

# Practical Course: Modeling, Simulation, Optimization

Week 12

**Daniël Veldman**

Chair in Dynamics, Control, and Numerics, Friedrich-Alexander-University Erlangen-Nürnberg

## Contents

**12.A** Deep (residual) neural networks and neural ODEs

**12.B** Sensitivity analysis with neural ODEs **12.C** An algorithm for the training of deep residual neural networks



## 12.A Deep (residual) neural networks and neural ODEs



## Classification problem

Classification problems: sort objects in different categories.

Example: MNIST data set.

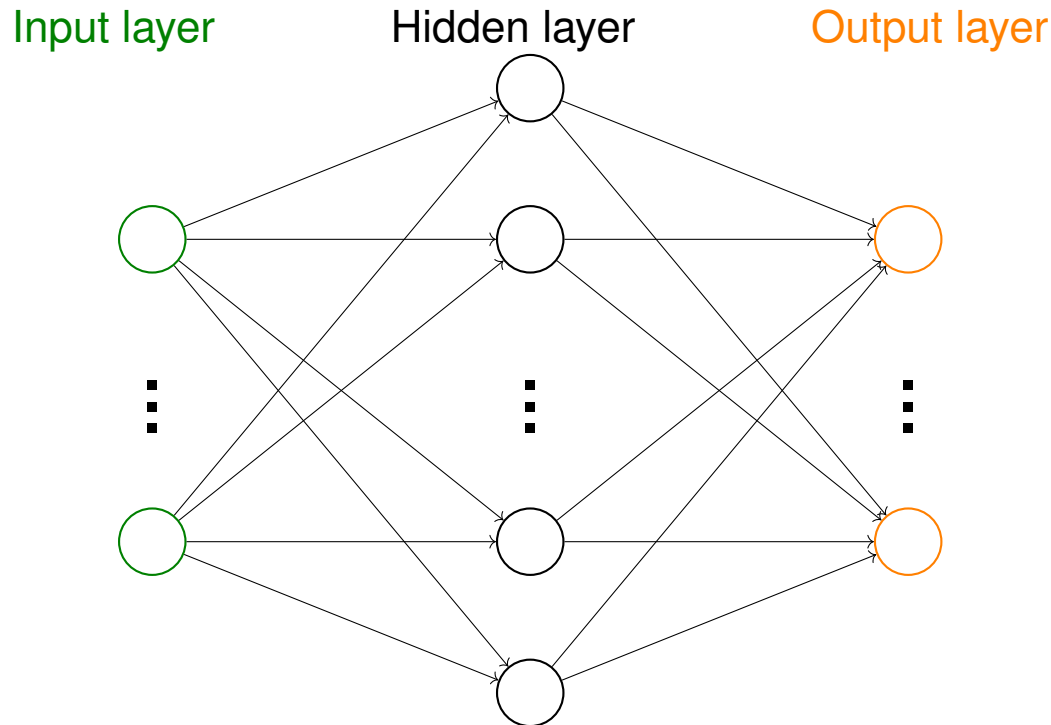


60,000 training images.

10,000 testing images.

Each image is  $28 \times 28$  pixels.

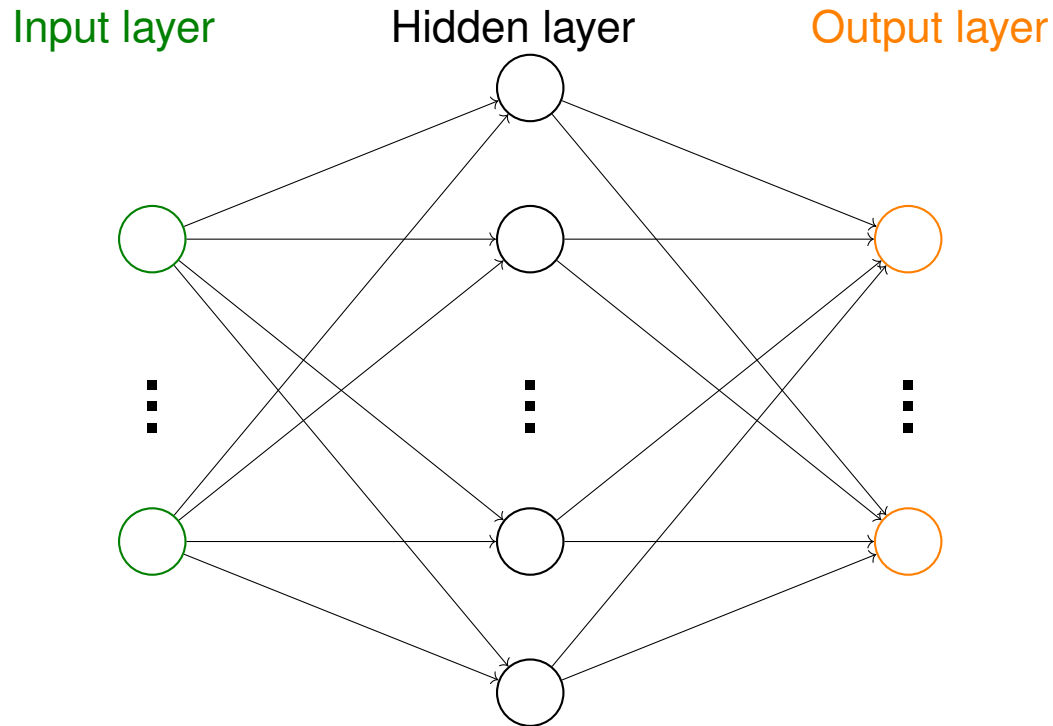
## Neural network



The number of nodes in the input layer is the number of inputs.  
Each arrow represents a linear map.  
Each node in the hidden layer represents a nonlinear operation  $\sigma(\cdot)$ .  
The number of nodes in the output layer is the number of outputs.

$$\mathbf{y} = \mathbf{V}\sigma(\mathbf{W}\mathbf{x} + \mathbf{c}) + \mathbf{b}$$

## Representation theorem



$$\mathbf{y} = \mathbf{V}\sigma(\mathbf{W}\mathbf{x} + \mathbf{c}) + \mathbf{b}$$

Any function  $\mathbf{y} = f(\mathbf{x})$  (e.g. in  $L^2$ ) can be approximated arbitrarily well by a sufficiently **wide** neural network  $\mathbf{y} = \mathbf{V}\sigma(\mathbf{W}\mathbf{x} + \mathbf{c}) + \mathbf{b}$ .

