

A decorative teal border with a wavy, scalloped edge runs vertically along the left side of the slide.

IMPLEMENTING AND TRAINING NEURAL NETWORKS(CONT.)

STFANO MELACCI

UNIVERSITY OF SIENA

DERIVATIVES: SCALARS, VECTORS, MATRICES, TENSORS

- Let's consider a scalar valued function (we will always assume values to be in the set of real numbers)

- $f(x)$, where x is a scalar

- $f(\mathbf{x})$, where \mathbf{x} is a vector (n elements)

- $f(X)$, where X is a matrix ($n \times m$)

- $f(X)$, where X is a “tensor” ($n \times m \times c$ – or any other)

$$\mathbf{x} = [x_1, \dots, x_n]^T$$

$$X = \begin{bmatrix} X_{11}, & \dots, & X_{1n} \\ \dots, & \dots, & \dots \\ X_{m1}, & \dots, & X_{mn} \end{bmatrix}$$

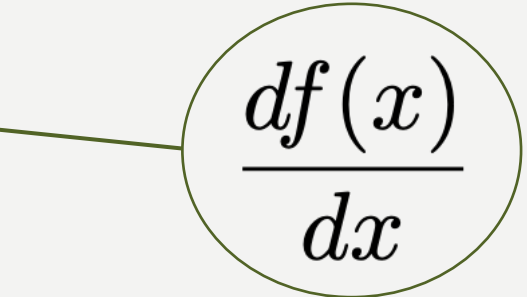
$$X = \begin{bmatrix} [X_{111}, \dots, X_{11c}] & \dots, & [X_{1n1}, \dots, X_{1nc}] \\ \dots, & \dots, & \dots \\ [X_{m11}, \dots, X_{m1c}], & \dots, & [X_{mn1}, \dots, X_{mnc}] \end{bmatrix}$$

Assumption: if we imagine to assing an unique numerical ID to each element of X , then the last dimension of X is the one on which elements are “contiguous” – explain this!

DERIVATIVES: SCALARS, VECTORS, MATRICES, TENSORS

- Let's consider a scalar valued function (we will always assume values to be in the set of real numbers)
 - $f(x)$, where x is a scalar
 - $f(\mathbf{x})$, where \mathbf{x} is a vector (n elements)
 - $f(X)$, where X is a matrix ($m \times n$)
 - $f(X)$, where X is a “tensor” ($m \times n \times c$)

most common
and straightforward
case


$$\frac{df(x)}{dx}$$

$$\frac{df(\mathbf{x})}{d\mathbf{x}}$$

$$\frac{df(X)}{dX}$$

DERIVATIVES W.R.T. VECTOR/MATIX

$$\frac{df(\mathbf{x})}{d\mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \quad \dots, \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right]^T$$

Denominator
Layout

$$\frac{df(X)}{dX} = \begin{bmatrix} \frac{\partial f(X)}{\partial X_{11}}, & \dots, & \frac{\partial f(X)}{\partial X_{1n}} \\ \frac{\partial f(X)}{\partial X_{21}}, & \dots, & \frac{\partial f(X)}{\partial X_{2n}} \\ \dots, & \dots, & \dots \\ \frac{\partial f(X)}{\partial X_{m1}}, & \dots, & \frac{\partial f(X)}{\partial X_{mn}} \end{bmatrix}$$

DERIVATIVES W.R.T. TENSOR: GENERAL CASE


- The derivative of a scalar function w.r.t. to any kind of input tensorial structure X is a tensor with the same dimensions of X , where each element is the partial derivative of the function w.r.t. to the corresponding element of X
 - WARNING: I will frequently relax the formality of several descriptions (and notation) with the purpose of making them more accessible

Denominator
Layout

$$\frac{df(X)}{dX} = \left[\frac{\partial f(X)}{\partial X_{ijk}}, \quad \begin{array}{l} i = 1, \dots, m \\ j = 1, \dots, n \\ k = 1, \dots, c \end{array} \right]$$

It generalizes all the cases considered so far!
(it holds also for 4D, 5D, whateverD tensors)

DERIVATIVES: LET'S RAISE THE BAR

- Let's consider a **vector-valued function** that outputs a vector with a elements
 - $f(X)$
 - What if we also consider a function that outputs an $a \times b$ matrix?
 - $F(X)$
 - What if we consider a function that returns a 3D structure of size $a \times b \times u$?
 - $F(X)$
 - *most generic case*
- 

DERIVATIVES: LET'S RAISE THE BAR

$$\frac{d\mathbf{f}(x)}{dx} = \left[\frac{\partial f_1(x)}{\partial x}, \quad \dots, \quad \frac{\partial f_a(x)}{\partial x} \right]$$

Denominator
Layout

$$\frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1}, & \dots, & \frac{\partial f_a(\mathbf{x})}{\partial x_1} \\ \dots, & \dots, & \dots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_n}, & \dots, & \frac{\partial f_a(\mathbf{x})}{\partial x_n} \end{bmatrix}$$

Jacobian Matrix

DERIVATIVES: LET'S RAISE THE BAR

Denominator
Layout

$$\frac{df(X)}{dX} = \begin{bmatrix} \frac{\partial f(X)}{\partial X_{11}}, & \cdots, & \frac{\partial f(X)}{\partial X_{1n}} \\ \cdots, & \cdots, & \cdots \\ \frac{\partial f(X)}{\partial X_{m1}}, & \cdots, & \frac{\partial f(X)}{\partial X_{mn}} \end{bmatrix}$$

vector

DERIVATIVES: LET'S RAISE THE BAR

- The derivative of a vectorial function w.r.t. to any kind of input tensorial structure X is a tensor whose “first” dimensions are the same dimensions of X
 - They index partial derivatives of the function w.r.t. to the corresponding elements of X
- The “last” dimensions are the ones of $F(X)$
 - They are used to index the output elements of $F(X)$

$$X : m \times n \times c$$

$$F(X) : a \times b \times u$$

$$\frac{dF(X)}{dX} : \underbrace{m \times n \times c}_{\text{shape of } X} \times \underbrace{a \times b \times u}_{\text{shape of } F(X)}$$

**Denominator
Layout**

$$\frac{\partial f(X)}{\partial X} = \left[\frac{\partial f_{stz}(X)}{\partial X_{ijk}}, \begin{matrix} i = 1, \dots, m \\ j = 1, \dots, n \\ k = 1, \dots, c \\ s = 1, \dots, a \\ t = 1, \dots, b \\ z = 1, \dots, u \end{matrix} \right]$$

CHAIN RULE

- Before diving into further tensor-related stuff, let's go back to a simple concept
- Whenever we are given a composite function $f(g(x))$, then the derivative of f with respect to x can be written as follows

$$z = g(x)$$

$$y = f(g(x)) = f(z)$$

$$\boxed{\frac{df(g(x))}{dx}} = \frac{dy}{dx} = \boxed{\frac{dy}{dz}} \boxed{\frac{dz}{dx}}$$

The whole derivative has been decomposed as the product of the derivatives of two “computational blocks/nodes”

CHAIN RULE

- Before diving into further tensor-related stuff, let's go back to a simple concept
- Whenever we are given a composite function $f(g(x))$, then the derivative of f with respect to x can be written as follows

$$z = g(x)$$

$$y = f(g(x)) = f(z)$$

$$\boxed{\frac{df(g(x))}{dx}} = \frac{dy}{dx} = \boxed{\frac{dy}{dz}} \boxed{\frac{dz}{dx}}$$

Computational block/node 1:
 $y = f(z)$

(we differentiate the block with respect to its input z)

CHAIN RULE

- Before diving into further tensor-related stuff, let's go back to a simple concept
- Whenever we are given a composite function $f(g(x))$, then the derivative of f with respect to x can be written as follows

$$z = g(x)$$

$$y = f(g(x)) = f(z)$$

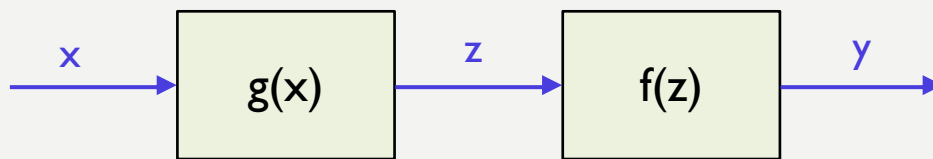
$$\boxed{\frac{df(g(x))}{dx}} = \frac{dy}{dx} = \boxed{\frac{dy}{dz}} \boxed{\frac{dz}{dx}}$$

