

# CONTROL AND MACHINE LEARNING

---

## Umberto Biccari

Chair of Computational Mathematics, Universidad de Deusto, Bilbao, Spain

[umberto.biccari@deusto.es](mailto:umberto.biccari@deusto.es)

[cmc.deusto.es](http://cmc.deusto.es)

*based on works by:*

**Enrique Zuazua** - Friedrich-Alexander-Universität, Erlangen

Universidad de Deusto, Bilbao

Universidad Autónoma de Madrid

**Rafael Orive-Illera** - Universidad Autónoma de Madrid and ICMAT

**Domènec Ruíz-Balet** - Université Paris Dauphine

**Antonio Álvarez-López** - Universidad Autónoma de Madrid

Control, Problemas Inversos y Aprendizaje Automático

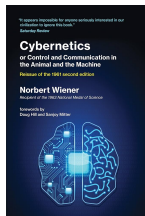
Almagro, December 2-3 2024



# Control and ML: two neighboring fields

**CYBERNETICS:** the science of communication and control in animals and machines.

Norbert Wiener, 1948



Two essential conceptual binomials involved in this definition

**control-communication:** sufficient and quality information about the state of a system is needed to make the right decisions, to reach a given objective or simply to avoid risky regimes to be eluded.

**animal-machine:** the human being, a rational animal, aims to build machines to carry out those tasks that prevent him from dedicating his time and energy to more sublime activities.

The close link between control and/or cybernetics and Machine Learning is built in Wiener's own definition.

# Control and ML: two neighboring fields

It is often hard to observe the interconnections between different mathematical disciplines since they are split by conceptual and technical mountain ranges, and often have evolved in different communities. Building the connecting paths requires an important effort of abstraction.

Exploring the notion of controllability helps disclosing one of those gateways.

**CONTROLLABILITY:** driving a dynamical system from an initial configuration to a final one, in a given time horizon, using skillfully designed and viable controls.

## LINEAR FINITE DIMENSIONAL-KALMAN

$$\dot{x}(t) = Ax(t) + Bu(t)$$

↓

$$\text{rank} \begin{bmatrix} B & AB & A^2B & \dots & A^{N-1}B \end{bmatrix} = N$$

The **size of the control** depends on the **length of the time horizon**: in short time horizons the control is enormous while in longer time horizons it has a smaller amplitude.

Are these ideas and methods relevant at all in Machine Learning?

# Supervised learning

## Goal: classification

Match points (images) to respective labels.



## Technique

Approximate a function  $f_{\theta} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , depending on some unknown parameters  $\theta \in \Theta$ , using a dataset

$$\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^d \times \mathbb{R}^d \quad \text{s.t. } \mathbf{y}_i = f_{\theta}(\mathbf{x}_i)$$

This is typically done by **training a neural network**. We will do it through **simultaneous or ensemble control of Neural Ordinary Differential Equations(NODEs)**.

# Supervised learning

## Residual Neural Network (ResNet)

For each item  $i \in \{1, \dots, N\}$

$$\mathbf{x}_i^{k+1} = \mathbf{x}_i^k + h\mathbf{W}^k \sigma(\mathbf{A}^k \mathbf{x}_i^k + \mathbf{b}^k), \quad k \in \{1, \dots, L-1\}$$

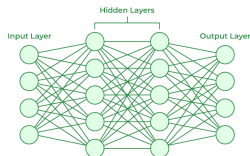
$(\mathbf{W}^k, \mathbf{A}^k, \mathbf{b}^k) \in \mathbb{R}^{d \times p} \times \mathbb{R}^{p \times d} \times \mathbb{R}^p$ : parameters

$\sigma$ : (nonlinear) **activation function**

$L > 1$ : number of layers (depth)

$p \geq 1$ : neurons per layer (width)

**GOAL**: find  $(\mathbf{W}^k, \mathbf{A}^k, \mathbf{b}^k)$  such that  $\mathbf{x}_i^L = \mathbf{y}_i$



In the limit  $L \rightarrow +\infty$ , this constitutes a simultaneous control problem for the **Neural Ordinary Differential Equation**

$$\dot{\mathbf{x}}_i(t) = \mathbf{W}(t) \sigma(\mathbf{A}(t) \mathbf{x}_i(t) + \mathbf{b}(t)), \quad t \in (0, T)$$

where  $(\mathbf{W}, \mathbf{A}, \mathbf{b}) \in L^\infty((0, T), \mathbb{R}^{d \times p} \times \mathbb{R}^{p \times d} \times \mathbb{R}^p)$  are control functions that must steer each data point to its corresponding label:  $\mathbf{x}_i(T) = \mathbf{y}_i$  for all  $i \in \{1, \dots, N\}$ .

