Some control aspects in deep learning

Domènec Ruiz-Balet¹

¹Imperial College London

25th August 2022

Benasque, IX Partial differential equations, optimal design and numerics, Special session for Enrique Zuazua's 60th birthday Joint works with Enrique Zuazua

Imperial College London



▲ロト ▲母ト ▲ヨト ▲ヨト 三国市 のへで

Classification

We have a large data set of pictures



How can we separate the dog pictures from the other types of pictures?

Domenec Ruiz-Balet (IC)	Domènec	Ruiz-Ba	let ((IC)	
-------------------------	---------	---------	-------	------	--

Control and DL

Learning a function through samples

Let $A \subset \mathbb{R}^d$ be the subset of dog pictures.

Classification

Can we recover the function $\mathbb{1}_A$ given a finite amount of samples of $\mathbb{1}_A$?

$$\{(x_i, y_i = \mathbb{1}_A(x_i))\}_{i=1}^N \subset \mathbb{R}^d \times \{0, 1\}$$

Approach

The approach would be to find an approximation of $\mathbbm{1}_A$ in a "large" family of functions

The family we will consider are Neural Networks

 Deep Residual networks (ResNets)¹ were implemented to ease the training of deep neural networks.

$$x^{k+1} = x^k + W^k \sigma (A^k x^k + b^k)$$

This formulation reminds an Euler discretization of an ODE²³.

$$\mathbf{x}' = \mathbf{W}(t)\sigma(\mathbf{A}(t)\mathbf{x} + \mathbf{b}(t))$$

Domènec Ruiz-Balet (IC)

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

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

³R. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud 2018. Neural ordinary differential equations - < 😑 - 🖳 😑 📀 🔍

A Neural ODE has the form of:

$$x' = W(t)\sigma(A(t)x + b(t))$$

where

Consider $\sigma : \mathbb{R} \to \mathbb{R}$ being the ReLU

$$\sigma(\mathbf{x}) = \max(\mathbf{x}, \mathbf{0})$$



Control problem

Let $f : \mathbb{R}^d \times \mathbb{R}^{d_u} \to \mathbb{R}^d$ be a Lipschitz function. **Controllability problem**: Let $x_0, x_T \in \mathbb{R}^d$, $\exists u \in L^{\infty}((0, T); \mathbb{R}^{d_u})$ s.t.

$$\begin{cases} x' = f(x, u) \\ x(0) = x_0, x(T) = x_T \end{cases}$$

is satisfied?

Simultaneous controllability problem: Does $u \in L^{\infty}((0, T); \mathbb{R}^{d_u})$ exist such that

$$\begin{cases} x' = f(x, u) \\ x(0) = x_{0,1}, x(T) = x_{T,1} \end{cases} \qquad \begin{cases} x' = f(x, u) \\ x(0) = x_{0,2}, x(T) = x_{T,2} \end{cases}$$

are satisfied?

Consider *N* samples $\{(x_i, y_i)\}_{i=1}^N \subset \mathbb{R}^d \times \{1, ..., M\}$ where x_i is the data associated with a class $y_i \in \{1, ..., M\}$.

Find a **control** strategy that brings **simultaneously** all points to their **prefixed locations**



$x' = W(t)\sigma(A(t)x + b(t))$



Domènec Ruiz-Balet (IC)

Aug'25 8/38

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □



Control and DL

・ロト・白ア・トロア・山下 シック

Aug'25 9/38

Classification

Let $d \ge 2$ and $M \ge 2$. Consider $\{x_i, y_i\}_{i=1}^N \subset \mathbb{R}^d \times \{1, ..., M\}$. Assume $x_i \ne x_j$ if $i \ne j$. Then, for every T > 0, $\exists A, W \in L^{\infty}$ ((0, T); $\mathbb{R}^{d \times d}$) and $b \in L^{\infty}$ ((0, T), \mathbb{R}^d) such that $\phi_T(x_i; A, W, b) \in S_{m_i}$.

 S_{m_i} being the subset corresponding to the label y_i .

Theorem RB-Zuazua (Simultaneous Control/Interpolation)

For $d \ge 2$ the system is approximately simultaneously controllable.

Simultaneous Control

Step 1: Apply the same arguments than in the classification



Simultaneous Control

Step 2: build, in an iterative manner piecewise controls that steer one by one each point to their target.



