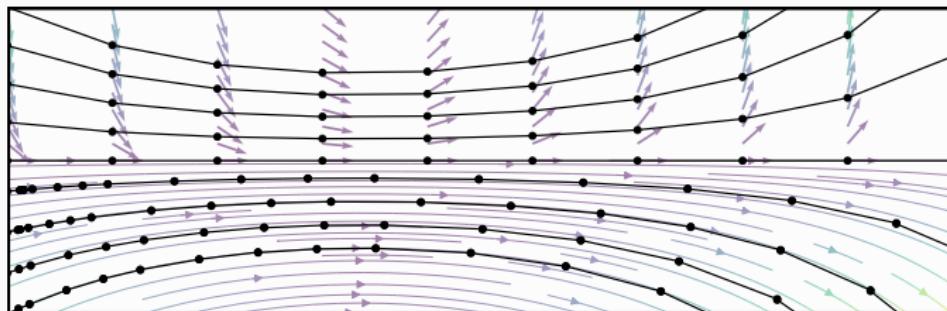


NEURAL ODES - SEMINAR TALK

RADA BERONOVÁ & GEORG TREDE



UNIVERSITY OF HEIDELBERG

JUNE 13, 2023

OVERVIEW

1. Introduction
 - 1.1 Neural Networks
 - 1.2 ResNets
 - 1.3 Limitations of classical ODEs
2. Applying Machine Learning to ODEs
 - 2.1 NNs
 - 2.2 ResNets
3. Neural ODEs
 - 3.1 Concept
 - 3.2 Training
 - 3.3 Applications and Examples
 - 3.4 Advantages and Disadvantages
 - 3.5 Outlook
4. Summary and Takeaways

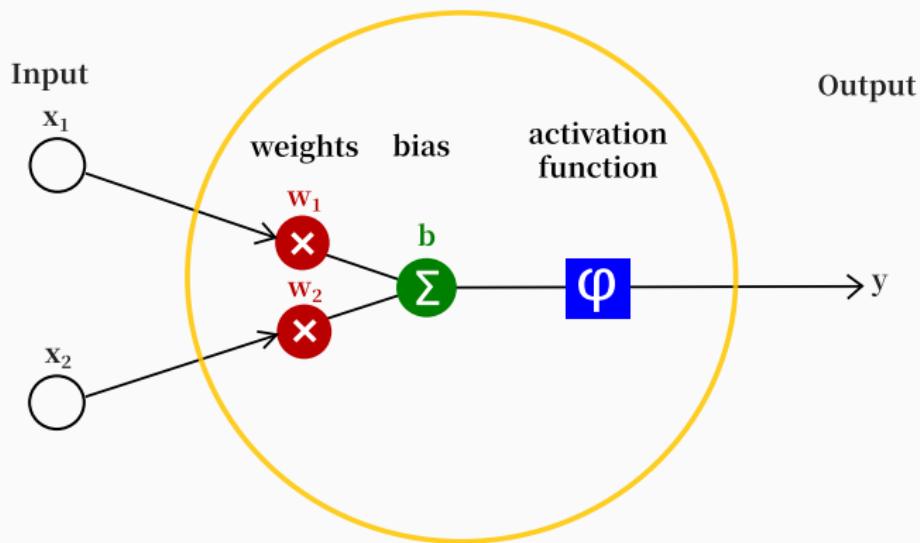
INTRODUCTION

NEURAL NETWORKS

Basic idea

Input data is processed through different layers of artificial neurons stacked together to produce the desired output.

NEURAL NETWORKS - NEURON



$$\varphi(x_1 \cdot w_1 + x_2 \cdot w_2 + b) = y$$

NEURAL NETWORKS - NEURON

What does feedforward mean?

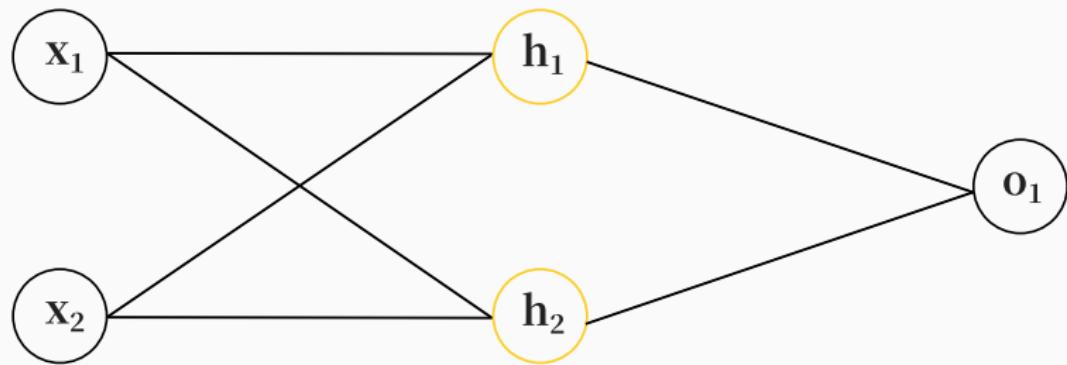
The process of passing inputs through the model to get an output is known as **feedforward**.

NEURAL NETWORKS - EXAMPLE: FEEDFORWARD

Input Layer

Hidden Layer

Output Layer



NEURAL NETWORKS - EXAMPLE: FEEDFORWARD

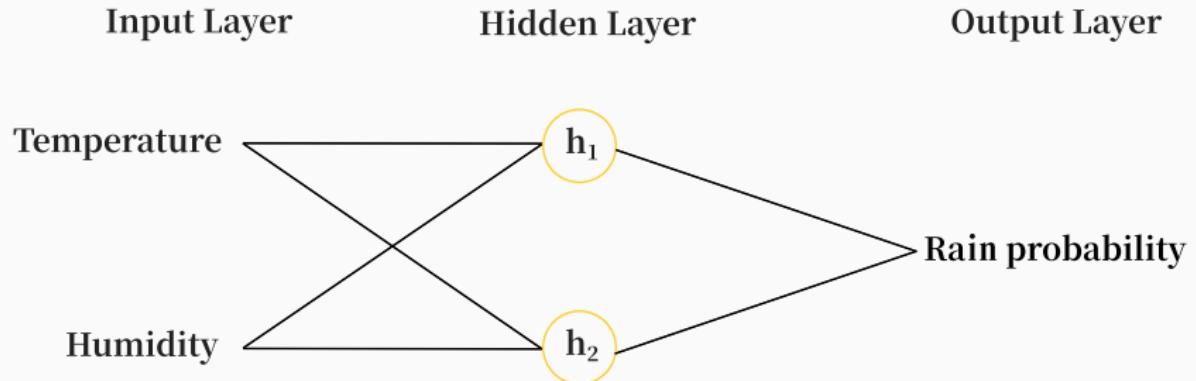
What happens when we pass in the input?

$$h_1 = \phi(x_1 \cdot w_{1,1} + x_2 \cdot w_{2,1} + b_1)$$

$$h_2 = \phi(x_1 \cdot w_{1,2} + x_2 \cdot w_{2,2} + b_2)$$

$$\Rightarrow o_1 = \phi(h_1 \cdot w_{h1} + h_2 \cdot w_{h2} + b_3)$$

NEURAL NETWORKS - TRAINING



NEURAL NETWORKS - TRAINING

What is a **loss function**?

