

# Introduction to Neural ODE

Della Bona Sarah, Dumez Erika

13 avril 2021

## Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Ordinary Differential Equations</b>	<b>3</b>
2.1	A reminder on ordinary differential equations . . . . .	3
2.2	A simple example . . . . .	4
2.3	Existence and uniqueness of a solution . . . . .	4
2.4	One-step methods . . . . .	4
2.5	Euler's method . . . . .	4
<b>3</b>	<b>Machine Learning</b>	<b>5</b>
3.1	Neural networks . . . . .	5
3.2	Back propagation . . . . .	6
3.3	Example . . . . .	6
3.4	Gradient descent . . . . .	8
3.5	Vanishing and exploding gradient . . . . .	8
3.6	Residual neural network . . . . .	8
3.7	Implicit Layers . . . . .	9
<b>4</b>	<b>Neural ODE</b>	<b>9</b>
4.1	Introduction . . . . .	9
4.2	Forward pass . . . . .	11
4.3	Backward pass : the Adjoint method . . . . .	11
4.4	Example . . . . .	13
4.5	Advantages and disadvantages of ODE-Nets . . . . .	14
<b>5</b>	<b>References</b>	<b>16</b>
<b>6</b>	<b>Annex</b>	<b>17</b>

# 1 Introduction

In this document, we introduce ODE-nets, which are deep neural networks models using ordinary differential equations. We focus in particular on the mathematical aspects of these neural networks. We will give definitions and properties for different notions such as ordinary differential equations, regular and residual neural networks, implicit layers, ...

At the end, we'll conclude with the advantages and disadvantages of ODE-nets.

## 2 Ordinary Differential Equations

### 2.1 A reminder on ordinary differential equations

An *ordinary differential equation* (ODE) is an equation that describes the changes of a function through time. The aim is to compute that function from the ODE which describes its derivative. In this setting, time is a continuous variable.

**Definition 1.** Let

$$\begin{aligned} f : \quad \mathbb{R}^n &\rightarrow \mathbb{R}^m \\ x = (x_1, \dots, x_n) &\mapsto f(x_1, \dots, x_n) \end{aligned}$$

be a function.

The *partial derivative* of  $f$  with respect to the variable  $x_i$  is denoted by

$$\frac{\partial f}{\partial x_i} : \mathbb{R}^n \rightarrow \mathbb{R}^m.$$

For  $a \in \mathbb{R}^n$ , the partial derivative of  $f$  with respect to  $x_i$ , if it exists, is defined as

$$\frac{\partial f}{\partial x_i}(a) = \lim_{h \rightarrow 0} \frac{f(a_1, \dots, a_{i-1}, a_i + h, a_{i+1}, \dots, a_n) - f(a_1, \dots, a_i, \dots, a_n)}{h}.$$

**Definition 2.** Let  $\Omega \subseteq \mathbb{R} \times \mathbb{R}^N$  an open set. Let  $f : \Omega \rightarrow \mathbb{R}^N$ .

A *first order ODE* takes the form

$$\frac{\partial u}{\partial t}(t) = f(t, u(t))$$

— A *solution* for this ODE is a function  $u : I \subseteq \mathbb{R} \rightarrow \mathbb{R}^N$ , where  $I$  is an interval, such that

- $u$  is differentiable on  $I$ ,
- $\forall t \in I, (t, u(t)) \in \Omega$ ,
- $\forall t \in I, \frac{\partial u}{\partial t}(t) = f(t, u(t))$

— An *initial condition* (IC) is a condition of the type

$$u(t_0) = u_0$$

where  $(t_0, u_0) \in \Omega$  is fixed.

— A *Cauchy problem* is an ODE with IC

$$\begin{cases} \frac{\partial u}{\partial t}(t) &= f(t, u(t)) \\ u(t_0) &= u_0 \end{cases}$$

## 2.2 A simple example

Let  $\frac{\partial x}{\partial t}(t) = x(t)$  an ODE. The solutions of this ODE are

$$\{x(t) = ae^t \mid a \in \mathbb{R}\}.$$

Indeed, for all  $a \in \mathbb{R}$  we have

$$\frac{\partial ae^t}{\partial t} = ae^t$$

If we add an initial condition  $x(0) = 1$ , we have a Cauchy problem and its solution is  $e^t$ , since  $e^0 = 1$  and  $\partial_t e^t = e^t$ .

