

# Hands-on derivation of tangent linear and adjoint codes

How to integrate and train neural networks inside 4D-Var

Marcin Chrust and Sébastien Massart

[marcin.chrust@ecmwf.int](mailto:marcin.chrust@ecmwf.int)

European Centre for Medium-Range Weather Forecasts

# Outline

- Bridging machine learning and data assimilation under the common 4D-Var framework
- Tangent linear and adjoint of a simple Multi Layer Perceptron (MLP) neural network
- General MLP neural network and its tangent linear and adjoint
- Link between the adjoint coding and the backpropagation algorithm

# Bridging machine learning and data assimilation under the common 4D-Var framework

- Let the dynamical model be parameterised by a set of parameters  $\mathbf{p}$  constant over the assimilation window:

$$\mathbf{x}_k = \mathcal{M}_{k:0}(\mathbf{p}, \mathbf{x}_0).$$

- The non-linear cost function is:

$$\begin{aligned} \mathcal{J}(\mathbf{p}, \mathbf{x}_0) &= \frac{1}{2} \left\| \mathbf{x}_0 - \mathbf{x}_0^b \right\|_{\mathbf{B}^{-1}}^2 + \frac{1}{2} \left\| \mathbf{p} - \mathbf{p}^b \right\|_{\mathbf{P}^{-1}}^2 \\ &+ \frac{1}{2} \sum_{k=0}^L \left\| \mathbf{y}_k - \mathcal{H}_k \circ \mathcal{M}_{k:0}(\mathbf{p}, \mathbf{x}_0) \right\|_{\mathbf{R}_k^{-1}}^2. \end{aligned}$$

- Consider  $\mathbf{p}$  to be a set of parameters (weights and biases) of a NN.

# Bridging machine learning and data assimilation under the common 4D-Var framework

Let's derive the quadratic cost function by making a change of variables  $(\delta\mathbf{p}, \delta\mathbf{x}_0) \triangleq (\mathbf{p} - \mathbf{p}^i, \mathbf{x}_0 - \mathbf{x}_0^i)$ , where  $(\mathbf{p}^i, \mathbf{x}_0^i)$  is the first guess:

$$\begin{aligned}\mathcal{J}(\mathbf{p}, \mathbf{x}_0) &= \mathcal{J}(\mathbf{p}^i + \delta\mathbf{p}, \mathbf{x}_0^i + \delta\mathbf{x}_0), \\ &= \frac{1}{2} \left\| \mathbf{x}_0^i - \mathbf{x}_0^b + \delta\mathbf{x}_0 \right\|_{\mathbf{B}^{-1}}^2 + \frac{1}{2} \left\| \mathbf{p}^i - \mathbf{p}^b + \delta\mathbf{p} \right\|_{\mathbf{P}^{-1}}^2 \\ &\quad + \frac{1}{2} \sum_{k=0}^L \left\| \mathbf{y}_k - \mathcal{H}_k \circ \mathcal{M}_{k:0}(\mathbf{p}^i + \delta\mathbf{p}, \mathbf{x}_0^i + \delta\mathbf{x}_0) \right\|_{\mathbf{R}_k^{-1}}^2, \\ &\approx \frac{1}{2} \left\| \mathbf{x}_0^i - \mathbf{x}_0^b + \delta\mathbf{x}_0 \right\|_{\mathbf{B}^{-1}}^2 + \frac{1}{2} \left\| \mathbf{p}^i - \mathbf{p}^b + \delta\mathbf{p} \right\|_{\mathbf{P}^{-1}}^2 \\ &\quad + \frac{1}{2} \sum_{k=0}^L \left\| \mathbf{d}_k - \mathbf{H}_k \mathbf{M}_{k:0} \begin{pmatrix} \delta\mathbf{p} \\ \delta\mathbf{x}_0 \end{pmatrix} \right\|_{\mathbf{R}_k^{-1}}^2, \\ &\triangleq \widehat{\mathcal{J}}(\delta\mathbf{p}, \delta\mathbf{x}_0).\end{aligned}$$

where:

