

Neural ODEs, Control and Machine Learning

Enrique Zuazua

Chair of Dynamics, Control and Numerics Department of Data Science FAU Erlangen, Germany https://caa-avh.nat.fau.eu

Universal approximation theorem I

Math. Control Signals Systems (1989) 2: 303-314

Mathematics of Control, Signals, and Systems

© 1989 Springer-Verlag New York Inc.

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

$$\sum_{j=1}^{N} \alpha_j \sigma(y_j^{\mathsf{T}} x + \theta_j), \qquad (1)$$

where $y_j \in \mathbb{R}^n$ and α_j , $\theta \in \mathbb{R}$ are fixed. (y^T is the transpose of y so that $y^T x$ is the inner product of y and x.) Here the univariate function σ depends heavily on the context of the application. Our major concern is with so-called sigmoidal σ 's:

$$\sigma(t) \rightarrow \begin{cases} 1 & \text{as } t \rightarrow +\infty, \\ 0 & \text{as } t \rightarrow -\infty. \end{cases}$$

Universal approximation theorem II







Supervised learning

Goal: Find an approximation of a function $f_{\rho} : \mathbb{R}^d \to \mathbb{R}^m$ from a dataset

$$\left\{\vec{x}_i, \vec{y}_i\right\}_{i=1}^N \subset \mathbb{R}^{d \times N} \times \mathbb{R}^{m \times N}$$

drawn from an unknown probability measure ρ on $\mathbb{R}^d \times \mathbb{R}^m$.

Classification: match points (images) to respective labels (cat, dog). \rightarrow Popular method: **training a neural network**.



Residual neural networks

[1] K He, X Zhang, S Ren, J Sun, 2016: Deep residual learning for image recognition

[2] E. Weinan, 2017. A proposal on machine learning via dynamical systems.

[3] R. Chen, Y. Rubanova, J. Bettencourt, D. Duvenaud, 2018. Neural ordinary differential equations.

[4] E. Sontag, H. Sussmann, 1997, Complete controllability of continuous-time recurrent neural networks.

ResNets

$$\begin{cases} \mathbf{x}_i^{k+1} = \mathbf{x}_i^k + \mathbf{h} W^k \boldsymbol{\sigma} (A^k \mathbf{x}_i^k + b^k), & k \in \{0, \dots, N_{layers} - 1 \\ \mathbf{x}_i^0 = \tilde{\mathbf{x}}_i, & i = 1, \dots, N \end{cases}$$

where h = 1, σ globally Lipschitz $\sigma(0) = 0$.

nODE Layer = timestep; $h = \frac{T}{N_{layers}}$ for given T > 0 $\begin{cases} \dot{\mathbf{x}}_i(t) = W(t)\sigma(A(t)\mathbf{x}_i(t) + b(t)) & \text{for } t \in (0, T) \\ \mathbf{x}_i(0) = \vec{\mathbf{x}}_i, \quad i = 1, ..., N \end{cases}$



The problem becomes then a giant simultaneous control problem in which each initial datum $x_i(0)$ needs to the driven to the corresponding destination for all i = 1, ..., N with the same controls:

• What happens when $T ightarrow \infty$, i.e. in the deep, high number of layers regime?¹ ²

¹C. Esteve, B. Geshkovski, D. Pighin, E. Zuazua, Large-time asymptotics in deep learning, arXiv:2008.02491

²D. Ruiz-Balet & Zuazua, Neural ODE control for classification, approximation and transport, arXiv:2104.05278

E. Zuazua (FAU - AvH)

Special features of the control of ResNets

- Nonlinearities are unusual in Mechanics: σ is flat in half of the phase space.
- We need to control many trajectories (one per item to be classified) with the same control! ³

The very nature of the activation function σ allows actually to achieve this monster simultaneous control goal. The fact that σ leaves half of the phase space invariant while deforming the other one, allows to build dynamics that are not encountered in the classical ODE systems in mechanics and for which such kind of simultaneous control property is unlikely or even impossible.



³This would be impossible for instance, for the standard linear system x' = Ax + Bu.

E. Zuazua (FAU - AvH)

Turnpike refers to the fact that, in long time-horizons, optimal controls and trajectories are exponentially close to the optimal steady-state control and state in most of the time-horizon.