A function that compares the target and predicted output values.

Example: **Mean-squared-error** loss

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_{true} - y_{pred})^2 = L(y_{pred})$$

⇒ **Training a neural network = Minimising its Loss**

NEURAL NETWORKS - TRAINING

How can we minimise the loss?

Backpropagation

How would loss change if we change $w_{1,1}$?

$$\frac{\partial L}{\partial w_{1,1}} = \frac{\partial L}{\partial y_{pred}} \cdot \frac{\partial y_{pred}}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{1,1}}$$

⇒ This process of calculating partial derivatives by working backwards is known as **backpropagation**.

NEURAL NETWORKS - TRAINING

How can we minimise the loss?

Problem

How does the neural network know how to change the weights and biases to minimise loss?

Solution

We can use an optimisation algorithm!

For example, **Gradient descent**:

$$w_{1,1} \leftarrow w_{1,1} - \eta \frac{\partial L}{\partial w_{1,1}} \quad \text{with } \eta: \text{learning rate}$$

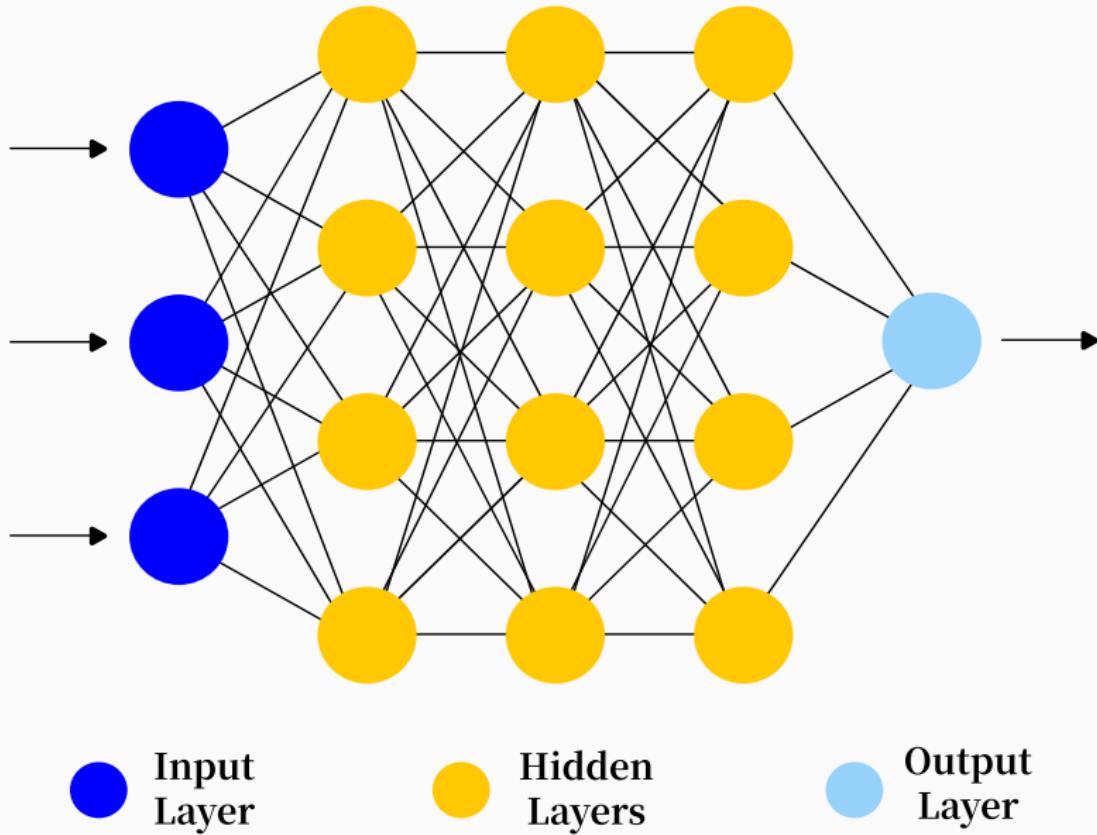
⇒ If we do this for **every weight and bias** in the network, the loss will slowly decrease, and our **network will improve**.

NEURAL NETWORKS - TRAINING

The whole training process:

1. Choose one sample from given dataset
2. Calculate the forward propagation
3. Compare the calculated output to the expected output (loss)
4. Calculate all the partial derivatives of loss with respect to weights and biases (backpropagation)
5. Use the update equation to adjust each weight and bias
6. Go back to step 1

DEEP NEURAL NETWORKS

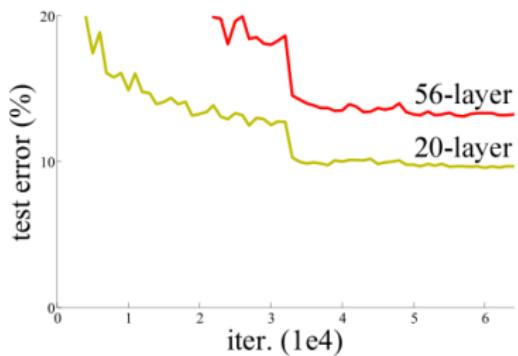
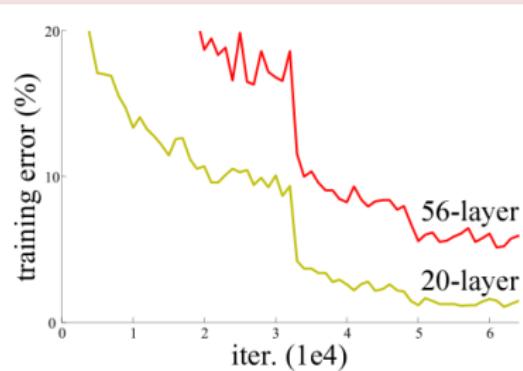


DEEP NEURAL NETWORKS

Problem: vanishing gradient

With **increasing** network **depth**:

- **Decreasing gradient** during backpropagation
- **Accuracy** gets saturated and then **degrades** rapidly

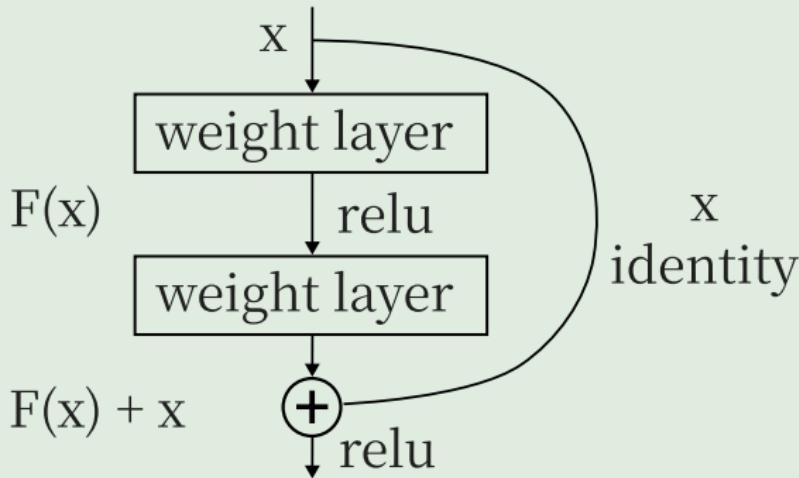


RESIDUAL NETWORKS (RESNETS)

Solution: ResNets

Idea: Breaking a deep NN down into small chunks of network that are connected through skip links.