## 2.3 Existence and uniqueness of a solution

If we want to find the solution to an ODE, we need to know the conditions under which this ODE has a solution. Thus, we define *Lipschitz continuous functions*. This notion is crucial for the following theorem which gives conditions for the existence and uniqueness of a solution to an ODE.

**Definition 3.** Let  $(X, d_X)$  and  $(Y, d_Y)$  be two metric spaces. A function  $f : X \rightarrow Y$  is called *Lipschitz continuous* if

$$\exists K \geq 0, \forall x_1, x_2 \in X, d_Y(f(x_1), f(x_2)) \leq K d_X(x_1, x_2).$$

### Theorem 1. Picard-Lindelöf theorem

Consider the Cauchy problem

$$\frac{\partial u}{\partial t}(t) = f(t, u(t)), \quad u(t_0) = u_0.$$

Suppose  $f$  is uniformly Lipschitz continuous in  $u$  and continuous in  $t$ . Then for some value  $T > 0$ , there exists a unique solution  $u(t)$  to the Cauchy problem on the interval  $[t_0, t_0 + T]$ .

## 2.4 One-step methods

Unfortunately, it is not always possible to explicitly find a solution to a Cauchy problem. However, let  $T > 0$  such that the solution  $u$  exists on  $[t_0, t_0 + T]$  and let  $n \geq 2$  be a natural. Let  $t_0 < \dots < t_n \in [t_0, t_0 + T]$  where  $t_n = t_0 + T$ . We can compute a finite number of points  $(u_1, \dots, u_n)$  such that :

$$\forall i \in \{0, \dots, n\}, u_i \approx u(t_i).$$

To compute those points, we use *one-step methods* which compute the points  $u_{i+1}$  from the previous point  $u_i$ , the time  $t_i$  and the *step*  $h_i := t_{i+1} - t_i$ .

## 2.5 Euler's method

Euler's method is a one-step method with a constant step  $h$ . It is similar to a Taylor development<sup>1</sup>, the idea is to compute  $u(t_{i+1})$  using the formula

$$u(t_{i+1}) \approx u(t_i) + h \frac{\partial u}{\partial t}(t_i)$$

---

1. See the annex.

where

$$\frac{\partial u}{\partial t}(t_i) = f(t_i, u(t_i)).$$

### 3 Machine Learning

In a typical machine learning problem, we have an output variable  $Y$  to  $p$  predictors  $X_1, \dots, X_p$ , also called input variable, where  $p \in \mathbb{N} \setminus \{0\}$ . The inputs belongs to an input space  $\mathcal{X}$  and usually  $\mathcal{X} \subset \mathbb{R}^p$ . The output belongs to a output space  $\mathcal{Y}$ . It depends on the problem, for example : if this is a regression problem,  $\mathcal{Y} \subset \mathbb{R}$ . But if we have a classification problem with  $K$  categories,  $\mathcal{Y} = \{1, 2, \dots, K\}$ .

Let's assume that there is some relationship between  $Y$  and  $X = (X_1, \dots, X_p)$ , which can be written in the general form

$$Y = f(X) + \epsilon.$$

Here  $f$  is some fixed but unknown function, called *target function*, of  $X_1, \dots, X_p$  and  $\epsilon$  is a random error term which is independent of  $X$  and has mean zero.

The goal of machine learning is to estimate this function  $f$  as precisely as possible. To do that, we need a *data set* to learn. The data is a set of  $n$  points in  $\mathcal{X} \times \mathcal{Y}$

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_n, y_n)\}.$$

We can predict  $Y$  using

$$\hat{Y} = \hat{f}(X),$$

where  $\hat{f}$  represents our estimate for  $f$ , and  $\hat{Y}$  represents the resulting prediction for  $Y$ .

To determine the precision of an estimation  $\hat{f}$ , we use a *loss function*, which is a function of a prediction and the output given by the target function. Some example of loss functions are

- Square error loss :  $L_1(y, \hat{y}) = (y - \hat{y})^2$ ;
- Absolute error loss :  $L_2(y, \hat{y}) = |y - \hat{y}|$ ;
- Zero-one loss :  $L_3(y, \hat{y}) = \mathbb{1}_{\{(y, \hat{y}) | y \neq \hat{y}\}}(y, \hat{y})$ .