- $\mathbf{d}_k \triangleq \mathbf{y}_k - \mathcal{H}_k \circ \mathcal{M}_{k:0}(\mathbf{p}^i, \mathbf{x}_0^i)$ ,
- $\mathbf{H}_k$  is the tangent linear (TL) operator of  $\mathcal{H}_k$  taken at  $\mathcal{M}_{k:0}(\mathbf{p}^i, \mathbf{x}_0^i)$ ,
- $\mathbf{M}_{k:0}$  is the TL operator of  $\mathcal{M}_{k:0}$  taken at  $(\mathbf{p}^i, \mathbf{x}_0^i)$
- $\mathbf{M}_{k:0}$  can be split into  $\mathbf{M}_{k:0}^x$  and  $\mathbf{M}_{k:0}^p$

# Bridging machine learning and data assimilation under the common 4D-Var framework

Let's derive the gradients of the quadratic cost function with respect to  $\delta \mathbf{x}_0$  and  $\delta \mathbf{p}$  respectively:

$$\begin{aligned}\nabla_{\delta \mathbf{x}_0} \widehat{\mathcal{J}}(\delta \mathbf{p}, \delta \mathbf{x}_0) &= \mathbf{B}^{-1} \left( \mathbf{x}_0^i - \mathbf{x}_0^b + \delta \mathbf{x}_0 \right) \\ &\quad - \sum_{k=0}^L \mathbf{M}_{k:0}^{\mathbf{x}}{}^T \mathbf{H}_k^T \mathbf{R}_k^{-1} \left( \mathbf{d}_k - \mathbf{H}_k \mathbf{M}_{k:0} \begin{pmatrix} \delta \mathbf{p} \\ \delta \mathbf{x}_0 \end{pmatrix} \right) \\ \nabla_{\delta \mathbf{p}} \widehat{\mathcal{J}}(\delta \mathbf{p}, \delta \mathbf{x}_0) &= \mathbf{P}^{-1} \left( \mathbf{p}^i - \mathbf{p}^b + \delta \mathbf{p} \right) \\ &\quad - \sum_{k=0}^L \mathbf{M}_{k:0}^{\mathbf{p}}{}^T \mathbf{H}_k^T \mathbf{R}_k^{-1} \left( \mathbf{d}_k - \mathbf{H}_k \mathbf{M}_{k:0} \begin{pmatrix} \delta \mathbf{p} \\ \delta \mathbf{x}_0 \end{pmatrix} \right).\end{aligned}$$

Remember the lecture yesterday on the 4DVar: we will change the loop on the observations to use  $\mathbf{M}_{k:k-1}$  instead of  $\mathbf{M}_{k:0}$

# Bridging machine learning and data assimilation under the common 4D-Var framework

Gradient of the incremental cost function  $\hat{\mathcal{J}}$ .

**Input:**  $\delta \mathbf{p}$  and  $\delta \mathbf{x}_0$

1:  $\mathbf{z}_0 \leftarrow \mathbf{R}_0^{-1} (\mathbf{H}_0 \delta \mathbf{x}_0 - \mathbf{d}_0)$

2: **for**  $k = 1$  **to**  $L$  **do**

3:      $\delta \mathbf{x}_k \leftarrow \mathbf{M}_{k:k-1} (\delta \mathbf{p}, \delta \mathbf{x}_{k-1})^\top$

4:      $\mathbf{z}_k \leftarrow \mathbf{R}_k^{-1} (\mathbf{H}_k \delta \mathbf{x}_k - \mathbf{d}_k)$

5: **end for**

▷ TL of the dynamical model  $\mathcal{M}_{k:k-1}$

# Bridging machine learning and data assimilation under the common 4D-Var framework

Gradient of the incremental cost function  $\hat{\mathcal{J}}$ .

**Input:**  $\delta \mathbf{p}$  and  $\delta \mathbf{x}_0$

1:  $\mathbf{z}_0 \leftarrow \mathbf{R}_0^{-1} (\mathbf{H}_0 \delta \mathbf{x}_0 - \mathbf{d}_0)$

2: **for**  $k = 1$  **to**  $L$  **do**

3:      $\delta \mathbf{x}_k \leftarrow \mathbf{M}_{k:k-1} (\delta \mathbf{p}, \delta \mathbf{x}_{k-1})^\top$

4:      $\mathbf{z}_k \leftarrow \mathbf{R}_k^{-1} (\mathbf{H}_k \delta \mathbf{x}_k - \mathbf{d}_k)$

5: **end for**

6:  $\delta \tilde{\mathbf{x}}_L \leftarrow \mathbf{0}$

7:  $\delta \tilde{\mathbf{p}}_L \leftarrow \mathbf{0}$

▷ TL of the dynamical model  $\mathcal{M}_{k:k-1}$

▷ AD variable for system state

▷ AD variable for model parameters

# Bridging machine learning and data assimilation under the common 4D-Var framework

Gradient of the incremental cost function  $\hat{\mathcal{J}}$ .