Computational block/node 2:
 $z = g(x)$

(we differentiate the block with respect to its input x)

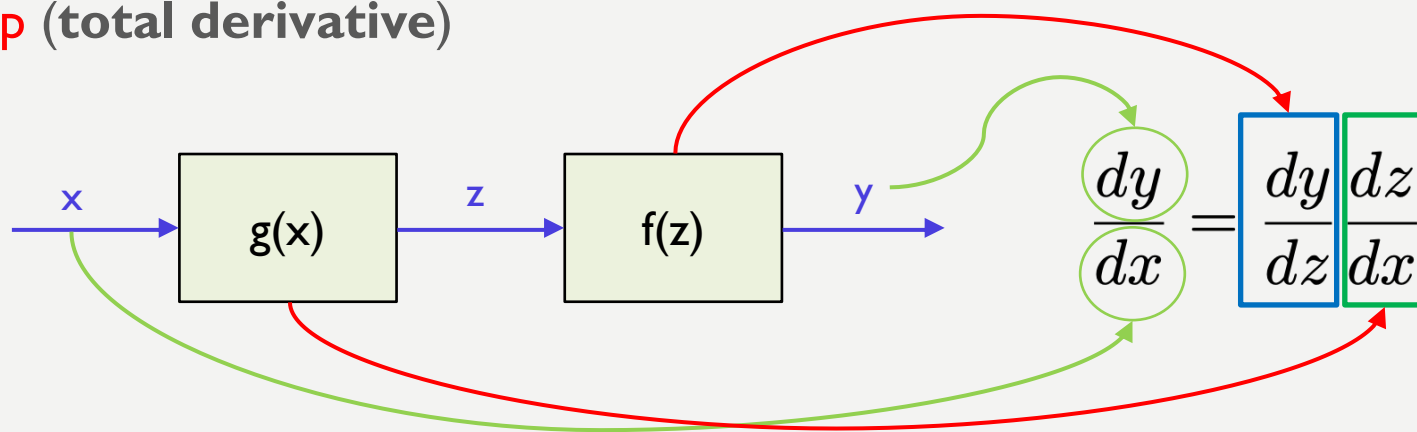
COMPUTATIONAL GRAPHS

- A computational graph describes a mathematical expression using a graph-based representation (directed graph) in which **nodes** are *operations* and **edges** are *data*
 - A node could be about a very low-level operation (a sum) or about a more "complicated" function - it is a design choice
 - In the case of the previous slides, $f(g(x))$, we get:



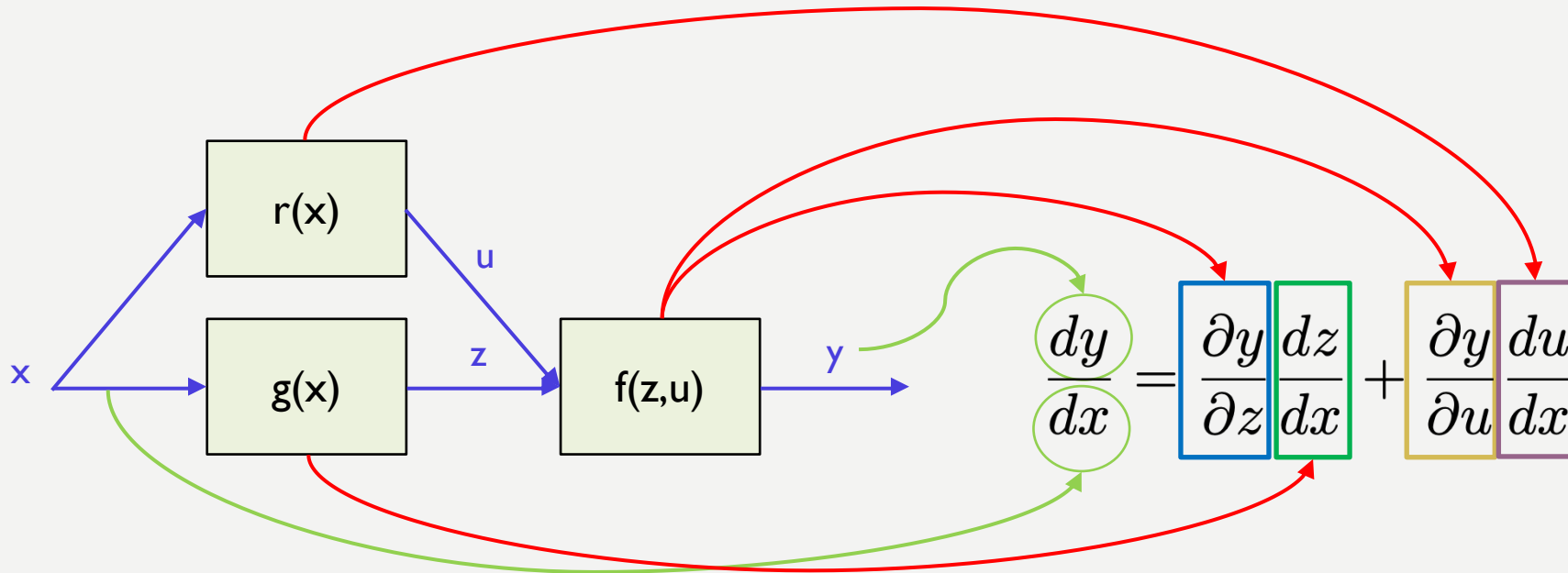
COMPUTATIONAL GRAPHS

- We can compute the derivative of some data **a** (edge) with respect to some other data **b** (edge) using the *chain rule*, following the path that connects **a** to **b**, and making products of the derivatives of each node we encounter, with respect to the data in the path
 - If there is not path, then the derivative is trivially zero
 - If there are **multiple paths** going from **a** to **b**, then the derivatives computed over all the paths **must be summed up** (total derivative)



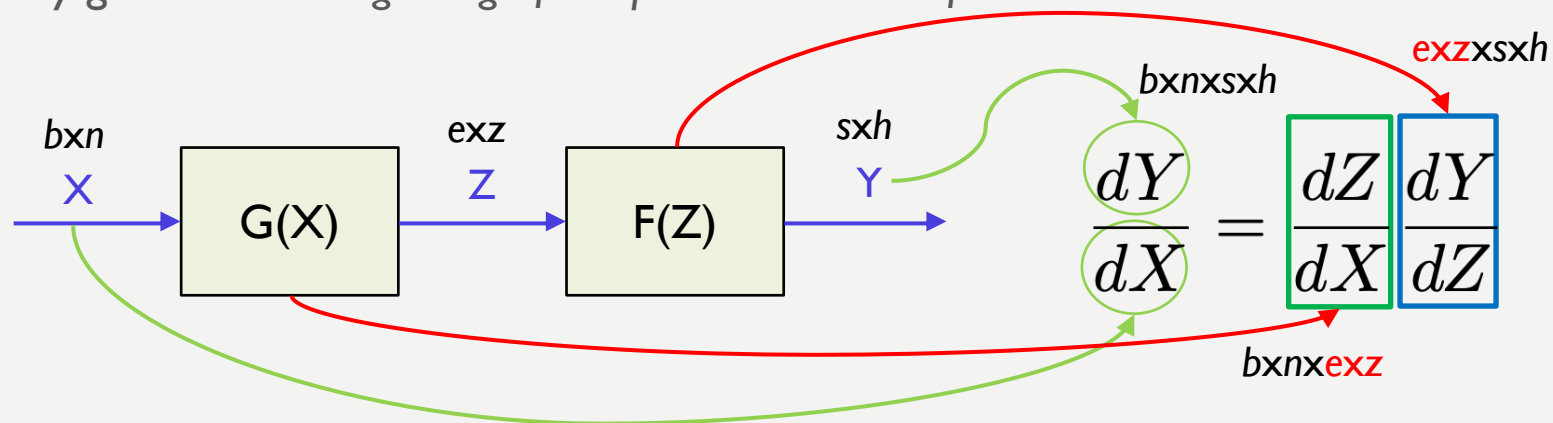
COMPUTATIONAL GRAPHS

- Another example (recall: If there are **multiple paths** going from **a** to **b**, then the derivatives computed over all the paths **must be summed up**)