#### 3.1 Neural networks

A *neural network* can be used to solve a machine learning problem. It consists of a series of layers. There are three types of layers :

- The *input* layer
- The *output* layer
- The *hidden* layers

Each layer consist of a certain number of neurons. We give an input  $x$  to the neurons of a layer, they do some calculus and give an output  $z$ . An *activation function* is then applied to this output and obtain a value  $h$  before transmitting it to the next layer thanks to the connections between the neurons of each layer. The simplest example of a neural network layer is

$$h = \sigma(wx + b)$$

where  $\sigma$  is an activation function,  $w$  is a weight matrix and  $b$  a bias vector.

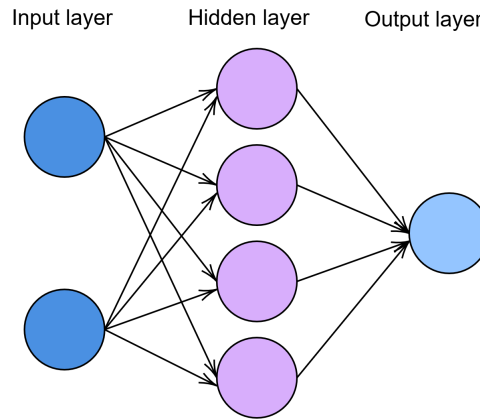


FIGURE 1 – Example of neural network

We begin by giving an input to the input layer, which transmits information to the first hidden layer<sup>2</sup>. In turn, it transmits information to the next layer and so on, until the output layer gives us the final output, the *prediction*. An example of neural network is given in figure 1.

The goal is to minimize the error for every input. To do that, we need to find the optimal parameters for the network which minimize the loss function.

### 3.2 Back propagation

Let  $\theta$  be the parameters of the network. We want to find  $\theta$  which minimize the loss function in order to have the error as small as possible. Therefore, we need to determine the partial derivative of the loss function with respect to the parameters,  $\frac{\partial L}{\partial \theta}$ . Indeed, we know that if the partial derivative of a function is 0 at a certain point, then this point is a local extremum.

Backpropagation is the process used to compute this derivative. It works by computing the gradient of the loss function with respect to each parameter by the chain rule, computing the gradient one layer at a time, iterating backward from the final layer to avoid redundant calculations of intermediate terms in the chain rule.

### 3.3 Example

Let's consider a neural network with one hidden layer that takes a two-dimensional input  $x = (x_1, x_2)$  and gives a 2-dimensional output  $\hat{y} = (\hat{y}_1, \hat{y}_2)$ . We can represent this network with the following equations :

---

2. There isn't always an hidden layer in a neural network

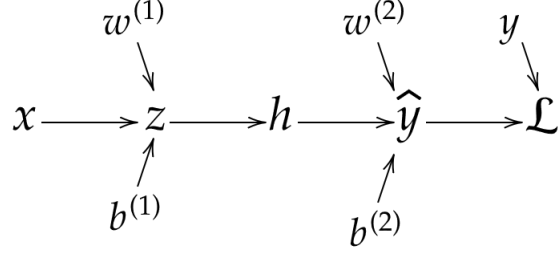


FIGURE 2 – Computation graph

$$\begin{aligned}
 z &= w^{(1)}x + b^{(1)} \\
 h &= \sigma(z) \\
 \hat{y} &= w^{(2)}h + b^{(2)} \\
 \mathcal{L} &= \frac{1}{2} \|\hat{y} - y\|_2^2
 \end{aligned}$$

where  $w^{(1)}, w^{(2)} \in \mathbb{R}^2 \times \mathbb{R}^2$  and  $b^{(1)}, b^{(2)} \in \mathbb{R}^2$  are parameters of the network.