**Input:**  $\delta \mathbf{p}$  and  $\delta \mathbf{x}_0$

1:  $\mathbf{z}_0 \leftarrow \mathbf{R}_0^{-1} (\mathbf{H}_0 \delta \mathbf{x}_0 - \mathbf{d}_0)$

2: **for**  $k = 1$  **to**  $L$  **do**

3:  $\delta \mathbf{x}_k \leftarrow \mathbf{M}_{k:k-1} (\delta \mathbf{p}, \delta \mathbf{x}_{k-1})^\top$

▷ TL of the dynamical model  $\mathcal{M}_{k:k-1}$

4:  $\mathbf{z}_k \leftarrow \mathbf{R}_k^{-1} (\mathbf{H}_k \delta \mathbf{x}_k - \mathbf{d}_k)$

5: **end for**

6:  $\delta \tilde{\mathbf{x}}_L \leftarrow \mathbf{0}$

▷ AD variable for system state

7:  $\delta \tilde{\mathbf{p}}_L \leftarrow \mathbf{0}$

▷ AD variable for model parameters

8: **for**  $k = L$  **to** 1 **do**

9:  $\delta \tilde{\mathbf{x}}_k \leftarrow \mathbf{H}_k^\top \mathbf{z}_k + \delta \tilde{\mathbf{x}}_k$

10:  $(\delta \tilde{\mathbf{p}}_{k-1}, \delta \tilde{\mathbf{x}}_{k-1})^\top \leftarrow (\mathbf{M}_{k:k-1})^\top \delta \tilde{\mathbf{x}}_k + (\delta \tilde{\mathbf{p}}_k, \mathbf{0})^\top$

▷ AD of the dyn. model  $\mathcal{M}_{k:k-1}$

11: **end for**

12:  $\delta \tilde{\mathbf{x}}_0 \leftarrow \mathbf{H}_0^\top \mathbf{z}_0 + \delta \tilde{\mathbf{x}}_0$

13:  $\delta \tilde{\mathbf{x}}_0 \leftarrow \mathbf{B}^{-1} (\mathbf{x}_0^i - \mathbf{x}_0^b + \delta \mathbf{x}_0) + \delta \tilde{\mathbf{x}}_0$

14:  $\delta \tilde{\mathbf{p}} \leftarrow \mathbf{P}^{-1} (\mathbf{p}^i - \mathbf{p}^b + \delta \mathbf{p}) + \delta \tilde{\mathbf{p}}_0$

**Output:**  $\nabla_{\delta \mathbf{p}} \hat{\mathcal{J}} = \delta \tilde{\mathbf{p}}$  and  $\nabla_{\delta \mathbf{x}_0} \hat{\mathcal{J}} = \delta \tilde{\mathbf{x}}_0$



# Tangent linear and adjoint of a simple Multi Layer Perceptron (MLP) neural network

Let's consider that our state vector  $\mathbf{x}$  is a scalar  $x$  and that the dynamical model  $\mathcal{M}$  is a Multi Layer Perceptron with a single hidden layer:

$$x_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, x_k) = f^1(w^1 x_k + b^1).$$

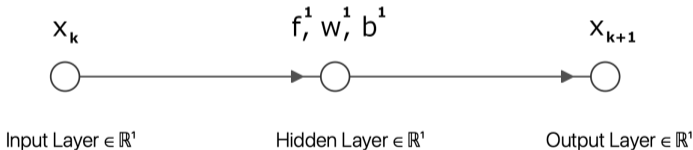


Figure: Simple multi-layer perceptron (MLP).