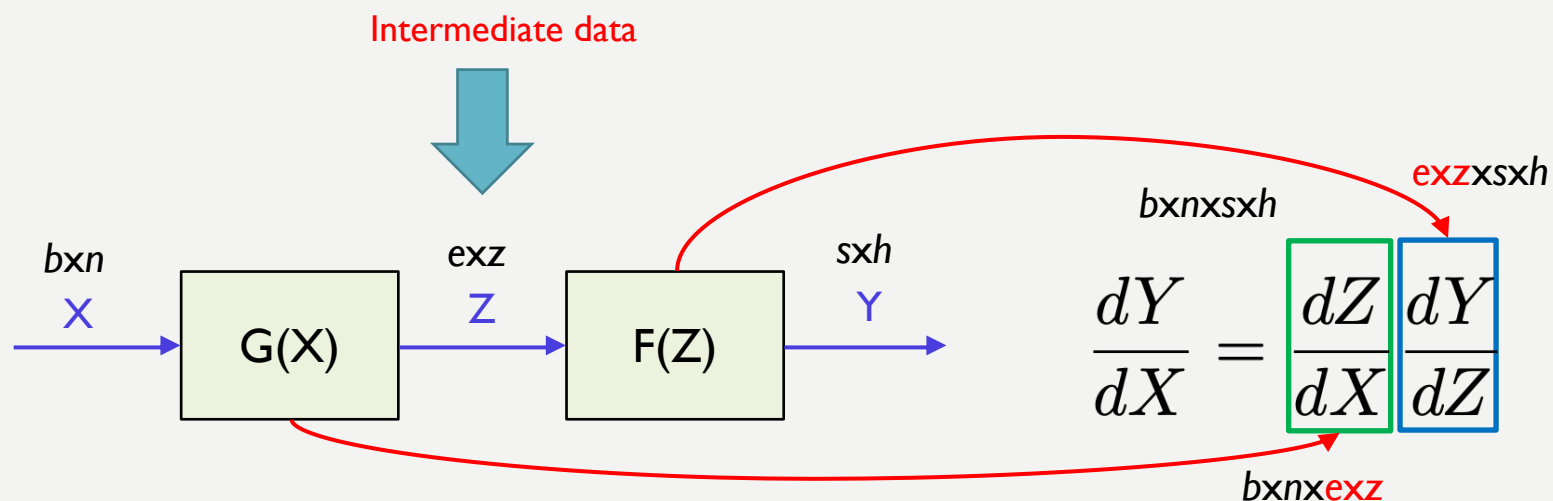
COMPUTATIONAL GRAPHS

- Let's consider computational nodes that are about vector functions
- Input and output data are now tensors
 - Check the following example: derivatives are combined in the **opposite order** with respect to the previous example(s), since they are now **tensors** and we assumed a **denominator layout**
 - They go from the *beginning of the path to the end of it*



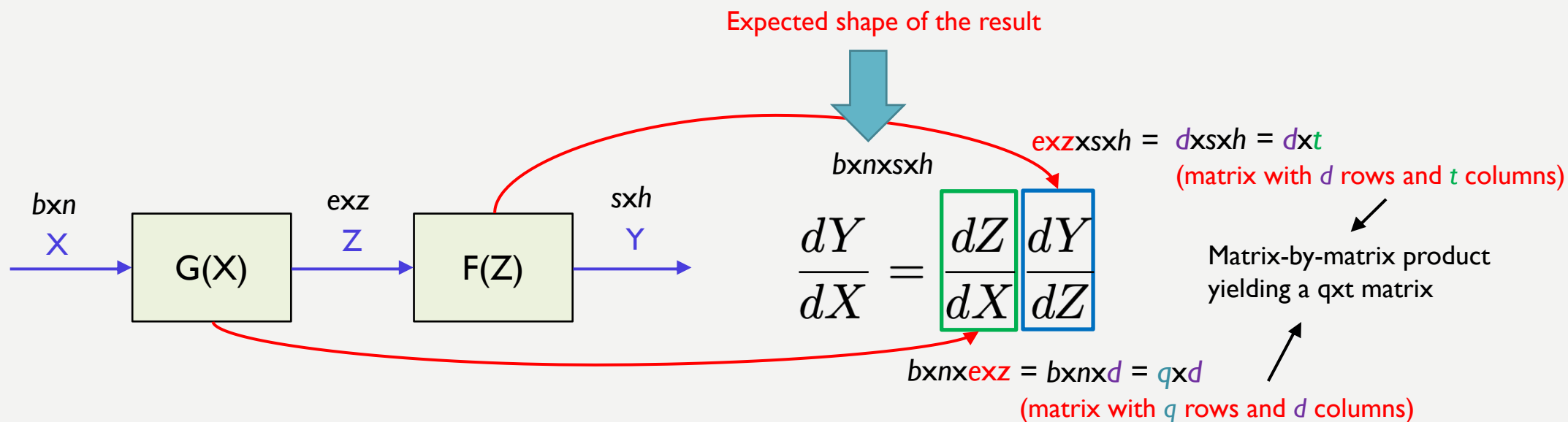
COMPUTATIONAL GRAPHS

- For **each pair** of nodes in path we have some **intermediate** data (Z in this example).
- The dimensions of these intermediate data are the ones on which the chain rule must “accumulate” the derivatives of the two nodes (so the dimensions of Z in this example, that are exz).



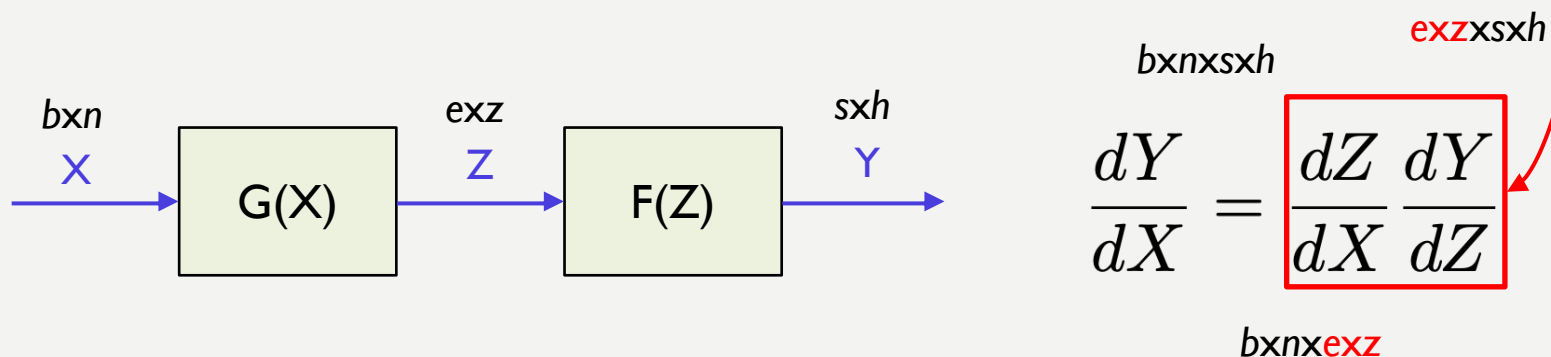
COMPUTATIONAL GRAPHS

- In order to qualitatively understand the way the derivatives are “accumulated”, reshape the tensor derivatives as shown in the example, then perform a matrix-by-matrix product
- Finally, we restore the expected shape of the result by another reshaping operation



COMPUTATIONAL GRAPHS

- Notice that when we apply the chain rule, we are performing **new operations**
- In turn, these new operations could be described with **their own** computational graph ☺



GOING BACK TO BACKPROPAGATION

- *The derivative of the objective function w.r.t. the weights W of the network is what is computed by means of Backpropagation*

$$\frac{\partial \text{obj}(f(X, W), Y)}{\partial W}$$

X, Y

Training set (inputs, targets)

$f(X, W)$

Output of the network on the training set

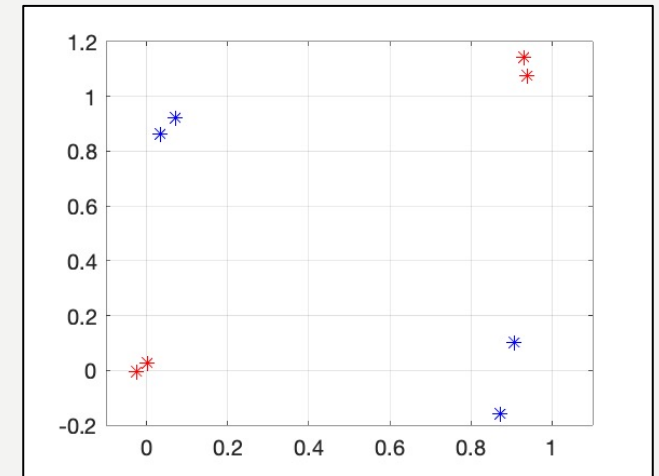
$\text{obj}(f(X, W), Y)$

Loss function (objective)

