 \quad (ab)_4], \quad \mathbf{a}\mathbf{B} = \begin{bmatrix} (ab)_1 & (ab)_2 \\ (ab)_3 & (ab)_4 \end{bmatrix}$$

We are finally ready. As detailed before, $\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1} = \begin{bmatrix} \frac{\partial a_1^1}{\partial \mathbf{W}^1} \\ \frac{\partial a_2^1}{\partial \mathbf{W}^1} \\ \frac{\partial a_3^1}{\partial \mathbf{W}^1} \end{bmatrix}$ is a tensor of shape $3 \times 3 \times 2$. Let's

use the flattening approach, obtaining $\overline{\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1}} \in \mathbb{R}^{3 \times 6}$. We can compute the flattened gradient $\overline{\frac{\partial J}{\partial \mathbf{W}^1}} \in \mathbb{R}^{1 \times 6}$ as:

$$\begin{aligned} \overline{\frac{\partial J}{\partial \mathbf{W}^1}} &= [-(t-z)][\sigma'(a^2)] [w_1^2 \sigma'(a_1^1) \quad w_2^2 \sigma'(a_2^1) \quad w_3^2 \sigma'(a_3^1)] \begin{bmatrix} x_1 & x_2 & 0 & 0 & 0 & 0 \\ 0 & 0 & x_1 & x_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & x_1 & x_2 \end{bmatrix} \\ &= [-(t-z)\sigma'(a^2)w_1^2 \sigma'(a_1^1)x_1 \quad \dots \quad \dots \quad \dots \quad -(t-z)\sigma'(a^2)w_3^2 \sigma'(a_3^1)x_2] \end{aligned}$$

We can now reshape the matrix to obtain our gradient $\frac{\partial J}{\partial \mathbf{W}^1} \in \mathbb{R}^{3 \times 2}$

$$\frac{\partial J}{\partial \mathbf{W}^1} = \begin{bmatrix} -(t-z)\sigma'(a^2)w_1^2 \sigma'(a_1^1)x_1 & -(t-z)\sigma'(a^2)w_1^2 \sigma'(a_1^1)x_2 \\ -(t-z)\sigma'(a^2)w_2^2 \sigma'(a_2^1)x_1 & -(t-z)\sigma'(a^2)w_2^2 \sigma'(a_2^1)x_2 \\ -(t-z)\sigma'(a^2)w_3^2 \sigma'(a_3^1)x_1 & -(t-z)\sigma'(a^2)w_3^2 \sigma'(a_3^1)x_2 \end{bmatrix}$$

Now that we know the result, we can notice that in our case, since the pre-activation a_i^1 only depends on the subset of weights connecting the input to that particular neuron, we have sparsity and regularities in $\frac{\partial \mathbf{a}^1}{\partial \mathbf{W}^1}$. We can avoid the computation with the tensor entirely, noticing that the result corresponds to a simpler outer product:

$$\frac{\partial J}{\partial \mathbf{W}^1} = ([-(t-z)][\sigma'(a^2)] (\mathbf{w}^2)^T \odot [\sigma'(a_1^1) \quad \sigma'(a_2^1) \quad \sigma'(a_3^1)])^T \mathbf{x}^T$$

Notice that we are multiplying the transpose of the first term that has shape 3×1 with the transpose of the input that has shape 1×2 , thus obtaining a matrix of shape 3×2 that is the correct shape.

Extension to multiple examples in matrix notation

Neural Network implementations work directly with tensors, even when using multiple examples to compute the gradient (as in the first part of this lecture notes).

It is possible to exploit the matrix notation also in this case.

$$J = \frac{1}{N} [J^{(1)} \quad \dots \quad J^{(N)}] \begin{bmatrix} 1 \\ \dots \\ 1 \end{bmatrix}$$

Where $J^{(i)}$ is defined as detailed above for a single example.

In practice, multiple examples are managed with an added dimension to all the tensors discussed before. Still, in this case, the mathematics becomes more difficult due to the newly added dimension that would generate tensors with 4 dimensions.

From Gradient computation to Backpropagation

Let's consider the backpropagation algorithm, in particular the backward pass. We will not consider the regularization terms for this exercise.

Algorithm 6.4 Backward computation for the deep neural network of algorithm 6.3, which uses in addition to the input \mathbf{x} a target \mathbf{y} . This computation yields the gradients on the activations $\mathbf{a}^{(k)}$ for each layer k , starting from the output layer and going backwards to the first hidden layer. From these gradients, which can be interpreted as an indication of how each layer's output should change to reduce error, one can obtain the gradient on the parameters of each layer. The gradients on weights and biases can be immediately used as part of a stochastic gradient update (performing the update right after the gradients have been computed) or used with other gradient-based optimization methods.

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l-1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g}$$

$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top}$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

end for

Exercise: show that backpropagation correctly computes the gradients for the 2-layer network considered before.

Solution:

In our notation, we referred to \hat{y} as z .

We have 2 layers in total (1 hidden layer and 1 output layer). Thus, $l=2$.

Before the for loop, we initialize $g \leftarrow \nabla_z J$

First iteration of for loop: $k=2$

We update g with the element-wise multiplication with $f'(a^2)$. In our case $f = \sigma$ and $a^2 \in \mathbb{R}$, thus the element-wise product corresponds to the classical product that we have in our derivation.

Now we have $g \leftarrow \nabla_z J \sigma'(a^2)$

This is exactly the update for the bias (we didn't compute it before, it was left as exercise).

The algorithm now computes $\nabla_{w^2} J = \frac{\partial J}{\partial w^2} = g (h^1)^T$

In our notation, we referred to h^1 simply as h .

Thus, expanding the equation, the backprop algorithm computes

$$\nabla_{w^2} J = \frac{\partial J}{\partial w^2} = g (h)^T = \nabla_z J \sigma'(a^2) (h)^T$$

That is exactly the gradient we computed in matrix notation.

Before finishing the for loop, the algorithm updates g multiplying it by the transpose of the weights of the current layer. Notice that, according to our computation, this is the third term of $\frac{dJ}{dw^1}$. Notice also that the backpropagation algorithm as presented in the book uses another convention for gradients and Jacobian matrices, that makes more complicated the operations (in fact, the multiplication with $(W^{(k)})^T$ is on the left, while our notation makes everything simpler and we can just compute the multiplications in order).

Thus, g is updated as $g \leftarrow g (W^{(k)})^T$. In our case w^2 is a vector, and with the aforementioned consideration we have that:

$$g \leftarrow g \cdot (w^2)^T = \nabla_z J \sigma'(a^2) (w^2)^T$$

Second iteration of for loop: $k=2$

Similarly as before, we update g with the element-wise multiplication with $f'(a^1)$, that in our case is the sigmoid function:

$$g \leftarrow g \odot f'(a^1) = \nabla_z J \sigma'(a^2) (w^2)^T \odot \sigma'(a^1)$$

Notice that this term is exactly the fourth term in our derivation of $\frac{dJ}{dw^1}$.

Again, similarly as before, this term is the correct update for the bias, while to compute $\nabla_{w^1} J$ we have to multiply g with the transpose of $h^{k-1} = h^{1-1} = x$.

Again for the different notations, the multiplication with x^T corresponds to our "trick" of using the outer product. This was the last missing term of $\frac{dJ}{dw^1}$, that is thus computed correctly.

Notice how the variable g allows not to re-compute the first terms of the gradients (this is where the computational efficiency arises).

However, we are required to store both h and a for all the layers (see forward part of backpropagation)