

# Neural Networks (Part III)

## Automatic differentiation

$$\rightarrow y = F(x), x \in \mathbb{R}^d, y$$

Setup: Evaluate the derivative of a function  $F: \mathbb{R}^d \rightarrow \mathbb{R}$ ,

e.g. evaluate the gradients of the loss of a neural network.

$\hookrightarrow x \in \mathbb{R}^d$  is a vector of containing all NN params.

$y \in \mathbb{R}$  is the value of loss, for a given  $x$ . (2.8)

$$F: \underbrace{[\dots]}_{x \in \mathbb{R}^d} \rightarrow \underbrace{\square}_{y \in \mathbb{R}}$$

Let's focus on cases where  $F$  has a compositional form:

$$F(x) = D(C(B(A(x)))) \quad F = D \circ C \circ B \circ A$$

$$y = D(c), c = C(b), b = B(a), a = A(x)$$



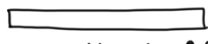
We want to evaluate the Jacobian of  $F$ , i.e.

$$\nabla_x F = F'(x) = \frac{\partial y}{\partial x} = \left[ \frac{\partial y}{\partial x_1}, \frac{\partial y}{\partial x_2}, \dots, \frac{\partial y}{\partial x_d} \right]$$


chain  
rule  $\Rightarrow$

$$\nabla_x F = \frac{\partial y}{\partial c} \cdot \frac{\partial c}{\partial b} \cdot \frac{\partial b}{\partial a} \cdot \frac{\partial a}{\partial x}$$

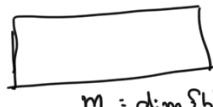
$$\frac{\partial y}{\partial c} = D'(c)$$

$\hookrightarrow$    $K := \dim\{c\}$

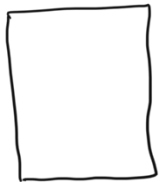
$$\frac{\partial c}{\partial b} = C'(b)$$

$K$    $L := \dim\{b\}$

$$\frac{\partial b}{\partial a} = B'(a)$$

$L$    $m := \dim\{a\}$

$$\frac{\partial a}{\partial x} = A'(x)$$

$m$    $d := \dim\{x\}$

Does it matter in which order we multiply these matrices?

i) Forward accumulation:  $K \times d$

$$\nabla_x F = \frac{\partial y}{\partial x} \left( \frac{\partial c}{\partial b} \left( \frac{\partial b}{\partial a} \cdot \frac{\partial a}{\partial x} \right) \right)$$

Jacobian vector product

$$\nabla_x F \cdot \vec{v} = \frac{\partial y}{\partial c} \left( \frac{\partial c}{\partial b} \left( \frac{\partial b}{\partial a} \left( \frac{\partial a}{\partial x} \right) \vec{v} \right) \right)$$

$\tau = [1, 0] \rightarrow$  row of  $\tau_1$

$$\begin{array}{c} 1 \times d \\ 1 \times k \end{array} \quad \begin{array}{c} 0 \dots 0 \\ 1 \times k \end{array} \quad \begin{array}{c} \underbrace{\hspace{2cm}}_{m \times d} \\ \frac{\partial b}{\partial x} = \begin{bmatrix} \frac{\partial b_1}{\partial x_1} & \dots & \frac{\partial b_1}{\partial x_d} \\ \vdots & & \vdots \\ \frac{\partial b_m}{\partial x_1} & \dots & \frac{\partial b_m}{\partial x_d} \end{bmatrix} \end{array}$$

$V = [0 \ 1] \Rightarrow \frac{\partial}{\partial x}$

Constructs the Jacobian one column at a time

ii) Reverse accumulation:

$$\nabla_x F = \left( \left( \frac{\partial y}{\partial c} \cdot \frac{\partial c}{\partial b} \right) \frac{\partial b}{\partial a} \right) \frac{\partial a}{\partial x}$$

$1 \times d$

$$\frac{\partial y}{\partial b} = \left[ \frac{\partial y}{\partial b_1}, \frac{\partial y}{\partial b_2}, \dots, \frac{\partial y}{\partial b_m} \right]$$

$1 \times m$

Vector jacobian product

$$\vec{v}^T \nabla_x F = \left( \left( \vec{v}^T \frac{\partial y}{\partial c} \right) \frac{\partial c}{\partial b} \right) \frac{\partial b}{\partial a}$$

Constructs the Jacobian matrix one row at a time

(\*) Use fwd mode auto-diff:

when  $\dim\{y\} \gg \dim\{x\}$

Use reverse mode auto-diff:

when  $\dim\{x\} \gg \dim\{y\}$

(this is the most common case in ML)

## Network initialization

### Glorot initialization:

Suppose we have an input  $X$  and a linear neuron with random weights that generates outputs  $Y$ . Also assume that  $\{X, Y\}$  are normalized to have zero mean and unit variance.

$$Y = w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$

(\*) Also assume that has a zero mean.

What's  $\text{Var}[Y]$ ?

$$\text{Var}[w_i x_i] = \cancel{\text{IE}[x_i]^2} \overset{0}{\text{Var}[w_i]} + \cancel{\text{IE}[w_i]^2} \overset{0}{\text{Var}[x_i]} + \text{Var}[w_i] \text{Var}[x_i]$$

$\dots$

$$= \text{Var}[W_i] \text{Var}[x_i]$$

$$\Rightarrow \text{Var}[y] = \underline{d} \cdot \underbrace{\text{Var}[W_i]}_{?} \text{Var}[x_i]$$

Choose  $\text{Var}[W_i]$  such :  $\text{Var}[W_i] = \frac{1}{d_{in}}$  ,  $d_{in} := \dim\{x\}$

If we repeat the same analysis for the back-propagated gradient we get a similar result :  $\text{Var}[W_i] = \frac{1}{d_{out}}$  ,  $d_{out} := \dim\{y\}$

These constraints can only be simultaneously satisfied if  $d_{in} = d_{out}$ .

In general we use an empirical rule to chose  $\text{Var}[W_i]$  as :

$$\text{Var}[W_i] = \frac{2}{d_{in} + d_{out}}$$

This leads to the so-called Glorot initialization scheme :

- Uniform :  $W \sim \mathcal{U}\left[-\frac{\sqrt{6}}{\sqrt{d_{in} + d_{out}}}, \frac{\sqrt{6}}{\sqrt{d_{in} + d_{out}}}\right]$
- Normal :  $W \sim \mathcal{N}\left(0, \frac{2}{d_{in} + d_{out}}\right)$