- If we describe the **forward** operations performed by a neural network, up to the loss computation, using a **computational graph**, then we can apply the **chain rule** **over all the possible paths connecting the weights of the network to the objective function**
 - Summing up the results, we get the outcome of Backpropagation 😊
 - Notice that what does matter is the use of the chain rule: the computational graph is just a formalism/tool we use to describe what is going on, it is not playing any other special roles

BACKPROPAGATION FROM SCRATCH

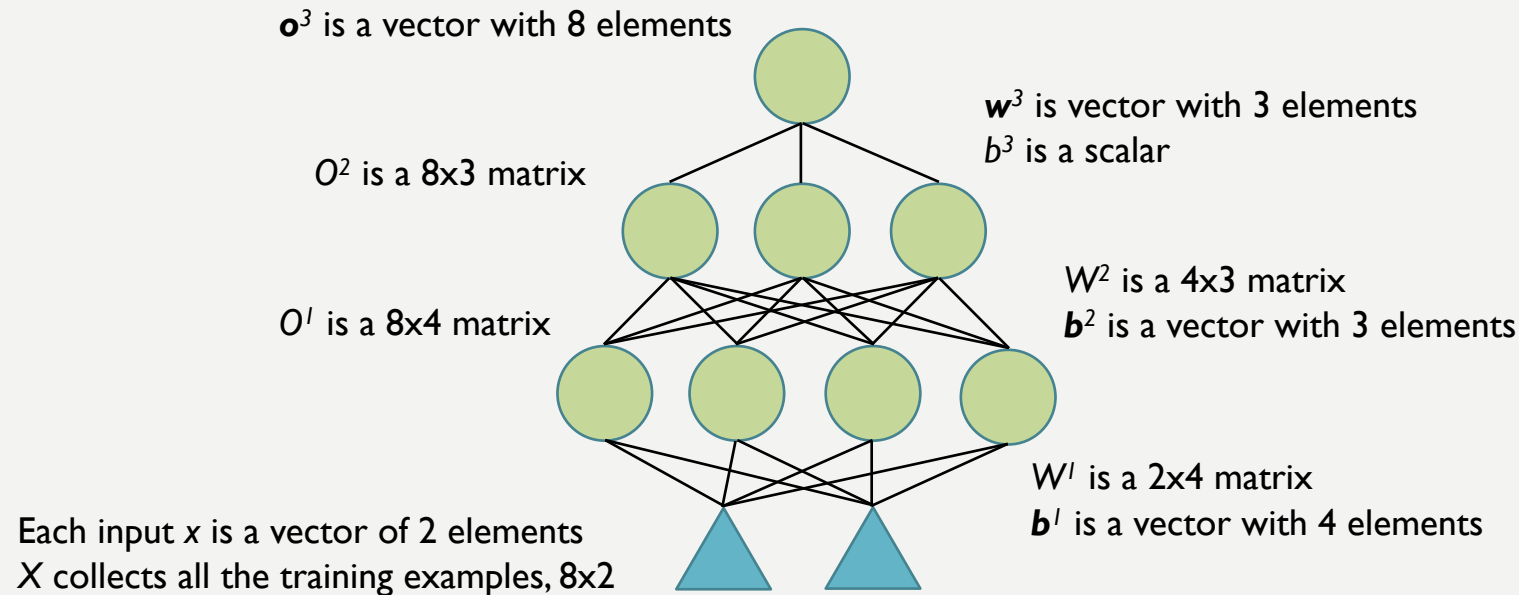
- We consider a training set composed of 8 bidimensional examples (XOR problem)
 - 4 examples of class label 0 (blue), 4 examples of class label 1 (red)
 - Training examples are collected into matrix X (one row per example, 8x2 matrix of floating-point numbers)
 - Binary targets of the training examples are collected in vector y (one target per example, 8 element-vector composed of zeros and ones)
- **Goal: train a neural network to learn from these data**
 - We will reimplement Backpropagation from scratch, using Numpy!
 - Batch-mode learning (forward on the whole training set, backward, ...)



X, y

BACKPROPAGATION FROM SCRATCH

- We consider the following feed-forward network with **sigmoidal** activation functions
 - **Weights** are collected in W^1, W^2, \mathbf{w}^3 , while **biases** (not shown in the picture) are collected in $\mathbf{b}^1, \mathbf{b}^2, b^3$



$$O^1 = \sigma(XW^1 + (\mathbf{b}^1 \mathbf{1}^T)^T)$$

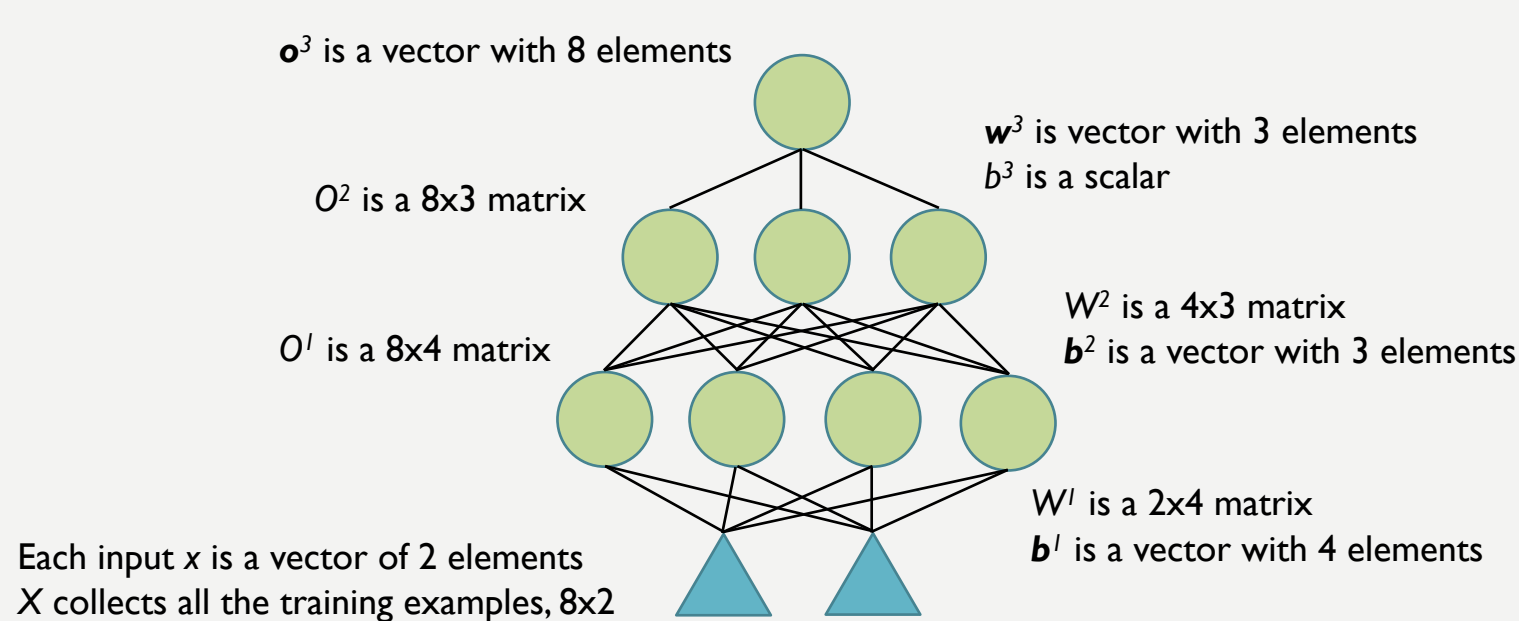
$$O^2 = \sigma(O^1 W^2 + (\mathbf{b}^2 \mathbf{1}^T)^T)$$

$$\mathbf{o}^3 = \sigma(O^2 \mathbf{w}^3 + \mathbf{1} b^3)$$

$$\text{loss} = \|\mathbf{o}^3 - \mathbf{y}\|^2$$

BACKPROPAGATION FROM SCRATCH

- We consider the following feed-forward network with **sigmoidal** activation functions
 - **Weights** are collected in W^1, W^2, w^3 , while **biases** (not shown in the picture) are collected in b^1, b^2, b^3