We can now use the backpropagation algorithm to easily compute  $\frac{\partial \mathcal{L}}{\partial w^{(1)}}, \frac{\partial \mathcal{L}}{\partial w^{(2)}}, \frac{\partial \mathcal{L}}{\partial b^{(1)}}, \frac{\partial \mathcal{L}}{\partial b^{(2)}}$ , the partial derivatives of the loss function with regards to the parameters.

$$\begin{aligned}
 \frac{\partial \mathcal{L}}{\partial \mathcal{L}} &= 1 \\
 \frac{\partial \mathcal{L}}{\partial \hat{y}} &= \frac{\partial \mathcal{L}}{\partial \mathcal{L}} (\hat{y} - y) \\
 \frac{\partial \mathcal{L}}{\partial w^{(2)}} &= \frac{\partial \mathcal{L}}{\partial \hat{y}} h^T \\
 \frac{\partial \mathcal{L}}{\partial b^{(2)}} &= \frac{\partial \mathcal{L}}{\partial \hat{y}} \\
 \frac{\partial \mathcal{L}}{\partial h} &= (w^{(2)})^T \frac{\partial \mathcal{L}}{\partial \hat{y}} \\
 \frac{\partial \mathcal{L}}{\partial z} &= \frac{\partial \mathcal{L}}{\partial h} \circ \sigma'(z) \\
 \frac{\partial \mathcal{L}}{\partial w^{(1)}} &= \frac{\partial \mathcal{L}}{\partial z} x^T \\
 \frac{\partial \mathcal{L}}{\partial b^{(1)}} &= \frac{\partial \mathcal{L}}{\partial z}
 \end{aligned}$$

### 3.4 Gradient descent

*Gradient descent* is a process used to find a local minimum of a differentiable function. It works as follow : at each step of the process, we take a step in the opposite direction of the gradient of the function at the current point, because this is the direction of the steepest descent.

More formally, if we have a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $n > 1$ , differentiable and a point  $x_0 \in \mathbb{R}^n$ , we have that if

$$x_{n+1} = x_n - \gamma_n \nabla f(x_n), n \geq 0$$

for  $\gamma_n \in \mathbb{R}^+$  small enough, then  $f(x_n) \geq f(x_{n+1})$ .

We get a sequence  $x_0, x_1, \dots$  that hopefully converges to the desired local minimum, such that

$$f(x_0) \geq f(x_1) \geq \dots$$

If the function  $f$  is convex, all local minima are also global minima, so the gradient descent can converge to the global minimum.

### 3.5 Vanishing and exploding gradient

The problem when using the gradient descent algorithm on neural network is that each weight is updated using the partial derivative of the loss function with regards to the current weight, and if this gradient is too small, it will prevent the weight from changing its value. In this case, the neural network will not be able to learn.

One example of this problem is when we use the hyperbolic tangent as activation function. Because this function has gradients in the range  $]0, 1[$  and backpropagation computes gradients by the chain rule, we multiply several of these small numbers which leads the gradient to decrease exponentially. The deeper is the neural network, the more likely this problem can occur.

The exploding gradient problem is the opposite, it happens when the derivatives take on larger values.

### 3.6 Residual neural network

A *residual neural network*, also called ResNet, is simply a regular neural network except that it has more connections. Not only do we feed the output of the previous layer to the next, but also the input of that layer. An example of the representation of a ResNet is given in Figure 3.

In these networks, the output of the  $k + 1$ th layer is given by

$$x_{k+1} = x_k + f_k(x_k)$$

where  $f_k$  is the function of the  $k$ th layer and its activation.

We can see that this simple formula is a special case of the formula

$$x_{k+1} = x_k + h f_k(x_k),$$

which is the formula for the Euler method for solving ODEs when  $h = 1$ . It is with this observation that we can later introduce neural ODE networks (Section 4).



## ResNet

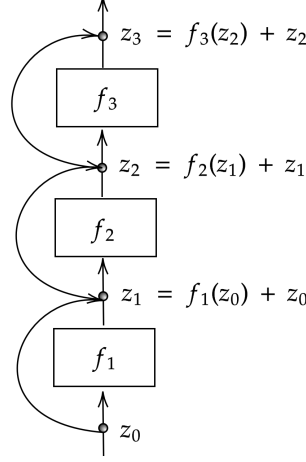


FIGURE 3 – Example of residual neural network

With these additional connections, we can avoid the problems of the *vanishing gradient* and the *exploding gradient* and thus have a better accuracy.