Universal approximation

Universal Approximation

Let $d \ge 2$ and T > 0. Let $\Omega \subset \mathbb{R}^d$ be a bounded set, and σ be the ReLU. Then, for any $f \in L^2(\Omega; \mathbb{R}^d)$ and $\epsilon > 0$ there exist $A, W \in L^{\infty}((0, T); \mathbb{R}^d)$ and $b \in L^{\infty}((0, T); \mathbb{R}^d)$ such that the flow generated by the Neural ODE, $\phi_T(\cdot; A, W, b)$, satisfies:

$$\|\phi_{\mathcal{T}}(\cdot; \boldsymbol{A}, \boldsymbol{W}, \boldsymbol{b}) - f(\cdot)\|_{L^{2}(\Omega)} < \epsilon.$$

In the NODE formulation, the result translates into the following each input $x \in \mathbb{R}^d$ has a target $f(x) \in \mathbb{R}^d$. The actual value of f(x). Ideally, we would like to find controls *A*, *W* and *b* such that for every $x \in \Omega$ we have that

$$f=\sum_{m=1}^{M}\alpha_m\chi_{\Omega_m},$$



$$\begin{cases} \dot{x} = W(t)\sigma(A(t)x(t) + b(t)) \\ x(0) = x \\ \phi_T(x; A, W, b) = f(x). \end{cases}$$

The proof is based on approximating the target function by a **simple function** and then trying to reduce the problem to a simultaneous control one.

1.Cover the boundary of the characteristic sets with small hypercubes and generate a mesh with them. 2. Remove a small strip around the mesh 3. The white areas in the figure above belong to Ω_h and the function will not be very well approximated there.





E 5 4 E

"Badly mixed data"

- The number of hyperrectangles depend on how many hypercubes one needs.
- If the interface between the characteristic sets is very irregular, one will need more hypercubes and therefore the control cost will be higher



"Badly mixed data"

Let *D* be the **box-counting dimension** of Γ^4 . Let $N_{\Gamma}(h)$ be the number of hypercubes of side *h* needed to cover the boundary Γ , the box-counting dimension is defined as

$$\mathcal{D} := \lim_{h o 0} rac{\log N_{\Gamma}(h)}{\log \left(rac{1}{h}
ight)}.$$

The bounds on the control cost are able to capture such complexity.

$$\|W\|_{L^{\infty}} \lesssim \epsilon^{-2(D+1)d}, \qquad \|b\|_{L^{\infty}} \lesssim \epsilon^{-2dD} \qquad ext{as } \epsilon o 0$$

The number of switches of *A*, *W*, *b* will be of the order of e^{-2dD} .

⁴K. Falconer 06: Fractal geometry: mathematical foundations and applications. 🕢 🗆 🕨 🖉 🖉 🗸 🗇 🔍 🖓

Curse of dimensionality

Constructive procedure: Meshing the boundary produces a high number of components and hence a high control cost

The number of switches of *A*, *W*, *b* will be, at most, of the order of e^{-2dD} . ONE CANNOT AVOID THE DEPENDENCE ON THE DIMENSION

For \mathbb{B}_d , the number of switches of A, W, b will be of the order of $e^{-\frac{d-1}{2}}$.



$$\min_{p_N \in P_N} |\mathbb{B}_d riangle p_N| \sim rac{c(d)}{N^{2/(d-1)}}$$

K. Börözky Jr, Polytopal approximation bounding the number of k-faces, Journal of Approximation Theory, (2000) and the second

Domènec Ruiz-Balet (IC)

Control and DI

Aug'25 18/38

Transport

Note that the set of differential equations

$$\begin{cases} x' = W(t)\sigma(A(t)x + b(t)) \\ x(0) = z \in \mathbb{R}^d \end{cases}$$

correspond to the projected characteristics of the transport equation:

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

《曰》《圖》《曰》《曰》 되는

Definition

Let $\mu, \nu \in \mathcal{P}_{c}(\mathbb{R}^{d})$ be probability measures. The Wasserstein-1 distance $\mathcal{W}_{1}(\mu, \nu)$ is defined by as:

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

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

⁵C. Villani 2008: Optimal transport: old and new

⁶L.V. Kantorovich 1942, On the transfer of masses.



2

x

3

1

The L^{p} -norm does not "see" the Euclidean distance between the supports

$$\|v_1 - v_2\|_{L^p} = 1$$

$$\|v_1 - \tilde{v}_2\|_{L^p} = 1$$

However, the Wasserstein distance does,

$$\mathcal{W}_1(\delta_{x_1},\delta_{x_2})=|x_1-x_2|$$

0

1

0.8

0.6 0.4 0.2 0 Consider target measures ρ^* in the form

$$\rho^* = \sum_{m=1}^M \beta_m \delta_{\alpha_m}, \qquad \sum_{m=1}^M \beta_m = \int_{\mathbb{R}^d} \rho^0 dx$$

where δ_{α_m} is the Dirac delta located at $\alpha_m \in \mathbb{R}^d$ and $\beta_m > 0$.