# Supervised learning

How do we achieve that? **Empirical risk minimization**: minimize under the ResNet dynamics

$$\underbrace{\frac{1}{N} \sum_{i=1}^N \text{loss}(\mathbf{x}_i, \mathbf{y}_i)}_{\text{empirical risk}} + \lambda \underbrace{\sum_{k=1}^L \|(\mathbf{W}^k, \mathbf{A}^k, \mathbf{b}^k)\|^2}_{\text{regularization}}$$

$\downarrow L \rightarrow +\infty$

$$\underbrace{\frac{1}{N} \sum_{i=1}^N \text{loss}(\mathbf{x}_i(T), \mathbf{y}_i)}_{\text{empirical risk}} + \lambda \underbrace{\int_0^T \|(\mathbf{W}(t), \mathbf{A}(t), \mathbf{b}(t))\|^2 dt}_{\text{regularization}}$$

Difficult problem due to

**computational complexity**:  $N$  is most-often very large

**lack of convexity** coming from the non-linearity of the neural network

**curse of dimensionality**: as the data dimension increases, the amount of data needed to properly minimize empirical risk grows exponentially.

# Supervised learning

Why does it work? The answer is **universal approximation**: finite combinations of rescaled and shifted activation functions

$$\sum_{k=1}^L W_k \sigma(A_k x + b_k) \leftarrow \text{Cybenko ansatz}$$

are dense in a variety of functional classes.

G. Cybenko, Approximation by superpositions of a sigmoidal function, 1989.

By universal approximation, whatever the data set is, it can be classified simply because we can approximate with the Cybenko ansatz any characteristic function taking value 1 in one set of items and 0 in the complementary one, allocating in this way the right label to each item.

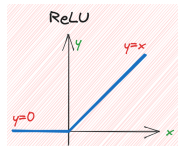
Given that universal approximation guarantees that all goals can be achieved by identifying the right parameters in Cybenko's ansatz, we may adopt the least squares point of view and look for those parameter values minimizing the distance to the needed function, in the so-called **training phase**.

# Classification by simultaneous control

## Theorem

Consider the NODE  $\dot{\mathbf{x}}(t) = \mathbf{w}(t)\sigma(\mathbf{a}(t) \cdot \mathbf{x}(t) + b(t))$  with  $t \in (0, T)$ ,  $\mathbf{w}, \mathbf{a} \in \mathbb{R}^d$ ,  $b \in \mathbb{R}$  and ReLU activation function

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



Then, in any time horizon  $[0, T]$ , a finite arbitrary number  $N$  of items can be driven to an arbitrary number of open subsets of the Euclidean space corresponding to its labels. Moreover

- Controls are piecewise constant with at most  $\mathcal{O}(N)$  switches.
- The control time  $T > 0$  can be made arbitrarily small (absorbed into  $\|\mathbf{w}\|$  by scaling).

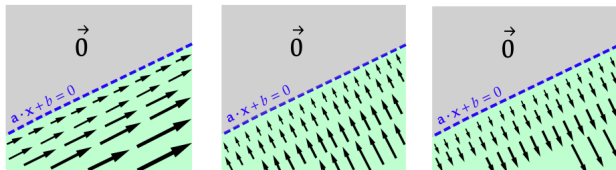
D. Ruiz-Balet and E. Zuazua, *Neural ODE control for classification, approximation, and transport*, SIAM Rev., 2023



# What is the ResNet doing? Basic control actions

$$\dot{\mathbf{x}}(t) = \mathbf{w}(t)\sigma(\mathbf{a}(t) \cdot \mathbf{x}(t) + b(t)) \rightarrow \text{flow map } \phi_T : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

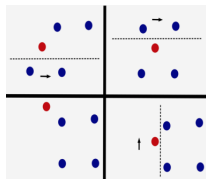
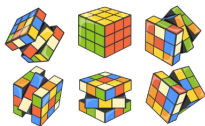
- Control functions  $(\mathbf{w}, \mathbf{a}, b)$   $\rightarrow$  piecewise constant
- Each time discontinuity  $\sim$  one change of layer
- $\mathbf{a}(t)$  and  $b(t)$  define a hyperplane  $H(x) = \mathbf{a}(t) \cdot \mathbf{x}(t) + b(t) = 0$  in  $\mathbb{R}^d$
- $\sigma(x) = \max\{x, 0\}$  activates  $H(x) > 0$  and freezes  $H(x) \leq 0$
- $\mathbf{w}(t)$  gives the direction of the field inside the half-space  $H(x) > 0$



# What is the ResNet doing? Basic control actions

The very features of the ReLU enable the realization of the simultaneous control goal. The fact that  $\sigma$  preserves one half-space invariant while deforming the other half-space introduces dynamics not found in mechanics, where such simultaneous control is unlikely.