Supervised Learning* \iff **minimize**⁴

$$\frac{1}{N}\int_0^T \|P\mathbf{x}(t)-\overline{y}\|^2 dt + \alpha \|u\|_{H^1(0,T;\mathbb{R}^{d_u})}^2.$$

$$(SL^*)$$

 $\overline{y} := [\vec{y_1}, \dots, \vec{y_N}], \ u := [A, W, b] \text{ in } (5) \text{ and } P : \mathbb{R}^d \to \mathbb{R}^m.$

Theorem (Turnpike): Under controllability assumptions, for any sufficiently large T, an optimal solution (u_T, \mathbf{x}_T) to (SL^*) –(5) satisfies

$$\|u_T\|_{H^1(0,T;\mathbb{R}^{d_u})} \leq C_1$$

and

$$\|P\mathbf{x}_T(t) - \overline{y}\| \leq C_2 e^{-\mu t} \qquad \forall t \in [0, T]$$

for some $C_1 > 0$, $C_2 > 0$ and $\mu > 0$, all independent of T.

E. Zuazua (FAU - AvH)

⁴Note that in this context we do not impose a perfect classification. We just expect that it will occur with high probability as an outcome of the optimar control problem

 $T o \infty \sim N_{\text{layers}} o \infty.$





E. Zuazua (FAU - AvH)

Classical SL problem?



Convergence of $\mathbf{x}(T)$ to $P^{-1}(\{\overline{y}\})$ when $T \to \infty$, but slow (no turnpike).

E. Zuazua (FAU - AvH)





E. Zuazua (FAU - AvH)



E. Zuazua (FAU - AvH)



E. Zuazua (FAU - AvH)

The classification problem is a relaxed version of the simultaneous control problem. We are given N points in \mathbb{R}^d and M classes $y_i \in \{1, ..., M\}$.

We then proceed as follows:

- We identify a region in the euclidean space corresponding to each class of data.
- Output: Look for a control strategy (A, W, b) bringing simultaneously all points to the location corresponding to its class.

 S_3

 S_2

 S_1

 $x^{(2)}$ -axis



$\dot{\mathbf{x}}(t) = W(t)\sigma(A(t)\mathbf{x}(t) + b(t)).$

- b(t) induces a time-dependent translation of the Euclidean space. It plays an important role to place the center of the action of the sigmoid.
- A(t) compresses, expands, and induces rotations in the euclidean space with different objectives:
 - Compression can help gathering data into clusters so that they might be manipulated simultaneously.
 - Expansion allows to separate data of different classes to better focus the action of the control on just one of them.
 - Rotations allow to better choose the hemisphere where the action will be focused.
- W(t) determines the direction and intensity with which the flow will evolve in the active hemisphere.

Some canonical flows induced by nODE



Controlling one datum



E. Zuazua (FAU - AvH)

Compression after classification



E. Zuazua (FAU - AvH)



Theorem (Classification, Domènec Ruiz-Balet EZ, 2021)

^a Let σ be the ReLU. Let $d \ge 2$, and $N, M \ge 2$. Let $\{x_i\}_{i=1}^N \subset \mathbb{R}^d$ be data to be classified into disjoint open non-empty subsets $S_m, m = 1, ..., M$ with labels m = m(i), i = 1, ..., N. Then, for every T > 0, there exist control functions $A, W \in L^{\infty}((0, T); \mathbb{R}^{d \times d})$ and $b \in L^{\infty}((0, T), \mathbb{R}^d)$ such that the flow associated to the Neural ODE, when applied to all initial data $\{x_i\}_{i=1}^N$, classifies them simulatenously, i.e.

$$\phi_T(x_i; A, W, b) \in S_{m_i}, \quad \forall i = 1, ..., N.$$

Furthermore,

- Controls are piecewise constant with a maximal finite number of switches of the order of O(N). They also lie in BV.
- The control time T > 0 can be made arbitrarily small (scaling).
- The complexity of controls diminishes when initial data are structured in clusters.
- The complexity of controls also diminishes when the control requirement is relaxed so that not all data need to be classified.

The targets S_m can be just N distinct points in the euclidean space.

^aRelated results for smooth sigmoids using Lie bracket control techniques: A. Agrachev and A. Sarychev, arXiv:2008.12702, (2020).

Complexity of the control strategy



E. Zuazua (FAU - AvH)

Neural transport equations

Note that the differential equation

$$\begin{cases} \dot{x} = W(t)\sigma(A(t)x + b(t)) \\ x(0) = x_0 \end{cases}$$

corresponds to the characteristics of the transport equation:

$$\begin{cases} \partial_t \rho + \operatorname{div}_x \left[(W(t)\sigma(A(t)x + b(t)))\rho \right] = 0\\ \rho(0) = \rho^0 \end{cases}$$

The results above can therefore be understood in terms of the controllability of the transport equation: "Atomic initial data can be driven to atomic final targets". This also allow for a more general interpretation in terms of approximation control in Wasserstein-1 distance. Or for systems of transport equations, so that each scalar component corresponds to the density within one of the classes of data.