Vectors of ones (8 elements)

$$O^1 = \sigma(XW^1 + (\mathbf{b}^1 \mathbf{1}^T)^T)$$

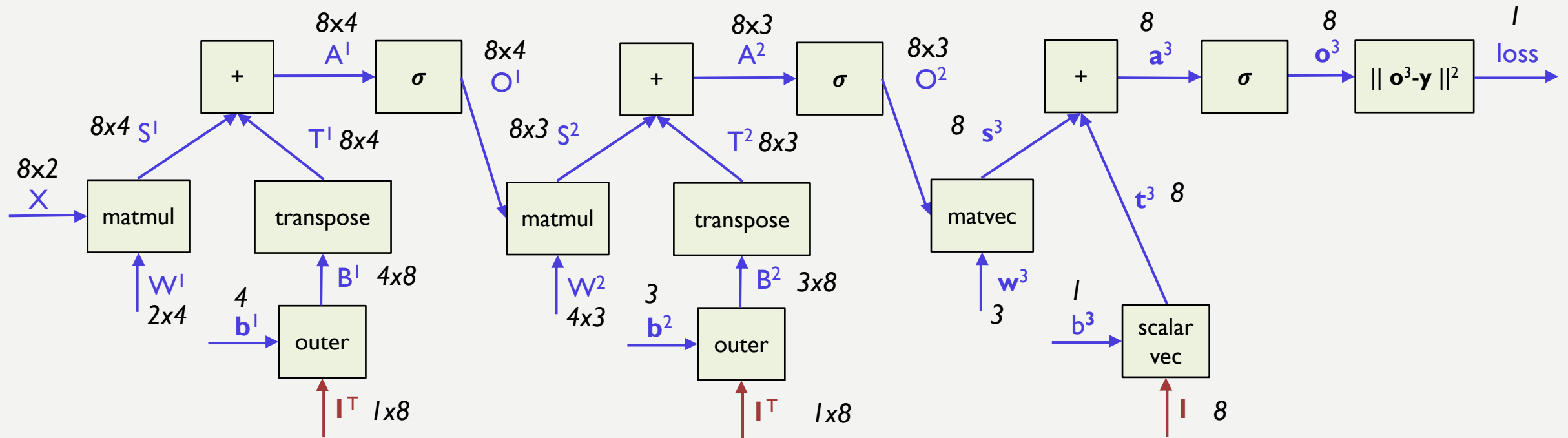
$$O^2 = \sigma(O^1 W^2 + (\mathbf{b}^2 \mathbf{1}^T)^T)$$

$$\mathbf{o}^3 = \sigma(O^2 \mathbf{w}^3 + \mathbf{1} b^3)$$

$$\text{loss} = \|\mathbf{o}^3 - \mathbf{y}\|^2$$

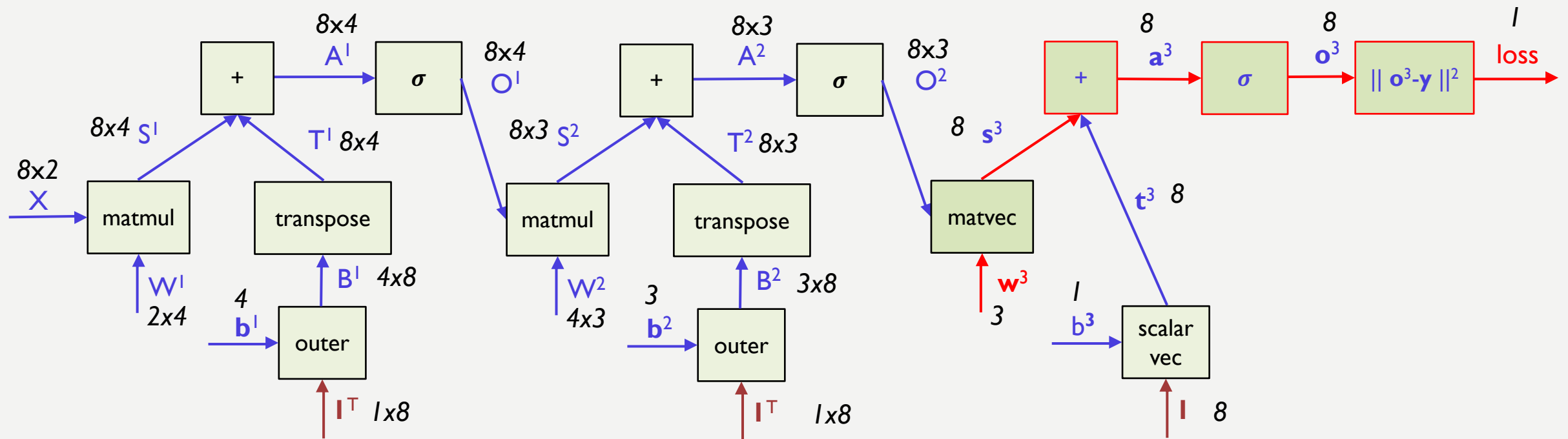
BACKPROPAGATION FROM SCRATCH

- Computational graph of the operations required to compute the loss (*forward*)



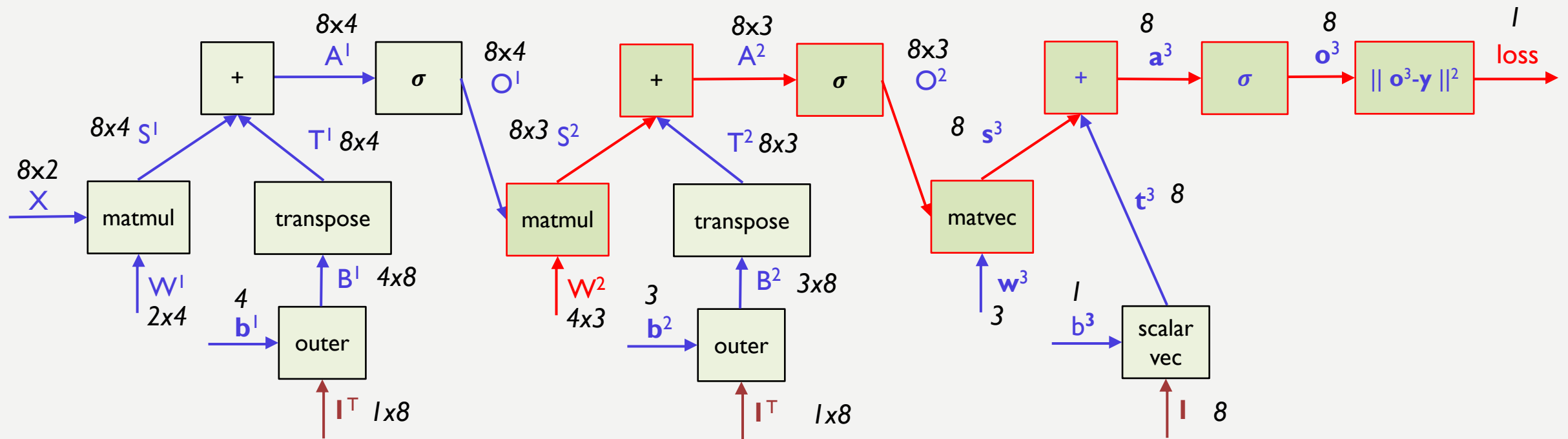
BACKPROPAGATION FROM SCRATCH

- Path that we must follow (in reverse order) to get the derivative of the loss w.r.t. to w^3



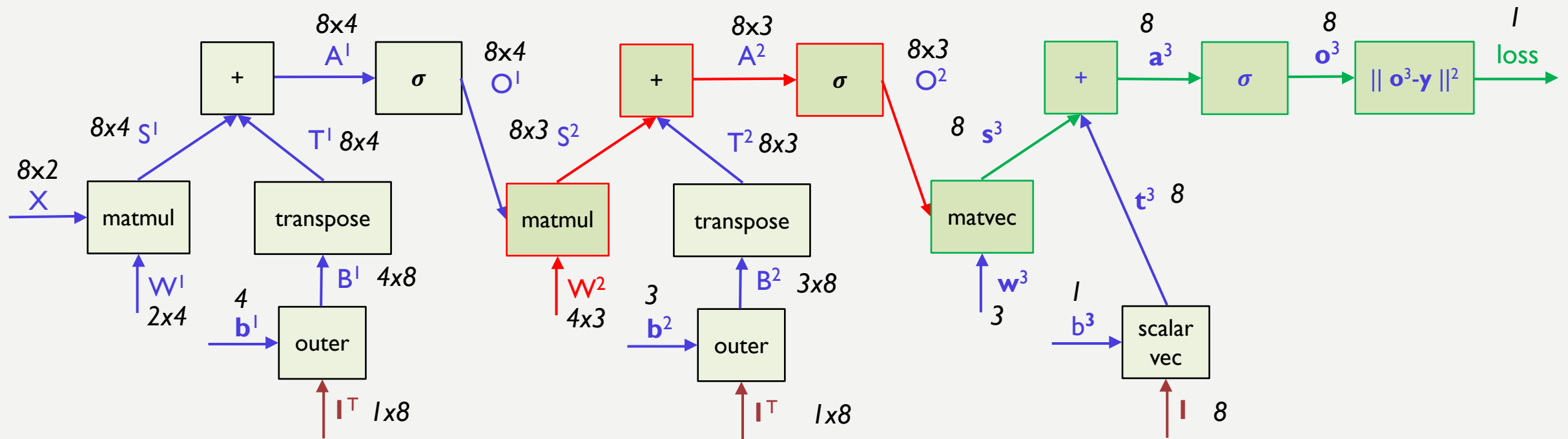
BACKPROPAGATION FROM SCRATCH

- Path that we must follow (in reverse order) to get the derivative of the loss w.r.t. to W^2