The NODE driven by the ReLU behaves like a Rubik Cube, **solvable in a finite number of smart operations** in which part of the cube is frozen, while the other one rotates in the appropriate direction and sense. The goal in the Rubik Cube game is to ensure that all faces are homogeneous in color, similar to the task that NODEs have to perform: drive each initial item to a given target.



A strategic **inductive choice** of the different hyperplanes (via the controls  $\mathbf{a}$  and  $\mathbf{b}$ ) and of the direction of the dynamics, through the control  $\mathbf{w}$ , guarantees classification in a finite number of steps

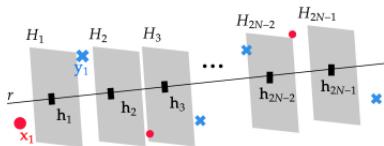
# Can we do better?

**Key idea for an optimal classification:** identifying a **minimal set of hyper-planes** that can separate the points in the dataset according to their labels. The required number of hyperplanes is a **measure of the complexity of the classifier**.

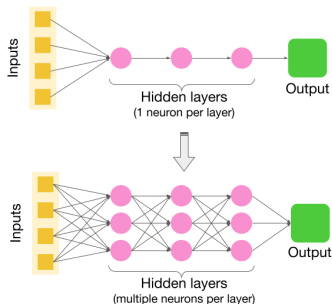
## Two possible way of reducing complexity:

Optimize the classification algorithm:

define sharply the hyperplanes and controls that accomplish the classification task



Increase the width  $p$ : use more neurons



# Optimize the classification algorithm

Classification through NODEs is intrinsically related with organizing the data points according to their labels using hyper-planes.

## WARNING!!

Not all hyper-planes families is adequate for classification. For example, if we were to divided our domain in polygonal cells separating the points of different classes, it would be not at all clear how to move them.

Classification via NODEs is **more restrictive** than separating the points and must be done properly. Actually, we can ensure the existence of an effective control  $\mathbf{w}$  only when the hyper-planes have an appropriate structure.

# Optimize the classification algorithm

## Theorem

Let  $2 \leq d \leq 2N$ . Consider a dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{2N} \subset \mathbb{R}^d \times \{1, 0\}$  such that  $\{\mathbf{x}_i\}_{i=1}^N \subset \mathbb{R}^d$  and  $\{\mathbf{x}_{i+N}\}_{i=1}^N \subset \mathbb{R}^d$  are in **general position**. Then

- For any  $T > 0$ , there exists piecewise constant controls

$$(\mathbf{w}, \mathbf{a}, \mathbf{b}) \in L^\infty((0, T); \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R})$$

such that

$$\phi_T(\mathbf{x}_i) > 1 \quad \text{and} \quad \phi_T(\mathbf{x}_{i+N}) < 1, \quad \text{for all } i \in \{1, \dots, N\}$$

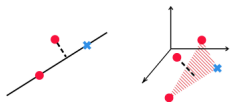
- The number of discontinuities of  $(\mathbf{w}, \mathbf{a}, \mathbf{b})$  is of the order  $\mathcal{O}\left(1 + \frac{N}{d}\right)$ .

A. Alvarez-López, R. Orive and E. Zuazua, Optimized classification with neural ODEs via separability, 2023

# Optimize the classification algorithm

## General position

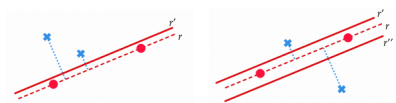
A set  $X$  of  $N$  points in  $\mathbb{R}^d$ , with  $N \geq d$ , is in **general position** if no hyperplane of dimension  $d - 1$  contains more than  $d$  points of  $X$ .