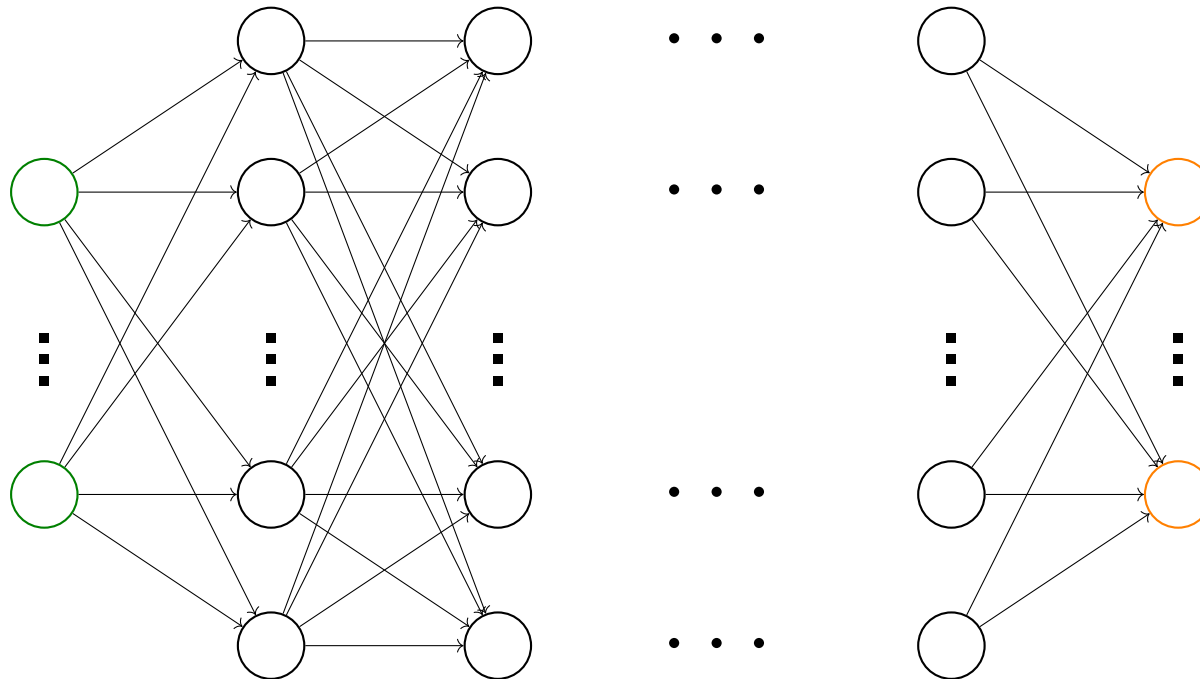
(Cybenko, 1989)

## Deep Neural Networks

Input layer

$K$  hidden layers

Output layer



There are now multiple hidden layers:

$$\mathbf{x}_0 = \mathbf{x}_{\text{in}},$$

$$\mathbf{x}_k = \mathbf{V}_k \sigma(\mathbf{x}_{k-1}) + \mathbf{b}_k,$$

$$\mathbf{y} = \mathbf{x}_K.$$

## Residual neural networks (ResNet) and neural ODEs

$$\mathbf{x}_0 = \mathbf{x}_{\text{in}}, \quad \mathbf{x}_k = \mathbf{V}_k \sigma(\mathbf{x}_{k-1}) + \mathbf{b}_k, \quad \mathbf{y} = \mathbf{x}_K.$$

When the number of hidden layers  $K$  is large,  
it is better to consider a **residual neural network (ResNet)**

$$\mathbf{x}_0 = \mathbf{x}_{\text{in}}, \quad \mathbf{x}_k = \mathbf{x}_{k-1} + \mathbf{V}_k \sigma(\mathbf{x}_{k-1}) + \mathbf{b}_k, \quad \mathbf{y} = \mathbf{x}_K.$$

We can view a ResNet as Forward Euler discretization of the **neural ODE**:

$$\mathbf{x}(0) = \mathbf{x}_{\text{in}}, \quad \dot{\mathbf{x}}(t) = \mathbf{V}(t) \sigma(\mathbf{x}(t)) + \mathbf{b}(t), \quad \mathbf{y} = \mathbf{x}(T).$$

How do we find the weights  $\mathbf{V}(t)$  and  $\mathbf{b}(t)$ ?

## Training a neural ODE

$$\mathbf{x}(0) = \mathbf{x}_{\text{in}}, \quad \dot{\mathbf{x}}(t) = \mathbf{V}(t)\sigma(\mathbf{x}(t)) + \mathbf{b}(t), \quad \mathbf{y} = \mathbf{x}(T).$$

Given training data: pairs  $(\mathbf{x}_{\text{in}}^i, \mathbf{y}_{\text{out}}^i)$ ,  $i = 1, 2, 3, \dots, I$ .

$\mathbf{y}_{\text{out}}^i$  is the desired output for the input  $\mathbf{x}_{\text{in}}^i$ .