BACKPROPAGATION FROM SCRATCH

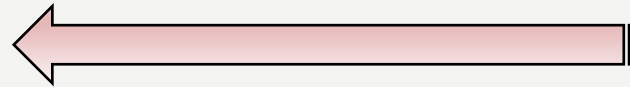
- The greenish portion of the path was already computed when computing the derivative w.r.t. w^3 , we can reuse some intermediate results. Oh, they are the already described Backprop **DELTA**S!



BACKPROPAGATION FROM SCRATCH

- Let's compute all the derivatives we need, in backward fashion, starting from \mathbf{w}^3
- We can recover the same elements/notation we introduced when describing Backpropagation

$$\begin{array}{ccccccc} & 3 & & 3 \times 8 & 8 \times 8 & 8 \times 8 & 8 \\ \frac{\partial \text{loss}}{\partial \mathbf{w}^3} & = & \frac{\partial \mathbf{s}^3}{\partial \mathbf{w}^3} & \frac{\partial \mathbf{a}^3}{\partial \mathbf{s}^3} & \frac{\partial \mathbf{o}^3}{\partial \mathbf{a}^3} & \frac{\partial \text{loss}}{\partial \mathbf{o}^3} \\ & & \underbrace{\hspace{1.5cm}}_{\frac{\partial \mathbf{a}^3}{\partial \mathbf{w}^3}} & & \underbrace{\hspace{1.5cm}}_{\delta^3} & & \\ & & 3 \times 8 & & 8 & & \end{array}$$



Better compute them in this order!

BACKPROPAGATION FROM SCRATCH

- We apply the tensor derivatives we previously described...

$$\begin{aligned}
 \frac{\partial \text{loss}}{\partial \mathbf{w}^3} &= \underbrace{\frac{\partial \mathbf{s}^3}{\partial \mathbf{w}^3}}_{\substack{3 \times 8 \\ \frac{\partial \mathbf{a}^3}{\partial \mathbf{w}^3}}} \underbrace{\frac{\partial \mathbf{a}^3}{\partial \mathbf{s}^3}}_{8 \times 8} \underbrace{\frac{\partial \mathbf{o}^3}{\partial \mathbf{a}^3}}_{\substack{8 \times 8 \\ \delta^3}} \underbrace{\frac{\partial \text{loss}}{\partial \mathbf{o}^3}}_{8 \times 1} \\
 &= \frac{\partial \mathbf{o}^3}{\partial \mathbf{a}^3} \frac{\partial \text{loss}}{\partial \mathbf{o}^3}
 \end{aligned}$$

$\frac{\partial \text{loss}}{\partial \mathbf{o}^3} = 2(\mathbf{o}^3 - \mathbf{y})$ **vector**

$\frac{\partial \mathbf{o}^3}{\partial \mathbf{a}^3} = \begin{bmatrix} \frac{\partial o_1^3}{\partial a_1^3} & \frac{\partial o_2^3}{\partial a_1^3} & \cdots & \frac{\partial o_8^3}{\partial a_1^3} \\ \frac{\partial o_1^3}{\partial a_2^3} & \frac{\partial o_2^3}{\partial a_2^3} & \cdots & \frac{\partial o_8^3}{\partial a_2^3} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial o_1^3}{\partial a_8^3} & \frac{\partial o_2^3}{\partial a_8^3} & \cdots & \frac{\partial o_8^3}{\partial a_8^3} \end{bmatrix} = \begin{bmatrix} \sigma'(a_1^3) & 0 & \cdots & 0 \\ 0 & \sigma'(a_2^3) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(a_8^3) \end{bmatrix}$

derivative of the activation function w.r.t. its argument

diagonal matrix!

BACKPROPAGATION FROM SCRATCH

- When we chain two derivatives, we can simplify the computations exploiting the sparseness of the (tensors that store the) derivatives

$$\frac{\partial \text{loss}}{\partial \mathbf{w}^3} = \underbrace{\frac{\partial \mathbf{s}^3}{\partial \mathbf{w}^3}}_{\substack{\frac{\partial \mathbf{a}^3}{\partial \mathbf{w}^3} \\ 3 \times 8}} \underbrace{\frac{\partial \mathbf{a}^3}{\partial \mathbf{s}^3}}_{8 \times 8} \underbrace{\frac{\partial \mathbf{o}^3}{\partial \mathbf{a}^3} \frac{\partial \text{loss}}{\partial \mathbf{o}^3}}_{\substack{\delta^3 \\ 8 \times 8}} \overset{8}{\frac{\partial \mathbf{o}^3}{\partial \mathbf{a}^3} \frac{\partial \text{loss}}{\partial \mathbf{o}^3} = \sigma'(\mathbf{a}^3) \odot 2(\mathbf{o}^3 - \mathbf{y})}$$

a diagonal matrix multiplied by a vector
can be written as an element-wise
vector multiplication:
Hadamard product

$$\overset{8}{\delta^3 = \sigma'(\mathbf{a}^3) \odot 2(\mathbf{o}^3 - \mathbf{y})}$$

The first operand of this product was originally supposed to be an 8x8 matrix, while it is now a vector with 8 elements (we exploited the structure/sparseness of the two operands)

Simplification like these will be frequently applied

BACKPROPAGATION FROM SCRATCH

- After having computed the last two terms, we can get the result we were looking for

$$\begin{array}{c}
 \begin{array}{cccc}
 3 & 3 \times 8 & 8 \times 8 & 8 \times 8 & 8 \\
 \frac{\partial \text{loss}}{\partial \mathbf{w}^3} = & \boxed{\frac{\partial \mathbf{s}^3}{\partial \mathbf{w}^3}} & \boxed{\frac{\partial \mathbf{a}^3}{\partial \mathbf{s}^3}} & \underbrace{\frac{\partial \mathbf{o}^3}{\partial \mathbf{a}^3} \frac{\partial \text{loss}}{\partial \mathbf{o}^3}}_{\delta^3} \\
 & \underbrace{\frac{\partial \mathbf{a}^3}{\partial \mathbf{w}^3}}_{3 \times 8} & & \underbrace{\delta^3}_{8} &
 \end{array}
 \end{array}$$

$\delta^3 = \sigma'(\mathbf{a}^3) \odot 2(\mathbf{o}^3 - \mathbf{y})$

$\frac{\partial \mathbf{a}^3}{\partial \mathbf{s}^3} = \frac{\partial(\mathbf{s}^3 + \mathbf{t}^3)}{\partial \mathbf{s}^3} = I$

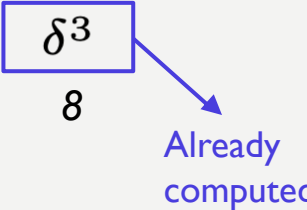
$\frac{\partial \mathbf{s}^3}{\partial \mathbf{w}^3} = \frac{\partial(O^2 \mathbf{w}^3)}{\partial \mathbf{w}^3} = (O^2)^T$

$$\frac{\partial \text{loss}}{\partial \mathbf{w}^3} = (O^2)^T \delta^3 = (O^2)^T (\sigma'(\mathbf{a}^3) \odot 2(\mathbf{o}^3 - \mathbf{y}))$$

BACKPROPAGATION FROM SCRATCH

- What about the bias term b^3 ?

$$\begin{array}{c}
 \frac{\partial \text{loss}}{\partial b^3} = \underbrace{\frac{\partial t^3}{\partial b^3} \frac{\partial a^3}{\partial t^3}}_{\substack{\frac{\partial a^3}{\partial b^3} \\ 1 \times 8}} \underbrace{\frac{\partial o^3}{\partial a^3} \frac{\partial \text{loss}}{\partial o^3}}_{\substack{\delta^3 \\ 8}}
 \end{array}$$


