

Introduction to Neural Ordinary Differential Equation

Della Bona Sarah, Dumez Erika

June 15, 2021

Contents

1	Introduction	3
2	Machine Learning	3
2.1	Neural Networks	4
2.2	Back propagation	5
2.3	Example	5
2.4	Gradient descent	6
2.5	Vanishing and exploding gradient	7
2.6	Residual neural network	7
2.7	Implicit Layers	9
3	Ordinary Differential Equations	9
3.1	A simple example	10
3.2	Existence and uniqueness of a solution	10
3.3	One-step methods	10
3.4	Euler's method	10
4	Neural ODE	11
4.1	Introduction	11
4.2	Forward pass	12
4.3	Backward pass: the Adjoint method	12
4.4	Simple Example	15
5	Example with real data	17
6	Advantages and disadvantages of ODE-Nets	18
7	Appendix	20

1 Introduction

In this document, we introduce ODE-Nets, which are deep neural networks models using ordinary differential equations. ODE-Nets are a fairly recent field of study. It uses ODEs which are, on the contrary, very old and have been studied for a long time, and thus are well-known.

This model is very promising since it has advantages over simple neural networks, such as a constant memory cost as well as better results on continuous time series data.

Here, we will 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, ...

At the end, we'll conclude with the advantages and disadvantages of ODE-nets. The code used to make the examples can be found at <https://github.com/DumezErika/ProjetMachineLearning>.

We can begin by the basis and describe what is a machine learning problem.

2 Machine Learning

In a typical machine learning problem [8], we have an output variable Y to p predictors X_1, \dots, X_p , also called input variables, 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} . 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 us assume that there is some relationship between Y and $X = (X_1, \dots, X_p)$, which can be written in the general form

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

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

There is multiple types of machine learning problem. One of them is *supervised learning* and it is the one we will consider in this document. The goal of supervised learning is to estimate the function f as precisely as possible thanks to a model that we will train. To do that, we need a *data set* in order for the model to learn. The data is a set of n points in $\mathcal{X} \times \mathcal{Y}$

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

a sample of the variables.

Let x be a point in the input space \mathcal{X} , then we can predict its output 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 measure the accuracy of a prediction given by \hat{f} , we use a *loss function* \mathcal{L} from \mathbb{R}^2 to \mathbb{R} , which is a function of a prediction and the output given by the target function. Some example of loss functions are

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

To train the model we use a data set called the *train set*. The model is constructed with different parameters θ . During the training we search for the parameters that will minimize a chosen error measure (loss function). Once the model is trained, we can test it on an other data set, called the *test set*, using an error measure which can be a loss function. The error on the test set is called *out-of-sample error* whereas the error on the training set is called the *in-sample error*.

We will focus on one type of model in particular which are neural networks.

2.1 Neural Networks

A *neural network* [1] is 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 and they 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 most used activation functions are :

- Sigmoid : $\text{sigmoid}(z) = \frac{1}{1+e^{-z}}$;
- Hyperbolic tangent : $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$;
- ReLU : $\text{ReLU}(z) = \max(0, z)$.

We use an activation function to add non-linearity to the network.

The simplest example of a neural network layer is

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

where σ is an activation function, W is a weight matrix and b a bias vector.

We begin by giving an input to the input layer, which transmits information to the first hidden layer¹. In turn, it transmit information to the next layer and so on, until the output layer gives us the final output, the *prediction*. It is called the *forward pass*, which is the process of getting a prediction from an input. An example of neural network is given in figure 1.

¹There isn't always a hidden layer in a neural network

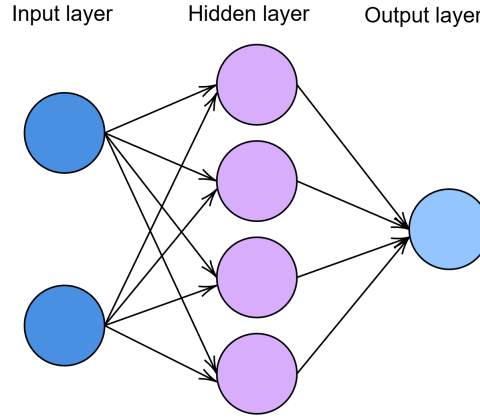


Figure 1: Example of neural network

The goal is to minimize the training error for every input of the training set. To do that, we need to find the optimal parameters for the network which minimize the loss function. We need to compute the derivatives of the loss with respect to the parameters. The process used to get these derivatives is the *backward pass*.