⇒ **Residual blocks**



⇒ **Skip connections**

RESIDUAL NETWORKS (RESNETS)

Residual blocks

Idea: Instead of letting the layers learn the underlying mapping, let the network fit the residual mapping.

$$\mathcal{F}(x) := \mathcal{H}(x) - x$$

$\mathcal{H}(x)$: underlying mapping, $\mathcal{F}(x)$: fitted mapping

$$\Rightarrow y = \mathcal{H}(x) = \mathcal{F}(x) + x$$

⇒ possibility to just learn identity function

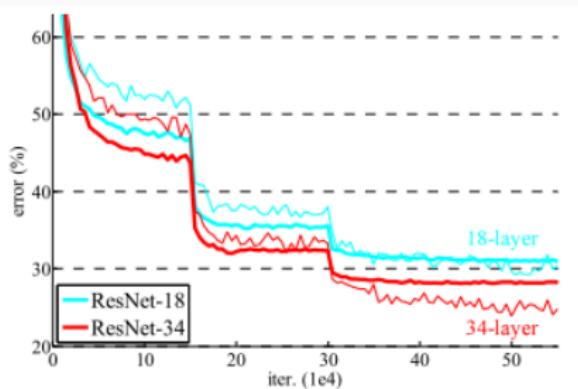
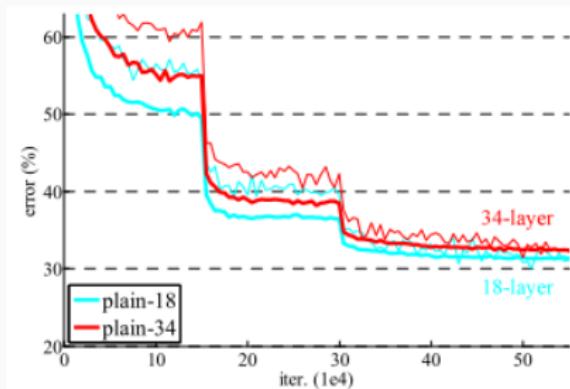
→ input = output

→ "skipping" layers without degrading network performance

RESIDUAL NETWORKS (RESNETS)

Main takeaways:

- By **stacking** many residual blocks, we can form a very deep NN with a lower error rate
- The **skip-connections** allow the network to skip any unneeded layers (regularisation)
 - ⇒ Help with vanishing gradient
 - ⇒ **Improved accuracy and faster convergence**



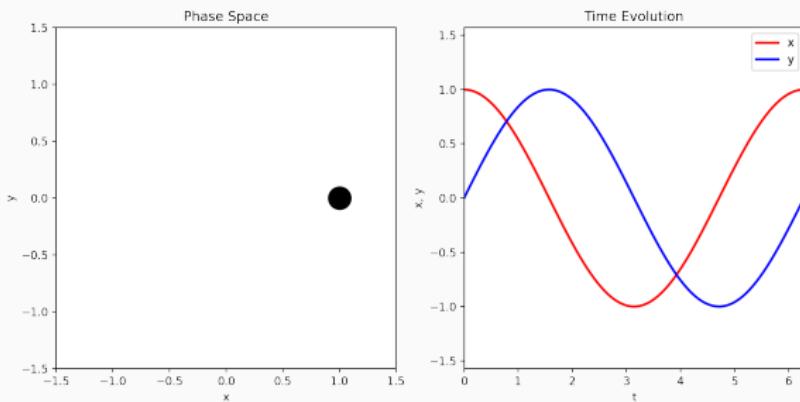
LIMITATIONS OF CLASSICAL ODES

You're asked!

Can you guess the underlying ODEs?

LIMITATIONS OF CLASSICAL ODES

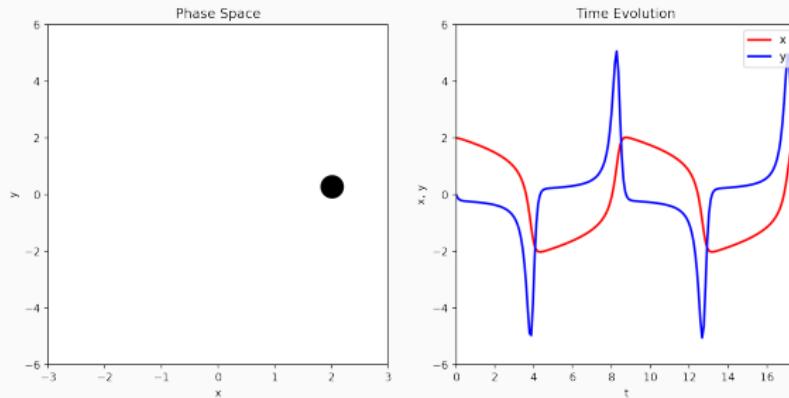
A few examples:



$$\begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} -y \\ x \end{pmatrix}$$

LIMITATIONS OF CLASSICAL ODES

A few examples:

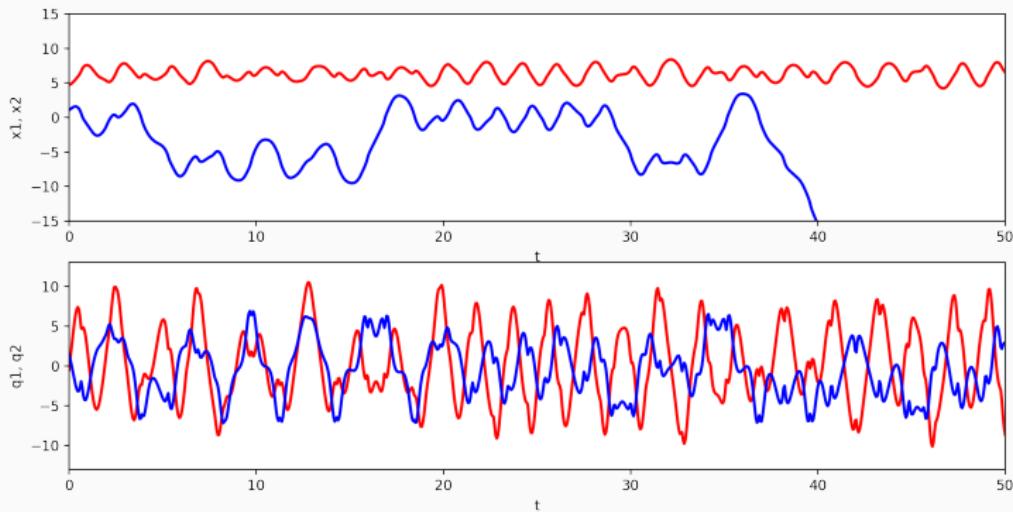


$$\ddot{x} - \epsilon(1-x^2)\dot{x} + x = 0$$

$$\Leftrightarrow \begin{pmatrix} \dot{x} \\ \dot{y} \end{pmatrix} = \begin{pmatrix} y \\ \epsilon(1-x^2)y - x \end{pmatrix}$$

