



HAL
open science

Optimal control for neural ODE in a long time horizon and applications to the classification and simultaneous controllability problems

Jon Asier Bárcena-Petisco

► **To cite this version:**

Jon Asier Bárcena-Petisco. Optimal control for neural ODE in a long time horizon and applications to the classification and simultaneous controllability problems. 2024. hal-03299270v4

HAL Id: hal-03299270

<https://hal.science/hal-03299270v4>

Preprint submitted on 4 Jan 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Optimal control for neural ODE in a long time horizon
2 and applications to the classification and ensemble
3 controllability problems

4 Jon Asier Bárcena-Petisco*

5 January 4, 2024

6 **Abstract:** We study the optimal control, in a long time horizon, of neural ordinary differ-
7 ential equations which are control-affine or whose activation function is homogeneous. When
8 considering the classical regularized empirical risk minimization problem we show that, in long
9 time and under structural assumption on the activation function, the final state of the optimal
10 trajectories has zero training error if the data can be interpolated and if the error can be taken
11 to zero with a cost proportional to the error. These hypotheses are fulfilled in the classification
12 and ensemble controllability problems for some relevant activation and loss functions. Finally,
13 we show the sharpness of our hypotheses by giving an example for which the error of the final
14 state of the optimal trajectory, even if it decays, is strictly positive for any time.

15 **Key words:** data classification, exact controllability, neural ODE, nonlinear systems, optimal
16 control, ensemble controllability

17 **Abbreviated title:** Optimal control for neural ODE

18 **AMS subject classification:** 34H05, 49N10, 93B05

*Department of Mathematics, University of the Basque Country UPV/EHU, Barrio Sarriena s/n, 48940, Leioa, <https://orcid.org/0000-0002-6583-866X> Spain. E-mail: jonasier.barcena@ehu.eus.

1 Introduction

In this paper we study the optimal control of neural ordinary differential equations for a long time horizon. Neural ODE have been used in Machine Learning in the last seven years, a trend started with [43, 16]. However, they date back to the 90s, when they were already used for the construction of controls (see the survey [35]) and when their controllability properties were first studied (see, for example, [47] and [34]). The control systems governed by neural ODE have considerably better controllability properties than linear control systems. In fact, as pointed out in [27], for a fixed $d \in \mathbb{N}$, if chosen the right neural ODE we can interpolate an arbitrarily large amount of data in \mathbb{R}^d , whereas in linear systems we can at most interpolate an amount of data equal to the dimension of the control. In this paper d denotes the dimension of the space where each element of the dataset is, and N the size of the dataset.

Roughly, the problem under study is the following: given a set of initial values $\mathbf{x} = (x^1, \dots, x^N) \in (\mathbb{R}^d)_*^N$, for:

$$(\mathbb{R}^d)_*^N := \{(x^1, \dots, x^N) \in (\mathbb{R}^d)^N : x^i \neq x^j \ \forall i, j \in \{1, \dots, N\} : i \neq j\},$$

we seek to take simultaneously the data set to some target points or regions in \mathbb{R}^d in a given time $T > 0$. This is usually called dataset as it is a set of values. The control problem is important in the context of ensemble controllability. The distance to those targets is measured with an error function (also known as *loss function*). The control is the minimizer of the risk minimization functional, which provides a balance between a small cost for the control and a small value for the loss function at the final state of the optimal trajectory. For a detailed introduction to the notation and its background, I recommend [8, 27].

We study the controllability on control-affine neural networks, which are given by the following equations:

$$\begin{cases} \dot{y}(t) = w(t)\sigma(y(t)) + b(t), \\ y(0) = x, \end{cases} \quad (1.1)$$

for $x \in \mathbb{R}^d$ the initial value, and $\sigma : \mathbb{R}^d \mapsto \mathbb{R}^d$ a Lipschitz function, which is called the *activation function*. The functions (w, b) are the controls and they belong to $L^2(0, T; \mathcal{U})$, for \mathcal{U} defined by:

$$\mathcal{U} := \mathbb{R}^{d \times d} \times \mathbb{R}^{d \times 1}.$$

If we want to emphasize the dependence of (1.1) to the initial value and the control, we write $y(\cdot; x, w, b)$. Similarly, we denote the sequence of solutions of (1.1) for some fixed control (w, b) applied simultaneously to a data set \mathbf{x} as:

$$y(\cdot; \mathbf{x}, w, b) := (y(\cdot; x^1, w, b), \dots, y(\cdot; x^N, w, b)). \quad (1.2)$$

Since σ is Lipschitz, (1.1) is well-posed by the Cauchy-Lipschitz Theorem.

1 In addition, we also study more compound neural networks, which are given by the equations:

$$\begin{cases} \dot{y}(t) = r(t)\sigma(w(t)y(t) + b(t)), \\ y(0) = x. \end{cases} \quad (1.3)$$

2 Here x is the initial value and (r, w, b) is the control, which belongs to $L^2(0, T; \tilde{\mathcal{U}})$, for:

$$\tilde{\mathcal{U}} := X \times \mathbb{R}^{d \times d} \times \mathbb{R}^{d \times 1},$$

3 for:

$$X \subseteq \{M \in \mathbb{R}^{d \times d} : M_{i,i} \in \{1, -1\}, \forall i = 1, \dots, d, M_{i,j} = 0, \forall i \neq j\}. \quad (1.4)$$

4 In fact, the intensity of the flow is modelled by (w, b) , and the direction of the flow, by r .
 5 We may take $X = \{I\}$, which makes sense when σ admits negative values. However, we have
 6 considered the general setting to have relevant results also for the case in which σ is a positive
 7 function; that is, in which $\sigma \geq 0$. We assume that the activation function σ is Lipschitz and
 8 homogeneous in the sense that:

$$\sigma(\lambda x) = \lambda \sigma(x), \quad \forall \lambda > 0, \quad \forall x \in \mathbb{R}^d. \quad (1.5)$$

This includes important *activation functions* such as rectified linear units, which are given by:

$$\sigma(x) = (\max\{x_1, 0\}, \dots, \max\{x_d, 0\}),$$

see [24]; and parametric rectified units, given by:

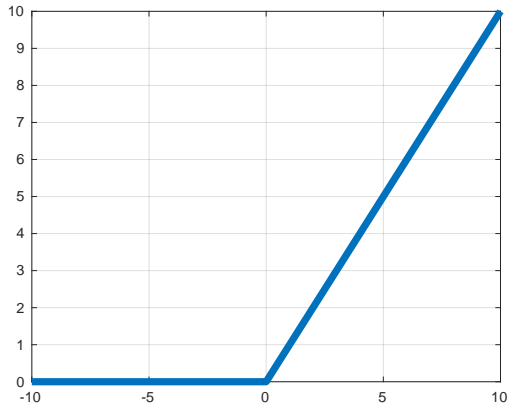
$$\sigma(x) = (\alpha x_1 \mathbf{1}_{x_1 < 0} + x_1 \mathbf{1}_{x_1 > 0}, \dots, \alpha x_d \mathbf{1}_{x_d < 0} + x_d \mathbf{1}_{x_d > 0}),$$

9 see [18] (see Figure 1 for the graphs of such activations functions in one dimension). As in the
 10 previous system:

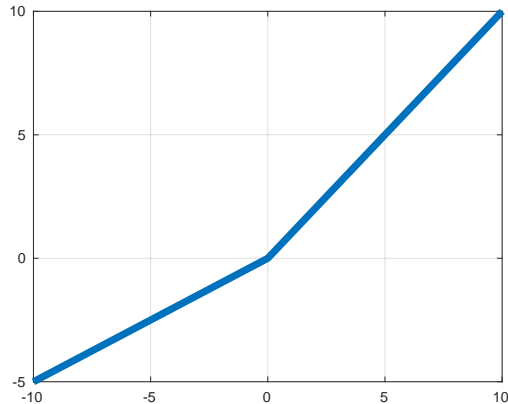
$$y(\cdot; \mathbf{x}, r, w, b) = (y(\cdot; x^1, r, w, b), \dots, y(\cdot; x^N, r, w, b)), \quad (1.6)$$

11 where $y(\cdot; x, r, w, b)$ denotes the solutions of (1.3), which is a well-posed system by the Cauchy-
 12 Lipschitz Theorem.

13 As stated in the first paragraph, we study the properties of any optimal control in a long
 14 time horizon. The main contribution of our paper is that, if the data can be interpolated
 15 and the error can be taken to 0 with a cost proportional to the current error, we improve the
 16 asymptotic bound $\mathcal{O}(1/T)$ for the error of the final state of the optimal trajectory obtained in
 17 [6] and prove that it is exactly 0 for a sufficiently large time. In fact, the paper is inspired in
 18 the simulations presented in [6, Examples 4.2 and 4.4] where the final errors seem to be 0, as we
 19 want to determine theoretically if, like their simulation suggests, the error is taken exactly to
 20 0. Even if approximate controllability is usually enough for practical purposes, obtaining null



(a) Rectified linear units.



(b) Parametric linear units.

Figure 1: Some usual activation functions for $d = 1$.

1 controllability is interesting to broaden the perspective of the field. We work in an abstract
 2 setting, though we give concrete examples of problems that satisfy our assumptions, notably
 3 the ensemble controllability and classification problems. In ensemble controllability we aim to
 4 control two or more independent equations by applying the same control. The study of ensemble
 5 controllability dates back to [28] and [20, Chapter 5], and relevant papers on this topic include
 6 [40, 33, 21, 44, 1, 27, 32, 26]. The main difference of this paper and [27, 26] with the previous
 7 ones is that in our papers the trajectories satisfy the same differential equation, whereas in
 8 the other papers they satisfy different differential equations (i.e. differential equations which
 9 at least do not have the same coefficients). As for the classification problem, it is a simplified
 10 version of the ensemble controllability problem, where the objective is to split the data into
 11 two sets, for instance, $\{x_1 \leq -1\}$ and $\{x_1 \geq 1\}$. An additional contribution of our paper is an
 12 example of neural ODE and loss functions where the error can be taken to 0, but for all time
 13 $T > 0$ the error at time T of the optimal trajectories is strictly positive. This illustrates that
 14 the results are far from being trivial.

15 This paper follows a well-established research line that studies the properties of the optimal
 16 control and trajectories in a long time horizon. This allows, for instance, that when doing
 17 numerical simulations, one may identify when a local minimum is not an optimal control. The
 18 *turnpike property* is a notion developed since the 1950s which means that when minimizing
 19 certain functionals all the optimal trajectories are most of the time near some specific state
 20 (the *turnpike*) independently of the initial value and the target (see, for instance, [22] and
 21 [5]). An important recent paper regarding the study of the turnpike property is [25], the first
 22 work which provides rigorous mathematical proof and a framework for the turnpike property
 23 for linear quadratic optimal control problems. Also, interesting recent studies on the turnpike
 24 property include discrete optimal control problems in [4] and [12], finite-dimensional nonlinear

1 control problems in [39] and [37], optimal control problems for hyperbolic systems in [15] and
 2 [30], general Hilbert spaces in [38] and [2], boundary optimal control problems in [13], Navier-
 3 Stokes equation in [46], nonlinear optimal control problems from a geometrical approach in
 4 [31], fractional parabolic equation in [42], the finite time turnpike phenomena in [14], hands-off
 5 controls in [29] and in deep neural networks in [10], and Lipschitz nonlinear functions in [7].
 6 Finally, the optimal control for neural ODE is also studied in [8], where the authors consider
 7 the cost of the L^1 -norm of the control instead of the L^2 -norm and obtain that for that norm the
 8 optimal control satisfies some sparsity properties and that the error of the final state belongs
 9 to $\mathcal{O}(1/T)$.