$X$  is in general position in  $\mathbb{R}^2$  if no three points lie on the same line

$X$  is in general position in  $\mathbb{R}^3$  if there are no four points on the same plane

A dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{2N} \subset \mathbb{R}^d \times \{1, 0\}$  such that  $\{\mathbf{x}_i\}_{i=1}^N \subset \mathbb{R}^d$  and  $\{\mathbf{x}_{i+N}\}_{i=1}^N \subset \mathbb{R}^d$  are in general position can be separated with at most  $2 \lceil N/2d \rceil$  hyperplanes.



# Optimize the classification algorithm

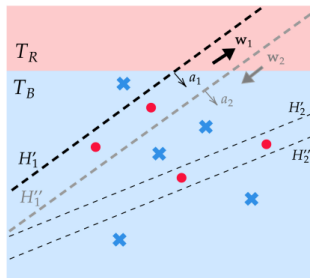
## Strategy

**Step 1 - isolate points according to their classes:** find a family of pairwise parallel hyperplanes

$$\mathcal{H} = \{H'_1, H''_2, \dots, H'_{\lceil N/d \rceil}, H''_{\lceil N/d \rceil}\}$$

so that the region bounded by the pair  $(H'_i, H''_i)$  contains only points of the same class. This corresponds to find suitable controls  $\mathbf{a}$  and  $\mathbf{b}$ .

**Step 2 - control:** move each group of point to its corresponding region by means of suitable designed controls  $\mathbf{w}$ .



To move the two red points between  $H'_1$  and  $H''_1$  to the red region  $T_R$ :

1. we move in direction  $\mathbf{w}_1$  the region defined by  $\mathbf{a}_1$ .
2. we move in direction  $\mathbf{w}_2$  the region defined by  $\mathbf{a}_2$ .

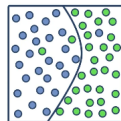
# Optimize the classification algorithm

The proposed methodology can classify any set of points, independently of their distribution. At the same time, it may not be useful in real-world classification problems.

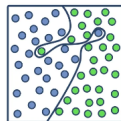
The algorithm characterizes the capacity of a NODE model to classify  $2N$  points in the worst-case scenario. It provides an upper bound of the controls complexity sufficient to induce overfitting. Said differently, the algorithm can serve to know when overfitting arises, depending on the model and dataset dimensions.

**Overfitting** happens when a model learns the training data too well, capturing not only the underlying patterns but also the noise and outliers. This leads to poor performance on new unseen data, because the model is too complex.

The good functioning of any real-life algorithm requires **avoiding overfitting**.



Optimal



Overfitting



# Increase the network's width

1 neuron ( $p=1$ )

$$\dot{\mathbf{x}}(t) = \mathbf{w}(t)\sigma(\mathbf{a}(t) \cdot \mathbf{x}(t) + b(t))$$

multiple neurons ( $p>1$ )

$$\dot{\mathbf{x}}(t) = \sum_{j=1}^p \mathbf{w}_j(t)\sigma(\mathbf{a}_j(t) \cdot \mathbf{x} + b_j) \quad (1)$$

$\mathbf{w}_j$  : columns of  $\mathbf{W} \in \mathbb{R}^{d \times p}$

$\mathbf{a}_j$  : rows of  $\mathbf{A} \in \mathbb{R}^{p \times d}$

$b_j$  : coordinates of  $\mathbf{b} \in \mathbb{R}^d$

## Theorem

Let  $d \geq 2$  and  $T > 0$  be fixed. Consider a dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \subset \mathbb{R}^d \times \mathbb{R}^d$  such that  $\mathbf{x}_i \neq \mathbf{x}_j$  and  $\mathbf{y}_i \neq \mathbf{y}_j$  for all  $i, j \in \{1, \dots, N\}$ . Then

- For any  $p \geq 1$ , there exists piecewise constant controls

$$(\mathbf{W}, \mathbf{A}, \mathbf{b}) \in L^\infty((0, T); \mathbb{R}^{d \times p} \times \mathbb{R}^{p \times d} \times \mathbb{R}^p)$$

such that the flow  $\phi_T$  generated by (1) **interpolates**  $\mathcal{D}$ , i.e.,  $\phi_T(\mathbf{x}_i) = \mathbf{y}_i$ , for all  $i \in \{1, \dots, N\}$ .

- The number of discontinuities of  $(\mathbf{W}, \mathbf{A}, \mathbf{b})$  is of the order  $\mathcal{O}\left(\frac{N}{p} - 1\right)$ .