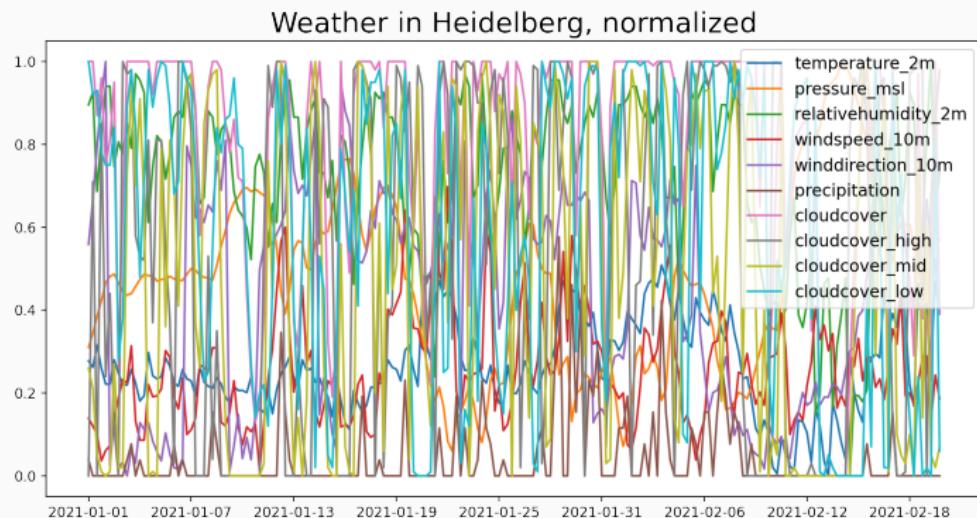
LIMITATIONS OF CLASSICAL ODES

A few examples:



LIMITATIONS OF CLASSICAL ODES

A few examples:



LIMITATIONS OF CLASSICAL ODES

Problem

Classical ODEs are a great concept, but often not practical with real world data

How can we solve this?

Blackbox mystery stuff Machine Learning!

APPLYING ML TO ODES

APPLYING MACHINE LEARNING TO ODEs - NNs

Idea

Approximate next state \mathbf{x}_{t+1} from current state \mathbf{x}_t using a neural network \mathbf{f} :

$$\mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{w})$$

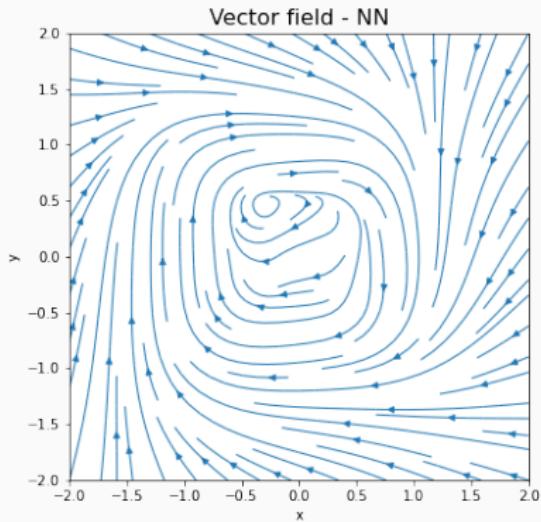
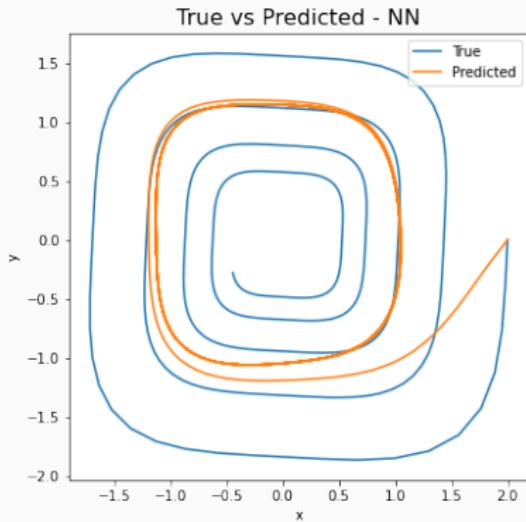
APPLYING MACHINE LEARNING TO ODEs - NNs

$$\mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{w})$$

```
1 from torch import nn
2
3 class NN(nn.Module):
4     def __init__(self, dim, n, hidden_layers):
5         layers = [nn.Linear(dim, n)]
6         for _ in range(hidden_layers):
7             layers.append(nn.Linear(n, n))
8         layers += [nn.Linear(n, dim)]
9         self.layers = nn.ModuleList(layers)
10        self.act = nn.Tanh()
11
12    def forward(self, x):
13        for l in self.layers:
14            x = self.act(l(x))
15        return x
```

APPLYING MACHINE LEARNING TO ODEs - NNs

Results:



APPLYING MACHINE LEARNING TO ODEs - NNs

Good

Learn time evolution from data, i.e. no prior knowledge needed

Bad

Inefficient, we have to cover a wide range of possible outcomes

APPLYING MACHINE LEARNING TO ODEs - RESNETS

Idea

Approximate only the residual for next state \mathbf{x}_{t+1} from current state \mathbf{x}_t using a neural network \mathbf{f} :

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \mathbf{f}(\mathbf{x}_t, \mathbf{w})$$

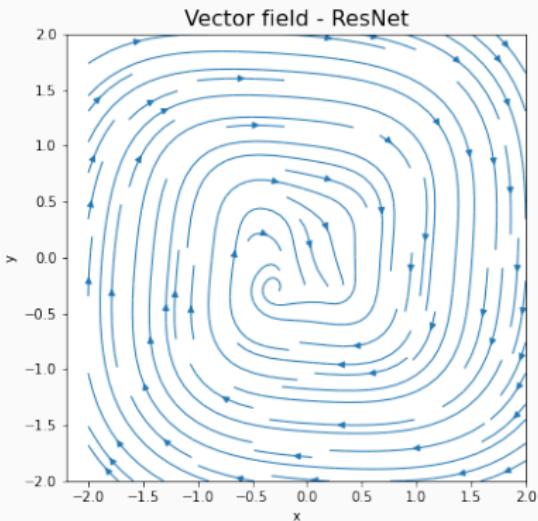
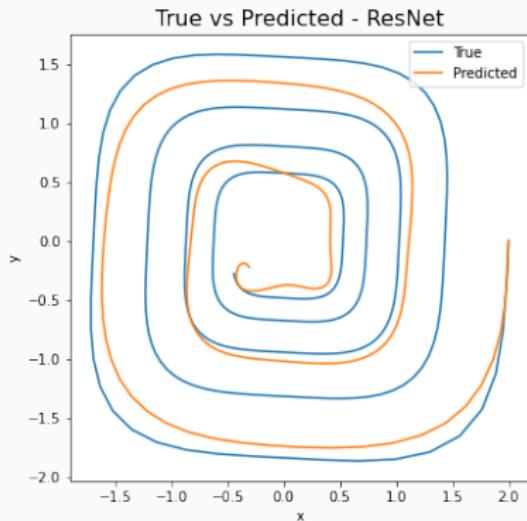
APPLYING MACHINE LEARNING TO ODEs - RESNETS