# Tangent linear and adjoint of a simple Multi Layer Perceptron (MLP) neural network

$$x_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, x_k) = f^1(w^1 x_k + b^1).$$

**NL:**

**Input:**  $x_k, w^1, b^1$

$$a^0 = x_k$$

$$z^1 = w^1 a^0 + b^1$$

$$a^1 = f^1(z^1)$$

$$x_{k+1} = a^1$$

**Output:**  $x_{k+1}$

we will use  $(f^1)'|_{z^1}$  as the derivative of  $f^1$  with regards to  $z^1$

**TL:**

**AD:**

# Tangent linear and adjoint of a simple Multi Layer Perceptron (MLP) neural network

$$x_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, x_k) = f^1(w^1 x_k + b^1).$$

**NL:**

**Input:**  $x_k, w^1, b^1$

$$a^0 = x_k$$

$$z^1 = w^1 a^0 + b^1$$

$$a^1 = f^1(z^1)$$

$$x_{k+1} = a^1$$

**Output:**  $x_{k+1}$

we will use  $(f^1)'|_{z^1}$  as the derivative of  $f^1$  with regards to  $z^1$

**TL:**

**Input:**  $\delta x_k, \delta w^1, \delta b^1, x_k, w^1, b^1$

$$\delta a^0 = \delta x_k$$

$$\delta z^1 = w^1 \delta a^0 + a^0 \delta w^1 + \delta b^1$$

$$\delta a^1 = (f^1)'|_{z^1} \delta z^1$$

$$\delta x_{k+1} = \delta a^1$$

**Output:**  $\delta x_{k+1}$

Note: we need to recompute  $z^1$  inside the TL/AD for  $(f^1)'|_{z^1}$  (not shown)

**AD:**

# Tangent linear and adjoint of a simple Multi Layer Perceptron (MLP) neural network

$$x_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, x_k) = f^1(w^1 x_k + b^1).$$

**NL:**

**Input:**  $x_k, w^1, b^1$

$$a^0 = x_k$$

$$z^1 = w^1 a^0 + b^1$$

$$a^1 = f^1(z^1)$$

$$x_{k+1} = a^1$$

**Output:**  $x_{k+1}$

we will use  $(f^1)'|_{z^1}$  as the derivative of  $f^1$  with regards to  $z^1$

**TL:**

**Input:**  $\delta x_k, \delta w^1, \delta b^1, x_k, w^1, b^1$

$$\delta a^0 = \delta x_k$$

$$\delta z^1 = w^1 \delta a^0 + a^0 \delta w^1 + \delta b^1$$

$$\delta a^1 = (f^1)'|_{z^1} \delta z^1$$

$$\delta x_{k+1} = \delta a^1$$

**Output:**  $\delta x_{k+1}$

Note: we need to recompute  $z^1$  inside the TL/AD for  $(f^1)'|_{z^1}$  (not shown)

**AD:**

**Input:**  $\delta \tilde{x}_{k+1}, x_k, w^1, b^1$

$$\delta \tilde{a}^1 = \delta \tilde{a}^0 = \delta \tilde{z}^1 = 0$$

$$\delta \tilde{a}^1 = \delta \tilde{a}^1 + \delta \tilde{x}_{k+1}$$

$$\delta \tilde{x}_{k+1} = 0$$

$$\delta \tilde{z}^1 = \delta \tilde{z}^1 + (f^1)'|_{z^1} \delta \tilde{a}^1$$

$$\delta \tilde{a}^1 = 0$$

$$\delta \tilde{a}^0 = \delta a^0 + w^1 \delta \tilde{z}^1$$

$$\delta \tilde{w}^1 = \delta \tilde{w}^1 + a^0 \delta \tilde{z}^1$$

$$\delta \tilde{b}^1 = \delta \tilde{b}^1 + \delta \tilde{z}^1$$

$$\delta \tilde{z}^1 = 0$$

$$\delta \tilde{x}_k = \delta \tilde{x}_k + \delta \tilde{a}^0$$

$$\delta \tilde{a}^0 = 0$$

**Output:**  $\delta \tilde{x}_k, \delta \tilde{w}^1, \delta \tilde{b}^1$

# General MLP neural network and its tangent linear and adjoint

Let's consider now that our state vector  $\mathbf{x}$  is a vector of dimension  $n$  and that the model  $\mathcal{M}$  is a Multi Layer Perceptron with  $L$  layers:

$$\mathbf{x}_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, \mathbf{x}_k) = \mathbf{f}^L \left( \mathbf{W}^L \mathbf{f}^{L-1} \left( \mathbf{W}^{L-1} \right) \dots \mathbf{f}^2 \left( \mathbf{W}^2 \mathbf{f}^1 \left( \mathbf{W}^1 \mathbf{a}^0 \right) \right) \dots \right)$$

where  $\mathbf{a}^0 = \begin{bmatrix} \mathbf{x}_k \\ 1 \end{bmatrix} = \mathbf{P}\mathbf{x}_k$  and  $\mathbf{W}^i = \begin{bmatrix} w_{1,1}^i & \dots & w_{1,n_{i-1}}^i & b_1^i \\ \vdots & \ddots & \vdots & \vdots \\ w_{n_i,1}^i & \dots & w_{n_i,n_{i-1}}^i & b_{n_i}^i \end{bmatrix}$

# General MLP neural network and its tangent linear and adjoint

$$\mathbf{x}_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, \mathbf{x}_k) = \mathbf{f}^L \left( \mathbf{W}^L \mathbf{f}^{L-1} \left( \mathbf{W}^{L-1} \right) \dots \mathbf{f}^2 \left( \mathbf{W}^2 \mathbf{f}^1 \left( \mathbf{W}^1 \mathbf{a}^0 \right) \right) \dots \right)$$

**TASK:** write NL, TL and ADJ pseudo-code

- write NL, TL and ADJ pseudo-code
- follow the notation where  $z^l$  is the input to the  $l - th$  hidden layer and  $a^l$  is the output of the  $l - th$  hidden layer
- make use of a loop structure

# General MLP neural network and its tangent linear and adjoint

$$\mathbf{x}_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, \mathbf{x}_k) = \mathbf{f}^L(\mathbf{W}^L \mathbf{f}^{L-1}(\mathbf{W}^{L-1}) \dots \mathbf{f}^2(\mathbf{W}^2 \mathbf{f}^1(\mathbf{W}^1 \mathbf{a}^0)) \dots)$$

NL:

```
Input:  $\mathbf{x}_k, \mathbf{W}^i$   
   $\mathbf{a}^0 = \mathbf{P}\mathbf{x}_k$   
  for  $k = 1$  to  $L$  do  
     $\mathbf{z}^i = \mathbf{W}^i \mathbf{a}^{i-1}$   
     $\mathbf{a}^i = \mathbf{f}^i(\mathbf{z}^i)$   
  end for  
   $\mathbf{x}_{k+1} = \mathbf{a}^L$   
Output:  $\mathbf{x}_{k+1}$ 
```

TL:

AD:

# General MLP neural network and its tangent linear and adjoint

$$\mathbf{x}_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, \mathbf{x}_k) = \mathbf{f}^L(\mathbf{W}^L \mathbf{f}^{L-1}(\mathbf{W}^{L-1}) \dots \mathbf{f}^2(\mathbf{W}^2 \mathbf{f}^1(\mathbf{W}^1 \mathbf{a}^0)) \dots)$$

**NL:**

**Input:**  $\mathbf{x}_k, \mathbf{W}^i$   
 $\mathbf{a}^0 = \mathbf{P}\mathbf{x}_k$   
**for**  $k = 1$  to  $L$  **do**  
     $\mathbf{z}^i = \mathbf{W}^i \mathbf{a}^{i-1}$   
     $\mathbf{a}^i = \mathbf{f}^i(\mathbf{z}^i)$   
**end for**  
 $\mathbf{x}_{k+1} = \mathbf{a}^L$   
**Output:**  $\mathbf{x}_{k+1}$

**TL:**

**Input:**  $\delta\mathbf{x}_k, \delta\mathbf{W}^i, \mathbf{x}_k, \mathbf{W}^i$   
 $\delta\mathbf{a}^0 = \mathbf{P}\delta\mathbf{x}_k$   
**for**  $k = 1$  to  $L$  **do**  
     $\delta\mathbf{z}^i = \mathbf{W}^i \delta\mathbf{a}^{i-1} + \mathbf{a}^{i-1} \delta\mathbf{W}^i$   
     $\delta\mathbf{a}^i = (\mathbf{f}^i)'|_{\mathbf{z}^i} \delta\mathbf{z}^i$   
**end for**  
 $\delta\mathbf{x}_{k+1} = \delta\mathbf{a}^L$   
**Output:**  $\delta\mathbf{x}_{k+1}$