10 2 Main results

11 2.1 Optimal trajectories for control-affine neural ODE

12 As stated in the introduction, we study the optimal control of a data set ruled by a neural
 13 ODE. To measure how far the data is from the objective we introduce the *error function* (also
 14 referred in the literature of Machine Learning as *loss function*) $\mathcal{E} : (\mathbb{R}^d)^N \mapsto \mathbb{R}^+ := [0, \infty)$. We
 15 assume that \mathcal{E} is continuous and satisfies the Hypothesis 1, which is later introduced in this
 16 section.

17 This allows to define the *empirical risk minimization functional for a target time T* :

$$J_T(w, b) := \mathcal{E}(y(T; \mathbf{x}, w, b)) + \int_0^T |(w(t), b(t))|^2 dt, \quad (2.1)$$

18 where y denotes a solution of (1.1) and $|\cdot|$ denotes the Frobenius norm. We denote any
 19 minimizer of J_T by (w_T, b_T) . Moreover, the trajectories induced by such minimizers, called
 20 *optimal trajectories*, are denoted by $y_T(t; \mathbf{x}) := y(t; \mathbf{x}, w_T, b_T)$.

21 *Example 2.1.* A usual definition for the error function is:

$$\mathcal{E}(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N E_i(x^i), \quad \forall \mathbf{x} \in (\mathbb{R}^d)^N, \quad (2.2)$$

22 for $E_i(x) = d(x, A_i)$, for d the euclidean distance and for given sets $A_i \subset \mathbb{R}^d$ (that might consist
 23 of a single element).

24 First of all, we recall that the functional J_T has at least a minimizer:

25 **Proposition 2.2** (Existence of minimizers). *Let $\mathcal{E} : (\mathbb{R}^d)^N \mapsto \mathbb{R}^+ := [0, \infty)$ a continuous*
 26 *function, σ a globally Lipschitz continuous function, $T > 0$ and $\mathbf{x} \in (\mathbb{R}^d)^N$. Then, the functional*
 27 *J_T given in (2.1) for y given by (1.2), where we consider the solution of (1.1), has at least one*
 28 *minimizer in $L^2(0, T; \mathcal{U})$.*

1 Proposition 2.2 is classical, and the proof can be found, for instance, in [36, Proposition 6.2.3].
 2 The main idea of the proof is that J_T is a sum of a positive weakly continuous functional and
 3 a positive continuous convex functional. For the sake of completeness, the proof is given in
 4 Appendix A.

5 Let us now present the hypotheses that we consider throughout the paper:

6 *Hypothesis 1.* Let $\mathbf{x} \in (\mathbb{R}^d)_*^N$, let $\mathcal{E} : (\mathbb{R}^d)^N \mapsto \mathbb{R}^+ := [0, \infty)$ be a continuous function, and let
 7 y denote (1.2), where we consider the solutions of (1.1). Then,

1. For the data set \mathbf{x} there are controls:

$$(w_*, b_*) \in L^2(0, 1; \mathcal{U}),$$

8 such that $\mathcal{E}(y(1; \mathbf{x}, w_*, b_*)) = 0$.

2. There are $C, \tilde{\varepsilon} > 0$ both just depending on \mathcal{E} such that for all $\bar{\mathbf{x}} = (\bar{x}^1, \dots, \bar{x}^N) \in (\mathbb{R}^d)_*^N$ satisfying $\mathcal{E}(\bar{\mathbf{x}}) < \tilde{\varepsilon}$, there are some controls (w, b) satisfying:

$$\|(w, b)\|_{L^\infty(0, 1; \mathcal{U})} < C\mathcal{E}(\bar{\mathbf{x}}),$$

such that:

$$\mathcal{E}(y(1; \bar{\mathbf{x}}, w, b)) = 0.$$

9 The first item of Hypothesis 1 is that the error can be taken to 0, a property known in
 10 Machine Learning as *interpolation* (see [6]), and the second one is a local controllability of the
 11 system.

12 *Remark 2.3.* The choice of the target time in Hypothesis 1 is arbitrary. Because of the linearity
 13 (see Lemma 3.1 below), if the system is controllable for some time, in this case $T = 1$, it is
 14 controllable for any time.

Example 2.4 (Application of Theorem 2.5 to the classification problem). Let us fix $M \in \mathbb{N}$ and consider:

$$\mathbf{x} = (x^1, \dots, x^M, x^{M+1}, \dots, x^N) \in (\mathbb{R}^d)_*^N,$$

the error function given by (2.2), for:

$$E_i(x) = \begin{cases} (x_1 + 1)1_{x_1 > -1}(x_1), & i = 1, \dots, M, \\ (x_1 - 1)1_{x_1 > 1}(x_1), & i = M + 1, \dots, N, \end{cases}$$

15 and any neural function σ of the type $\sigma(x) = (\tilde{\sigma}(x_1), \dots, \tilde{\sigma}(x_d))$ such that there is $c > 0$
 16 such that $cs \leq \tilde{\sigma}(s)$ for all $s \geq 0$ and $\tilde{\sigma}(s) \leq cs$ for all $s \leq 0$. The second item of Hy-
 17 pothesis 1 is clearly satisfied, as it suffices to consider $\tilde{\varepsilon} = 1/(2N)$, $b = 0$ and $w(t)x =$

1 $(2Nc^{-1}\mathcal{E}(\bar{\mathbf{x}})x_1, 0, \dots, 0)$. Thus, Theorem 2.5 implies that if the data can be classified (i.e.
 2 if the first item of Hypothesis 1 is satisfied), then by computing the optimal control for a suffi-
 3 ciently large time, the data is sent to the sets $\{x_1 \leq -1\}$ and $\{x_1 \geq 1\}$. More detailed examples
 4 can be found in [6] and [27].

5 Now we have all the tools to state the first main result of this paper:

6 **Theorem 2.5** (Annihilation of the error in a long time horizon). *Let $\mathbf{x} \in (\mathbb{R}^d)_*^N$, σ be a
 7 Lipschitz activation function, \mathcal{E} be an error function such that Hypotheses 1 is satisfied and J_T
 8 given in (2.1). Then, for $T > 0$ large enough depending on σ , \mathbf{x} and \mathcal{E} , and for all $\varepsilon > 0$ there
 9 is $\delta > 0$ such that $J_T(w, b) < \inf J_T + \delta$ implies:*

$$\mathcal{E}(y(T; \mathbf{x}, w, b)) < \varepsilon. \quad (2.3)$$

10 Moreover, for $T > 0$ large enough the following equality holds for any optimal trajectory:

$$\mathcal{E}(y_T(T; \mathbf{x})) = 0. \quad (2.4)$$

11 Here, y is given by (1.2), where we consider the solution of (1.1).

Theorem 2.5 is proved by showing that if T is sufficiently large and if $\mathcal{E}(y(T; \mathbf{x}, w, b))$ is small
 and strictly positive, we can construct with the second item of Hypothesis 1 a control (\tilde{w}, \tilde{b})
 such that:

$$J_T(\tilde{w}, \tilde{b}) \leq J_T(w, b) - \frac{1}{2}\mathcal{E}(y(T; \mathbf{x}, w, b)).$$

12 For that, we show in Lemma 3.1 that the trajectories may be preserved when we perform a
 13 diffeomorphism in the time variable. Then, in Lemma 3.4 given a control with a non-constant
 14 norm we construct a more efficient one and in Proposition 3.5 we use this to construct a control
 15 for which the value of the empirical risk minimization functional is smaller for all controls with
 16 a non-constant norm.

17 The construction of such control is far from trivial and, as illustrated in Appendix B, the
 18 hypotheses are rather sharp. As explained in the first part of the introduction, Theorem 2.5
 19 improves the results presented in [6], where the authors prove that the error of the final state
 20 of the optimal trajectory is of size $\mathcal{O}(1/T)$.

21 **2.2 Optimal trajectories for neural ODE with a homogeneous acti- 22 vation function**

23 In this section we present the analogous results to those in Section 2.1 for the neural ODE
 24 (1.3) with activation functions which satisfy (1.5). Let us reformulate Hypothesis 1 in the
 25 context of (1.3):

1 *Hypothesis 2.* Let $\mathbf{x} \in (\mathbb{R}^d)_*^N$, let $\mathcal{E} : (\mathbb{R}^d)^N \mapsto \mathbb{R}^+ := [0, \infty)$ be a continuous function, and let
 2 y denote (1.6), where we consider the solutions of (1.3). Then:

1. For the data set \mathbf{x} there are controls:

$$(r_*, w_*, b_*) \in L^2(0, 1; \mathcal{U}),$$

3 such that $\mathcal{E}(y(1; \mathbf{x}, r_*, w_*, b_*)) = 0$.

2. There are $C, \tilde{\varepsilon} > 0$ both just depending on \mathcal{E} such that for all $\bar{\mathbf{x}} = (\bar{x}^1, \dots, \bar{x}^N) \in (\mathbb{R}^d)_*^N$
 satisfying $\mathcal{E}(\bar{\mathbf{x}}) < \tilde{\varepsilon}$, there are some controls (r, w, b) satisfying:

$$\|(w, b)\|_{L^\infty(0, 1; \mathcal{U})} < C\mathcal{E}(\bar{\mathbf{x}}),$$

such that:

$$\mathcal{E}(y(1; \bar{\mathbf{x}}, r, w, b)) = 0.$$

4 *Example 2.6* (Hypothesis 2 in a context of ensemble controllability). Hypothesis 2 can be
 5 considered in an ensemble controllability problem. Let $\mathbf{x} \in (\mathbb{R}^d)_*^N$ for $d \geq 2$, X given in (1.4):

$$\sigma(x) = (\max\{x_1, 0\}, \dots, \max\{x_d, 0\}), \quad (2.5)$$

6 the activation function, $\mathbf{z} = (z^1, \dots, z^N) \in (\mathbb{R}^d)_*^N$ the targets, and \mathcal{E} given by (2.2) for $E_i(x) =$
 7 $|x - z^i|$ the error function. Note that σ satisfies:

$$|\sigma(u)| \leq |u| \quad \forall u \in \mathbb{R}^d. \quad (2.6)$$

8 Then, it is proved in [27, Theorem 2] that the first item of Hypothesis 2 is satisfied. Moreover,
 9 as we prove in Appendix D, the second item of Hypothesis 2 also holds. We present the proof
 10 because the bounds for the cost of the control is not a straight consequence of the computations
 11 in [27]. Consequently, Theorem 2.7 below (and all the auxiliary results and corollaries) can be
 12 applied to this neural problem.