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \mathbf{f}(\mathbf{x}_t, \mathbf{w})$$

```
1 from torch import nn
2
3 class ResNet(nn.Module):
4     def __init__(self, dim, n, hidden_layers):
5         layers = [nn.Linear(dim, n)]
6         for _ in range(hidden_layers):
7             layers.append(nn.Linear(n, n))
8         layers += [nn.Linear(n, dim)]
9         self.layers = nn.ModuleList(layers)
10        self.act = nn.Tanh()
11
12    def forward(self, x):
13        y = x
14        for l in self.layers:
15            x = self.act(l(x))
16        return x + y
```

APPLYING MACHINE LEARNING TO ODEs - RESNETS

Results:



APPLYING MACHINE LEARNING TO ODEs - RESNETS

Good

More accurate representation with same number of weights
(compared to plain NNs)

Bad

Time step is fixed

- What about states in between?
- What about data points at irregular times?
- We basically created maps/applied the Euler method

NEURAL ODES

NEURAL ODEs - CONCEPT

Let's go to the limit

Time step $\Delta t \rightarrow 0$

$$\mathbf{x}(t+1) = \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), \mathbf{w})$$

$$\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), \mathbf{w}) \cdot \Delta t$$

$$\mathbf{f}(\mathbf{x}(t), \mathbf{w}) = \frac{\mathbf{x}(t + \Delta t) - \mathbf{x}(t)}{\Delta t}$$

$$\xrightarrow{\Delta t \rightarrow 0} \mathbf{f}(\mathbf{x}(t), \mathbf{w}) = \frac{d\mathbf{x}(t)}{dt}$$

NEURAL ODEs - CONCEPT

Let's go to the limit

Time step $\Delta t \rightarrow 0$

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{f}(\mathbf{x}(t), \mathbf{w})$$

This is amazing, we've got an ODE with our neural network ($\mathbf{f}(\mathbf{x}, \mathbf{w})$) in it!

NEURAL ODEs - CONCEPT

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}(t), \mathbf{w})$$

```
1 from torch import nn, tensor
2
3 class ODEFunc(nn.Module):
4     def __init__(self, dim, n, hidden_layers):
5         layers = [nn.Linear(dim, n)]
6         # for _ in range(hidden_layers):
7         #     layers.append(nn.Linear(n,n))
8         layers += [nn.Linear(n, dim)]
9         self.layers = nn.ModuleList(layers)
10        self.act = nn.Tanh()
11
12    def forward(self, x):
13        for l in self.layers:
14            x = self.act(l(x))
15        return x
```

NEURAL ODEs - CONCEPT

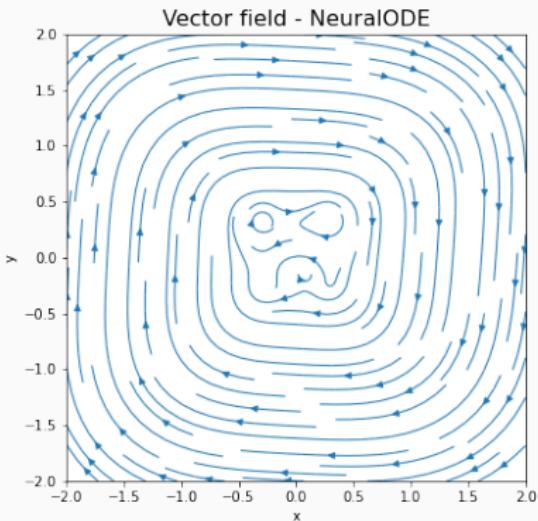
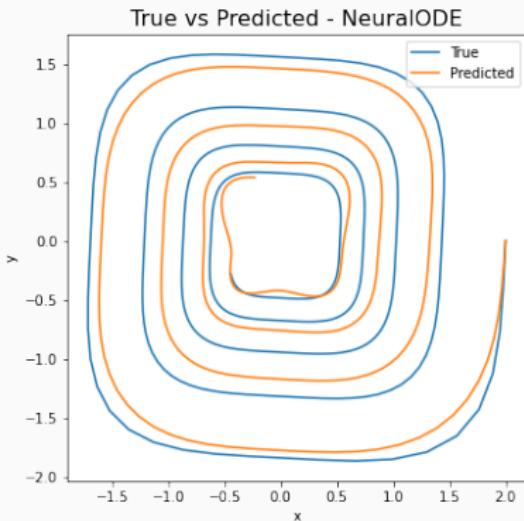
$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}(t), \mathbf{w})$$

```
1 from torchdiffeq import odeint
2
3 x0 = torch.tensor([2.0, 0.0])
4 t = torch.tensor([0.0, 0.1, 0.3, 0.31, 0.8, 1.0])
5 x = odeint(ODEFunc, x0, t)
```

'odeint' could be any ODE solver

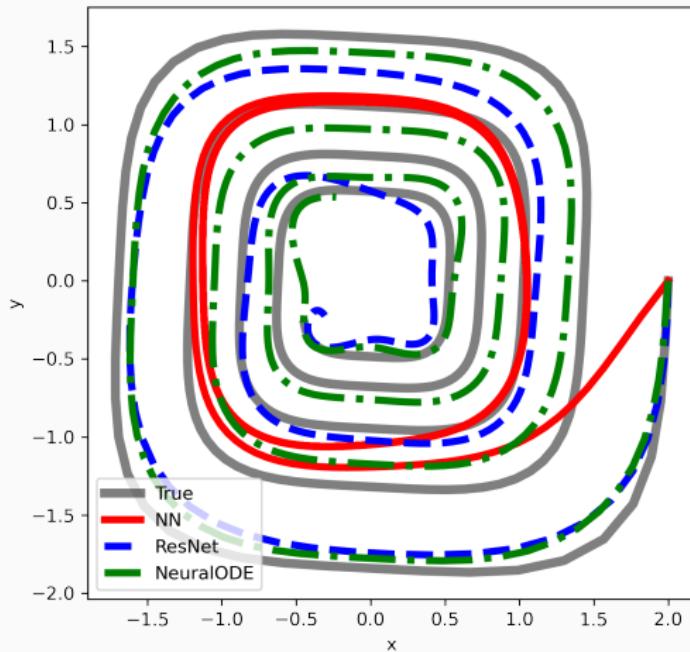
NEURAL ODEs - CONCEPT

Results:



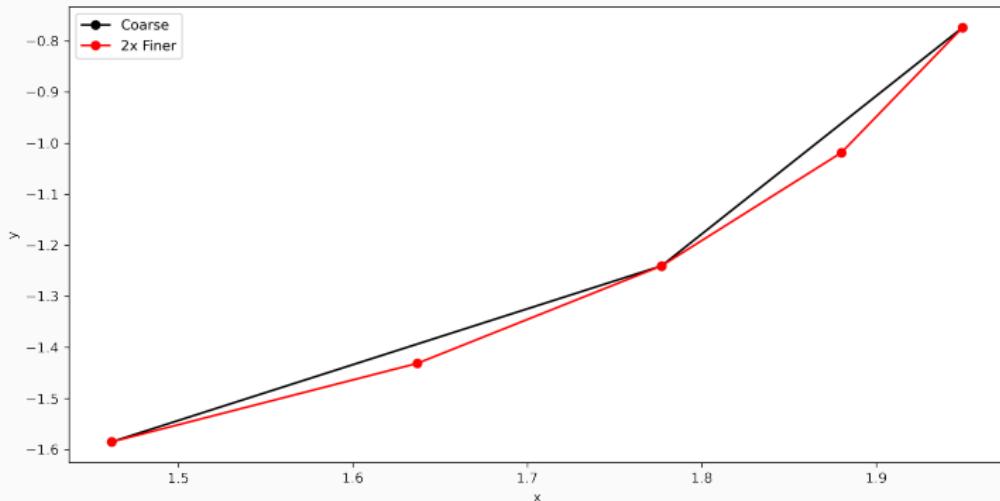
NEURAL ODEs - CONCEPT

Comparison:



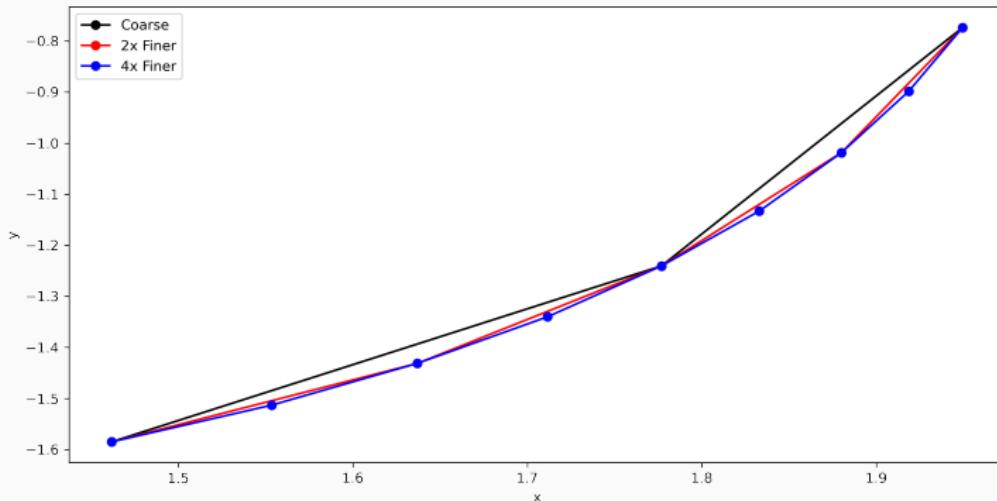
NEURAL ODEs - CONCEPT

Exploiting continuous time:



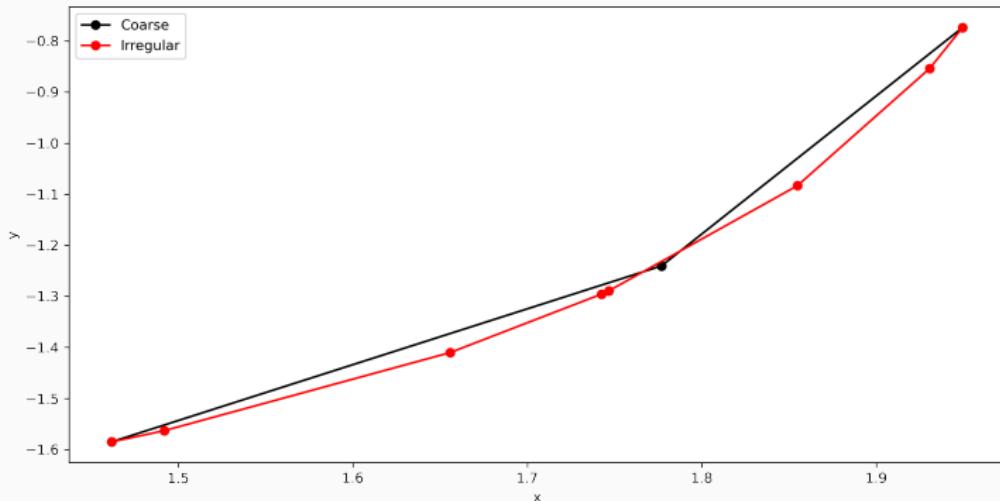
NEURAL ODEs - CONCEPT

Exploiting continuous time:



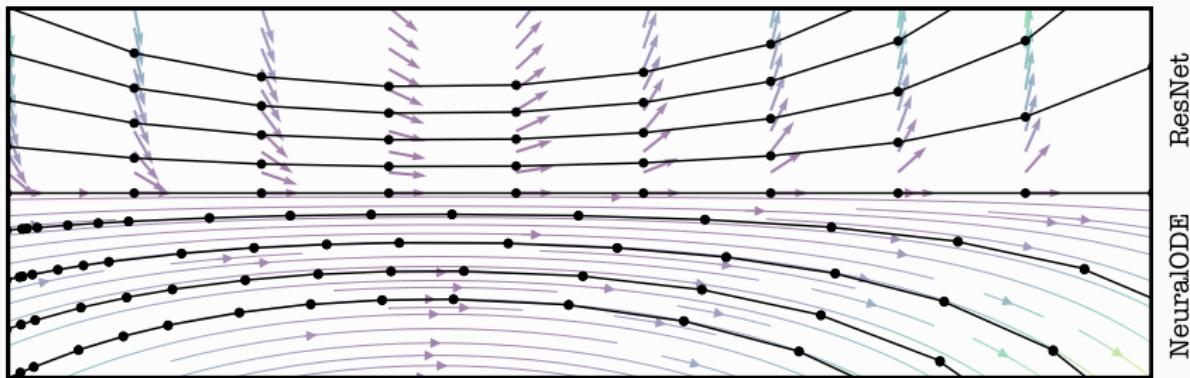
NEURAL ODEs - CONCEPT

Exploiting continuous time:



NEURAL ODES - CONCEPT

Exploiting continuous time:



NEURAL ODEs - TRAINING

So where's the catch?

- Training/backpropagation is not trivial in continuous time.
- Even with finite computations from ODE solver, it is intractable

Adjoint method

For loss function $L(\mathbf{x})$, define *adjoint state* $\mathbf{a}(t) = \partial L / \partial \dot{\mathbf{x}}(t)$.
Loss gradient is then given by

$$\frac{dL}{d\mathbf{w}} = - \int_{t_1}^{t_0} \mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{x}(t), t, \mathbf{w})}{\partial \mathbf{w}} dt$$

NEURAL ODEs - TRAINING

$$\frac{dL}{d\mathbf{w}} = - \int_{t_1}^{t_0} \mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{x}(t), t, \mathbf{w})}{\partial \mathbf{w}} dt$$

Components

$\mathbf{a}(t) = \frac{\partial L}{\partial \mathbf{x}(t)}$: sensitivity of loss to changes in $\mathbf{x}(t)$

$\frac{\partial \mathbf{f}(\mathbf{x}, \mathbf{w})}{\partial \mathbf{w}}$: sensitivity of change of $\mathbf{x}(t)$ to changes in \mathbf{w}

$-\int_{t_1}^{t_0} dt$: backpropagation through time/summing up

NEURAL ODES - TRAINING

$$\frac{dL}{dw} = - \int_{t_1}^{t_0} \mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{x}(t), t, w)}{\partial w} dt$$

Table 1

Comparison of the different gradient computation methods. The flow eqn. is either linear or non-linear, with P parameters and N state variables.