For certain weights  $\mathbf{V}(t)$  and  $\mathbf{b}(t)$ , the output resulting from the input  $\mathbf{x}_{\text{in}}^i$  is  $\mathbf{x}^i(T)$

We thus want to minimize

$$J(\mathbf{V}, \mathbf{b}) = \frac{1}{2} \sum_{i=1}^I |\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i|^2 + \frac{w_1}{2} \sum_{i=1}^I \int_0^T |\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i|^2 dt + \frac{w_2}{2} \int_0^T (\|\mathbf{V}(t)\|_F^2 + |\mathbf{b}(t)|^2) dt.$$

subject to the dynamics (for  $i = 1, 2, 3, \dots, I$ )

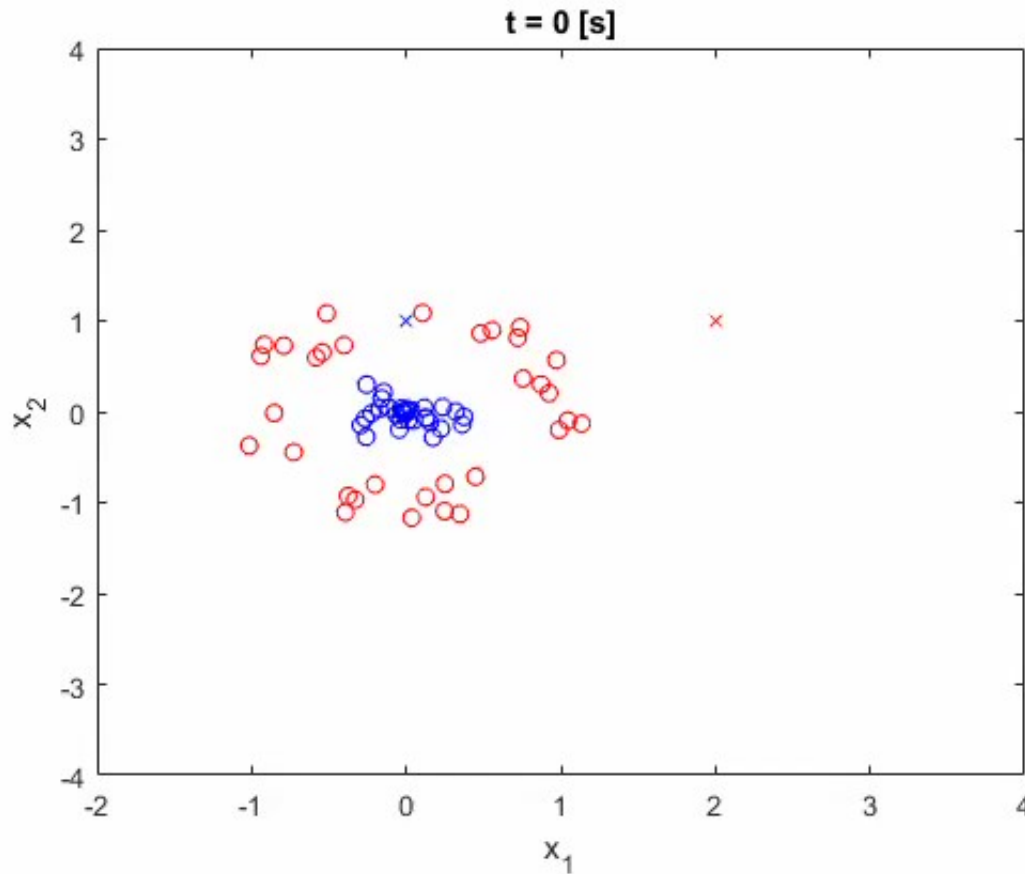
$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i, \quad \dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t).$$

Note:  $\mathbf{b}(t)$  is a vector, but  $\mathbf{V}(t)$  is a (square) matrix.

Note:  $\|\cdot\|_F$  denotes the Frobenius norm,  $\|\mathbf{V}\|_F := \sqrt{\text{trace}(\mathbf{V}^\top \mathbf{V})}$ .



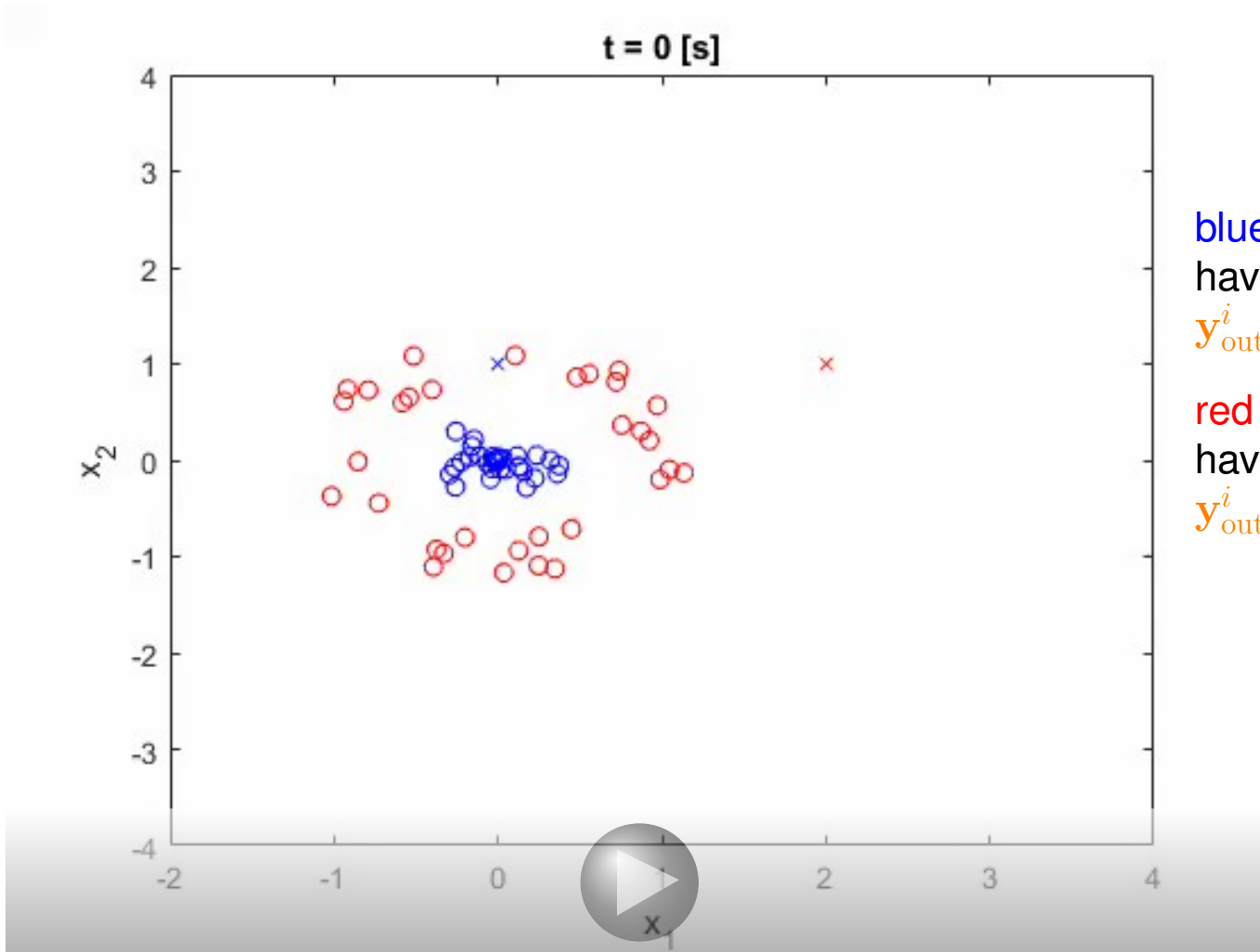
## An example: initial data $x_{\text{in}}^i$