Residual networks avoid the problem of vanishing gradient by introducing short paths which can carry a gradient over the entire extent of very deep networks. This is because adding the information from the previous layer will make these activations larger, so to some extent, they will prevent these activations from becoming exponentially small.

### 3.7 Implicit Layers

There are two different ways to define a layer : *explicitly* or *implicitly*. When we define a layer explicitly, we specify the exact sequence of operations to do from the input to the output layer like in the example of the section 3.3.

However, when we add some functionality to the layers, it can become complex to define them explicitly. Instead, we can define them implicitly : we specify the condition we want the layer's output to satisfy.

An *explicit layer* is defined by a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$ . For an implicit layer, we give a condition that a function  $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^n$  should satisfy. For example we can search for a  $y$  such that  $g(x, y) = 0$ .

## 4 Neural ODE

### 4.1 Introduction

In a residual neural network, the output for an input  $x$  is a composition of functions. We want to extract all these individual layers to only have one "shared" layer.

### ODE-Net

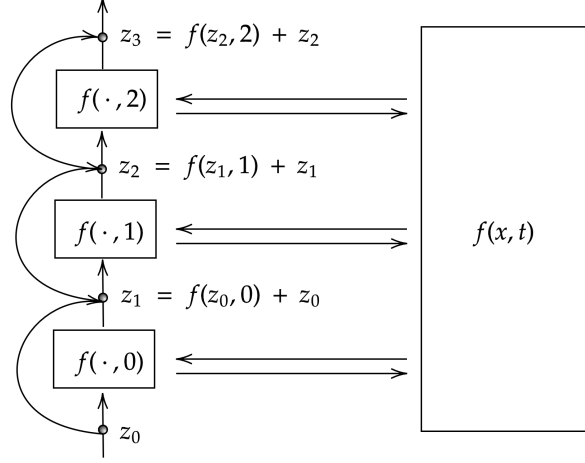


FIGURE 4 – Representation of an ODE-Net

A *neural ODE network* (or ODE-Net) takes a simple layer as a building block. This “base layer” is going to specify the dynamics of an ODE. ODE-Net enable us to replace layers of neural networks with a continuous-depth model. This means that we do not need to specify the number of layers beforehand.

Let us return to ResNets to give intuition behind this definition. We know that any output of the  $k^{th}$  layer of a residual network can be computed with the function

$$F(z_t, t; \theta) = f(z_t, t) + z_t$$

where  $t = k - 1$ .

Thus, in the ResNet, the output for the input  $z_0 = x$  is a composition of the functions  $F(z_t, t; \theta)$  where  $\theta$  represents the parameters of the layers.

We can then view  $z$  as a function of  $t$ . For example,

$$z(1) = f(x, 0) + x.$$

With that, we can write  $F(z_t, t, \theta) = F(z(t), t, \theta)$ . However, we need to give it the initial value of  $z$ , which is  $z(0) = x$  (the input).

We saw that in ResNets, the outputs of each layer are the solutions of an ODE using Euler’s method (cf Section 2.5). The ODE from which it is a solution is  $\frac{\partial z}{\partial t}(t) = f(z(t), t; \theta)$ . But here we want to use a more precise method and then use a more complex ODE solver such as linear multistep methods. With what we’ve just shown, it is possible!

If we consider that the value given by  $f(z(t), t, \theta)$  is the derivative of  $z(t)$ , we obtain the following Cauchy problem :

$$\begin{cases} \frac{\partial z}{\partial t}(t) = f(z(t), t; \theta) \\ z(0) = x \end{cases} \quad (1)$$

## 4.2 Forward pass

Layers in an ODE-Net are implicit. The output  $z(t)$  of a layer in an ODE-Net is defined by the Cauchy problem (1).

But how do we find the solution to this ODE, i.e. the output ? We can simply use an ODE Solver. In the case of the Euler method, the result is equivalent to a residual neural network, as we saw in Section 3.6.

To be able to use an ODE solver we have to make sure that the function satisfies the hypotheses in the theorem of existence and uniqueness (cf Section 2.3). For example, if the activation function used in the network is ReLu, we can't apply the theorem since it is not differentiable at 0.

## 4.3 Backward pass : the Adjoint method

Now that we know how to calculate the output from the input and the parameter  $\theta$ , we need a method to find the optimal  $\theta$  that minimize the loss function.