	Finite differences	Forward sensitivities	Adjoint
Suitability	Arbitrary	$N \gg P$	$P \gg N$
Cost	$(1 + P)$ flow eqns.	P non-linear sensitivity eqns. + 1 flow eqn.	1 linear adjoint eqn. + 1 flow eqn.
Key steps	1. Integrate flow eqn. 2. Parametrically perturb flow P times	1. Integrate the coupled flow and sensitivity eqns.	1. Integrate flow eqn. 2. Integrate adjoint eqn.

Figure: Comparison of gradient computing methods [Sen+14]

NEURAL ODEs - TRAINING

General remarks on training with time series data

- Predicting how far?
 - ▶ One step ahead?
 - ▶ Multiple steps ahead?
 - ▶ Scheduled learning?
- How to handle missing data?
 - Not a problem for continuous time, but for discrete time
- How to measure quality of prediction?
 - Not trivial for chaotic systems
 - ▶ Error in N-step prediction
 - ▶ Kullback-Leibler divergence
 - ▶ Power Spectrum analysis
- Regularisation?
 - ▶ Jacobian norm penalty
 - ▶ Kinetic energy penalty

NEURAL ODES - APPLICATIONS AND EXAMPLES

Where can we use Neural ODEs?

- **Time-series analysis**
- Dynamical system modelling
- Physics-informed learning
- Image and signal processing
- Generative modelling → **Continuous normalising flows**

NEURAL ODES - EXAMPLE: WEATHER SIMULATION

Time-series: Look at weather data of Heidelberg over the span of 3 days in July 2021

- **Parameters:** temperature, pressure, relative humidity, precipitation, cloud cover, shortwave radiation, wind speed, hour of the day

NEURAL ODEs - EXAMPLE: CNF

Continuous normalising flows (CNF):

Concept

Using Neural ODEs to parameterise and compute the transformations.

⇒ **Continuous-time dynamics** described by ODEs

Advantages:

- More **flexibility** and **adaptivity** in modelling complex data distributions
- Ability to capture even **intricate transformations**

NEURAL ODES - ADVANTAGES

What are the benefits?

- Continuous-time modeling
- Memory efficiency
- Adaptive computation
- Interpolation and extrapolation
- Flexibility in model design
- End-to-end learning
- Scalable and invertible normalising flows

NEURAL ODEs - DISADVANTAGES

What are the limitations?

- Computational cost
- Limited control over the ODE-Solver
- Limited generalisation
- Sensitivity to training data
- Difficulty in capturing discontinuities

NEURAL ODEs - OUTLOOK

Expanding to more complicated Differential Equations:

Partial Differential Equations (PDEs)

⇒ Controlled Differential Equations (CDEs)

→ adaptive Framework for studying the dynamics and optimising control policies in systems governed by PDEs

Stochastic Differential Equations (SDEs)

⇒ Neural SDEs

→ Framework for learning and modelling complex stochastic dynamics

SUMMARY AND TAKEAWAYS

SUMMARY AND TAKEAWAYS

- Introduction of Neural Networks and ResNets
 - universal function approximators
- Application of ResNets (and NNs) to ODEs
 - Easy to implement and good for fixed time steps
- Transition to continuous time (Neural ODEs)
 - Good for adaptive time steps
- Applications/Examples of Neural ODEs
 - Used for time-series analysis
 - Many advantages due to continuous-time modelling

SUMMARY AND TAKEAWAYS

Takeaways for working with time series data

- Know your data!
 - Make sure it's covering phase space as much as possible
- For regular data, use ResNets!
 - They're easier/faster to train
- For irregular data or intermediate states, use Neural ODEs!
 - Tweak the training length to improve your results

THANK YOU!

REFERENCES

REFERENCES

Neural Ordinary Differential Equations

Ricky T. Q. Chen*, Yulia Rubanova*, Jesse Bettencourt*, David Duvenaud
University of Toronto, rubanova, jessebett, duvenaud@cs.toronto.edu

Abstract

We introduce a new family of deep neural network models. Instead of specifying a discrete sequence of hidden layers, we parameterize the derivative of the hidden state using a neural network. The output of the network is computed using a black-box numerical integrator. This allows the model to bypass expensive implicit differentiation costs, and their evaluation strategy to each input, and can significantly trade numerical precision for speed. We demonstrate these properties in continuous-depth residual networks and continuous-time latent variable models. We also construct continuous normalizing flows, a generative model that can train by maximum likelihood, without ever evaluating its log-density. Finally, we show how to scalably backpropagate through any ODE solver, without access to its internal operations. This allows end-to-end training of ODEs within larger models.

1 Introduction

Models such as residual networks, recurrent neural network decoders, and normalizing flows build complicated transformations by composing a sequence of transformations to a hidden state:

$$\mathbf{h}_{t+1} = \mathbf{h}_t + f(\mathbf{h}_t, \theta_t) \quad (1)$$

where $t \in \{0, \dots, T\}$ and $\mathbf{h}_0 \in \mathbb{R}^D$. These iterative updates can be seen as an Euler discretization of a continuous transformation (Lu et al., 2017; Haber and Helman, 2017).

What happens as we add more layers and take smaller steps? In the limit, we parameterize the continuous dynamics of hidden units using an ordinary differential equation (ODE) specified by a neural network:

$$\frac{dh(t)}{dt} = f(\mathbf{h}(t), t, \theta) \quad (2)$$

Starting from the input \mathbf{h}_0 at time 0, we can define the output layer $\mathbf{h}(T)$ to be the solution to this ODE initial value problem at some time T . This value can be computed by a black-box differential equation solver, which evaluates the hidden unit dynamics f , whenever necessary to determine the solution with the desired accuracy. Figure 1 contrasts these two approaches.

Defining and evaluating models using ODE solvers has several benefits:

Memory efficiency. In Section 3, we show how to compute gradients of a scalar-valued loss with respect to all inputs of any ODE solver, without backpropagating through the operations of the solver. Not storing any intermediate quantities of the forward pass allows us to train our models with constant memory cost as a function of depth, a major bottleneck of training deep models.

32nd Conference on Neural Information Processing Systems (NeurIPS 2018), Montréal, Canada.

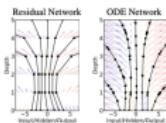


Figure 1: Left: A Residual network defines a discrete sequence of finite transformations.

Right: An ODE network defines a vector field, which continuously transforms the state.

Bottom: Circles represent evaluation locations.