blue data points  
have target  
 $y_{\text{out}}^i = (0, 1)$ .

red data points  
have target  
 $y_{\text{out}}^i = (2, 1)$ .

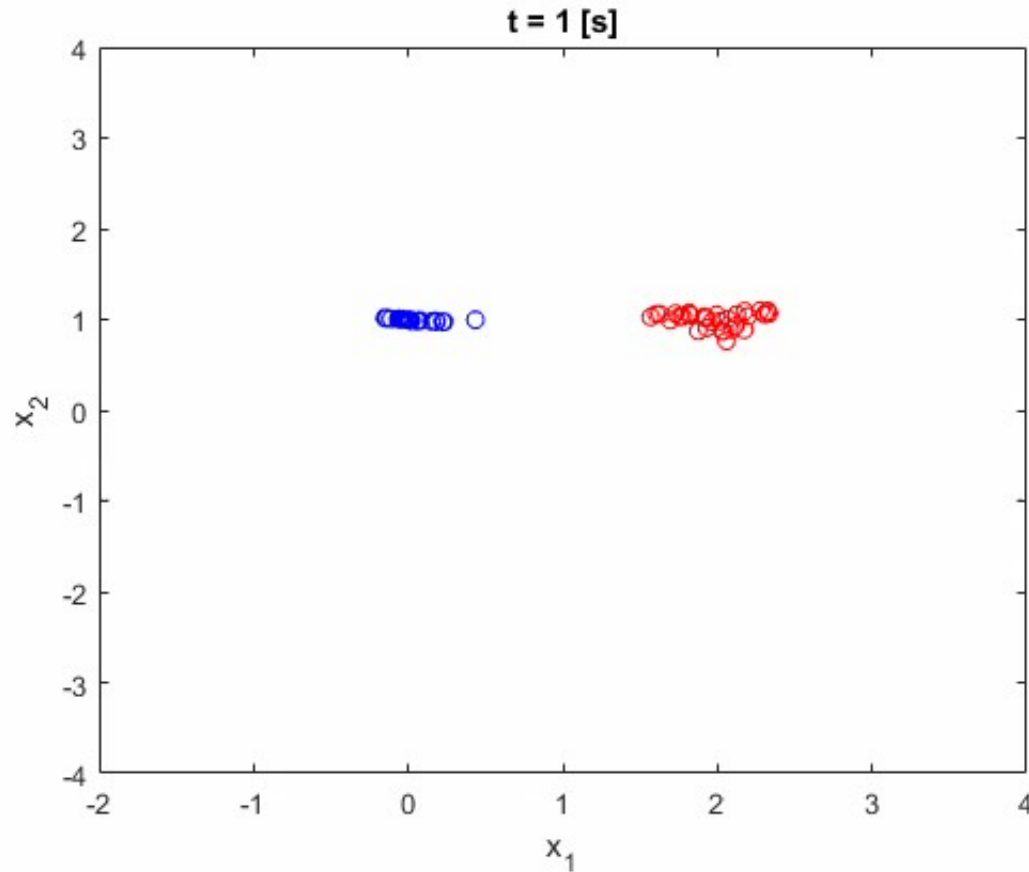
## An example: evolution of $x^i(t)$



blue data points  
have target  
 $y_{\text{out}}^i = (0, 1)$ .

red data points  
have target  
 $y_{\text{out}}^i = (2, 1)$ .

## An example: final data $x^i(T)$



blue data points  
have target  
 $y_{\text{out}}^i = (0, 1)$ .

red data points  
have target  
 $y_{\text{out}}^i = (2, 1)$ .