 Already computed

$$\frac{\partial \mathbf{a}^3}{\partial \mathbf{t}^3} = \frac{\partial (\mathbf{s}^3 + \mathbf{t}^3)}{\partial \mathbf{t}^3} = \mathbf{I}^{8 \times 8}$$

$$\frac{\partial \mathbf{t}^3}{\partial b^3} = \frac{\partial (b^3 \mathbf{1})}{\partial b^3} = \mathbf{1}^T \quad 1 \times 8$$

$$\frac{\partial \text{loss}}{\partial b^3} = \mathbf{1}^T \delta^3 = \mathbf{1}^T (\sigma'(\mathbf{a}^3) \odot 2(\mathbf{o}^3 - \mathbf{y}))$$

BACKPROPAGATION FROM SCRATCH

- Let's keep going backward, derivative w.r.t. W^2

$$\frac{\partial \text{loss}}{\partial W^2} = \underbrace{\frac{\partial S^2}{\partial W^2}}_{\substack{4 \times 3 \\ \frac{\partial A^2}{\partial W^2} \\ 4 \times 3 \times 8 \times 3}} \underbrace{\frac{\partial A^2}{\partial S^2} \frac{\partial O^2}{\partial A^2} \frac{\partial s^3}{\partial O^2}}_{\substack{8 \times 3 \times 8 \times 3 \\ 8 \times 3 \times 8 \\ \Delta^2 \\ 8 \times 3}} \underbrace{\frac{\partial a^3}{\partial s^3}}_{\substack{8 \times 8 \\ 8 \\ \delta^3}}$$

$\delta^3 = \sigma'(a^3) \odot 2(o^3 - y)$
 Already computed 😊

$\frac{\partial a^3}{\partial s^3} = \frac{\partial (s^3 + t^3)}{\partial s^3} = I$

It is just the **identity matrix**, we can forget about it / remove it

BACKPROPAGATION FROM SCRATCH

- Let's keep going backward, derivative w.r.t. of W^2

$$\frac{\partial \text{loss}}{\partial W^2} = \underbrace{\frac{\partial S^2}{\partial W^2}}_{\substack{4 \times 3 \\ \frac{\partial A^2}{\partial W^2} \\ 4 \times 3 \times 8 \times 3}} \underbrace{\frac{\partial A^2}{\partial S^2} \frac{\partial O^2}{\partial A^2} \frac{\partial s^3}{\partial O^2}}_{\substack{8 \times 3 \times 8 \times 3 \\ \Delta^2 \\ 8 \times 3}} \underbrace{\frac{\partial a^3}{\partial s^3}}_{\substack{8 \times 8 \\ \delta^3}}$$

$\frac{\partial s^3}{\partial O^2} = \frac{\partial (O^2 \mathbf{w}^3)}{\partial O^2} =$

$$\begin{bmatrix} \frac{\partial (O^2 \mathbf{w}^3)}{\partial O_{11}} & \frac{\partial (O^2 \mathbf{w}^3)}{\partial O_{12}} & \frac{\partial (O^2 \mathbf{w}^3)}{\partial O_{13}} \\ \vdots & \vdots & \vdots \\ \frac{\partial (O^2 \mathbf{w}^3)}{\partial O_{81}} & \frac{\partial (O^2 \mathbf{w}^3)}{\partial O_{82}} & \frac{\partial (O^2 \mathbf{w}^3)}{\partial O_{83}} \end{bmatrix}$$

8x3x8

You can read this result (tensor) as follows:
each element of the 8x3 matrix above is a vector
with 8 elements

Let's see the whole tensor in the next slide!
(be ready)

BACKPROPAGATION FROM SCRATCH

- The tensor is actually very **sparse**!

$$\frac{\partial \mathbf{s}^3}{\partial O^2} = \frac{\partial (O^2 \mathbf{w}^3)}{\partial O^2} = \begin{bmatrix} \boxed{[w_1^3, 0, 0, 0, 0, 0, 0, 0]^T} & [w_2^3, 0, 0, 0, 0, 0, 0, 0]^T & [w_3^3, 0, 0, 0, 0, 0, 0, 0]^T \\ [0, w_1^3, 0, 0, 0, 0, 0, 0]^T & [0, w_2^3, 0, 0, 0, 0, 0, 0]^T & [0, w_3^3, 0, 0, 0, 0, 0, 0]^T \\ \vdots & \vdots & \vdots \\ [0, 0, 0, 0, 0, 0, 0, w_1^3]^T & [0, 0, 0, 0, 0, 0, 0, w_2^3]^T & [0, 0, 0, 0, 0, 0, 0, w_3^3]^T \end{bmatrix}$$

sparse vector

- Looking back at the previous slide, notice that we are now going to multiply it by the vector ‘delta’

$$\dots \frac{\partial \mathbf{s}^3}{\partial O^2} I \delta^3 = \dots \frac{\partial \mathbf{s}^3}{\partial O^2} \delta^3$$

This is the formula in the previous slide, where some elements on the left portion of the formula were replaced by dots, for better readability

- This means that each of the sparse vectors above will be multiplied (dot product) by ‘delta’
- As a result, **only one component of ‘delta’ will be preserved after each dot product**
- We get a matrix in which the element in position i,j is $\delta_i^3 w_j^3$

BACKPROPAGATION FROM SCRATCH

- We can write the result of the just described computation in a more compact manner!
 - No more big-sparse tensors 😊

$$\begin{aligned}
 \frac{\partial \text{loss}}{\partial W^2} &= \underbrace{\frac{\partial S^2}{\partial W^2}}_{\substack{4 \times 3 \\ \frac{\partial A^2}{\partial W^2} \\ 4 \times 3 \times 8 \times 3}} \underbrace{\frac{\partial A^2}{\partial S^2}}_{8 \times 3 \times 8 \times 3} \underbrace{\frac{\partial O^2}{\partial A^2}}_{8 \times 3 \times 8} \underbrace{\frac{\partial \mathbf{s}^3}{\partial O^2} \frac{\partial \mathbf{a}^3}{\partial \mathbf{s}^3} \delta^3}_{\substack{8 \times 3 \times 8 \quad 8 \times 8 \quad 8 \\ \Delta^2 \\ 8 \times 3}} \\
 &\rightarrow \frac{\partial \mathbf{s}^3}{\partial O^2} I \delta^3 = \frac{\partial \mathbf{s}^3}{\partial O^2} \delta^3 = \boxed{\delta^3 (\mathbf{w}^3)^T}
 \end{aligned}$$

'delta' was moved on the left hand side of the product

BACKPROPAGATION FROM SCRATCH

- The remaining terms yield a similar structure to the one we already experienced in the layer above
 - Now we have matrices, a bit more generic case

$$\frac{\partial \text{loss}}{\partial W^2} = \underbrace{\frac{\partial S^2}{\partial W^2}}_{4 \times 3} \underbrace{\frac{\partial A^2}{\partial S^2}}_{4 \times 3 \times 8 \times 3} \underbrace{\frac{\partial O^2}{\partial A^2}}_{8 \times 3 \times 8 \times 3} \underbrace{\frac{\partial s^3}{\partial O^2}}_{8 \times 3 \times 8} \underbrace{\frac{\partial a^3}{\partial s^3}}_{8 \times 8} \delta^3$$

$\frac{\partial O^2}{\partial A^2} = \frac{\partial (\sigma(A^2))}{\partial A^2}$

It is always zero except for those elements whose indices are (i,j,i,j) .