2.2 Back propagation

Let θ be the parameters of the network. We want to find θ 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 [5] 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.

2.3 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:

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

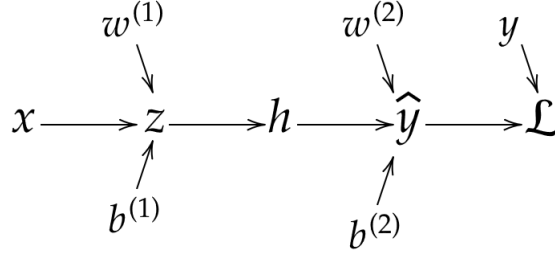


Figure 2: Computation graph

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.

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

2.4 Gradient descent

Gradient descent [8] 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 converges to the desired local minimum under some conditions (q.v. Annex: theorem 2), 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.

2.5 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.

2.6 Residual neural network

A *residual neural network* [4], 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

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

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

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.

We can implement a simple ResNet to approximate the function

$$h(x) = x^3 + 0.1x.$$

ResNet

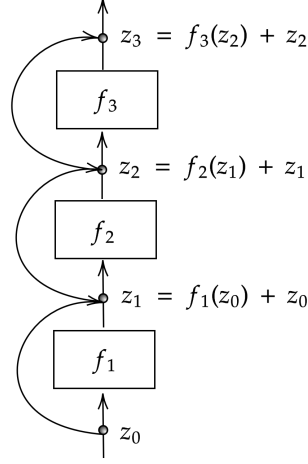


Figure 3: Example of residual neural network

To do that, we generate 10 points between -2.5 and 2.5 . Their associated output comes from the function

$$h(x) + \varepsilon,$$

where ε is a noise variable with mean 0 and standard deviation 1.

We train a ResNet with 3 layers, the hidden one having 20 neurons. Each layer has the hyperbolic tangent as activation function. The learning rate used during the training was 0.01. After 1000 iterations, we get the function given in Figure 4.

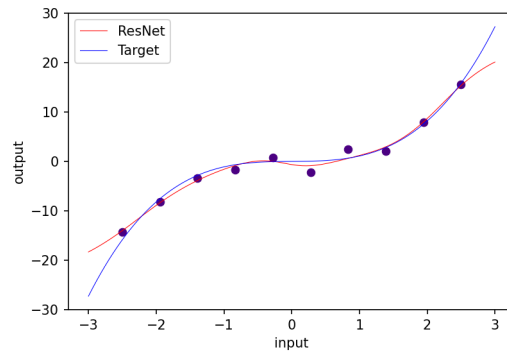


Figure 4: Result of the training for the ResNet

The purple points represent the data used for the training, the blue line is the function we want to approximate and the red line is the function represented by the ResNet. We tested the function returned by the ResNet on the points

we used to trace the blue line. The out-of-sample error for these points is 4.673. The loss function was the mean squared error.

2.7 Implicit Layers

There is two different ways to define a layer : *explicitly* or *implicitly* [7]. 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 2.3.

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 $g(x, y) = 0$.

3 Ordinary Differential Equations

An *ordinary differential equation* (ODE) [9] 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 $\Omega \subseteq \mathbb{R} \times \mathbb{R}^N$ 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

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

3.1 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 $\partial_t e^t = e^t$.

3.2 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 2. 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 $[t_0, t_0 + T]$.

3.3 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* $h_i := t_{i+1} - t_i$.

3.4 Euler's method

Euler's method is a one-step method with a constant step h . It is similar to a Taylor development (q.v. Section 7), the idea is to compute $u(t_{i+1})$ using the formula

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

where

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

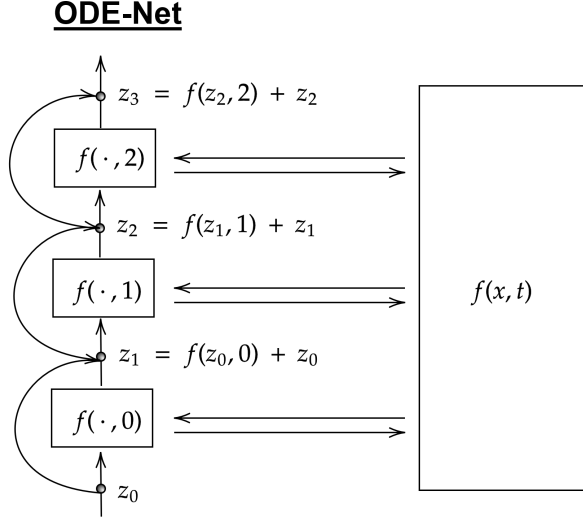


Figure 5: Representation of an ODE-Net

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.