## 12.B Sensitivity analysis with neural ODEs



## The directional derivative w.r.t. $\mathbf{b}(t)$

Consider a perturbation  $\tilde{\mathbf{b}}(t)$  and compute:

$$\tilde{\mathbf{x}}^i(t) := \lim_{h \rightarrow 0} \frac{\mathbf{x}^{ih}(t) - \mathbf{x}^i(t)}{h}$$

where

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i,$$

$$\mathbf{x}^{ih}(0) = \mathbf{x}_{\text{in}}^i,$$

$$\dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t),$$

$$\dot{\mathbf{x}}^{ih}(t) = \mathbf{V}(t)\sigma(\mathbf{x}^{ih}(t)) + \mathbf{b}(t) + h\tilde{\mathbf{b}}(t).$$

## The directional derivative w.r.t. $\mathbf{b}(t)$

Consider a perturbation  $\tilde{\mathbf{b}}(t)$  and compute:

$$\tilde{\mathbf{x}}^i(t) := \lim_{h \rightarrow 0} \frac{\mathbf{x}^{ih}(t) - \mathbf{x}^i(t)}{h}$$

where

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i,$$

$$\mathbf{x}^{ih}(0) = \mathbf{x}_{\text{in}}^i,$$

$$\dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t),$$

$$\dot{\mathbf{x}}^{ih}(t) = \mathbf{V}(t)\sigma(\mathbf{x}^{ih}(t)) + \mathbf{b}(t) + h\tilde{\mathbf{b}}(t).$$

Because

$$\frac{\dot{\mathbf{x}}^{ih} - \dot{\mathbf{x}}^i}{h} = \frac{\mathbf{V}\sigma(\mathbf{x}^{ih}) + \mathbf{b} + h\tilde{\mathbf{b}} - \mathbf{V}\sigma(\mathbf{x}^i) - \mathbf{b}}{h} = \mathbf{V} \frac{\sigma(\mathbf{x}^{ih}) - \sigma(\mathbf{x}^i)}{h} + \tilde{\mathbf{b}},$$

taking the limit  $h \rightarrow 0$  shows that  $\tilde{\mathbf{x}}^i(t)$  satisfies

$$\tilde{\mathbf{x}}^i(0) = 0, \quad \dot{\tilde{\mathbf{x}}}^i(t) = \mathbf{V}(t) \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}^i(t)) \right) \tilde{\mathbf{x}}^i(t) + \tilde{\mathbf{b}}(t).$$

## The directional derivative w.r.t. $\mathbf{b}(t)$

Consider a perturbation  $\tilde{\mathbf{b}}(t)$  and compute:

$$\tilde{\mathbf{x}}^i(t) := \lim_{h \rightarrow 0} \frac{\mathbf{x}^{ih}(t) - \mathbf{x}^i(t)}{h}$$

where

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i,$$

$$\mathbf{x}^{ih}(0) = \mathbf{x}_{\text{in}}^i,$$

$$\dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t),$$

$$\dot{\mathbf{x}}^{ih}(t) = \mathbf{V}(t)\sigma(\mathbf{x}^{ih}(t)) + \mathbf{b}(t) + h\tilde{\mathbf{b}}(t).$$

Because

$$\frac{\dot{\mathbf{x}}^{ih} - \dot{\mathbf{x}}^i}{h} = \frac{\mathbf{V}\sigma(\mathbf{x}^{ih}) + \mathbf{b} + h\tilde{\mathbf{b}} - \mathbf{V}\sigma(\mathbf{x}^i) - \mathbf{b}}{h} = \mathbf{V} \frac{\sigma(\mathbf{x}^{ih}) - \sigma(\mathbf{x}^i)}{h} + \tilde{\mathbf{b}},$$

taking the limit  $h \rightarrow 0$  shows that  $\tilde{\mathbf{x}}^i(t)$  satisfies

$$\tilde{\mathbf{x}}^i(0) = 0, \quad \dot{\tilde{\mathbf{x}}}^i(t) = \mathbf{V}(t) \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}^i(t)) \right) \tilde{\mathbf{x}}^i(t) + \tilde{\mathbf{b}}(t).$$

We then obtain

$$\begin{aligned} \langle \nabla_{\mathbf{b}} J, \tilde{\mathbf{b}} \rangle &= \lim_{h \rightarrow 0} \frac{J(\mathbf{V}, \mathbf{b} + h\tilde{\mathbf{b}}) - J(\mathbf{V}, \mathbf{b})}{h} = \sum_{i=1}^I (\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(T) + \\ &\quad w_1 \sum_{i=1}^I \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(t) dt + w_2 \int_0^T (\mathbf{b}(t))^\top \tilde{\mathbf{b}}(t) dt. \end{aligned}$$

## The directional derivative w.r.t. $\mathbf{V}(t)$

Consider a perturbation  $\hat{\mathbf{V}}(t)$  and compute:

$$\hat{\mathbf{x}}^i(t) := \lim_{h \rightarrow 0} \frac{\mathbf{x}^{ih}(t) - \mathbf{x}^i(t)}{h}$$

where

$$\begin{aligned}\mathbf{x}^i(0) &= \mathbf{x}_{\text{in}}^i, \\ \mathbf{x}^{ih}(0) &= \mathbf{x}_{\text{in}}^i,\end{aligned}$$

$$\begin{aligned}\dot{\mathbf{x}}^i(t) &= \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t), \\ \dot{\mathbf{x}}^{ih}(t) &= (\mathbf{V}(t) + h\hat{\mathbf{V}}(t))\sigma(\mathbf{x}^{ih}(t)) + \mathbf{b}(t).\end{aligned}$$



## The directional derivative w.r.t. $\mathbf{V}(t)$

Consider a perturbation  $\hat{\mathbf{V}}(t)$  and compute:

$$\hat{\mathbf{x}}^i(t) := \lim_{h \rightarrow 0} \frac{\mathbf{x}^{ih}(t) - \mathbf{x}^i(t)}{h}$$

where

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i,$$

$$\dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t),$$

$$\mathbf{x}^{ih}(0) = \mathbf{x}_{\text{in}}^i,$$

$$\dot{\mathbf{x}}^{ih}(t) = (\mathbf{V}(t) + h\hat{\mathbf{V}}(t))\sigma(\mathbf{x}^{ih}(t)) + \mathbf{b}(t).$$