13 Again, we seek to get sufficient conditions so that the optimal trajectories induced by:

$$\tilde{J}_T(r, w, b) := \mathcal{E}(y(T; \mathbf{x}, r, w, b)) + \int_0^T |(w(t), b(t))|^2 dt, \quad (2.7)$$

14 satisfy $\mathcal{E}(y_T(T; \mathbf{x})) = 0$. Since $|r|$ is constant (see (1.4)), it makes no sense to include it in the
 15 definition of \tilde{J}_T . For the functional \tilde{J}_T the following result holds:

16 **Theorem 2.7** (Annihilation of the error for a sufficiently large time). *Let σ be a Lipschitz*
 17 *activation function satisfying (1.5) and \mathcal{E} an error function satisfying Hypothesis 2. Then,*
 18 *for $T > 0$ large enough depending on σ , \mathbf{x} and \mathcal{E} , and all $\varepsilon > 0$ there is $\delta > 0$ such that if*
 19 $J_T(r, w, b) < \inf J_T + \delta$:

$$\mathcal{E}(y(T; \mathbf{x}, r, w, b)) < \varepsilon. \quad (2.8)$$

1 Moreover, if T is large enough and if \tilde{J}_T has an optimal trajectory:

$$\mathcal{E}(y_T(T; \mathbf{x})) = 0. \quad (2.9)$$

2 Here y is given by (1.2), where we consider the solution of (1.1).

3 The proof of Theorem 2.7 is analogous to that of Theorem 2.5, so we just give some brief
 4 explanations in the first comment of Section 4. As with Theorem 2.5, Theorem 2.7 improves
 5 the results presented in [6], where the authors prove that the error of the optimal trajectory
 6 at a final time T is of magnitude $\mathcal{O}(1/T)$ also for the solutions of (1.3) with an activation
 7 functions satisfying (1.5).

8 *Remark 2.8* (Existence of minimizers of \tilde{J}_T). We have stated “if \tilde{J}_T has an optimal trajectory” in
 9 Theorem 2.7 because, as far as we know, it is an open question to see if \tilde{J}_T admits a minimizer.
 10 The main obstacle to adapt the proof of Proposition 2.2 is that nonlinear functions and weak
 11 limits may not commute. However, as we see in the first comment of Section 4, we can improve
 12 Theorem 2.7 and obtain that for T large enough and all $\varepsilon > 0$ there are controls (r, w, b) such
 13 that $J_T(r, w, b) < \inf J_T + \varepsilon$ and $\mathcal{E}(y(T; \mathbf{x}, r, w, b)) = 0$.

14 *Remark 2.9* (Functionals allowing expensive controls). As in [6], we can consider the functional:

$$J_{T,\delta}(w, b) := \mathcal{E}(y(T; \mathbf{x}, w, b)) + \delta \int_0^T |(w(t), b(t))|^2 dt,$$

15 instead of J_T for (1.1), and:

$$J_{T,\delta}(r, w, b) := \mathcal{E}(y(T; \mathbf{x}, r, w, b)) + \delta \int_0^T |(w(t), b(t))|^2 dt,$$

16 instead of J_T for (1.3)-(1.5). By linearity (see Remark 3.2) it holds that:

$$J_{T,\delta}(w, b) = J_{T\delta^{-1},1}(\delta w(t\delta), \delta b(t\delta)),$$

17 and:

$$\tilde{J}_{T,\delta}(r, w, b) = \tilde{J}_{T\delta^{-1},1}(r(t\delta), \delta w(t\delta), \delta b(t\delta)),$$

18 respectively. A straight consequence is that (w, b) is a minimizer of $J_{T,\delta}$ if and only if $(\delta w(t\delta), \delta b(t\delta))$
 19 is a minimizer of $J_{T\delta^{-1},1}$. Similarly, (r, w, b) is a minimizer of $\tilde{J}_{T,\delta}$ if and only if $(r(t\delta), \delta w(t\delta), \delta b(t\delta))$
 20 is a minimizer of $\tilde{J}_{T\delta^{-1},1}$. Thus, analogous results to Theorems 2.5 and 2.7 and all the auxiliary
 21 results hold true for $J_{T,\delta}$ and $\tilde{J}_{T,\delta}$ when T is fixed and $\delta > 0$ is small enough depending on σ ,
 22 \mathcal{E} , \mathbf{x} and T .

23 3 Optimal control for control-affine neural ODE

24 In this section we work in the control problem described by (1.1) and the risk minimization
 25 functional J_T given by (2.1). In this section $C > 0$ denotes an arbitrary constant that may

1 change from line to line and depends only on σ , \mathcal{E} and \mathbf{x} . Similarly, when we assume that T is
2 large enough we mean with respect to σ , \mathcal{E} and \mathbf{x} . We first present some technical results in
3 Section 3.1, then conclude the proof of Theorem 2.5 in Section 3.2 by a proof by contradiction,
4 and finally provide additional properties of the optimal controls in Section 3.3.

5 3.1 Preliminaries

6 We first construct controls to follow the same trajectory in the state space but with a different
7 velocity by reparametrization thanks to the structure of the controls:

Lemma 3.1 (A technical result regarding the time variable). *Let $\mathbf{x} \in (\mathbb{R}^d)^N$ and:*

$$\phi \in L^1_{loc}(0, \infty; \mathbb{R}^+).$$

8 *Then:*

$$y(\Gamma(t); \mathbf{x}, w, b) = y\left(t; \mathbf{x}, \phi(\Gamma(s))w(\Gamma(s)), \phi(\Gamma(s))b(\Gamma(s))\right), \quad \forall t \in [0, T^*], \quad (3.1)$$

9 *for $T^* > 0$, y given by (1.2), where we consider the solutions of (1.1), and Γ any solution of:*

$$\begin{cases} \dot{\Gamma}(s) = \phi(\Gamma(s)), & s \in [0, T^*], \\ \Gamma(0) = 0. \end{cases} \quad (3.2)$$

10 *Remark 3.2* (Invariance of trajectories when ϕ is constant). An important application of Lemma
11 3.1 is the case $\phi(t) = \lambda \in \mathbb{R}^+$; that is, when ϕ is constant. Then, (3.1) becomes:

$$y(\lambda t; \mathbf{x}, w, b) = y(t; \mathbf{x}, \lambda w(\lambda s), \lambda b(\lambda s)). \quad (3.3)$$

Proof of Lemma 3.1. It suffices to see that for all i the function $t \mapsto y(\Gamma(t); x^i, w, b)$ is a solution
of (1.1) with initial value x^i and controls $\phi(\Gamma(t))w(\Gamma(t))$ and $\phi(\Gamma(t))b(\Gamma(t))$, since (1.1) has a
unique solution by the Cauchy-Lipschitz Theorem. From the initial condition on (3.2) we obtain
that:

$$y(\Gamma(0); x^i, w, b) = y(0; x^i, w, b) = x^i.$$

12 Moreover, from the first equation of (3.2) and the chain rule:

$$\begin{aligned} \frac{d}{dt} \left(y(\Gamma(t); x^i, w, b) \right) &= \phi(\Gamma(t)) \dot{y}(\Gamma(t); x^i, w, b) \\ &= \phi(\Gamma(t)) \left(w(\Gamma(t)) \sigma(y(\Gamma(t); x^i, w, b)) + b(\Gamma(t)) \right) \\ &= \left(\phi(\Gamma(t)) w(\Gamma(t)) \right) \sigma \left(y(\Gamma(t); x^i, w, b) \right) + \left(\phi(\Gamma(t)) b(\Gamma(t)) \right). \end{aligned} \quad (3.4)$$

13 □

1 Next, we recall that the first item of Hypothesis 1 implies that the error is at most of size
 2 $\mathcal{O}(1/T)$:

3 **Lemma 3.3** (Boundedness of the error with respect to T). *Let $T > 0$, σ be an activation
 4 function and \mathcal{E} an error function satisfying Hypothesis 1. Then:*

$$\mathcal{E}(y_T(T; \mathbf{x})) \leq \frac{C}{T}, \quad (3.5)$$

5 for y_T given by (1.2), where we consider the solutions with the optimal control of (1.1).

6 Lemma 3.3 is proven in [6, Theorem 3.1]. Note that in [6] they only require the first item of
 7 Hypothesis 1 (as well as the positivity and continuity of \mathcal{E}), so we can use their result. Briefly,
 8 it is a consequence of the definition of y_T as the optimal trajectory and that, by Remark 3.2,
 9 $\frac{1}{T}w_*(\frac{\cdot}{T})$ and $\frac{1}{T}b_*(\frac{\cdot}{T})$ are controls that take the error to 0 (see Hypothesis 1 for the definition of
 10 (w_*, b_*)). In fact:

$$\begin{aligned} \mathcal{E}(y_T(T; \mathbf{x})) &\leq J_T(w_T, b_T) \leq J_T\left(\frac{1}{T}w_*\left(\frac{\cdot}{T}\right), \frac{1}{T}b_*\left(\frac{\cdot}{T}\right)\right) \\ &= \frac{1}{T} \int_0^1 |(w_*(t), b_*(t))|^2 dt. \end{aligned} \quad (3.6)$$

11 Finally, we show how to construct a more efficient control when the norm is not constant
 12 and they do not take the null value:

Lemma 3.4 (Construction of more efficient controls). *Let:*

$$(w, b) \in C^1([0, T]; \mathcal{U} \setminus \{(0_{\mathbb{R}^{d \times d}}, 0_{\mathbb{R}^d})\}),$$

13 *be such that $t \mapsto |(w(t), b(t))|$ is not constant. Then, there is a control (\tilde{w}, \tilde{b}) such that $t \mapsto$
 14 $|(\tilde{w}(t), \tilde{b}(t))|$ is constant, such that:*

$$|(\tilde{w}(t), \tilde{b}(t))| \in \left(\min_{[0, T]} |(w, b)|, \max_{[0, T]} |(w, b)| \right) \text{ in } [0, T], \quad (3.7)$$

15 *and such that:*

$$y(T; \mathbf{x}, \tilde{w}, \tilde{b}) = y(T; \mathbf{x}, w, b), \quad \int_0^T |(\tilde{w}(t), \tilde{b}(t))|^2 dt < \int_0^T |(w(t), b(t))|^2 dt, \quad (3.8)$$

16 for y given by (1.2), where we consider the solutions of (1.1).