A *neural ODE network* (or ODE-Net) [2–4, 7] 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 θ 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 $z(t_0) = x$ (the input).

We saw that in ResNets, the outputs of each layer are the solutions of an ODE using Euler’s method (cf Section 3.4). 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:

$$\begin{cases} \frac{\partial z}{\partial t}(t) = f(z(t), t; \theta) \\ z(t_0) = x \end{cases} \quad (2)$$

4.2 Forward pass

The layer in an ODE-Net is implicit. The output $z(t_N)$ of an ODE-Net with the input $z(t_0)$ is defined by the Cauchy problem (2). We see that the Cauchy problem depends on the parameters $z(t_0), t_0, t_N, \theta$.

But how do we solve this problem? We can simply use an ODE Solver with the parameters given above. In the case of the Euler method, the result is equivalent to a residual neural network, as we saw in Section 2.6.

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 3.2). 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 of an ODE-Net, we need a method to find the optimal parameters 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* [2]. This method computes the gradient by solving a second ODE backwards and is applicable to all ODE solvers.

Let $\mathcal{L} : \mathbb{R}^p \rightarrow \mathbb{R}$ be a loss function. To minimize this loss function \mathcal{L} , we need gradients with respect to the parameters $z(t_0), t_0, t_N, \theta$. To achieve that, we can determine how the gradient of the loss depends on the hidden state $z(t)$ for each t , which is

$$a(t) = \frac{\partial \mathcal{L}}{\partial z(t)} \quad (3)$$

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

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

$$z(t + \varepsilon) = \int_t^{t+\varepsilon} f(z(t), t, \theta) dt + z(t) \quad (4)$$

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

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

In other words

$$a(t) = a(t + \varepsilon) \frac{\partial G(\varepsilon)}{\partial z(t)} \quad (5)$$

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

$$\begin{aligned} \frac{\partial a}{\partial t}(t) &= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t + \varepsilon) - a(t)}{\varepsilon} \text{ by definition of the derivative.} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t + \varepsilon) - a(t + \varepsilon) \frac{\partial G(\varepsilon)}{\partial z(t)}}{\varepsilon} \text{ by (5).} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t + \varepsilon) - a(t + \varepsilon) \frac{\partial z(t) + \varepsilon f(z(t), t, \theta) + O(\varepsilon^2)}{\partial z(t)}}{\varepsilon} \text{ by Taylor's development of G in 0.} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{a(t + \varepsilon) - a(t + \varepsilon) \left(\mathbf{1} + \varepsilon \frac{\partial f(z(t), t, \theta)}{\partial z(t)} + O(\varepsilon^2) \right)}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{-\varepsilon a(t + \varepsilon) \frac{\partial f(z(t), t, \theta)}{\partial z(t)} + O(\varepsilon^2)}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} -a(t + \varepsilon) \frac{\partial f(z(t), t, \theta)}{\partial z(t)} + O(\varepsilon) \\ &= -a(t) \frac{\partial f(z(t), t, \theta)}{\partial z(t)} \end{aligned}$$

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

$$\frac{\partial a(t)}{\partial t} = -a(t) \frac{\partial f(z(t), t, \theta)}{\partial z(t)} \quad (6)$$

As we are searching for $a(t_0) = \frac{\partial L}{\partial z(t_0)}$, we need to solve an ODE for the adjoint backwards in time because the value for $a(t_N)$ is already known. The constraint on the last time point, which is simply the gradient of the loss with respect to $z(t_N)$,

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

has to be specified to the ODE solver. Then, the gradients with respect to the hidden state can be calculated at any time, including the initial value. We have

$$\begin{aligned} a(t_0) &= a(t_N) + \int_{t_N}^{t_0} \frac{\partial a(t)}{\partial t} dt \text{ by the fundamental theorem of calculus} \\ &= a(t_N) - \int_{t_N}^{t_0} a(t) \frac{\partial f(z(t), t, \theta)}{\partial z(t)} dt \text{ par (6).} \end{aligned}$$