Because

$$\frac{\dot{\mathbf{x}}^{ih} - \dot{\mathbf{x}}^i}{h} = \frac{(\mathbf{V} + h\hat{\mathbf{V}})\sigma(\mathbf{x}^{ih}) + \mathbf{b} - \mathbf{V}\sigma(\mathbf{x}^{ih}) + \mathbf{V}\sigma(\mathbf{x}^{ih}) - \mathbf{V}\sigma(\mathbf{x}^i) - \mathbf{b}}{h},$$

taking the limit  $h \rightarrow 0$  shows that  $\hat{\mathbf{x}}^i(t)$  satisfies

$$\hat{\mathbf{x}}^i(0) = 0, \quad \dot{\hat{\mathbf{x}}}^i(t) = \hat{\mathbf{V}}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{V}(t)\text{diag}\left(\frac{d\sigma}{dx}(\mathbf{x}^i(t))\right)\hat{\mathbf{x}}^i(t).$$

## The directional derivative w.r.t. $\mathbf{V}(t)$

Consider a perturbation  $\hat{\mathbf{V}}(t)$  and compute:

$$\hat{\mathbf{x}}^i(t) := \lim_{h \rightarrow 0} \frac{\mathbf{x}^{ih}(t) - \mathbf{x}^i(t)}{h}$$

where

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i,$$

$$\mathbf{x}^{ih}(0) = \mathbf{x}_{\text{in}}^i,$$

$$\dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t),$$

$$\dot{\mathbf{x}}^{ih}(t) = (\mathbf{V}(t) + h\hat{\mathbf{V}}(t))\sigma(\mathbf{x}^{ih}(t)) + \mathbf{b}(t).$$

Because

$$\frac{\dot{\mathbf{x}}^{ih} - \dot{\mathbf{x}}^i}{h} = \frac{(\mathbf{V} + h\hat{\mathbf{V}})\sigma(\mathbf{x}^{ih}) + \mathbf{b} - \mathbf{V}\sigma(\mathbf{x}^{ih}) + \mathbf{V}\sigma(\mathbf{x}^{ih}) - \mathbf{V}\sigma(\mathbf{x}^i) - \mathbf{b}}{h},$$

taking the limit  $h \rightarrow 0$  shows that  $\hat{\mathbf{x}}^i(t)$  satisfies

$$\hat{\mathbf{x}}^i(0) = 0, \quad \dot{\hat{\mathbf{x}}}^i(t) = \hat{\mathbf{V}}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{V}(t)\text{diag}\left(\frac{d\sigma}{dx}(\mathbf{x}^i(t))\right) \hat{\mathbf{x}}^i(t).$$

We then obtain

$$\begin{aligned} \langle \nabla_{\mathbf{V}} J, \hat{\mathbf{V}} \rangle_F &= \lim_{h \rightarrow 0} \frac{J(\mathbf{V} + h\hat{\mathbf{V}}, \mathbf{b}) - J(\mathbf{V}, \mathbf{b})}{h} = \sum_{i=1}^I (\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \hat{\mathbf{x}}^i(T) + \\ &\quad w_1 \sum_{i=1}^I \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \hat{\mathbf{x}}^i(t) dt + w_2 \int_0^T \langle \mathbf{V}(t), \hat{\mathbf{V}}(t) \rangle_F dt. \end{aligned}$$

## The adjoint state

Just as in the previous lecture, we need the adjoint state to find the gradients.

We now define the adjoint states  $\varphi^i(t)$ :

$$\varphi^i(T) = \mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i, \quad -\dot{\varphi}^i(t) = \left( \mathbf{V}(t) \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}^i(t)) \right) \right)^\top \varphi^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i),$$

for  $i = 1, 2, 3, \dots, I$ .

Note: the final condition is now nonzero because the state  $\mathbf{x}^i(T)$  at the final time appears in the cost functional.

## Question 1

Write:

$$\mathbf{A}(t) = \mathbf{V}(t) \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}^i(t)) \right).$$

From the previous slides

$$\tilde{\mathbf{x}}^i(0) = 0, \quad \dot{\tilde{\mathbf{x}}}^i(t) = \mathbf{A}(t)\tilde{\mathbf{x}}^i(t) + \tilde{\mathbf{b}}(t).$$

$$\boldsymbol{\varphi}^i(T) = \mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i, \quad -\dot{\boldsymbol{\varphi}}^i(t) = (\mathbf{A}(t))^\top \boldsymbol{\varphi}^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i).$$

What is

$$\int_0^T \frac{d}{dt} \left( (\boldsymbol{\varphi}^i(t))^\top \tilde{\mathbf{x}}^i(t) \right) dt = (\boldsymbol{\varphi}^i(T))^\top \tilde{\mathbf{x}}^i(T) - (\boldsymbol{\varphi}^i(0))^\top \tilde{\mathbf{x}}^i(0)?$$

A) 0

B)  $\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i$

C)  $(\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(T)$

D)  $-(\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(T)$

E) None of the above

## Question 2

We have that:

$$\tilde{\mathbf{x}}^i(0) = 0, \quad \dot{\tilde{\mathbf{x}}}^i(t) = \mathbf{A}(t)\tilde{\mathbf{x}}^i(t) + \tilde{\mathbf{b}}(t).$$

$$\boldsymbol{\varphi}^i(T) = \mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i, \quad -\dot{\boldsymbol{\varphi}}^i(t) = (\mathbf{A}(t))^\top \boldsymbol{\varphi}^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i).$$

What is

$$\int_0^T \frac{d}{dt} \left( (\boldsymbol{\varphi}^i(t))^\top \tilde{\mathbf{x}}^i(t) \right) dt = \int_0^T (\dot{\boldsymbol{\varphi}}^i(t))^\top \tilde{\mathbf{x}}^i(t) dt + \int_0^T (\boldsymbol{\varphi}^i(t))^\top \dot{\tilde{\mathbf{x}}}^i(t) dt?$$

A)  $\int_0^T \left( (\mathbf{A}(t))^\top \boldsymbol{\varphi}^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i) \right)^\top \tilde{\mathbf{x}}^i(t) dt + \int_0^T (\boldsymbol{\varphi}^i(t))^\top \left( \mathbf{A}(t)\tilde{\mathbf{x}}^i(t) + \tilde{\mathbf{b}}(t) \right) dt$