17 The proof consists on constructing new controls with a time-transformation that allows us
 18 to arrive to the same target with a smaller cost:

1 *Proof of Lemma 3.4.* Let us consider the auxiliary function:

$$\phi_\gamma(t) = \frac{\gamma}{|(w(t), b(t))|} \mathbf{1}_{[0, T]}(t), \quad (3.9)$$

2 for $\gamma > 0$ to be fixed later, and Γ_γ given by:

$$\begin{cases} \dot{\Gamma}_\gamma(s) = \phi_\gamma(\Gamma_\gamma(s)), & s \in [0, T_\gamma], \\ \Gamma_\gamma(0) = 0, \end{cases} \quad (3.10)$$

3 for:

$$T_\gamma := \sup\{t : \Gamma_\gamma(t) < T\}.$$

4 Note that from the definition of T_γ it follows that:

$$\Gamma_\gamma(T_\gamma) = T. \quad (3.11)$$

5 Since ϕ_γ is C^1 (as $(w, b) \neq 0$, by compactness $\min_{[0, T]} |(w, b)| > 0$), (3.10) has a unique
6 solution by the Cauchy-Lipschitz Theorem. Moreover, $\gamma \mapsto T_\gamma$ is continuous and decreasing,
7 $\lim_{\gamma \rightarrow 0} T_\gamma = \infty$, and $\lim_{\gamma \rightarrow \infty} T_\gamma = 0$, so there is $\gamma^* > 0$ such that:

$$T_{\gamma^*} = T. \quad (3.12)$$

8 For that value γ^* it holds:

$$\gamma^* = \frac{1}{\frac{1}{T} \int_0^T \frac{dt}{|(w(\Gamma_{\gamma^*}(t)), b(\Gamma_{\gamma^*}(t)))|}}, \quad (3.13)$$

9 because by (3.11) and (3.12):

$$\int_0^T \frac{\gamma^* dt}{|(w(\Gamma_{\gamma^*}(t)), b(\Gamma_{\gamma^*}(t)))|} = \int_0^T \dot{\Gamma}_{\gamma^*}(t) dt = \Gamma_{\gamma^*}(T) = \Gamma_{\gamma^*}(T_{\gamma^*}) = T.$$

10 Considering (3.1) we obtain that:

$$\begin{aligned} y(T; \mathbf{x}, \phi_{\gamma^*}(\Gamma_{\gamma^*}(s))w(\Gamma_{\gamma^*}(s)), \phi_{\gamma^*}(\Gamma_{\gamma^*}(s))b(\Gamma_{\gamma^*}(s))) &= y(\Gamma_{\gamma^*}(T); \mathbf{x}, w, b) \\ &= y(T; \mathbf{x}, w, b). \end{aligned} \quad (3.14)$$

11 Moreover, considering the strict inequality between the harmonic and arithmetic means (see,
12 for instance, [19]), (3.10)-(3.13), and the Change of Variables Theorem we get that:

$$\begin{aligned} \int_0^T \phi_{\gamma^*}^2(\Gamma_{\gamma^*}(t)) |(w(\Gamma_{\gamma^*}(t)), b(\Gamma_{\gamma^*}(t)))|^2 dt &= \int_0^T (\gamma^*)^2 dt = (\gamma^*)^2 T \\ &= \frac{\gamma^* T}{\frac{1}{T} \int_0^T \frac{dt}{|(w(\Gamma_{\gamma^*}(t)), b(\Gamma_{\gamma^*}(t)))|}} \\ &< \frac{1}{T} \int_0^T \gamma^* T |(w(\Gamma_{\gamma^*}(t)), b(\Gamma_{\gamma^*}(t)))| dt \\ &= \int_0^T \phi_{\gamma^*}(\Gamma_{\gamma^*}(t)) |(w(\Gamma_{\gamma^*}(t)), b(\Gamma_{\gamma^*}(t)))|^2 dt \\ &= \int_0^T |(w(t), b(t))|^2 dt. \end{aligned} \quad (3.15)$$

Therefore, combining (3.14) and (3.15) we obtain (3.8) for:

$$\tilde{w}(t) = \phi_{\gamma^*}(\Gamma_{\gamma^*}(t))w(\Gamma_{\gamma^*}(t)), \quad \tilde{b}(t) = \phi_{\gamma^*}(\Gamma_{\gamma^*}(t))b(\Gamma_{\gamma^*}(t)).$$

Finally, since γ^* is the harmonic mean of values in:

$$\left[\min_{[0,T]} |(w, b)|, \max_{[0,T]} |(w, b)| \right],$$

1 we obtain (3.7). □

2 3.2 Construction of controls which take the error to zero

3 Let us state the properties of the controls that we construct in this section:

4 **Proposition 3.5.** *Let σ be an activation function and E an error function that satisfy Hypoth-*
5 *esis 1. Then, there is $T_0 > 0$ such that if $T \geq T_0$ and (w, b) are such that:*

$$J_T(w, b) \leq 2 \inf J_T, \tag{3.16}$$

6 *there is a control (\hat{w}, \hat{b}) such that:*

$$\mathcal{E}(y(T; \mathbf{x}, \hat{w}, \hat{b})) = 0, \tag{3.17}$$

7 *and:*

$$J_T(\hat{w}, \hat{b}) \leq J_T(w, b) - \frac{1}{2} \mathcal{E}(y(T; \mathbf{x}, w, b)), \tag{3.18}$$

8 *for y given by (1.2), where we consider the solutions of (1.1).*

9 The first step is to remark that $\mathcal{E}(y(T; \mathbf{x}, w, b))$ is small for T large enough and all (w, b)
10 satisfying (3.16) by Lemma 3.3. The second step is to approximate (w, b) by some control (\tilde{w}, \tilde{b})
11 with a constant norm thanks to regularization and Lemma 3.4. The third step is to show that
12 if (2.4) is false, we may prolong for some $\tau > 0$ the controls \tilde{w} and \tilde{b} in $[T, T + \tau]$ so that:

$$\tilde{y}(T + \tau; \mathbf{x}, \tilde{w}, \tilde{b}) = 0. \tag{3.19}$$

13 The fourth step is to take those trajectories to $[0, T]$ with (3.3), which we recall is a consequence
14 of Lemma 3.1. The fifth and last step is to check that the new control satisfies (3.18).

15 *Proof of Proposition 3.5. Step 1: estimate of $\mathcal{E}(y(T; \mathbf{x}, w, b))$.* If $\mathcal{E}(y(T; \mathbf{x}, w, b)) = 0$, then
16 it suffices to consider $(\hat{w}, \hat{b}) = (w, b)$, so we suppose from now on that $\mathcal{E}(y(T; \mathbf{x}, w, b)) > 0$. Let
17 $\tilde{\varepsilon}$ the value in Hypothesis 1 and let $T_0 = \frac{3 \int_0^1 |(w_*(t), b_*(t))|^2}{\tilde{\varepsilon}}$, for (w_*, b_*) given in Hypothesis 1.
18 Then, from (3.6) and (3.16) we obtain for $T \geq T_0$ and all (w, b) satisfying (3.16) that:

$$\mathcal{E}(y(T; \mathbf{x}, w, b)) \in (0, \tilde{\varepsilon}/2). \tag{3.20}$$

Step 2: approximating the control. Clearly, $C^1([0, T]; \mathcal{U} \setminus \{0\})$ is dense in $L^2(0, T; \mathcal{U})$. Moreover:

$$(w, b) \mapsto \mathcal{E}(y(T; \mathbf{x}, w, b)),$$

1 is continuous from $L^2(0, T; \mathcal{U})$ to \mathbb{R} . Thus, there is $(\tilde{w}, \tilde{b}) \in C^1([0, T]; \mathcal{U} \setminus \{0\})$ such that:

$$\|(\tilde{w}, \tilde{b})\|_{L^2(0, T; \mathcal{U})} \leq \|(w, b)\|_{L^2(0, T; \mathcal{U})}, \quad (3.21)$$

2 and:

$$\mathcal{E}(y(T; \mathbf{x}, \tilde{w}, \tilde{b})) \leq 2\mathcal{E}(y(T; \mathbf{x}, w, b)). \quad (3.22)$$

3 Moreover, by Lemma 3.4 we can suppose that $t \mapsto |(\tilde{w}(t), \tilde{b}(t))|$ is constant. In addition, from
4 (3.6), (3.16) and (3.21) it follows that:

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} \leq \frac{C}{\sqrt{T}}. \quad (3.23)$$

Step 3: taking the error to 0. From Hypothesis 1, (3.20) and (3.22) we obtain a control $(\bar{w}, \bar{b}) \in L^\infty(0, T; \mathcal{U})$ that takes the solution from $y(T; \mathbf{x}, \tilde{w}, \tilde{b})$ to a state $\tilde{\mathbf{x}}$ such that $\mathcal{E}(\tilde{\mathbf{x}}) = 0$. Moreover:

$$\|(\bar{w}, \bar{b})\|_{L^\infty(0, 1; \mathcal{U})} \leq C\mathcal{E}(y(T; \mathbf{x}, \tilde{w}, \tilde{b})) \leq C\mathcal{E}(y(T; \mathbf{x}, w, b)).$$

5 Consequently, by Remark 3.2, for some:

$$\tau \leq C \frac{\mathcal{E}(y(T; \mathbf{x}, \tilde{w}, \tilde{b}))}{\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})}}, \quad (3.24)$$

6 the control (\tilde{w}, \tilde{b}) can be prolonged to $[0, T + \tau]$ so that both:

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T + \tau; \mathcal{U})} = \|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})}, \quad (3.25)$$

7 and (3.19) are satisfied.

8 **Step 4: taking the trajectory to $[0, T]$.** We consider:

$$\begin{aligned} \hat{w}(t) &:= \frac{T + \tau}{T} \tilde{w} \left(\frac{T + \tau}{T} t \right), \\ \hat{b}(t) &:= \frac{T + \tau}{T} \tilde{b} \left(\frac{T + \tau}{T} t \right). \end{aligned} \quad (3.26)$$

9 Then, (3.17) is true. In fact, the equation (3.3) with $\lambda = \frac{T + \tau}{T}$ implies:

$$y(T; \mathbf{x}, \hat{w}, \hat{b}) = y(T + \tau; \mathbf{x}, \tilde{w}, \tilde{b}).$$

Step 5: efficiency of the new control. First, we realize that:

$$J_T(w, b) - J_T(\hat{w}, \hat{b}) = \mathcal{E}(y(T; \mathbf{x}, w, b)) + \int_0^T |(w(t), b(t))|^2 dt - \left(\frac{T+\tau}{T}\right)^2 \int_0^T \left| \left(\tilde{w}\left(\frac{T+\tau}{T}t\right), \tilde{b}\left(\frac{T+\tau}{T}t\right) \right) \right|^2 dt. \quad (3.27)$$

1 Considering that $t \mapsto (\tilde{w}(t), \tilde{b}(t))$ is constant in $[0, T]$, and that (3.21) and (3.25) are satisfied
2 we deduce that:

$$\int_0^T |(w(t), b(t))|^2 dt - \int_0^T \left| \left(\tilde{w}\left(\frac{T+\tau}{T}t\right), \tilde{b}\left(\frac{T+\tau}{T}t\right) \right) \right|^2 dt \geq 0. \quad (3.28)$$