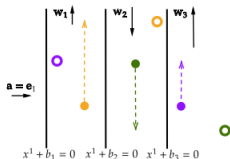
# Increase the network's width

## Strategy

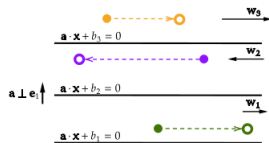
**Step 1 - control  $d - 1$  coordinates:** for all  $j \in \{1, \dots, p\}$  we take  $\mathbf{a}_j = \mathbf{e}_1$  and  $b_j \in \mathbb{R}$  satisfying  $-b_1 < x_1^{(1)} < -b_2 < x_2^{(1)} < \dots < -b_p < x_p^{(1)}$ . This selection is possible up to a change of coordinates.

These controls confine the trajectory of each data point  $\mathbf{x}_i$  within the hyperplane  $x_j^{(1)} = b_j$ . By properly choosing the velocities  $\mathbf{w}_j$ , we can control  $d - 1$  coordinates of these points using  $\mathcal{O}\left(\frac{N}{p}\right)$  switches.

**Step 2 - control the last coordinate:** use the controlled coordinates and suitable designed velocities  $\mathbf{w}_j$  to adjust  $x_i^{(1)}$ . This also can be done with  $\mathcal{O}\left(\frac{N}{p}\right)$  switches.



Control  $d - 1$  coordinates



Control the remaining coordinate

# Increase the network's width

## Remark

If the target points  $\{\mathbf{y}_i\}_{i=1}^N$  are not distinct, interpolation is not achievable due to the uniqueness of solutions in the NODE. In such cases, we can relax the statement from **exact** to **approximate** controllability by applying the theorem to an  $\varepsilon$ -perturbation of the targets.

# Additional bibliography

- D. Ruiz-Balet and E. Zuazua, *Control of neural transport for normalizing flows*, J. Math. Pres et Appl., 2024.
- Z. Li, K. Liu, L. Liverani and E. Zuazua, *Universal approximation of dynamical systems by semi-autonomous Neural ODEs and applications*, preprint, 2024
- M. Hernández and E. Zuazua, *Deep neural networks: multi-classification and universal approximation*, preprint, 2024
- Y. Song, Z. Wang and E. Zuazua, *Approximate and weighted data reconstruction attack in Federated Learning*, preprint, 2023
- Y. Song, Z. Wang, E. Zuazua, *FedADMM-InSa: an Inexact and self-Adaptive ADMM for Federated Learning*, Neural Networks, 2024
- B. Geshkovski and E. Zuazua, *Turnpike in optimal control of PDEs, ResNets, and beyond*, Acta Num., 2022
- C. Esteve and B. Geshkovski, *Sparsity in long-time control of neural ODEs*, Syst. Control. Lett., 2023
- B. Geshkovski, C. Letrouit, Y. Polyanskiy and P. Rigollet, *The emergence of clusters in self-attention dynamics*, Adv. Neur. Inf. Process. Syst., 2024.
- B. Geshkovski, C. Letrouit, Y. Polyanskiy and P. Rigollet, *A mathematical perspective on transformers*, preprint, 2024
- A. Alcalde, G. Fantuzzi and E. Zuazua, *Clustering in pure-attention hardmax transformers and its role in sentiment analysis*, preprint, 2024

ERC Advanced Grant **CoDeFeL**

PI: Enrique Zuazua

Universidad de Deusto

Friedrich-Alexander-Universität

Universidad Autónoma de Madrid



**European Research Council**

Established by the European Commission

Develop new tools with a solid mathematical foundation suitable for data-aware/physics-inspired modeling, combining ML and control theoretical ideas.

Four Work Packages:

**WP1:** Asymptotics and Turnpike for Deep Neural Networks

**WP2:** Complexity of ResNet dynamics

**WP3:** Federated Learning

**WP4:** Modelling through Control and Machine Learning

**Pre-doc** and **postdoc** students applications are welcome.

# THANKS FOR YOUR ATTENTION!!

The **CoDeFeL** project (ERC-2022-ADG) has received funding from the European Union's Horizon ERC Grants programme under grant agreement No. 101096251. The views and opinions expressed are solely those of the author(s) and do not necessarily reflect those of the European Research Council Executive Agency (ERCEA), the European Union or the granting authority who cannot be held responsible for them.

## Additional funding

- PID2020-112617GB-C22-**KiLearn** (MINECO)
- PID2023-146872OB-I00-**DyCMaMod** (MINECO)
- TED2021-131390B-I00-**DasEl** (MINECO)



**European Research Council**

Established by the European Commission