B)  $-\int_0^T \left( (\mathbf{A}(t))^\top \boldsymbol{\varphi}^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i) \right)^\top \tilde{\mathbf{x}}^i(t) dt + \int_0^T (\boldsymbol{\varphi}^i(t))^\top \left( \mathbf{A}(t)\tilde{\mathbf{x}}^i(t) + \tilde{\mathbf{b}}(t) \right) dt$

C)  $w_1 \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(t) dt + \int_0^T (\boldsymbol{\varphi}^i(t))^\top \tilde{\mathbf{b}}(t) dt$

D)  $-w_1 \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(t) dt + \int_0^T (\boldsymbol{\varphi}^i(t))^\top \tilde{\mathbf{b}}(t) dt$

## The gradient w.r.t. $\mathbf{b}(t)$

$$\langle \nabla_{\mathbf{b}} J, \tilde{\mathbf{b}} \rangle = \lim_{h \rightarrow 0} \frac{J(\mathbf{V}, \mathbf{b} + h\tilde{\mathbf{b}}) - J(\mathbf{V}, \mathbf{b})}{h} = \sum_{i=1}^I (\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(T) +$$

$$w_1 \sum_{i=1}^I \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(t) dt + w_2 \int_0^T (\mathbf{b}(t))^\top \tilde{\mathbf{b}}(t) dt.$$

From the previous questions:

$$(\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(T) = -w_1 \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \tilde{\mathbf{x}}^i(t) dt + \int_0^T (\boldsymbol{\varphi}^i(t))^\top \tilde{\mathbf{b}}(t) dt$$

Resulting expression:

$$\langle \nabla_{\mathbf{b}} J, \tilde{\mathbf{b}} \rangle = \sum_{i=1}^I \int_0^T (\boldsymbol{\varphi}^i(t))^\top \tilde{\mathbf{b}}(t) dt + w_2 \int_0^T (\mathbf{b}(t))^\top \tilde{\mathbf{b}}(t) dt$$

Resulting gradient (w.r.t. the standard  $L^2$ -innerproduct):

$$(\nabla_{\mathbf{b}} J)(t) = \sum_{i=1}^I \boldsymbol{\varphi}^i(t) + w_2 \mathbf{b}(t)$$

## The gradient w.r.t. $\mathbf{V}(t)$

$$\begin{aligned} \langle \nabla_{\mathbf{V}} J, \hat{\mathbf{V}} \rangle_F &= \lim_{h \rightarrow 0} \frac{J(\mathbf{V} + h\hat{\mathbf{V}}, \mathbf{b}) - J(\mathbf{V}, \mathbf{b})}{h} = \sum_{i=1}^I (\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \hat{\mathbf{x}}^i(T) + \\ &\quad w_1 \sum_{i=1}^I \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \hat{\mathbf{x}}^i(t) \, dt + w_2 \int_0^T \langle \mathbf{V}(t), \hat{\mathbf{V}}(t) \rangle_F \, dt. \end{aligned}$$

We can now verify similarly as before that

$$(\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i)^\top \hat{\mathbf{x}}^i(T) = -w_1 \int_0^T (\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i)^\top \hat{\mathbf{x}}^i(t) \, dt + \int_0^T (\boldsymbol{\varphi}^i(t))^\top \hat{\mathbf{V}}(t) \sigma(\mathbf{x}^i(t)) \, dt$$

In a similar way, we can show that the gradient w.r.t.  $\mathbf{V}(t)$  (w.r.t. the Frobenius inner product) is

$$(\nabla_{\mathbf{V}} J)(t) = \sum_{i=1}^I \sigma(\mathbf{x}^i(t)) (\boldsymbol{\varphi}^i(t))^\top + w_2 \mathbf{V}(t).$$

## An algorithm for the computation of the gradients

Computation of the gradients  $\nabla_{\mathbf{V}} J(\mathbf{V}, \mathbf{b})$  and  $\nabla_{\mathbf{b}} J(\mathbf{V}, \mathbf{b})$  (gradient in the point  $(\mathbf{V}, \mathbf{b})$ )

- Compute for  $i = 1, 2, 3, \dots, I$  the solutions of

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i, \quad \dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t).$$

- Compute for  $i = 1, 2, 3, \dots, I$  the solutions of

$$\varphi^i(T) = \mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i, \quad -\dot{\varphi}^i(t) = \left( \mathbf{V}(t) \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}^i(t)) \right) \right)^\top \varphi^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i),$$

- The gradients are now given by

$$(\nabla_{\mathbf{V}} J)(t) = \sum_{i=1}^I \sigma(\mathbf{x}^i(t))(\varphi^i(t))^\top + w_2 \mathbf{V}(t), \quad (\nabla_{\mathbf{b}} J)(t) = \sum_{i=1}^I \varphi^i(t) + w_2 \mathbf{b}(t)$$



## An algorithm for the computation of the gradients

Computation of the gradients  $\nabla_{\mathbf{V}} J(\mathbf{V}, \mathbf{b})$  and  $\nabla_{\mathbf{b}} J(\mathbf{V}, \mathbf{b})$  (gradient in the point  $(\mathbf{V}, \mathbf{b})$ )

- Compute for  $i = 1, 2, 3, \dots, I$  the solutions of

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i, \quad \dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t).$$

- Compute for  $i = 1, 2, 3, \dots, I$  the solutions of

$$\varphi^i(T) = \mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i, \quad -\dot{\varphi}^i(t) = \left( \mathbf{V}(t) \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}^i(t)) \right) \right)^\top \varphi^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i),$$

- The gradients are now given by

$$(\nabla_{\mathbf{V}} J)(t) = \sum_{i=1}^I \sigma(\mathbf{x}^i(t)) (\varphi^i(t))^\top + w_2 \mathbf{V}(t), \quad (\nabla_{\mathbf{b}} J)(t) = \sum_{i=1}^I \varphi^i(t) + w_2 \mathbf{b}(t)$$

Remark: when  $I$  is large, we need to solve many ODEs at each iteration. As we also need many iterations (e.g. 10,000) this can lead to a huge computational cost.