Consequently, we obtain from (3.20), (3.23)-(3.25) and (3.27)-(3.28) that:

$$\begin{aligned} J_T(w, b) - J_T(\hat{w}, \hat{b}) &\geq \mathcal{E}(y_T(T; \mathbf{x}, w, b), \mathbf{x}) \\ &\quad - \left(\frac{2\tau}{T} + \frac{\tau^2}{T^2}\right) \int_0^T \left| \left(\tilde{w}\left(\frac{T+\tau}{T}t\right), \tilde{b}\left(\frac{T+\tau}{T}t\right) \right) \right|^2 dt \\ &\geq \left(1 - C\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} - C\tilde{\varepsilon}T^{-1}\right) \mathcal{E}(y(T; \mathbf{x}, w, b)) \\ &\geq (1 - CT^{-1/2}) \mathcal{E}(y(T; \mathbf{x}, w, b)), \end{aligned}$$

3 which implies (3.18) if $T > T_0$ for T_0 large enough depending only on σ , \mathcal{E} and \mathbf{x} . □

4 Now we may conclude the proof of Theorem 2.5 by a proof by contradiction:

Conclusion of the proof of Theorem 2.5. Let $\varepsilon > 0$. It suffices to consider $\delta = \varepsilon/3$. If (w, b) are such that $J_T(w, b) \leq \inf J_T + \varepsilon/3$, then $\mathcal{E}(y(T; \mathbf{x}, w, b)) < \varepsilon$. Otherwise, by Proposition 3.5 there are (\hat{w}, \hat{b}) such that:

$$J_T(\hat{w}, \hat{b}) \leq J_T(w, b) - \frac{\varepsilon}{2} \leq \inf J_T - \frac{\varepsilon}{6},$$

5 which is absurd. Similarly, if (w_T, b_T) is a minimizer of J_T and (2.4) is not satisfied, then the
6 control (\hat{w}, \hat{b}) of Proposition 3.5 satisfies $J_T(\hat{w}, \hat{b}) < J_T(w_T, b_T)$, contradicting the definition of
7 minimizer. □

8 **3.3 Additional properties of the optimal control**

As a consequence of Remark 3.2, we can easily prove that, assuming Hypothesis 1, for a sufficiently large time the optimal controls are of the form:

$$\left(\frac{1}{T} w_* \left(\frac{t}{T} \right), \frac{1}{T} b_* \left(\frac{t}{T} \right) \right),$$

1 for (w_*, b_*) the minimizers of the functional:

$$(w, b) \mapsto \int_0^T |(w(t), b(t))|^2 dt,$$

2 considered in the domain:

$$\{(w, b) : \mathcal{E}(y(1; \mathbf{x}, w, b)) = 0\}.$$

3 In addition, we can prove that such minimizers belong to $L^\infty(0, T)$ and satisfy that $t \mapsto$
 4 $|(w(t), b(t))|$ is constant, which follows from:

Proposition 3.6 (A more efficient control). *Let (w, b) a control in $L^2(0, T)$ such that $t \mapsto$
 5 $|(w(t), b(t))|$ is not constant. Then, there is a control (\tilde{w}, \tilde{b}) such that:*

$$y(T; \mathbf{x}, \tilde{w}, \tilde{b}) = y(T; \mathbf{x}, w, b),$$

$$\|(\tilde{w}, \tilde{b})\|_{L^2(0, T; \mathcal{U})} < \|(w, b)\|_{L^2(0, T; \mathcal{U})},$$

6 and, if $(w, b) \in L^\infty(0, T; \mathcal{U})$,

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} \leq 2\|(w, b)\|_{L^\infty(0, T; \mathcal{U})}. \quad (3.29)$$

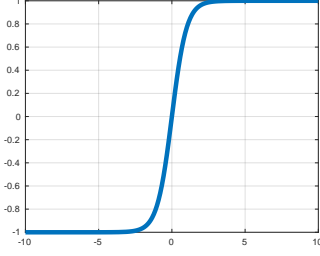
7 Here, y is given by (1.2), where we consider the solutions of (1.1).

8 The proof of Proposition 3.6 is based on classical results from Measure Theory and is postponed
 9 to Appendix C. Note that, opposed to the L^1 case proved in [8], where the norm is constant up
 10 to some time $T^* \leq T$ and then is null, in our model the norm is constant in the whole interval
 11 $(0, T)$. Proposition 3.6, compared to Lemma 3.4, has the advantage of having a less restrictive
 12 hypothesis. However, it has the disadvantage that we do not obtain neither a contraction for
 13 the L^∞ norm (see Remark C.2) nor a control with constant norm, which is needed for proving
 14 Proposition 3.5.

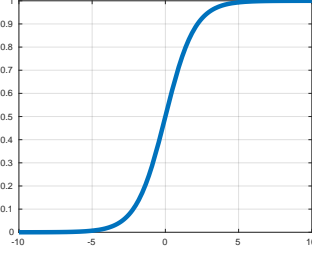
15 4 Further comments and open problems

- **Analogous results for neural ODE whose dynamics are described by (1.3).** Clearly Lemmas 3.1, 3.3, 3.4, and Propositions 3.5 and 3.6 can be proved for system (1.3) with σ satisfying (1.5) as in Section 3. The key lemma is Lemma 3.1, since the other results use the homogeneity of the system via Lemma 3.1. The analogous of Lemma 3.1 can be proved by replacing (3.4) by:

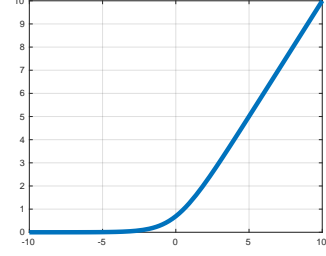
$$\begin{aligned} \frac{d}{dt}(y(\Gamma(t); x^i, w, b, r)) &= \phi(\Gamma(t))r(\Gamma(t))\dot{y}(\Gamma(t); x^i, w, b) \\ &= \phi(\Gamma(t))r(\Gamma(t))\sigma\left(w(\Gamma(t))y(\Gamma(t); x^i, w, b) + b(\Gamma(t))\right) \\ &= r(\Gamma(t))\sigma\left(\phi(\Gamma(t))w(\Gamma(t))y(\Gamma(t); x^i, w, b) + \phi(\Gamma(t))b(\Gamma(t))\right). \end{aligned}$$



(a) Hyperbolic tangent.



(b) Sigmoid.



(c) Softplus.

Figure 2: Some activation functions for $d = 1$.

1 The last equality follows from (1.5). Finally, Theorem 2.7 and the analogous of Proposi-
 2 tion 3.5 imply that for all $\delta > 0$ there is a control (r, w, b) such that $\tilde{J}_T(r, w, b) < \inf \tilde{J}_T - \delta$
 3 and $\mathcal{E}(y(T; \mathbf{x}, r, w, b)) = 0$.

- **Optimal control for non-homogenous activation functions.** It remains an open problem to determine if similar results to Theorem 2.7 hold for non-homogeneous activation functions satisfying $\sigma(0) = 0$ such as the hyperbolic tangent:

$$\sigma(x) = (\tanh(x_1), \dots, \tanh(x_d)),$$

see [9]. We may wonder whether similar results hold with more general activation functions if we replace X (see (1.4)) by the unitary matrices or by $\mathbb{R}^{d \times d}$ (of course, the cost of r must also be included in the risk minimization functional). This would include, for instance, sigmoid:

$$\sigma(x) = ((1 + e^{-x_1})^{-1}, \dots, (1 + e^{-x_d})^{-1}),$$

see [23]; softplus:

$$\sigma(x) = (\log(1 + e^{x_1}), \dots, \log(1 + e^{x_d})),$$

4 see [11] (see Figure 2 for their graph in one dimension), and others like logistic and cross-
 5 entropy functions. The main difficulty is that the analogue of Lemma 3.1 cease to be
 6 true, so another tool is needed to prove the main result, probably a local inverse theorem
 7 result.

- 8 • **Optimal control with the H^1 norm.** It is a relevant problem to determine if similar
 9 results to Theorems 2.5 and 2.7 hold for any other Lebesgue or Sobolev penalty. In par-
 10 ticular, an interesting scenario is to replace both in J_T and \tilde{J}_T the terms $\|(w, b)\|_{L^2(0, T; \mathcal{U})}^2$
 11 by $\|(w, b)\|_{H^1(0, T; \mathcal{U})}^2$ and adding the restriction that the component of r can only change
 12 signs if $(w, b) = 0$ or to measure the H^1 norm of r if the space X is connected. The
 13 interest of this is double: thinking in potential applications it makes sense to also try
 14 to bound the variations in the time variable, which can be obtained by minimizing the
 15 time derivative. Moreover, if we consider the H^1 -norm we can prove as in Proposition

2.2 that \tilde{J}_T admits a minimizer. The main difficulties when studying these norms are that Lemmas 3.4 and D.1 and Proposition 3.6 may not be proved as easily (if they are true) because we need to keep track of the time derivative and because we cannot define the control on $[T, T + \tau]$ independently to the controls on $[0, T]$ due to the necessity of bounding the time derivative.

- **Optimal control with the BV norm.** It is also a relevant problem to determine if results similar to Theorems 2.5 and 2.7 hold when we consider a BV penalty. The existence of minimizers, as shown in [6, Section 4], follows from the fact that any minimizing sequence in BV converges strongly in L^1 . However, the main difficulty when studying these penalties, as before, is that Lemmas 3.4 and D.1 and Proposition 3.6 may not be proved as easily (if they are true) because we need to keep track of the jumps and because we cannot define the control on $[T, T + \tau]$ independently to the controls on $[0, T]$ due to the time derivative.

A Proof of Proposition 2.2

Let us consider (w_n, b_n) a minimizing sequence; that is, a sequence such that:

$$\lim_{n \rightarrow \infty} J_T(w_n, b_n) = \inf J_T.$$

Since $\mathcal{E} \geq 0$, (w_n, b_n) is bounded in $L^2(0, T; \mathcal{U})$. From now on:

$$y_n(t) := y_n(t; \mathbf{x}, w_n, b_n) = (y_n^1(t; x^1, w_n, b_n), \dots, y_n^N(t; x^N, w_n, b_n)).$$

We recall that each y_n^i is a solution of:

$$\begin{cases} \dot{y}_n^i(t) = w_n(t)\sigma(y_n^i(t)) + b_n(t), \\ y_n^i(0) = x^i, \end{cases} \quad (\text{A.1})$$

Multiplying (A.1) by $2y_n^i$ and adding up, we obtain that:

$$\frac{d}{dt} (|y_n|^2) = \sum_{i=1}^N 2(w_n(t)\sigma(y_n^i(t))) \cdot y_n^i(t) + 2b_n(t) \cdot y_n^i(t)$$

Thus, with Cauchy-Schwarz inequality and using that σ is Lipschitz, we obtain that:

$$\frac{d}{dt} (|y_n(t)|^2) \leq C(|w_n(t)| + 1)|y_n(t)|^2 + |b_n(t)|^2,$$

which implies:

$$|y_n(t)|^2 \leq \int_0^t C(|w_n(s)| + 1)|y_n(s)|^2 ds + \int_0^t |b_n(s)|^2 ds + |\mathbf{x}|^2.$$

1 Consequently, from Grönwall inequality in its integral form (see [41, Appendix E]), and the
 2 continuous inclusion $L^2(0, T) \subset L^1(0, T)$ we obtain that:

$$|y_n(t)|^2 \leq \left(|\mathbf{x}|^2 + \int_0^t |b_n(s)|^2 ds \right) \exp \left(\int_0^t C(|w_n(s)| + 1) ds \right).$$

3 Thus, the sequence (y_n) is uniformly bounded in $C^0([0, T]; (\mathbb{R}^d)^N)$. Consequently, considering
 4 (1.1) and that σ is Lipschitz, we obtain that (y_n) is uniformly bounded in $H^1(0, T; (\mathbb{R}^d)^N)$. Thus,
 5 considering the compact inclusion $H^1(0, T; (\mathbb{R}^d)^N) \subset C^0([0, T]; (\mathbb{R}^d)^N)$, there are $(w^*, u^*) \in$
 6 $L^2(0, T; \mathcal{U})$ and $y^* \in H^1(0, T; (\mathbb{R}^d)^N)$ and subsequences $(w_{n_k}, b_{n_k}) \rightharpoonup (w^*, b^*)$ in $L^2(0, T; \mathcal{U})$,
 7 $y_{n_k} \rightarrow y^*$ in $C^0([0, T]; (\mathbb{R}^d)^N)$ and $y_{n_k} \rightharpoonup y^*$ in $H^1(0, T; (\mathbb{R}^d)^N)$.