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

$$\frac{\partial L}{\partial \theta} = - \int_{t_N}^{t_0} a(t) \frac{\partial f(z(t), t, \theta)}{\partial \theta} dt. \quad (7)$$

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

$$\frac{\partial}{\partial t} \begin{bmatrix} z \\ \theta \\ t \end{bmatrix} (t) = f_{aug}([z(t), \theta, t]) := \begin{bmatrix} f([z(t), \theta, t]) \\ 0 \\ 1 \end{bmatrix},$$

$$a_{aug}(t) := \begin{bmatrix} a \\ a_\theta \\ a_t \end{bmatrix} (t), \quad a(t) = \frac{\partial \mathcal{L}}{\partial z(t)}, \quad a_\theta(t) = \frac{\partial \mathcal{L}}{\partial \theta(t)}, \quad a_t(t) := \frac{\partial \mathcal{L}}{\partial t(t)}.$$

The jacobian of f has the form

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

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

We can inject a_{aug} in (6) and we get

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

We can see that the first component, $-a(t) \frac{\partial f(z(t), t, \theta)}{\partial z(t)}$, is the adjoint differential equation that we calculated previously in (6). The total gradient with respect to the parameters is given by integrating the second component, $-a(t) \frac{\partial f(z(t), t, \theta)}{\partial \theta(t)}$ over the full interval and by setting $a_\theta(t_N) = \mathbf{0}$. We obtain

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

We can also get gradients with respect to t_0 and t_N by integrating the last component, $-a(t) \frac{\partial f(z(t), t, \theta)}{\partial t(t)}$, and by the Chain rule respectively. We have

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

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.

We can represent this process with the Figure 6. As we can see, during the forward pass, the loss function is evaluated at each time to be able to determine the smallest error.

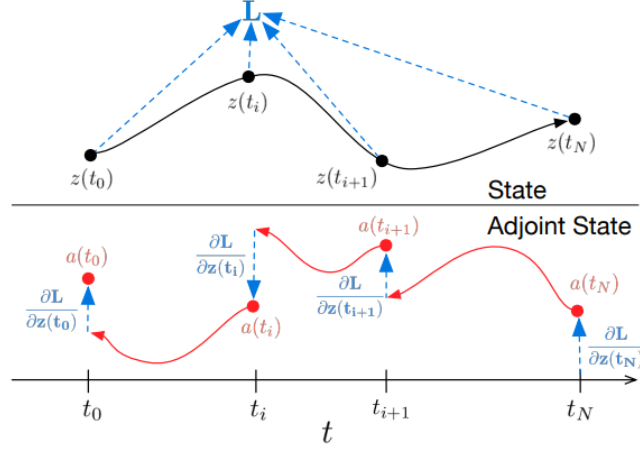


Figure 6: Graphical representation of the forward and backward pass for the ODE-Net. The figure has been taken from the paper [2].

4.4 Simple Example

We can now use the same example than for ResNets. We have the function

$$h(x) = x^3 + 0.1x$$

that we wish to approximate. We use the same training data as for the ResNet.

The dynamics of the ODE-Net is specified by a layer of size 20. The learning rate in this case was also 0.005. Each layer has the hyperbolic tangent as activation function. After 1500 iteration, we get the function given in Figure 7.

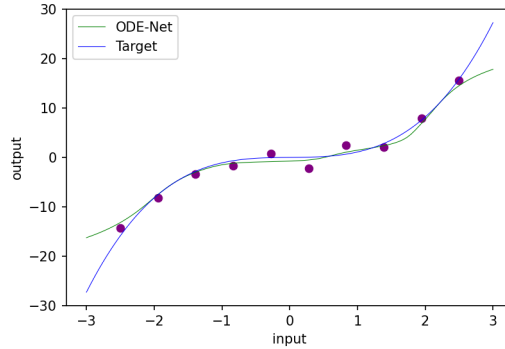


Figure 7: Result of the training for the ODE-Net

The purple points represent the data used for the training, the blue line is the function we want to approximate and the red line is the function represented by the ODE-Net. With the mean squared error, we had an out-of-sample error of 7.72.

We can compare these results with those we had for the ResNet, we can see that the ResNets is slightly better with these parameters. The comparison

graph is in the Figure 8.