**Stochastic gradient descent methods** reduce this cost by considering a randomly selected subset of indices  $i \in \{1, 2, \dots, I\}$  at each time step.

We will not go into this further in this lecture.

## 12.C An algorithm for the training of deep residual neural networks



## Cost functional and dynamics

$$\min_{\mathbf{V}, \mathbf{b}} J(\mathbf{V}, \mathbf{b}) = \frac{1}{2} \sum_{i=1}^I |\mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i|^2 +$$

$$\frac{w_1}{2} \sum_{i=1}^I \int_0^T |\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i|^2 dt + \frac{w_2}{2} \int_0^T (\|\mathbf{V}(t)\|_F^2 + |\mathbf{b}(t)|^2) dt.$$

subject to the dynamics (for  $i = 1, 2, 3, \dots, I$ )

$$\mathbf{x}^i(0) = \mathbf{x}_{\text{in}}^i, \quad \dot{\mathbf{x}}^i(t) = \mathbf{V}(t)\sigma(\mathbf{x}^i(t)) + \mathbf{b}(t).$$

We consider a uniform grid  $t_k = (k - 1)\Delta t$  ( $k = 1, 2, \dots, N_T$ ), so  $\Delta t = T/(N_T - 1)$ .

We denote  $\mathbf{x}_k^i \approx \mathbf{x}^i(t_k)$  with  $k = 1, 2, \dots, N_T$ ,

$\mathbf{V}_k = \mathbf{V}(t_k)$  and  $\mathbf{b}_k = \mathbf{b}(t_k)$  with  $k = 1, 2, \dots, N_T - 1$ .

We discretize with forward Euler:

$$\min_{\mathbf{V}_k, \mathbf{b}_k} J(\mathbf{V}, \mathbf{b}) = \frac{1}{2} \sum_{i=1}^I |\mathbf{x}_{N_T}^i - \mathbf{y}_{\text{out}}^i|^2 +$$

$$\frac{w_1 \Delta t}{2} \sum_{i=1}^I \sum_{k=2}^{N_T-1} |\mathbf{x}_k^i - \mathbf{y}_{\text{out}}^i|^2 + \frac{w_2 \Delta t}{2} \sum_{k=1}^{N_T-1} (\|\mathbf{V}_k\|_F^2 + |\mathbf{b}_k|^2).$$

$$\mathbf{x}_1^i = \mathbf{x}_{\text{in}}^i, \quad \mathbf{x}_{k+1}^i = \mathbf{x}_k^i + \Delta t (\mathbf{V}_k \sigma(\mathbf{x}_k^i) + \mathbf{b}_k).$$

Note: Forward Euler gives us precisely the structure of a ResNet.

## Adjoint state

In the continuous time setting, we could compute the gradient from the adjoint state:

$$\varphi^i(T) = \mathbf{x}^i(T) - \mathbf{y}_{\text{out}}^i, \quad -\dot{\varphi}^i(t) = \left( \mathbf{V}(t) \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}^i(t)) \right) \right)^\top \varphi^i(t) + w_1(\mathbf{x}^i(t) - \mathbf{y}_{\text{out}}^i),$$

Adjoint variables are  $\varphi_k \approx \varphi(t_k)$ ,  $k = 1, 2, 3, \dots, N_T - 1$ .

Compute the adjoint variables starting from

$$\varphi_{N_T-1}^i = \mathbf{x}_{N_T}^i - \mathbf{y}_{\text{out}}^i$$

and then backward in time according to

$$\varphi_{k-1}^i = \varphi_k^i + \Delta t \left( \mathbf{V}_k \text{diag} \left( \frac{d\sigma}{dx}(\mathbf{x}_k^i) \right) \right)^\top \varphi_k^i + \Delta t w_1(\mathbf{x}_k^i - \mathbf{y}_{\text{out}}^i).$$

## Gradient computation

In the continuous time setting:

$$(\nabla_{\mathbf{V}} J)(t) = \sum_{i=1}^I \sigma(\mathbf{x}^i(t)) (\boldsymbol{\varphi}^i(t))^\top + w_2 \mathbf{V}(t),$$

$$(\nabla_{\mathbf{b}} J)(t) = \sum_{i=1}^I \boldsymbol{\varphi}^i(t) + w_2 \mathbf{b}(t)$$

After discretization:

$$(\nabla_{\mathbf{V}} J)_k = \Delta t \sum_{i=1}^I \boldsymbol{\varphi}_k^i \sigma(\mathbf{x}_k^i)^\top + w_2 \Delta t \mathbf{V}_k, \quad k = 1, 2, \dots, N_T - 1,$$

$$(\nabla_{\mathbf{b}} J)_k = \Delta t \sum_{i=1}^I \boldsymbol{\varphi}_k^i + w_2 \Delta t \mathbf{b}_k, \quad k = 1, 2, \dots, N_T - 1.$$

## Optimization algorithm

We can now just use a basic gradient descent algorithm to minimize  $J$ .  
In every iteration we thus use the updates

$$\mathbf{V}_k^{j+1} = \mathbf{V}_k^j - \gamma_j (\nabla_{\mathbf{V}} J)_k, \quad \mathbf{b}_k^{j+1} = \mathbf{b}_k^j - \gamma_j (\nabla_{\mathbf{b}} J)_k.$$

The stepsize  $\gamma_j$  is also called the learning rate.

But the problem is now very much nonconvex:

- ▶ We cannot guarantee the uniqueness of the (global) minimizer.
- ▶ We do not know whether the algorithm converges to a global minimizer.
- ▶ The convergence rate is generally slow.  
We need many iterations (e.g. 10,000) to obtain good results.