8 Let us now show that:

$$y^*(t) = y(t; \mathbf{x}, w^*, b^*). \quad (\text{A.2})$$

Indeed, for all $i = 1, \dots, N$, by taking weak limit in both sides of:

$$\dot{y}_{n_k}^i = w_{n_k} \sigma(y_{n_k}^i) + b_{n_k},$$

9 we obtain that:

$$(\dot{y}^*)^i = w^* \sigma((y^*)^i) + b^*. \quad (\text{A.3})$$

Indeed, in $w_{n_k} \sigma(y_{n_k}^i)$ we have the product of a weak limit times a strong limit in the continuous
 space. Thus, with (A.3) and:

$$(y^*)^i(0) = \lim_{k \rightarrow \infty} y_{n_k}^i(0) = x^i,$$

10 we get (A.2).

11 Finally, consider that the norm of the control is weakly-lower semi-continuous:

$$\begin{aligned} J_T(w^*, b^*) &= \mathcal{E}(y^*(T)) + \int_0^T |(w^*(t), b^*(t))|^2 dt \\ &= \lim_{k \rightarrow \infty} \mathcal{E}(y_{n_k}(T)) + \int_0^T |(w^*(t), b^*(t))|^2 dt \\ &\leq \liminf_{k \rightarrow \infty} \mathcal{E}(y_{n_k}(T)) + \liminf_{k \rightarrow \infty} \int_0^T |(w_{n_k}(t), b_{n_k}(t))|^2 dt \\ &\leq \liminf_{k \rightarrow \infty} J_T(w_{n_k}, b_{n_k}) \\ &= \inf J_T, \end{aligned}$$

12 so J_T attains its minimum on (w^*, b^*) \square .

13 B A pathological case

14 In this section we prove that without the second item of Hypothesis 1 the error may not be
 15 taken exactly to 0 if the ratio between the cost of correcting the error and the error explodes

1 as the error vanishes. We present an example for the sake of simplicity, though the proof can
 2 be replicated whenever the gradient of the error is null on all the points where the error is null,
 3 which is the key impediment for taking the error exactly to zero.

4 **Proposition B.1** (Necessity of local controllability). *Let us consider $d = 1$, $\mathbf{x} = x_1 = 1$,
 5 $\mathcal{E}(x) = x^2$, $\sigma(s) = s$ and J_T given by (2.1). Then, $y_T(T) > 0$ for all $T > 0$. Here y_T is the
 6 solution of (1.1) with the optimal control.*

7 *Proof of Proposition B.1.* Let (w_T, b_T) be a minimizer of J_T . Clearly $w_T, b_T \leq 0$. Let us prove
 8 by contradiction that $y_T(T) > 0$. For that, we suppose that $y_T(T) = 0$. By Proposition 3.6,
 9 $t \mapsto |(w_T(t), b_T(t))|$ is a constant function equal to some constant \mathfrak{c} . In particular, for $\delta > 0$
 10 small enough the following inequality is satisfied:

$$(y_T(T - \delta))^2 - \int_{T-\delta}^T |(w_T(t), b_T(t))|^2 = (y_T(T - \delta))^2 - \mathfrak{c}^2\delta \leq (C\delta)^2 - \mathfrak{c}^2\delta < 0. \quad (\text{B.1})$$

The estimate $|y_T(T - \delta)| \leq C\delta$ follows from the formula:

$$y_T(T - \delta) = - \int_{T-\delta}^T b_T(s) \exp\left(- \int_{T-\delta}^s w_T(z) dz\right) ds,$$

11 which follows from $y_T(T) = 0$. Consequently, we obtain from (B.1) that:

$$J_T(w_T 1_{(0, T-\delta)}, b_T 1_{(0, T-\delta)}) - J_T(w_T, b_T) < 0, \quad (\text{B.2})$$

12 which contradicts that (w_T, b_T) is a minimizer of J_T . □

13 *Remark B.2* (On the first item of Hypothesis 1). It is trivial that the first item of Hypothesis
 14 1 is satisfied by the activation and error function introduced in Propositions B.1.

15 C Proof of Proposition 3.6

16 In this section we prove Proposition 3.6. Here μ denotes the Lebesgue measure. In order to
 17 prove Proposition 3.6 we need the following classical result of measure theory, whose proof can
 18 be found in [45, Theorem 3.25]:

19 **Lemma C.1** (Comparison between sets of positive measure and open sets). *Let $S \subset [0, T]$ be
 20 a measurable set such that $\mu(S) > 0$. Then, for all $\varepsilon > 0$ there is an open set $\mathcal{O}^\varepsilon = \bigcup_{i=1}^{n^\varepsilon} (a_i^\varepsilon, b_i^\varepsilon)$
 21 such that $\mu(\mathcal{O}^\varepsilon \Delta S) < \varepsilon$.*

1 *Proof of Proposition 3.6.* Since $|(w, b)|$ is not constant, there are some sets S_1 and S_2 and some
 2 constants $C_1, C_2 > 0$ such that $C_1 < C_2$, $|(w, b)| < C_1$ on S_1 , $|(w, b)| > C_2$ on S_2 and:

$$\inf\{|x^2 - x^1| : x^1 \in S_1, x^2 \in S_2\} > 0. \quad (\text{C.1})$$

3 From Lemma C.1 we get that for $\varepsilon > 0$ small enough there are two sets $\mathcal{O}_1^\varepsilon = \bigcup_{i=1}^{n_1^\varepsilon} (a_{1,i}^\varepsilon, b_{1,i}^\varepsilon)$
 4 and $\mathcal{O}_2^\varepsilon = \bigcup_{i=1}^{n_2^\varepsilon} (a_{2,i}^\varepsilon, b_{2,i}^\varepsilon)$ satisfying:

$$\mu(\mathcal{O}_1^\varepsilon \setminus S_1) < \varepsilon, \quad \mu(\mathcal{O}_2^\varepsilon \setminus S_2) < \varepsilon, \quad (\text{C.2})$$

5 and:

$$\mu(\mathcal{O}_1^\varepsilon) = \mu(\mathcal{O}_2^\varepsilon) = \frac{\min\{\mu(S_1), \mu(S_2)\}}{2}. \quad (\text{C.3})$$

6 If ε is small enough, because of (C.1) we may also assume that:

$$\mathcal{O}_1^\varepsilon \cap \mathcal{O}_2^\varepsilon = \emptyset. \quad (\text{C.4})$$

7 Let us consider the auxiliary function:

$$\phi_\gamma(t) = \begin{cases} 1 & t \in [0, T] \setminus (\mathcal{O}_1^\varepsilon \cap \mathcal{O}_2^\varepsilon), \\ 1 + \gamma & t \in \mathcal{O}_1^\varepsilon, \\ \frac{1+\gamma}{1+2\gamma} & t \in \mathcal{O}_2^\varepsilon, \\ 0 & t \geq T, \end{cases} \quad (\text{C.5})$$

8 for $\gamma > 0$ to be fixed later, and Γ_γ given by:

$$\begin{cases} \dot{\Gamma}_\gamma(s) = \phi_\gamma(\Gamma_\gamma(s)), & \forall s \geq 0, \\ \Gamma_\gamma(0) = 0. \end{cases} \quad (\text{C.6})$$

9 We remark that:

$$\Gamma_\gamma(T) = T. \quad (\text{C.7})$$

Indeed, it can be proved that if $\Gamma_\gamma(T_*) = a$ and $\phi_\gamma(t) = c$ on $[a, b]$, then $\Gamma_\gamma(T_* + \frac{b-a}{c}) = b$.
 Hence:

$$\Gamma_\gamma \left(\mu([0, T] \setminus (\mathcal{O}_1^\varepsilon \cap \mathcal{O}_2^\varepsilon)) + \frac{1}{1+\gamma} \mu(\mathcal{O}_1^\varepsilon) + \frac{1+2\gamma}{1+\gamma} \mu(\mathcal{O}_2^\varepsilon) \right) = T,$$

10 which considering (C.3), (C.4) and (C.5) implies (C.7).

11 Consequently, the following controls satisfy the conclusions of Proposition 3.6:

$$(\tilde{w}, \tilde{b}) = \phi_\gamma(\Gamma_\gamma(t))(w(\Gamma_\gamma(t)), b(\Gamma_\gamma(t))). \quad (\text{C.8})$$

Indeed, from (3.1) and (C.1) it holds that:

$$\begin{aligned} y(T; \mathbf{x}, \phi_\gamma(s)w(\Gamma_\gamma(s)), \phi_\gamma(s)b(\Gamma_\gamma(s))) &= y(\Gamma_\gamma(T); \mathbf{x}, w, b) \\ &= y(T; \mathbf{x}, w, b). \end{aligned}$$

In addition, if γ and ε are small enough:

$$\begin{aligned} &\int_0^T |(w(t), b(t))|^2 dt - \int_0^T |\phi_\gamma(\Gamma_\gamma(t))(w(\Gamma_\gamma(t)), b(\Gamma_\gamma(t)))|^2 dt \\ &= \int_{\mathcal{O}_1^\varepsilon \cup \mathcal{O}_2^\varepsilon} |(w(t), b(t))|^2 dt \\ &\quad - \int_{\Gamma_\gamma^{-1}(\mathcal{O}_1^\varepsilon) \cup \Gamma_\gamma^{-1}(\mathcal{O}_2^\varepsilon)} \phi_\gamma^2(\Gamma_\gamma(t)) |(w(\Gamma_\gamma(t)), b(\Gamma_\gamma(t)))|^2 dt \\ &= -\gamma \int_{\mathcal{O}_1^\varepsilon} |(w(t), b(t))|^2 dt + \frac{\gamma}{1+2\gamma} \int_{\mathcal{O}_2^\varepsilon} |(w(t), b(t))|^2 dt \\ &\geq \frac{\gamma}{1+2\gamma} C_2 \left(\frac{\min\{\mu(S_1), \mu(S_2)\}}{2} - \varepsilon \right) \\ &\quad - \gamma C_1 \frac{\min\{\mu(S_1), \mu(S_2)\}}{2} - \|(w, b)\|_{L^2(\mathcal{O}_1^\varepsilon \setminus S_1)}^2 > 0. \end{aligned}$$