In regular neural networks, we usually use the gradient descent. However in our case, it is more difficult because we used an ODE solver in the forward pass which is some sort of black box. This is why we are introducing the *adjoint method*. This method computes the gradient by solving a second ODE backwards and is applicable to all ODE solvers.

Let  $L$  be a loss function. To minimize this loss function  $L$ , we need gradients with respect to the parameters  $\theta$ . To achieve that, we first need to determine how the gradient of the loss depends on the hidden state  $z(t)$  for each  $t$ , which is

$$a(t) = \frac{\partial L}{\partial z(t)} \quad (2)$$

This quantity is called the *adjoint*. We would like to determine its dynamics, so we need to compute its derivative with respect to  $t$ .<sup>3</sup>

With a continuous hidden state, we can write the transformation after an  $\varepsilon$  change in time as :

$$z(t + \varepsilon) = \int_t^{t+\varepsilon} f(z(t), t, \theta) dt + z(t) \quad (3)$$

Let  $G : \varepsilon \mapsto z(t + \varepsilon)$ . We can apply the Chain rule and we have

$$\frac{\partial L}{\partial z(t)} = \frac{\partial L}{\partial z(t + \varepsilon)} \frac{\partial z(t + \varepsilon)}{\partial z(t)}.$$

In other words

$$a(t) = a(t + \varepsilon) \frac{\partial G(\varepsilon)}{\partial z(t)} \quad (4)$$

---

3. Here we see the vectors as row vectors.

We can now compute the derivative of  $a(t)$  :

$$\begin{aligned}
\frac{\partial a}{\partial t}(t) &= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t+\varepsilon) - a(t)}{\varepsilon} \text{ by definition of the derivative.} \\
&= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t+\varepsilon) - a(t+\varepsilon) \frac{\partial G(\varepsilon)}{\partial z(t)}}{\varepsilon} \text{ by (4).} \\
&= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t+\varepsilon) - a(t+\varepsilon) \frac{\partial z(t) + \varepsilon f(z(t), t, \theta) + O(\varepsilon^2)}{\partial z(t)}}{\varepsilon} \text{ by Taylor's development of G in 0.} \\
&= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t+\varepsilon) - a(t+\varepsilon) \left( \mathbb{1} + \varepsilon \frac{\partial f(z(t), t, \theta)}{\partial z(t)} + O(\varepsilon^2) \right)}{\varepsilon} \\
&= \lim_{\varepsilon \rightarrow 0^+} \frac{-\varepsilon a(t+\varepsilon) \frac{\partial f(z(t), t, \theta)}{\partial z(t)} + O(\varepsilon^2)}{\varepsilon} \\
&= \lim_{\varepsilon \rightarrow 0^+} -a(t+\varepsilon) \frac{\partial f(z(t), t, \theta)}{\partial z(t)} + O(\varepsilon) \\
&= -a(t) \frac{\partial f(z(t), t, \theta)}{\partial z(t)}
\end{aligned}$$

We now have the dynamics of  $a(t)$

$$\frac{\partial a(t)}{\partial t} = -a(t)^T \frac{\partial f(z(t), t, \theta)}{\partial z(t)} \quad (5)$$

We need to solve an ODE for the adjoint backwards in time. The constraint on the last time point, which is simply the gradient of the loss with respect to this point, has to be specified :

$$a(t_N) = \frac{\partial L}{\partial z(t_N)}$$

Then, the gradients with respect to the hidden state can be calculated at any time, including the initial value :

$$\begin{aligned}
a(t_0) &= a(t_N) + \int_{t_N}^{t_0} \frac{\partial a(t)}{\partial t} dt \text{ by the fundamental theorem of calculus} \\
&= a(t_N) - \int_{t_N}^{t_0} a(t)^T \frac{\partial f(z(t), t, \theta)}{\partial z(t)} dt \text{ par (5)}
\end{aligned}$$

If we want to compute the gradients with respect to the parameters  $\theta$ , we have to evaluate another integral, which depends on both  $z(t)$  and  $a(t)$  :

$$\frac{\partial L}{\partial \theta} = - \int_{t_N}^{t_0} a(t)^T \frac{\partial f(z(t), t, \theta)}{\partial \theta} dt \quad (6)$$