Theorem (RB-Zuazua 2021 Control transport)

Let T > 0, $d \ge 2$. Then, for every $\epsilon > 0$, $\exists W, A \in L^{\infty}((0, T); \mathbb{R}^{d \times d})$ and $b \in L^{\infty}((0, T); \mathbb{R}^{d})$ s.t. $\mathcal{W}_{1}(\rho(T), \rho^{*}) < \epsilon.$



Aug'25 23/38

<ロ> < @> < E> < E> EE のQの

Simultaneous Control of NTEs

Remark

The scalar transport equation does not allow distinguishing among labels. One should take a vectorial structure in order to have a transport formulation for classification.

- Let us consider *M* classes, and *M* compactly supported probability densities ρ_m for m = 1, ..., M.
- Assume that for every *x* there exists, at most, a unique label *y*. This implies that the supports of the probability measures ρ_m are disjoint:

$$\operatorname{supp}(\rho_m) \cap \operatorname{supp}(\rho_{m'}) = \varnothing, \quad \text{if } m \neq m'.$$

The system reads:

$$\begin{cases} \partial_t \rho_m + \operatorname{div}_x \left[(W(t)\sigma(A(t)x + b(t))\rho_m) \right] = 0, & m \in \{1, ..., M\} \\ \rho_m(0) = \rho_m^0 \in C_c(\mathbb{R}^d), & m \in \{1, ..., M\} \\ \operatorname{supp}(\rho_m^0) \cap \operatorname{supp}(\rho_{m'}^0) = \varnothing & \text{if } m \neq m'. \end{cases}$$

Assume that the target functions satisfy:

$$\operatorname{supp}(\rho_m^*) \cap \operatorname{supp}(\rho_{m'}^*) = \varnothing$$
 if $m \neq m'$, $\int \rho_m^* = \int \rho_m^0 dx = 1$ n

Theorem RB-Zuazua 2021: Simultaneous control transport

Let T > 0. Then, for any $\epsilon > 0$, $\exists W, A \in L^{\infty}((0, T), \mathbb{R}^{d \times d})$ and $\exists b \in L^{\infty}((0, T), \mathbb{R}^{d})$ such the solution satisfies:

$$\mathcal{W}_1(\rho_m(T),\rho_m^*) < \epsilon \quad m \in \{1,...,M\}$$

イロト イポト イモト イモト・モー

Other type of dynamics: Momentum ResNet

 $x'' + x' = W\sigma(Ax + b)$



Neural ODE



Domènec Ruiz-Balet (IC)

Control and D

Zorionak!

- D Ruiz-Balet, E Zuazua, Neural ODE control for classification, approximation and transport,arXiv preprint arXiv:2104.05278 (Accepted in SIAM Review)
- D Ruiz-Balet, E Affili, E Zuazua, Interpolation and approximation via momentum ResNets and neural ODEs, Systems & Control Letters 162, 2022



Enrique Zuazua



Elisa Affili

- Freeze. The fact that the vector field can be zero in half-space, implies that any vector field generated by a hyperplane leaves half being critical points for the dynamics.
- **2** Allocate. One can allocate the hyperplane at any place in \mathbb{R}^d .
- Setting an appropiate w one can generate
 - an expansion.
 - a compression.
 - a translation.

This allows us to consider fundamental vector fields that will be key for controlling the dynamics.

Classification Problem

First of all, before applying the main argument, we will have to prepeare our data set. We need that to find controls such that can justify without loss of generality the following assumption

$$x_j^{(1)}
eq x_i^{(1)}$$
 if $i \neq j$







Domènec Ruiz-Balet (IC)

Control and DL

Aug'25 30/38



Control and DL

・

Aug'25 31/38



Control and DL

・

Aug'25 32/38



Control and DL



Control and DL

・

Aug'25 34/38



< ロ > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < < 回 > < < ○ Aug'25

35/38



Control and DL

きょう きょう キャー・ 小学・ トー・

Aug'25 36/38



Aug'25 37/38

 The restriction of *d* ≥ 2 comes from the limitation already pointed out before. However, the transport equation

 $\partial_t \rho + \operatorname{div}_{\boldsymbol{X}}[\boldsymbol{V}(\boldsymbol{X},t)\rho] = 0$

with V as a control, can be approximately controlled in the one-dimensional case.

- It would be enough to design appropriate locations of attractors and repulsors depending on the mass distribution of ρ⁰.
- This would allow concentrating the mass in a finite number of points.
- Later, one needs to design a dynamical system that brings each mass approximately on the approximation of the target. Neural Transport Equations can achieve this in dimension d ≥ 2 just with the controls A, W and b.