- 1 The second equality follows from the change of variable $s = \Gamma_\gamma(t)$, the first inequality from the
2 definitions of S_1 , S_2 , $\mathcal{O}_1^\varepsilon$, $\mathcal{O}_2^\varepsilon$ and (C.3), and the last inequality from $C_2 > C_1$, (C.2), being γ
3 and ε small enough, and the well known identity:

$$\lim_{c \rightarrow 0} \sup_{\mu(A)=c} \|g\|_{L^2(A, dx)} = 0, \quad \forall g \in L^2(0, T).$$

- 4 Finally, if $(w, b) \in L^\infty(0, T; \mathcal{U})$ the estimate (3.29) follows from (C.5) and (C.8) by taking
5 $\gamma \leq 1$. □

- 6 *Remark C.2* (Sharpness of the estimate (3.29)). The construction provided in the previous
7 proof may not ensure us that:

$$\|(\tilde{w}, \tilde{b})\|_{L^\infty(0, T; \mathcal{U})} \leq \|(w, b)\|_{L^\infty(0, T; \mathcal{U})};$$

- 8 for instance if $|(w, b)| = 1_\Omega$, for $\Omega \subset [0, T]$ a set such that $\mu(\Omega) \in (0, T)$ and which contains an
9 open neighbourhood of every rational number in $[0, T]$. However, we can replace in the estimate
10 (3.29) the constant 2 by any constant strictly greater than 1.

11 D Local ensemble controllability

- 12 In this section we prove the following result:

1 **Lemma D.1** (Local ensemble controllability result). *Let σ be the activation function defined*
 2 *by (2.5) and \mathcal{E} defined in Example 2.6. Then σ and \mathcal{E} satisfy the second item of Hypothesis 2.*

3 The main contribution with respect to [27, Theorem 3.1] is that we keep track of the cost
 4 and continuity of the control. The controls that we construct are different to those in [27], in
 5 which w and b have a single non-zero component at any time, since we do not search for a
 6 sparse property, but to obtain the continuity of the controls with respect to the initial data.
 7 Here, the constants C_i are positive constants sufficiently large which depend on the target set
 8 \mathbf{z} . Finally, throughout the proof, we consider the space \mathbb{R}^d endowed with the euclidean norm
 9 and $\mathbb{R}^{d \times d}$ endowed with the norm:

$$|A|_{\mathbb{R}^{d \times d}} = \sup_{|u|=1} |Au|.$$

10 This can be done because in finite dimensional spaces all norms are equivalent.

11 *Proof of Lemma D.1.* In order to simplify the notation we prove Lemma D.1 for the case $d = 2$,
 12 though the proof is analogous for any $d \geq 2$. We prove Lemma D.1 by induction on N .

Step 1: the base case. Let us begin with the case $N = 1$. We may take $x^1 = (x_1^1, x_2^1)$ to
 $z^1 = (z_1^1, z_2^1)$ with a force proportional to $|z^1 - x^1|$ by applying the controls:

$$r = \begin{pmatrix} \text{sign}(z_1^1 - x_1^1) & 0 \\ 0 & \text{sign}(z_2^1 - x_2^1) \end{pmatrix}, \quad w = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \quad \text{and} \quad b = \begin{pmatrix} |z_1^1 - x_1^1| \\ |z_2^1 - x_2^1| \end{pmatrix}.$$

13 **Step 2: the inductive case. Step 2.1: rearranging the points.** We may suppose by
 14 rearranging the indexes that $|z^N| = \max_{i=1, \dots, N} |z^i|$. For the rest of the proof we define:

$$\delta := \min \left\{ |z^N| - \max_{i=1, \dots, N-1} \frac{|z^i \cdot z^N|}{|z^N|}, 1 \right\}. \quad (\text{D.1})$$

15 Then, $\delta > 0$ since, for $i = 1, \dots, N - 1$, either $|z^i| < |z^N|$ or $|z^i| = |z^N|$ but $z^i \neq z^N$, so
 16 $z^i \cdot z^N < |z^N|^2$.

17 **Step 2.2: controlling (x^1, \dots, x^{N-1}) in $[0, 1/2]$.** By the induction hypothesis (we may
 18 apply it to $T = 1/2$ instead of $T = 1$ by linearity), we know that for $\hat{\varepsilon}$ small enough, if
 19 $\sum_{i=1}^{N-1} |z^i - x^i| < \hat{\varepsilon}$ there are some controls (r, w, b) defined in $[0, \frac{1}{2}]$ and a constant $C_1 > 0$
 20 satisfying:

$$\|w\|_{L^\infty(0, 1/2; \mathbb{R}^{2 \times 2})} + \|b\|_{L^\infty(0, 1/2; \mathbb{R}^2)} < C_1 \sum_{i=1}^{N-1} |z^i - x^i|, \quad (\text{D.2})$$

21 and such that:

$$y(1/2; (x^1, \dots, x^{N-1}), r, w, b) = (z^1, \dots, z^{N-1}). \quad (\text{D.3})$$

We fix:

$$\tilde{\varepsilon} = \min \left\{ \hat{\varepsilon}, \frac{\delta}{2C_1(|z^N| + 1)}, \frac{\delta}{4} \right\}.$$

1 If:

$$\sum_{i=1}^N |z^i - x^i| < \tilde{\varepsilon}, \quad (\text{D.4})$$

2 then:

$$|y(t; x^N, r, w, b) - z^N| < \frac{\delta}{2} \quad \forall t \in \left[0, \frac{1}{2}\right]. \quad (\text{D.5})$$

Indeed,

$$|y(0; x^N, r, w, b) - z^N| = |x^N - z^N| < \frac{\delta}{4},$$

3 by (D.4). Moreover, if $|y(\cdot; x^N, w, b) - z^N| < \frac{\delta}{2}$ on $[0, t]$, for $t \leq 1/2$, then, considering (2.6),

4 (D.2) and that $\tilde{\varepsilon} < \frac{\delta}{C(|z^N|+1)}$:

$$\int_0^t |\sigma(w(s)y(s; x^N, r, w, b) + b(s))| ds \leq \|w\|_{L^\infty(0,t;\mathbb{R}^{2 \times 2})} \left(|z^N| + \frac{\delta}{2}\right) + \|b\|_{L^\infty(0,t;\mathbb{R}^2)} < \frac{\delta}{4}.$$

5 In a similar way, we can prove that for $C_2 > 0$ large enough:

$$|y(1/2; x^N, r, w, b) - z^N| \leq C_2 \sum_{i=1}^N |z^i - x^i|. \quad (\text{D.6})$$

6 Indeed, by (2.6) and (D.2):

$$\begin{aligned} \int_0^{1/2} |\sigma(w(t)y(t; x^N, r, w, b) + b(t))| dt &\leq \|w\|_{L^\infty(0,1/2;\mathbb{R}^{2 \times 2})} (|z^N| + \delta) + \|b\|_{L^\infty(0,1/2;\mathbb{R}^2)} \\ &\leq C_1(|z^N| + 1) \sum_{i=1}^{N-1} |z^i - x^i|. \end{aligned}$$

7 **Step 2.3: controlling $y(1/2; x^N, r, w, b)_1$ in $[1/2, 3/4]$.** We seek to obtain that:

$$y_1(3/4; x^N, r, w, b) = z_1^N. \quad (\text{D.7})$$

If $y(1/2; x^N, w, b)_1 = z_1^N$, it suffices to consider:

$$r = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad w = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}, \quad \text{and} \quad b = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

8 in $t \in [1/2, 3/4]$, so we may restrict to the case $y_1(1/2; x^N, r, w, b) \neq z_1^N$. To obtain (D.7) we

9 consider the controls:

- 1 • $r = \begin{pmatrix} \text{sign}(z_1^N - y_1(1/2; x^N, r, w, b)) & 0 \\ 0 & 1 \end{pmatrix},$
- 2 • $w = \mathbf{c}_1 \sum_{i=1}^N |z^i - x^i| \begin{pmatrix} \frac{z_1^N}{|z^N|} & \frac{z_2^N}{|z^N|} \\ 0 & 0 \end{pmatrix},$
- 3 • $b = \begin{pmatrix} \mathbf{c}_1 \sum_{i=1}^N |z^i - x^i| (-|z^N| + \delta) \\ 0 \end{pmatrix},$

in $[1/2, 3/4]$, for \mathbf{c}_1 to be fixed later on, as under that hypothesis $(w \cdot x)_1 + b_1 \leq 0$. First, we remark that:

$$\sigma(w \cdot x + b) = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$

- 4 for all x such that $x \cdot \frac{z^N}{|z^N|} \leq |z^N| - \delta$. In particular, from (D.1) and (D.3) we derive:

$$y(3/4; (x^1, \dots, x^{N-1}), r, w, b) = y(1/2; (x^1, \dots, x^{N-1}), r, w, b) = (z^1, \dots, z^{N-1}). \quad (\text{D.8})$$

- 5 Moreover, as $\dot{y}_2 = 0$ in $[1/2, 3/4]$, we obtain that:

$$y_2(3/4; x^N, r, w, b) = y_2(1/2; x^N, r, w, b). \quad (\text{D.9})$$

Next, $|y_1(t; x^N, r, w, b) - z_1^N|$ is decreasing on $[1/2, T_*]$, for:

$$T_* := \min\{\inf\{T_* \geq 1/2 : y(T_*; x^N, r, w, b)_1 = z_1^N\}, 3/4\}.$$

In addition, thanks to (D.5) in $[1/2, T_*]$ the following inequality is satisfied:

$$|\dot{y}_1(t; x^N, r, w, b)| = |(w \cdot y(t; x^N, r, w, b))_1 - b_1| \geq \mathbf{c}_1 \sum_{i=1}^N |z^i - x^i| \frac{\delta}{2} \quad \forall t \in \left[\frac{1}{2}, T_*\right].$$

Combining this with (D.6) we obtain that there is $C_4 > 0$ such that $T_* = 3/4$ for some $\mathbf{c}_1 < C_4$, if \mathbf{c}_1 is sufficiently large just with respect to \mathbf{z} (recall that δ is a fixed parameter depending only on \mathbf{z}). In particular, there are controls (r, w, b) defined in $[0, \frac{3}{4}]$ and a constant $C_5 > 0$ such that (D.7), (D.8) and (D.9) hold, and such that:

$$\|w\|_{L^\infty(0, 3/4; \mathbb{R}^{2 \times 2})} + \|b\|_{L^\infty(0, 3/4; \mathbb{R}^2)} < C_5 \sum_{i=1}^N |z^i - x^i|.$$

Step 2.4: controlling $y(\cdot; x^N, r, w, b)$ in $[3/4, 1]$. In a similar way, we can prolong the controls (w, b) in $[3/4, 1]$ so that $y(1; \mathbf{x}, r, w, b) = \mathbf{z}$ and so that there is $C_6 > 0$ such that:

$$\|w\|_{L^\infty(0, 1; \mathbb{R}^{2 \times 2})} + \|b\|_{L^\infty(0, 1; \mathbb{R}^2)} < C_6 \sum_{i=1}^N |z^i - x^i|.$$

- 6 This can be proved as in Step 2.3 by fixing:

$$\bullet r = \begin{pmatrix} 1 & 0 \\ 0 & \text{sign}(z_2^N - y_2(1/2; x^N, r, w, b)) \end{pmatrix},$$

$$\bullet w = \mathfrak{c}_2 \sum_{i=1}^N |z^i - x^i| \begin{pmatrix} 0 & 0 \\ \frac{z_1^N}{|z^N|} & \frac{z_2^N}{|z^N|} \end{pmatrix},$$

$$\bullet b = \begin{pmatrix} 0 \\ \mathfrak{c}_2 \sum_{i=1}^N |z^i - x^i| (-|z^N| + \delta) \end{pmatrix},$$

in $[3/4, 1]$, for some constant $\mathfrak{c}_2 > 0$ with an upper bound depending only on \mathbf{z} . □

Remark D.2. It is an open problem whether we can obtain a result similar to Lemma D.1 for system (1.1). In order to prove it with an inductive approach, the main difficulty is the obtention of (D.8), as it is essential that the activation function is applied to $wy + b$. Indeed, ensuring that $w\sigma(y) + b = 0$ does not seem straightforward.