To avoid computing each ODE on its own, we can do all of them at the same time. To do that we can generalize the ODE to :

$$\begin{aligned}
\frac{\partial}{\partial t} \begin{bmatrix} z \\ \theta \\ t \end{bmatrix} (t) &= f_{aug}([z, \theta, t]) := \begin{bmatrix} f([z, \theta, t]) \\ 0 \\ 1 \end{bmatrix}, \\
a_{aug} &:= \begin{bmatrix} a \\ a_\theta \\ a_t \end{bmatrix}, \quad a(t) = \frac{\partial L}{\partial z(t)}, \quad a_\theta(t) = \frac{\partial L}{\partial \theta(t)}, \quad a_t(t) := \frac{\partial L}{\partial t(t)}.
\end{aligned}$$

The jacobian of  $f$  has the form :

$$\frac{\partial f_{aug}}{\partial [z, \theta, t]} = \begin{bmatrix} \frac{\partial f}{\partial z} & \frac{\partial f}{\partial \theta} & \frac{\partial f}{\partial t} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} (t)$$

where each  $\mathbf{0}$  is a matrix of zeros with the corresponding dimensions.

We can use  $a_{aug}$  in (5) and we get :

$$\begin{aligned} \frac{\partial a_{aug}(t)}{\partial t} &= -[a(t) \ a_{\theta}(t) \ a_t(t)] \frac{\partial f_{aug}}{\partial [z, \theta, t]}(t) \\ &= -\left[ a \frac{\partial f}{\partial z} \ a \frac{\partial f}{\partial \theta} \ a \frac{\partial f}{\partial t} \right](t) \end{aligned}$$

We can see that the first element is the adjoint differential equation that we calculated previously in (5). The total gradient with respect to the parameters is given by integrating the second element over the full interval and by setting  $a_{\theta}(t_N) = \mathbf{0}$ . We obtain :

$$\frac{\partial L}{\partial \theta} = a_{\theta}(t_0) = - \int_{t_N}^{t_0} a(t) \frac{\partial f(z(t), t, \theta)}{\partial \theta} dt$$

We can also get gradients with respect to  $t_0$  and  $t_N$  by integrating the last element and by the Chain rule respectively.

$$\begin{aligned} \frac{\partial L}{\partial t_0} &= a_t(t_0) = a_t(t_N) - \int_{t_N}^{t_0} a(t) \frac{\partial f(z(t), t, \theta)}{\partial t} dt \\ \frac{\partial L}{\partial t_N} &= \frac{\partial L}{\partial z(t_N)} \frac{\partial z(t_N)}{\partial t_N} = a(t_N) f(z(t_N), t_N, \theta) \end{aligned}$$

With this generalized method, we have gradients for all possible inputs to a Cauchy problem solver.

In the development above, we assumed that the loss function  $L$  depends only on the last time point  $t_N$ . If function  $L$  depends also on intermediate time points  $t_1, t_2, \dots, t_{N-1}$ , we can repeat the adjoint step for each of the intervals  $[t_{N-1}, t_N], [t_{N-2}, t_{N-1}], \dots, [t_0, t_1]$  in the backward order and sum up the obtained gradients.

In practice, most ODE solvers have the option to output the state  $z(t)$  at multiple times. When the loss depends on these intermediate states, the reverse-mode derivative must be broken into a sequence of separate solves, one between each consecutive pair of output times. At each observation, the adjoint must be adjusted in the direction of the corresponding partial derivative  $\frac{\partial L}{\partial z(t_i)}$ .

#### 4.4 Example

ODE-Net are very useful but they can't approximate any function. For example, the simple function  $f(x) = x^2$  can't be approximated by an ODE-Net because it's not a bijective function.

As we can see in figure 4.4, the network have a problem with negative numbers : the lines can't cross and that causes the negative values to be mapped to 0.

On the other hand, ResNets have no problem approximating this type of function.