**AD:**



# General MLP neural network and its tangent linear and adjoint

$$\mathbf{x}_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, \mathbf{x}_k) = \mathbf{f}^L(\mathbf{W}^L \mathbf{f}^{L-1}(\mathbf{W}^{L-1}) \dots \mathbf{f}^2(\mathbf{W}^2 \mathbf{f}^1(\mathbf{W}^1 \mathbf{a}^0)) \dots)$$

**NL:**

**Input:**  $\mathbf{x}_k, \mathbf{W}^i$   
 $\mathbf{a}^0 = \mathbf{P}\mathbf{x}_k$   
**for**  $k = 1$  to  $L$  **do**  
 $\mathbf{z}^i = \mathbf{W}^i \mathbf{a}^{i-1}$   
 $\mathbf{a}^i = \mathbf{f}^i(\mathbf{z}^i)$   
**end for**  
 $\mathbf{x}_{k+1} = \mathbf{a}^L$   
**Output:**  $\mathbf{x}_{k+1}$

**TL:**

**Input:**  $\delta\mathbf{x}_k, \delta\mathbf{W}^i, \mathbf{x}_k, \mathbf{W}^i$   
 $\delta\mathbf{a}^0 = \mathbf{P}\delta\mathbf{x}_k$   
**for**  $k = 1$  to  $L$  **do**  
 $\delta\mathbf{z}^i = \mathbf{W}^i \delta\mathbf{a}^{i-1} + \mathbf{a}^{i-1} \delta\mathbf{W}^i$   
 $\delta\mathbf{a}^i = (\mathbf{f}^i)'|_{\mathbf{z}^i} \delta\mathbf{z}^i$   
**end for**  
 $\delta\mathbf{x}_{k+1} = \delta\mathbf{a}^L$   
**Output:**  $\delta\mathbf{x}_{k+1}$

**AD:**

**Input:**  $\delta\tilde{\mathbf{x}}_{k+1}, \mathbf{x}_k, \mathbf{W}^i$   
 $\delta\tilde{\mathbf{a}}^L = \delta\tilde{\mathbf{a}}^L + \delta\tilde{\mathbf{x}}_{k+1}$   
 $\delta\tilde{\mathbf{x}}_{k+1} = 0$   
**for**  $k = L$  to  $1$  **do**  
 $\delta\tilde{\mathbf{z}}^i = \delta\tilde{\mathbf{z}}^i + (\mathbf{f}^i)'|_{\mathbf{z}^i} \delta\tilde{\mathbf{a}}^i$   
 $\delta\tilde{\mathbf{a}}^i = 0$   
 $\delta\tilde{\mathbf{a}}^{i-1} = \delta\mathbf{a}^{i-1} + \mathbf{W}^{iT} \delta\tilde{\mathbf{z}}^i$   
 $\delta\tilde{\mathbf{W}}^i = \delta\tilde{\mathbf{W}}^i + \mathbf{a}^{i-1,T} \delta\tilde{\mathbf{z}}^i$   
 $\delta\tilde{\mathbf{z}}^i = 0$   
**end for**  
 $\delta\tilde{\mathbf{x}}_k = \delta\tilde{\mathbf{x}}_k + \mathbf{P}^T \delta\tilde{\mathbf{a}}^0$   
 $\delta\tilde{\mathbf{a}}^0 = 0$   
**Output:**  $\delta\tilde{\mathbf{a}}^0, \delta\tilde{\mathbf{W}}^i$

# Link between the adjoint coding and the backpropagation algorithm

$$\mathbf{x}_{k+1} = \mathcal{M}_{k+1:k}(\mathbf{p}, \mathbf{x}_k) = \mathbf{f}^L(\mathbf{W}^L \mathbf{f}^{L-1}(\mathbf{W}^{L-1}) \dots \mathbf{f}^2(\mathbf{W}^2 \mathbf{f}^1(\mathbf{W}^1 \mathbf{a}^0)) \dots)$$