Acknowledgements

This article has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement NO: 694126-DyCon). It is also supported by the Grant PID2021-126813NB-I00 funded by MCIN/AEI/10.13039/501100011033 and by “ERDF A way of making Europe”, and by the grant IT1615-22 funded the Basque Government. I would like to thank Carlos Esteve, Borjan Geshkovski and Domènec Ruiz i Balet for fruitful discussion. In particular, I would like to thank Carlos Esteve for proposing Proposition B.1. I would also like to thank two anonymous reviewers for their useful remarks.

References

- [1] J. B. Amara and E. Beldi. Simultaneous controllability of two vibrating strings with variable coefficients. *Evol. Equ. Control The.*, 8(4):687–694, 2019.
- [2] T. Breiten and L. Pfeiffer. On the turnpike property and the receding-horizon method for linear-quadratic optimal control problems. *SIAM J. Control Optim.*, 58(2):1077–1102, 2020.
- [3] L. Cesari. *Optimization—theory and applications: problems with ordinary differential equations*, volume 17. Springer Science & Business Media, 2012.

- 1 [4] T. Damm, L. Grüne, M. Stieler, and K. Worthmann. An exponential turnpike theorem for
2 dissipative discrete time optimal control problems. *SIAM J. Control Optim.*, 52(3):1935–
3 1957, 2014.
- 4 [5] R. Dorfman, P. A. Samuelson, and R. M. Solow. *Linear programming and economic*
5 *analysis*. Courier Corporation, 1987.
- 6 [6] C. Esteve, B. Geshkovski, D. Pighin, and E. Zuazua. Large-time asymptotics in deep
7 learning. *arXiv preprint arXiv:2008.02491v2*, 2021.
- 8 [7] C. Esteve, B. Geshkovski, D. Pighin, and E. Zuazua. Turnpike in lipschitz-nonlinear
9 optimal control. *Nonlinearity*, 35(4):1652, 2022.
- 10 [8] C. Esteve-Yagüe and B. Geshkovski. Sparsity in long-time control of neural odes. *Syst.*
11 *Control Lett.*, 172:105452, 2023.
- 12 [9] E. Fathi and B. M. Shoja. Deep neural networks for natural language processing. In
13 *Handbook of statistics*, volume 38, pages 229–316. Elsevier, 2018.
- 14 [10] T. Faulwasser, A.-J. Hempel, and S. Streif. On the turnpike to design of deep neural nets:
15 Explicit depth bounds. *arXiv preprint arXiv:2101.03000*, 2021.
- 16 [11] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In *Proceedings*
17 *of the fourteenth international conference on artificial intelligence and statistics*, pages
18 315–323. JMLR Workshop and Conference Proceedings, 2011.
- 19 [12] L. Gruüne and R. Guglielmi. Turnpike properties and strict dissipativity for discrete time
20 linear quadratic optimal control problems. *SIAM J. Control Optim.*, 56(2):1282–1302,
21 2018.
- 22 [13] M. Gugat and F. M. Hante. On the turnpike phenomenon for optimal boundary control
23 problems with hyperbolic systems. *SIAM J. Control Optim.*, 57(1):264–289, 2019.
- 24 [14] M. Gugat, M. Schuster, and E. Zuazua. The Finite-Time Turnpike Phenomenon for Op-
25 timal Control Problems: Stabilization by Non-smooth Tracking Terms. In G. Sklyar and
26 A. Zuyev, editors, *Stabilization of Distributed Parameter Systems: Design Methods and*
27 *Applications*, pages 17–41, Cham, 2021. Springer International Publishing.
- 28 [15] M. Gugat, R. Trélat, and E. Zuazua. Optimal Neumann control for the 1D wave equation:
29 Finite horizon, infinite horizon, boundary tracking terms and the turnpike property. *Syst.*
30 *Control Lett.*, 90:61–70, 2016.
- 31 [16] E. Haber and L. Ruthotto. Stable architectures for deep neural networks. *Inverse Probl.*,
32 34(1):014004, 2017.

- 1 [17] R. F. Hartl, S. P. Sethi, and R. G. Vickson. A survey of the maximum principles for
2 optimal control problems with state constraints. *SIAM Rev.*, 37(2):181–218, 1995.
- 3 [18] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-
4 level performance on imagenet classification. In *Proceedings of the IEEE international*
5 *conference on computer vision*, pages 1026–1034, 2015.
- 6 [19] J. Komić. *International Encyclopedia of Statistical Science. Harmonic Mean.*, pages 622–
7 624. Springer, Heidelberg, 2011.
- 8 [20] J.-L. Lions. Contrôlabilité exacte, perturbations et stabilisation de systèmes distribués.
9 Tome 1. *RMA*, 8, 1988.
- 10 [21] J. Lohéac and E. Zuazua. From averaged to simultaneous controllability. In *Annales de la*
11 *Faculté des sciences de Toulouse: Mathématiques*, volume 25, pages 785–828, 2016.
- 12 [22] L. W. McKenzie. Turnpike theory. *Econometrica: Journal of the Econometric Society*,
13 pages 841–865, 1976.
- 14 [23] J. Mira and F. Sandoval. *From Natural to Artificial Neural Computation: International*
15 *Workshop on Artificial Neural Networks, Malaga-Torremolinos, Spain, June 7-9, 1995:*
16 *Proceedings*, volume 930. Springer Science & Business Media, 1995.
- 17 [24] V. Nair and G. E. Hinton. Rectified linear units improve restricted boltzmann machines. In
18 *27th International Conference on International Conference on Machine Learning, ICML*
19 *10*, pages 807–814, 2010.
- 20 [25] A. Porretta and E. Zuazua. Long time versus steady state optimal control. *SIAM J.*
21 *Control Optim.*, 51(6):4242–4273, 2013.
- 22 [26] D. Ruiz-Balet, E. Affili, and E. Zuazua. Interpolation and approximation via momentum
23 resnets and neural odes. *Syst. Control Lett.*, 162:105182, 2022.
- 24 [27] D. Ruiz-Balet and E. Zuazua. Neural ode control for classification, approximation, and
25 transport. *SIAM Review*, 65(3):735–773, 2023.
- 26 [28] D. L. Russell. The Dirichlet–Neumann boundary control problem associated with
27 Maxwell’s equations in a cylindrical region. *SIAM J. Control Optim.*, 24(2):199–229, 1986.
- 28 [29] N. Sakamoto and M. Nagahara. The turnpike property in the maximum hands-off control.
29 In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 2350–2355. IEEE,
30 2020.
- 31 [30] N. Sakamoto, D. Pighin, and E. Zuazua. The turnpike property in nonlinear optimal
32 control—a geometric approach. In *2019 IEEE 58th Conference on Decision and Control*
33 *(CDC)*, pages 2422–2427. IEEE, 2019.

- 1 [31] N. Sakamoto and E. Zuazua. The turnpike property in nonlinear optimal control—A
2 geometric approach. *Automatica*, 134:109939, 2021.
- 3 [32] M. Schönlein. Computation of open-loop inputs for uniformly ensemble controllable sys-
4 tems. *Math. Control Relat. F.*, 12:813–829, 2022.
- 5 [33] M. Schönlein and U. Helmke. Controllability of ensembles of linear dynamical systems.
6 *Math. Comput. Simulat.*, 125:3–14, 2016.
- 7 [34] E. Sontag and H. Sussmann. Complete controllability of continuous-time recurrent neural
8 networks. *Syst. Control Lett.*, 30(4):177–183, 1997.
- 9 [35] E. D. Sontag. Neural nets as systems models and controllers. In *Proc. Seventh Yale*
10 *Workshop on Adaptive and Learning Systems*, pages 73–79, 1992.
- 11 [36] E. Trélat. *Contrôle optimal: théorie & applications*. Vuibert Paris, 2005.
- 12 [37] E. Trélat. Linear turnpike theorem. *Math. Control Signal*, 35:685–739, 2023.
- 13 [38] E. Trélat, C. Zhang, and E. Zuazua. Steady-state and periodic exponential turnpike prop-
14 erty for optimal control problems in Hilbert spaces. *SIAM J. Control Optim.*, 56(2):1222–
15 1252, 2018.
- 16 [39] E. Trélat and E. Zuazua. The turnpike property in finite-dimensional nonlinear optimal
17 control. *J. Differ. Equations*, 258(1):81–114, 2015.
- 18 [40] M. Tucsnak and G. Weiss. Simultaneous exact controllability and some applications. *SIAM*
19 *J. Control Optim.*, 38(5):1408–1427, 2000.
- 20 [41] A. Valli. *A compact course on linear PDEs*. Springer, 2020.
- 21 [42] M. Warma and S. Zamorano. Exponential turnpike property for fractional parabolic equa-
22 tions with non-zero exterior data. *ESAIM:COCV*, 27(1):1–35, 2021.
- 23 [43] E. Weinan. A proposal on machine learning via dynamical systems. *Commun. Math. Stat.*,
24 5(1):1–11, 2017.
- 25 [44] J. Wu, X. Zhu, and S. Li. Simultaneous controllability of damped wave equations. *Math.*
26 *Method Appl. Sci.*, 40(1):319–324, 2017.
- 27 [45] J. Yeh. *Real analysis: theory of measure and integration second edition*. World Scientific
28 Publishing Company, 2006.
- 29 [46] S. Zamorano. Turnpike property for two-dimensional Navier–Stokes equations. *J. Math.*
30 *Fluid Mech.*, 20(3):869–888, 2018.
- 31 [47] R. Zbikowski. Lie algebra of recurrent neural networks and identifiability. In *1993 American*
32 *Control Conference*, pages 2900–2901. IEEE, 1993.