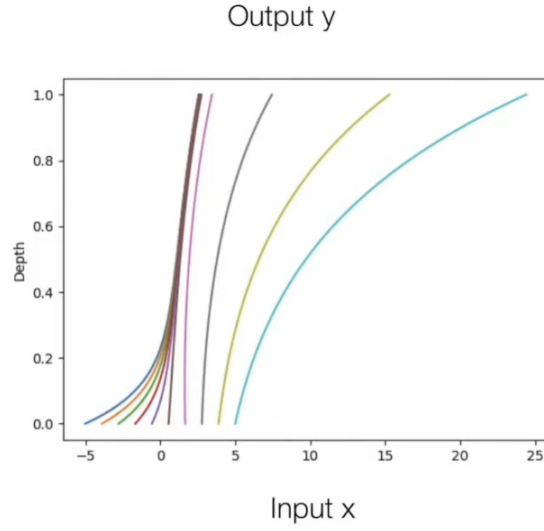


FIGURE 5 – Evolution of the output from the ODE-Net w.r.t. the depth

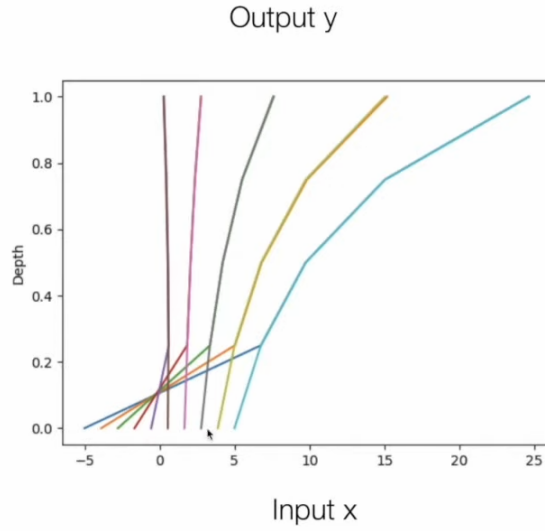


FIGURE 6 – Evolution of the output from the ResNet w.r.t. the depth

#### 4.5 Advantages and disadvantages of ODE-Nets

- The biggest advantage of ODE-Nets is that they have more accurate results for time series predictions. Regular neural network have discrete layers, which means they expect the intervals for these time series data sets to be fixed. Therefore, they are bad at predicting output for time series data that is irregular.
- They have a faster testing time than regular networks, but a slower train-

ning time. There is a trade-off between precision and speed. However, regular neural networks can be evaluated with a fixed amount of computation, and are typically faster to train. In this case, we don't have to choose an error tolerance for a solver.

- We can use ordinary differential equations solvers instead of gradient descent. These solvers have more than a hundred years of theory behind them.
- Lastly, there's a constant memory cost, instead of increasing the cost linearly with each layer in a regular network.

## 5 References

- *Neural Ordinary Differential Equations*, Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, David Duvenaud.  
<https://arxiv.org/pdf/1806.07366.pdf>
- *Deep Implicit Layers - Neural ODEs, Deep Equilibrium Models, and Beyond*, Zico Kolter, David Duvenaud, and Matt Johnson.  
<https://implicit-layers-tutorial.org/>
- *Stanford CS229 : Machine Learning | Autumn 2018*  
<https://www.youtube.com/playlist?list=PLoROMvodv4rMiGQp3WXShtMGgzqpfVfbU>
- *Neural Ordinary Differential Equation (Neural ODE)*, Ayan Das  
<https://ayandas.me/blog-tut/2020/03/20/neural-ode.html>
- Neural ODEs Introduction  
<https://www.youtube.com/watch?v=uPd0B0WhH5w>



## 6 Annex

**Definition 4.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be a function and  $a \in \mathbb{R}^n$ . We can write the *first-order Taylor's development* for  $f$  at  $x$  as :

$$f(x + a) = f(x) + a \cdot \partial f(x) + O(\|a\|^2).$$

**Definition 5.** If  $f$  and  $g$  are differentiable functions, then the *chain rule* expresses the derivative of their composite  $f \circ g$  in terms of the derivatives of  $f$  and  $g$  and the product of functions as follows :

$$\frac{\partial f \circ g}{\partial x} = \left( \frac{\partial f}{\partial x} \circ g \right) \frac{\partial g}{\partial x}$$

**Definition 6.** Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be function,  $n \geq 2$ . Then  $f$  is *convex* if and only if

$$\forall 0 \leq t \leq 1, \forall x_1, x_2 \in \mathbb{R}^n, \quad f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$$