We can apply again the 'trick' of the Hadamard product with the matrix $\sigma'(A^2)_{8 \times 3}$

$\frac{\partial S^2}{\partial W^2} = \frac{\partial (O^1 W^2)}{\partial W^2} = \dots \text{sparse tensor} \dots = (O^1)^T_{4 \times 8}$

It has no effects, we can forget about it (it is a tensor composed of identity matrices)

$$\frac{\partial \text{loss}}{\partial W^2} = (O^1)^T \Delta^2 = (O^1)^T (\sigma'(A^2) \odot \delta^3 (\mathbf{w}^3)^T)$$

BACKPROPAGATION FROM SCRATCH

- What about the bias vector \mathbf{b}^2 ?

$$\frac{\partial \text{loss}}{\partial \mathbf{b}^2} = \underbrace{\frac{\partial B^2}{\partial \mathbf{b}^2} \frac{\partial T^2}{\partial B^2} \frac{\partial A^2}{\partial T^2}}_{\substack{\frac{\partial A^2}{\partial \mathbf{b}^2} \\ 3 \times 8 \times 3}} \underbrace{\frac{\partial O^2}{\partial A^2} \frac{\partial \mathbf{s}^3}{\partial O^2} \frac{\partial \mathbf{a}^3}{\partial \mathbf{s}^3} \delta^3}_{\substack{\Delta^2 \\ 8 \times 3}}$$

It has no effects, we can forget about it (it is a tensor composed of identity matrices)

Already computed

$$\frac{\partial T^2}{\partial B^2} = \frac{\partial (B^2)^T}{\partial B^2}$$

It is a sparse tensor with a few elements equal to one. The effect of multiplying this tensor with the following product of terms in the chain rule is of transposing the outcome of such product

$$\frac{\partial T^2}{\partial B^2} \times Z = Z^T$$

$$\frac{\partial B^2}{\partial \mathbf{b}^2} = \frac{\partial (\mathbf{b}^2 \mathbf{1}^T)}{\partial \mathbf{b}^2}$$

The effect of multiplying this tensor with the following product of terms in the chain rule is of post-multiplying by a vector of ones

$$\frac{\partial B^2}{\partial \mathbf{b}^2} \times U = U \mathbf{1}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{b}^2} = (\Delta^2)^T \mathbf{1}$$

BACKPROPAGATION FROM SCRATCH

- Let's keep going backward, derivative w.r.t. W^l (and \mathbf{b}^l)

$$\frac{\partial \text{loss}}{\partial W^1} = \underbrace{\frac{\partial S^1}{\partial W^1} \frac{\partial A^1}{\partial S^1}}_{\substack{2 \times 4 \\ \frac{\partial A^1}{\partial W^1} \\ 2 \times 4 \times 8 \times 4}} \underbrace{\frac{\partial O^1}{\partial A^1} \frac{\partial S^2}{\partial O^1} \frac{\partial A^2}{\partial S^2} \Delta^2}_{\substack{8 \times 4 \times 8 \times 4 \\ 8 \times 4 \times 8 \times 3 \\ 8 \times 3 \times 8 \times 3 \\ \Delta^1 \\ 8 \times 4}} \frac{\partial S^2}{\partial O^1} = \frac{\partial (O^1 W^2)}{\partial O^1} =$$

$$\begin{bmatrix} W_{11}^2 & W_{12}^2 & W_{13}^2 \\ 0 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \end{bmatrix}$$

8x3

$$\begin{bmatrix} W_{41}^2 & W_{42}^2 & W_{43}^2 \\ 0 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \end{bmatrix}$$

8x3

$$\begin{bmatrix} 0 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \\ W_{41}^2 & W_{42}^2 & W_{43}^2 \end{bmatrix}$$

8x3

It has no effects, we can forget about it (it is a tensor composed of identity matrices)

We get the 'matrix version' of what we already described in the case in which we had a vector \mathbf{s}^3 instead of a matrix S^2

BACKPROPAGATION FROM SCRATCH

- Let's keep going backward, derivative w.r.t. W^1

$$\frac{\partial \text{loss}}{\partial W^1} = \underbrace{\frac{\partial S^1}{\partial W^1} \frac{\partial A^1}{\partial S^1}}_{\frac{\partial A^1}{\partial W^1}} \underbrace{\frac{\partial O^1}{\partial A^1} \frac{\partial S^2}{\partial O^1} \frac{\partial A^2}{\partial S^2} \Delta^2}_{\Delta^1}$$

2×4 $2 \times 4 \times 8 \times 4$ 8×4

$$\frac{\partial S^2}{\partial O^1} \Delta^2 = \Delta^2 (W^2)^T \quad 8 \times 4$$

'Delta' was moved on the left hand side of the product

We already discussed how to compute the other derivatives in the case of the layer above, so we can go straight to the final result (the case of the bias term is fully equivalent to what we already discussed)

$$\frac{\partial \text{loss}}{\partial W^1} = X^T \Delta^1 = X^T \left(\sigma' (A^1) \odot \Delta^2 (W^2)^T \right)$$

$$\frac{\partial \text{loss}}{\partial \mathbf{b}^1} = (\Delta^1)^T \mathbf{1}$$

BACKPROPAGATION FROM SCRATCH

- Extended summary

- Copy & paste of the previous formulas

$$\frac{\partial \text{loss}}{\partial \mathbf{w}^3} = (O^2)^T \delta^3 = (O^2)^T (\sigma'(\mathbf{a}^3) \odot 2(\mathbf{o}^3 - \mathbf{y}))$$

$$\frac{\partial \text{loss}}{\partial W^2} = (O^1)^T \Delta^2 = (O^1)^T (\sigma'(A^2) \odot \delta^3 (\mathbf{w}^3)^T)$$

$$\frac{\partial \text{loss}}{\partial W^1} = X^T \Delta^1 = X^T (\sigma'(A^1) \odot \Delta^2 (W^2)^T)$$

$$\frac{\partial \text{loss}}{\partial b^3} = \mathbf{1}^T \delta^3 = \mathbf{1}^T (\sigma'(\mathbf{a}^3) \odot 2(\mathbf{o}^3 - \mathbf{y}))$$

$$\frac{\partial \text{loss}}{\partial \mathbf{b}^2} = (\Delta^2)^T \mathbf{1}$$

$$\frac{\partial \text{loss}}{\partial \mathbf{b}^1} = (\Delta^1)^T \mathbf{1}$$

In our example we considered a single neuron in the last layer.

We can make this general, and consider multiple output neurons...

BACKPROPAGATION FROM SCRATCH

- **Compact summary**

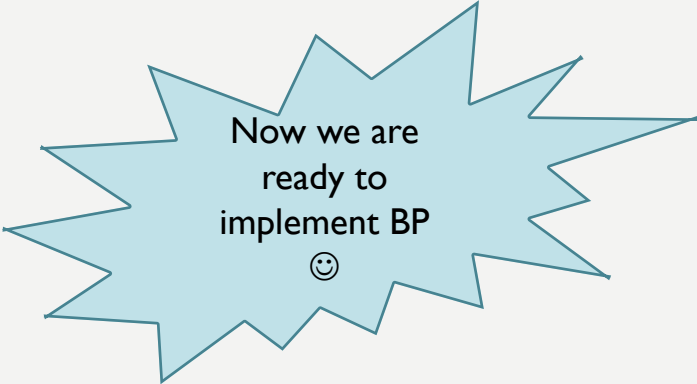
- Compute/update deltas
- Compute derivatives w.r.t. weights and biases

$$\text{num_examples} \times \text{num_neurons_layer_L} \quad \Delta^L = \sigma'(A^L) \odot \frac{\partial \text{loss}}{\partial O^L}$$

$$\text{num_examples} \times \text{num_neurons_layer_l} \quad \Delta^\ell = \sigma'(A^\ell) \odot \Delta^{\ell+1} (W^{\ell+1})^T, \quad 1 \leq \ell < L$$

$$\text{num_neurons_layer_l-1} \times \text{num_neurons_layer_l} \quad \frac{\partial \text{loss}}{\partial W^\ell} = (O^{\ell-1})^T \Delta^\ell, \quad 1 \leq \ell \leq L, \quad O^0 = X$$

$$\text{num_neurons_layer_l} \quad \frac{\partial \text{loss}}{\partial \mathbf{b}^\ell} = (\Delta^\ell)^T \mathbf{1}, \quad 1 \leq \ell \leq L$$



Now we are
ready to
implement BP



BACKPROPAGATION FROM SCRATCH

- **Some take-home hints**


- *If $f(Z)$ is a function that operates element-wise on its argument Z (activation functions)*

$$\frac{\partial f(Z)}{\partial Z} \frac{\partial \text{loss}}{f(Z)} = f'(Z) \odot \frac{\partial \text{loss}}{f(Z)}$$

- *If we have a tensor times another tensor, differentiated w.r.t. the second one*

$$\frac{\partial (OW)}{\partial W} \frac{\partial \text{loss}}{\partial W} = O^T \frac{\partial \text{loss}}{\partial W}$$

- *If we have a tensor times another tensor differentiated w.r.t. to the first one (**swap order!**)*

$$\frac{\partial (OW)}{\partial O} \frac{\partial \text{loss}}{\partial (OW)} = \frac{\partial \text{loss}}{\partial (OW)} W^T$$


CODE BREAK: IMPLEMENTING NEURAL NETWORKS & BACKPROP

- Let's have a look at the code!
 - NumPy-based implementation of
 - Feed-forward neural networks
 - Backpropagation