Let's take the derivative of the neural network w.r.t. the input  $\mathbf{x}_k$  and its parameters  $\mathbf{W}^i$  of the  $i$ -th layer:

$$\frac{d\mathbf{f}^L(\dots)}{d\mathbf{x}_k} = \frac{d\mathbf{f}^L}{d\mathbf{a}^L} \cdot \frac{d\mathbf{a}^L}{d\mathbf{z}^L} \cdot \frac{d\mathbf{z}^L}{d\mathbf{a}^{L-1}} \cdot \frac{d\mathbf{a}^{L-1}}{d\mathbf{z}^{L-1}} \dots \frac{d\mathbf{a}^1}{d\mathbf{z}^1} \cdot \frac{\partial \mathbf{z}^1}{\partial \mathbf{a}^0} \cdot \frac{d\mathbf{a}^0}{d\mathbf{x}_k}$$

$$\frac{d\mathbf{f}^L(\dots)}{d\mathbf{W}^i} = \frac{d\mathbf{f}^L}{d\mathbf{a}^L} \cdot \frac{d\mathbf{a}^L}{d\mathbf{z}^L} \cdot \frac{d\mathbf{z}^L}{d\mathbf{a}^{L-1}} \cdot \frac{d\mathbf{a}^{L-1}}{d\mathbf{z}^{L-1}} \dots \frac{d\mathbf{a}^i}{d\mathbf{z}^i} \cdot \frac{\partial \mathbf{z}^i}{\partial \mathbf{W}^i}$$

And evaluate the derivatives:

$$\frac{d\mathbf{f}^L(\dots)}{d\mathbf{x}_k} = (\mathbf{f}^L)'|_{z^L} \cdot \mathbf{W}^L \cdot (\mathbf{f}^{L-1})'|_{z^{L-1}} \cdot \mathbf{W}^{L-1} \dots (\mathbf{f}^1)'|_{z^1} \cdot \mathbf{W}^1 \cdot \mathbf{P}$$

$$\frac{d\mathbf{f}^L(\dots)}{d\mathbf{W}^i} = (\mathbf{f}^L)'|_{z^L} \cdot \mathbf{W}^L \cdot (\mathbf{f}^{L-1})'|_{z^{L-1}} \cdot \mathbf{W}^{L-1} \dots (\mathbf{f}^i)'|_{z^i} \cdot \mathbf{W}^i \cdot \mathbf{a}^{i-1}$$

# Link between the adjoint coding and the backpropagation algorithm

Note that we can write the Tangent Linear model as:

$$\delta \mathbf{x}_{k+1} = \mathbf{M}_{k+1:k}(\delta \mathbf{p}, \delta \mathbf{x}_k) = \frac{d\mathbf{f}^L(\dots)}{d\mathbf{x}_k} \delta \mathbf{x}_k + \sum_{i=1}^L \frac{d\mathbf{f}^L(\dots)}{d\mathbf{W}^i} \delta \mathbf{W}^i$$

Recall our TL model derived using the line by line approach:

TL:

**Input:**  $\delta \mathbf{x}_k, \delta \mathbf{W}^i$

$$\delta \mathbf{a}^0 = \mathbf{P} \delta \mathbf{x}_k$$

**for**  $k = 1$  to  $L$  **do**

$$\delta \mathbf{z}^i = \mathbf{W}^i \delta \mathbf{a}^{i-1} + \mathbf{a}^{i-1} \delta \mathbf{W}^i$$

$$\delta \mathbf{a}^i = (\mathbf{f}^i)'|_{\mathbf{z}^i} \delta \mathbf{z}^i$$

**end for**

$$\delta \mathbf{x}_{k+1} = \delta \mathbf{a}^L$$

# Link between the adjoint coding and the backpropagation algorithm

Recall that the gradient is the transpose of the derivative:

$$\begin{aligned}\nabla_{\mathbf{x}_k} \mathbf{f}^L(\dots) &= \mathbf{P}^T \cdot (\mathbf{W}^1)^T \cdot (\mathbf{f}^1)'|_{\mathbf{z}^1} \cdot \dots \cdot (\mathbf{W}^{L-1})^T \cdot (\mathbf{f}^{L-1})'|_{\mathbf{z}^{L-1}} \cdot (\mathbf{W}^L)^T \cdot (\mathbf{f}^L)'|_{\mathbf{z}^L} \\ \nabla_{\mathbf{w}^i} \mathbf{f}^L(\dots) &= (\mathbf{a}^{i-1})^T (\mathbf{W}^i)^T \cdot (\mathbf{f}^i)'|_{\mathbf{z}^i} \cdot \dots \cdot (\mathbf{W}^{L-1})^T \cdot (\mathbf{f}^{L-1})'|_{\mathbf{z}^{L-1}} \cdot (\mathbf{W}^L)^T \cdot (\mathbf{f}^L)'|_{\mathbf{z}^L} \\ &= (\mathbf{a}^{i-1})^T \delta_i\end{aligned}$$