On Neural Differential Equations

arXiv:2202.02435v1 [cs.LG] 4 Feb 2022



Patrick Kidger

Mathematical Institute

University of Oxford

A thesis submitted for the degree of

Doctor of Philosophy

Trinity 2021

REFERENCES

arXiv:1512.03385v1 [cs.CV] 10 Dec 2015

Deep Residual Learning for Image Recognition

Kaiming He¹ Xiangyu Zhang¹ Shaoqing Ren¹ Jian Sun²
Microsoft Research
{kaiming, v-xiangz, v-sher, junsun}@microsoft.com

Abstract

Deeper neural networks are more difficult to train. We present a residual learning framework to ease the training of networks that are substantially deeper than those used previously. We explicitly reformulate the layers as learning residual functions with reference to the layer inputs, instead of feature maps or channels. This simple but conceptually empirical evidence showing that these residual networks are easier to optimize, and can gain accuracy from considerably increased depth. On the ImageNet dataset we evaluate residual nets with a depth of up to 152 layers—8x deeper than VGG [16]. An ensemble of these residual networks achieves a 5.3% error on the ImageNet test set. This results won the 1st place in the ILSVRC 2015 classification task. We also present analysis on CIFAR-10 with 100 and 1000 layers.

The depth of representations is of critical importance for learning visual features. In this paper, after our extremely deep networks, we obtain a 28% relative improvement on the COCO object detection dataset. Deep residual nets are foundations of our submissions to ILSVRC & COCO 2015 competitions¹, where we also won the 1st places on the tasks of ImageNet detection, ImageNet localization, COCO detection, and COCO segmentation.

1. Introduction

Deep convolutional neural networks [22, 21] have led to a series of breakthroughs for image classification [21, 50, 40]. Deep networks naturally integrate low-level features [50] and classifiers in an end-to-end multi-layer fashion, and the “levels” of features can be enriched by the number of stacked layers (depth). Recent evidence [41] suggests that the depth of a network is not necessarily correlated with its performance. The leading results [41, 14, 13, 16] in the challenging ImageNet dataset [36] all exploit “very deep” [41] models with a depth of sixteen [41] to thirty [16]. Many other non-trivial visual recognition tasks [8, 12, 7, 32, 27] have also

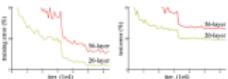


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 30-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

greatly benefited from very deep models.

Driven by the significance of depth, a question arises: Is learning better networks as easy as stacking more layers? An obstacle to answering this question was the notorious problem of vanishing/exploding gradients [1, 9], which hampered learning from the beginning. This problem, however, has been largely resolved by normalization [23, 9, 37, 13] and intermediate normalization layers [16], which enable networks with tens of layers to start converging for stochastic gradient descent (SGD) [22] with back-propagation [22].

What deeper networks are able to start converging, a deeper network does not necessarily benefit from it. As depth increasing, accuracy gets saturated (which might be unsatisfying) and then degrades rapidly. Unexpectedly, such degradation is not caused by overfitting, and adding more layers to a suitably deep model leads to *higher training error*, as reported in [11, 42] and thoroughly verified by our experiments in Fig. 1 and Fig. 4.

The degradation (of training accuracy) indicates that not all terms are similarly easy to optimize. Let us consider a shallower architecture and its deeper counterpart that adds more layers onto it. There exists a solution by construction to the deeper model: the added layers are identity mapping, and the overall function is the same as the shallower model. The existence of this constructed solution indicates that a deeper model should produce higher training error than its shallower counterpart. But experiments show that our current solvers on hand are unable to find solutions that

¹<http://image-net.org/challenges/LSVRC/2015/> and <http://msrcvse.org/dataset/detections-challenge2015/>.

How to Train Your Neural ODE: the World of Jacobian and Kinetic Regularization

Chris Finlay¹ Jiri-Henrik Jacobsen² Levon Nurbekyan³ Adam M Oberman¹

Abstract

Training neural ODEs on large datasets has not been tractable due to the necessity of allowing the adaptive numerical ODE solver to refine its step to very small values. This practice leads to dynamics equivalent to many hundreds or even thousands of layers. In this paper, we overcome this apparent difficulty by introducing a theoretically-grounded combination of both optimal transport and kinetic regularization, which encourage neural ODEs to prefer simple dynamics out of all the dynamics that solve a problem well. Simpler dynamics lead to faster convergence and to fewer discretizations of the solver, considerably decreasing wall-clock time without loss in performance. Our approach allows us to train neural ODEs to the same performance as the same architecture with the steepest dynamics, with significant reductions in training time. This brings neural ODEs closer to practical relevance in large-scale applications.

1. Introduction

Recent research has bridged dynamical systems, a workhorse of mathematical modeling, with the field of neural networks, the default function approximator for high dimensional data. The great promise of this pairing is that the vast mathematical machinery stemming from dynamical systems can be leveraged for modeling high dimensional problems in a dimension-independent fashion.

Connections between neural networks and ordinary differential equations (ODEs) were almost immediately noted after residual networks (He et al., 2015) were first proposed.

¹Department of Mathematics, Simon Fraser University, Burnaby, BC, Canada. ²NeuroSpin, Institute of Research in Biomedicine, Gif-sur-Yvette, France. ³Department of Mathematics, UCL, California, USA. Correspondence to: Chris Finlay <chrisfin@maths.ucl.ac.uk>.

Proceedings of the 27th International Conference on Machine Learning, Vienna, Austria. PMLR 119, 2020. Copyright 2020 by the author(s).

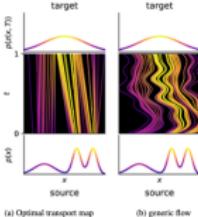


Figure 2. Optimal transport map and a generative normalizing flow.

Indeed, it was observed that there is a striking similarity between ResNets and the numerical solution of ordinary differential equations (El, 2017; Haber & Ruthotto, 2017; Ruthotto & Haber, 2018; Chen et al., 2018; 2019). In these works, deep dynamics are interpreted as discretized ordinary differential dynamics, where time indexes the “depth” of the network and the latent space indexes the “width” of the network. An alternate viewpoint was taken by neural ODEs (Chen et al., 2018), where the dynamics of the neural network are generated by an adaptive ODE solver on a fixed time interval. This is particularly interesting as it does not require specifying the number of layers of the network beforehand. Furthermore, it allows the learning of homeomorphisms without any structural constraints on the function computed by the residual block.

Neural ODEs have shown great promise in the physical sciences (Köhler et al., 2019), in modeling irregular time series (Ruthotto et al., 2019), in motion planning (Ruthotto et al., 2019), control and time modeling (Vlahis et al., 2019; Kambou et al., 2019), and in generative modeling through normalizing flows with free-form facades (Grindrod et al., 2019).

arXiv:2002.02798v3 [stat.ML] 23 Jun 2020

RECOMMENDATIONS

- **YouTube: Machine Learning and Simulation for adjoint method**
<https://youtube.com/watch?v=k6s2G5MZv-I>
- **Website of Patrick Kidger** for software overview
<https://kidger.site/>
- **Neural ODEs: breakdown of another deep learning breakthrough**
[https://towardsdatascience.com/
neural-odes-breakdown-of-another-deep-learning-
breakthrough-3e78c7213795](https://towardsdatascience.com/neural-odes-breakdown-of-another-deep-learning-breakthrough-3e78c7213795)
- **Lecture: Dynamical Systems Theory in Machine Learning**
[https://moodle.uni-heidelberg.de/course/view.php?
id=14482](https://moodle.uni-heidelberg.de/course/view.php?id=14482)