This establishes a link to the Theory of Optimal Transport: Neural Transport? ⁵

$$\mathcal{W}_1(\mu,
u) = \sup_{\textit{Lip}(g) \leq 1} \left\{ \int_{\mathbb{R}^d} g d\mu - \int_{\mathbb{R}^d} g d
u
ight\}.$$

where $Lip(g) \leq 1$ stands for the class of Lipschitz functions with Lipschitz constant less or equal than 1.

E. Zuazua (FAU - AvH)

Neural transport equations

The simultaneous control of the nODE

$$egin{aligned} \dot{x} &= W(t) \boldsymbol{\sigma}(A(t)x + b(t)) \ x(0) &= x_i, \quad i = 1, ..., N \end{aligned}$$

to arbitrary terminal states

$$x(T) = y_i, \quad i = 1, ..., N$$

in terms of the transport equation, leads to the control of an atomic initial datum from

$$\rho(x,0) = \sum_{i=1}^{N} m_i \delta_{x_i}$$

to the terminal one

$$\rho(x,T) = \sum_{i=1}^{N} m_i \delta_{y_i}.$$

But note that, even if the locations of the masses are transported, the amplitude of the masses do not vary.



Neural transport equations

We can enrich the strategy above to also regulate the amplitude of the masses. But this requires to relax the control statement into an ε -approximate one.

For that to be done we need to split initial masses so that

$$m_i = \sum_{j=1}^{J_i} m_{i,j}, \quad i = 1, ..., N$$

they are dispersed from the center x_i into the neighboring points $x_{i,j}$. This allows to enrich the transport diagram.



Neural Transport



Universal approximation

- As a corollary we can achieve Universal Approximation.
- By density, it is sufficient to consider targets that are simple piecewise constant functions.
- We can proceed making a partition of the departure and arrival spaces so that the problem becomes a countable version of simultaneous control of the nODE.



• The complexity of the needed controls depends on the nature of functions one aims to approximate.⁶



 $^{6}N_{\Gamma}(h)$ being the number of hypercubes of side h needed to cover the boundary Γ , the box-counting dimension is

$$D := \lim_{h \to 0} \left[\log N_{\Gamma}(h) / \log \left(h^{-1} \right) \right]$$

Then

$$\|W\|_{L^{\infty}} \lesssim \epsilon^{-\frac{4Dd}{d-D}}, \qquad \|b\|_{L^{\infty}} \lesssim \epsilon^{-\frac{-2dD}{d-D}} \qquad \text{as } \epsilon \to 0$$

and the number of switches of A, W, b will be of the order of $e^{-\overline{d-D}}$

E. Zuazua (FAU - AvH)

Universal approximation

Let us approximate a piece-weise constant function taking two different values P and Q on the sets represented by colors blue and red.

We aim to build a nODE so that the solution of

$$\dot{\varphi}(t) = W(t)\boldsymbol{\sigma}(A(t)\varphi(t) + b(t))$$

 $\varphi(0) = x,$

is such that

$$\varphi(T, x) = P$$
, when $x \in$ Blue Set

and

 $\varphi(T, x) = Q$, when $x \in \text{Red Set}$.

The same control inspired strategies allow to achieve the result up to an ε - error.



An extraordinary and fertile field in the interplay between Dynamical Systems, Control, Machine Learning and applications

- Control and dynamical systems tools allow to explain the amazing efficiency of Neural Networks (NN) in some specific applications.
- Long-time / Turnpike control arise naturally in Deep Learning
- Interesting open questions:
 - How to deal with Neural ODEs that switch in dimension of the Euclidean phase space.
 - Are there results explaining how the clustering of data (number of separating interfaces needed) diminishes in higher dimensions?
 - How close is our piecewise constant control strategy from the optimal one (in the Pontryagin sense?)
 - How does our control strategy compare to those obtained in a purely NN setting?
 - How does the complexity of the controls diminish when we relax the classification criteria?
 - Links with Optimal Transport.
 - Other objectives: Unsupervised learning?
- My sincere thanks to collaborators: A. Porretta (Roma 2), E. Trélat (Paris Sorbonne), M. Gugat (FAU), D. Pighin (Innovalia), B. Geskhovski (UAM & Deusto), C. Esteve (UAM & Deusto), M. Schuster (FAU), M. Lazar (Dubrovnik), V. Hernández-Santamaria (Mx), N. Sakamoto (Nanzan), J. Heiland (Magdeburg), H. Kouhkouh (Padova), D. Ruiz-Balet (UAM & Deusto).
- Funded by the ERC Advanced Grant DyCon and an Alexander von Humboldt Professorship and Marie-Sklodowska Curie ITN "ConFlex"



This project has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie Sklodowska-Curie grant agreemen: No 765579.