Note that  $\delta_i$  can be computed **recursively** going from the end to the beginning:

$$\begin{aligned}\delta_L &= (\mathbf{f}^L)'|_{\mathbf{z}^L} \\ \delta_{i-1} &= (\mathbf{f}^{i-1})'|_{\mathbf{z}^{i-1}} \cdot (\mathbf{W}^i)^T \cdot \delta_i\end{aligned}$$

This allows for an efficient computation of the gradient of the neural network and is known as the **backpropagation** algorithm.

# Link between the adjoint coding and the backpropagation algorithm

Finally note that we can write the adjoint model as:

$$\begin{aligned}\delta\tilde{\mathbf{x}}_k &= \nabla_{\mathbf{x}_k} \mathbf{f}^L(\dots) \delta\tilde{\mathbf{x}}_{k+1} \\ \delta\tilde{\mathbf{W}}_i &= \nabla_{\mathbf{W}^i} \mathbf{f}^L(\dots) \delta\tilde{\mathbf{x}}_{k+1}\end{aligned}$$

where

$$\begin{aligned}\nabla_{\mathbf{x}_k} \mathbf{f}^L(\dots) &= \mathbf{P}^T \cdot (\mathbf{W}^1)^T \cdot (\mathbf{f}^1)'|_{\mathbf{z}^1} \dots \dots \\ &\quad (\mathbf{W}^{L-1})^T \cdot (\mathbf{f}^{L-1})'|_{\mathbf{z}^{L-1}} \cdot \\ &\quad (\mathbf{W}^L)^T \cdot (\mathbf{f}^L)'|_{\mathbf{z}^L} \\ \nabla_{\mathbf{W}^i} \mathbf{f}^L(\dots) &= (\mathbf{a}^{i-1})^T (\mathbf{W}^i)^T \cdot (\mathbf{f}^i)'|_{\mathbf{z}^i} \dots \dots \\ &\quad (\mathbf{W}^{L-1})^T \cdot (\mathbf{f}^{L-1})'|_{\mathbf{z}^{L-1}} \cdot \\ &\quad (\mathbf{W}^L)^T \cdot (\mathbf{f}^L)'|_{\mathbf{z}^L}\end{aligned}$$

Recall our implementation of the adjoint code: AD:

```
Input:  $\delta\tilde{\mathbf{x}}_{k+1}$   
 $\delta\tilde{\mathbf{a}}^L = \delta\tilde{\mathbf{a}}^L + \delta\tilde{\mathbf{x}}_{k+1}$   
 $\delta\tilde{\mathbf{x}}_{k+1} = 0$   
for  $k = L$  to  $1$  do  
   $\delta\tilde{\mathbf{z}}^i = \delta\tilde{\mathbf{z}}^i + (\mathbf{f}^i)'|_{\mathbf{z}^i} \delta\tilde{\mathbf{a}}^i$   
   $\delta\tilde{\mathbf{a}}^i = 0$   
   $\delta\tilde{\mathbf{a}}^{i-1} = \delta\mathbf{a}^{i-1} + \mathbf{W}^{iT} \delta\tilde{\mathbf{z}}^i$   
   $\delta\tilde{\mathbf{W}}^i = \delta\tilde{\mathbf{W}}^i + \mathbf{a}^{i-1,T} \delta\tilde{\mathbf{z}}^i$   
   $\delta\tilde{\mathbf{z}}^i = 0$   
end for  
 $\delta\tilde{\mathbf{x}}_k = \delta\tilde{\mathbf{x}}_k + \mathbf{P}^T \delta\tilde{\mathbf{a}}^0$   
 $\delta\tilde{\mathbf{a}}^0 = 0$ 
```

Can you verify that we have effectively implemented the **backpropagation** algorithm?