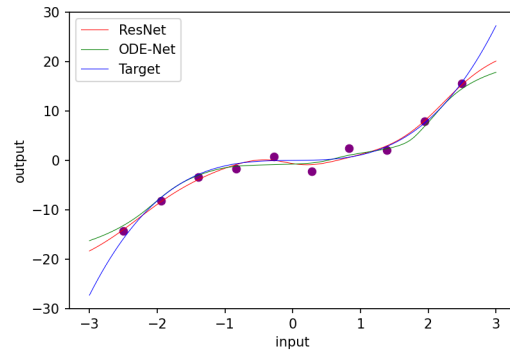


Figure 8: Comparison of ResNet and ODE-Net

5 Example with real data

Now that we have tested the ODE-Net architecture on different small examples, we can test it on real data to see how it performs.

The data we chose is the MIT Beth Israel Hospital (BIH) electrocardiogram dataset that can be found at <https://www.kaggle.com/shayanfazeli/heartbeat>. It contains two data sets: one for training and one for testing.

This dataset consists of around 110,000 samples, classified as either

- 0: normal;
- 1: supraventricular premature beat;
- 2: premature ventricular contraction;
- 3: fusion of ventricular and normal beat;
- 4: unclassified beat.

We will build a ResNet and an ODE-Net as similar as possible so that we can compare them better.

- The Resnet is composed of 3 first linear layers with ReLU as activation function, followed by 6 residual block also with ReLU. The output layer is a simple linear layer (with the input flattened).
- The ODE-Net is composed similarly of 3 first layers, each followed by the activation function ReLU. The dynamics are defined by ODEfunc, which is a linear layer taking time into account. The output layer is the same as for the ResNet, i.e. a simple linear layer with the input flattened.

To optimize our model we use the cross entropy error. The parameters are then updated using stochastic gradient descent with a learning rate of 0.1 and a momentum of 0.9. We trained each models for 20 epochs with batches of size 128.

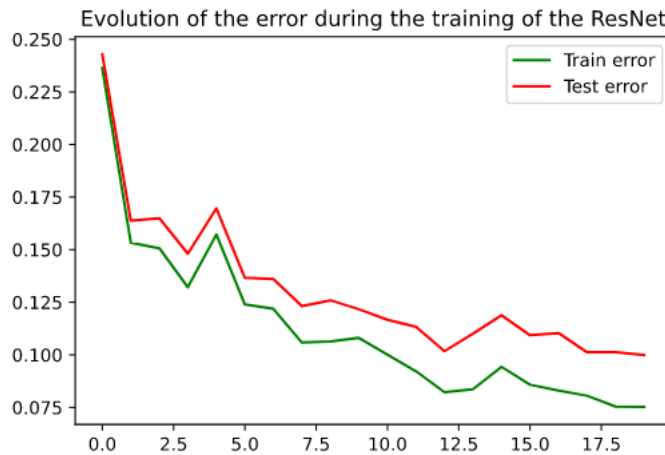


Figure 9: Evolution of the error during training of the ResNet

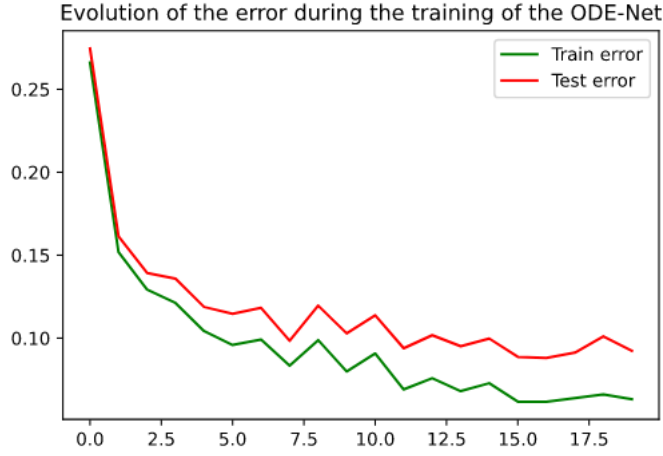


Figure 10: Evolution of the error during training of the ODE-Net

In Figure 9, we can see the evolution of the training and test loss. We tested our models on test data and used the cross entropy error to evaluate it. At the end of the training, the Resnet had a test error of 0.1.

For the ODE-Net, we can see the evolution of the loss on Figure 10. At the end of the training, the model had a test error of 0.09. Thus the ODE-Net is slightly better than the ResNet. Moreover, the figures show that the ODE-Net learns faster since the curve decreases in less epochs.

We also computed the accuracy of each models. To do so we first used softmax to compute predictions for the test data, then we used the zero-one loss on those predictions. Both the ResNet and the ODE-Net have an accuracy of 97,5%, which is a good accuracy.

6 Advantages and disadvantages of ODE-Nets

Advantages

- *Continuous time series predictions*

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 whereas ODE-Net have a continuous layer which means we can evaluate the hidden states at every point t in time. Therefore, regular neural networks are bad at predicting output for time series data that is irregular.

- *ODE solvers*

We can use ordinary differential equations solvers instead of gradient descent. These solvers have more than a hundred years of theory behind them which is a great advantage against gradient descent.

- *Robustness* [10]

After experimenting, it was proved that ODE-Net are very robust against perturbed data compared to regular neural network. Two experiments were conducted: in the first one they trained an ODE-Net and a convolutional neural network² on real images without perturbations. They tested these models on the original images and the ODE-Net outperformed the CNN. In the second experiment, they trained these networks on the original and perturbed images. Again, the ODE-Net was much better.

- *Constant memory cost*

Lastly, there's a constant memory cost, instead of increasing the cost linearly with each layer in a regular network. In ODE-Net, we know the state at every time t . Because of that, we can always reconstruct the entire trajectory of an ODE forwards and backwards in time only by knowing this point. This means that ODE-Nets can be trained with a memory cost constant in the number of evaluations of f . There is a trade-off between memory cost and computational time: ResNets are faster but use more memory and ODE-Net are slower but use less memory.

Disadvantages

- *Slower training time*

ODE-Net have a slower training time. Indeed, during training, the dynamics we want to learn tend to become expensive to solve since the network becomes deeper. 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.

There is then a trade-off between accuracy and computational time: if we choose a small error tolerance, then the computational time be bigger.

- *More Hyperparameters*

In ODE-Nets we need to choose a solver and its error tolerance, which induces more choices to find the parameters which works better.

- *Restriction on activation functions*

To ensure that the ODE has a solution we have to make sure the dynamics are uniformly continuous Lipschitz (q.v. Theorem 3.2). This is why we mostly use \tanh as an activation function.

²A neural network that is usually good with images.

7 Appendix

Definition 3. 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 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. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be function, $n \geq 2$. Then f is *convex* if and only if

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

Definition 6. 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 7. A smooth function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L-Lipschitz if for all $x, y \in \mathbb{R}^d$,

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2.$$

Theorem 2. If $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a function L-Lipschitz convex and if $x^* = \operatorname{argmin}_x f(x)$, then the gradient descent algorithm with step-size $\eta \leq \frac{1}{L}$ satisfy

$$f(x_k) \leq f(x^*) + \frac{\|x_0 - x^*\|_2}{2\eta k}.$$

References

- [1] Stanford cs229: Machine learning | autumn 2018, lecture 11, 12, 13. <https://www.youtube.com/playlist?list=PLoROMvody4rMiGQp3WXShMGgzqpfVfbU>.
- [2] Ricky T. Q. Chen, Yulia Rubanova, Jesse Bettencourt, and David Duvenaud. Neural ordinary differential equations. <https://arxiv.org/pdf/1806.07366.pdf>.
- [3] Ayan Das. Neural ordinary differential equation (neural ode). <https://ayandas.me/blog-tut/2020/03/20/neural-ode.html>.
- [4] Andriy Drozdyuk. Neural odes introduction. <https://www.youtube.com/watch?v=uPd0B0WhH5w>. Carleton University.
- [5] Roger Grosse. Csc321 lecture 6: Backpropagation. http://www.cs.toronto.edu/~rgrosse/courses/csc321_2017/slides/lec6.pdf.
- [6] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. An introduction to statistical learning with applications in r. Springer.
- [7] Zico Kolter, David Duvenaud, and Matt Johnson. Deep implicit layers - neural odes, deep equilibrium models, and beyond. <https://implicit-layers-tutorial.org/>.
- [8] Souhaib Ben Taieb. Machine learning. Université de MONS, 2019-2020.
- [9] Christophe Troestler. Introduction à l'analyse numérique. Université de MONS, 2018-2019.
- [10] Hanshu Yan, Jiawei Du, Vincent Tan, and Jiashi Feng. On robustness of neural ordinary differential equations. <https://openreview.net/forum?id=B1e9Y2NYvS>.
- [11] Amer Zayegh and Nizar Al Bassam. Neural network principles and applications. https://www.researchgate.net/publication/329264107_Neural_Network_Principles_and_Applications.