

Large-time asymptotics in Deep Learning

Borjan Geshkovski

Seminario de Estadística, UAM
October 23rd, 2020

●
○○○○○○○○○

○○○○○○○○○

○○○○○○○○○○○

○○○○

○○○○

Supervised learning

Supervised learning

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

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

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

- **Classification:** match points (images) to respective labels (cat, dog).
- High-dimensional interpolation problem.



Classification

Suppose $\vec{x}_i \in \mathbb{R}^2$ and $\vec{y}_i \in \{-1, 1\}$ for $i \leq N$.

- A simple idea:

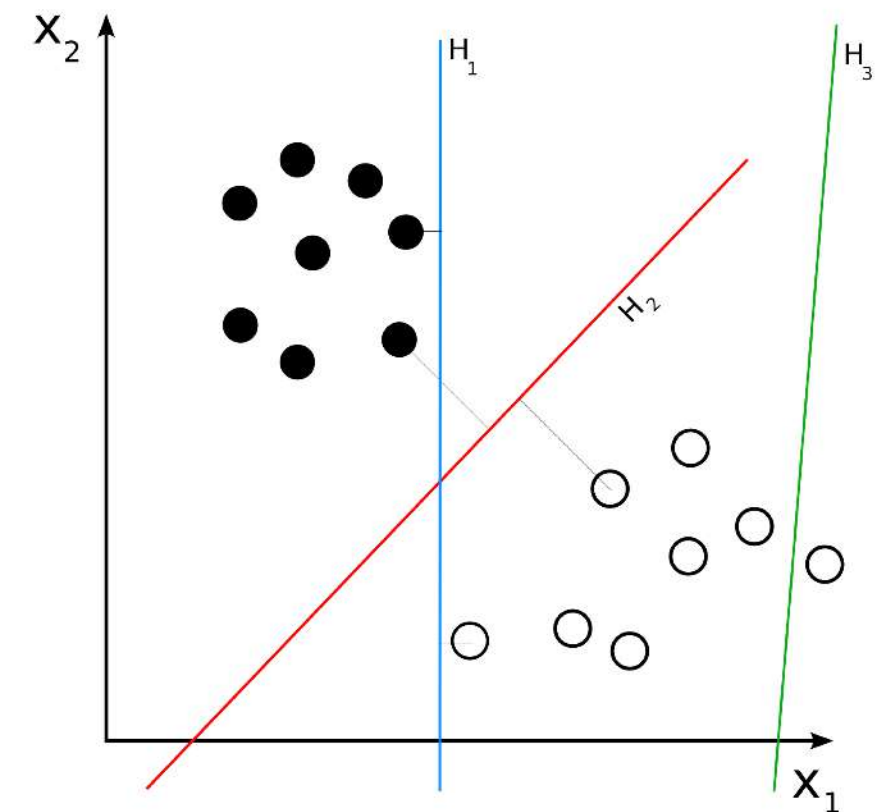
$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^N \left\| w^\top \vec{x}_i - \vec{y}_i \right\|^2 + \|w\|_{\ell^p}^p$$

$p = 1, 2$ and set

$$f(x) = \begin{cases} 1 & \text{if } w^\top x > 0 \\ -1 & \text{else.} \end{cases}$$

→ discriminative model (other ex: Support Vector Machine)

- Compared to generative models (e.g. Bayes classifier)
- But data is not linearly separable in general.



Neural networks

Neural network: for any $i \leq N$

$$\begin{cases} \mathbf{x}_i^{k+1} = \sigma(w^k \mathbf{x}_i^k + b^k) & \text{for } k \in \{0, \dots, N_{\text{layers}} - 1\} \\ \mathbf{x}_i^0 = \vec{x}_i \in \mathbb{R}^d, \end{cases} \quad (\text{NN}_1)$$

- $w^k \in \mathbb{R}^{d_{k+1} \times d_k}$ (**weights**) and $b^k \in \mathbb{R}^{d_k}$ (**biases**) are *controls*.
- $N_{\text{layers}} \geq 1$ given **depth**;
- d_k called **widths** with $d_0 = d$ and $d_{N_{\text{layers}}} = m$.
- $\sigma \in \text{Lip}(\mathbb{R})$ & $\sigma(0) = 0$ defined componentwise

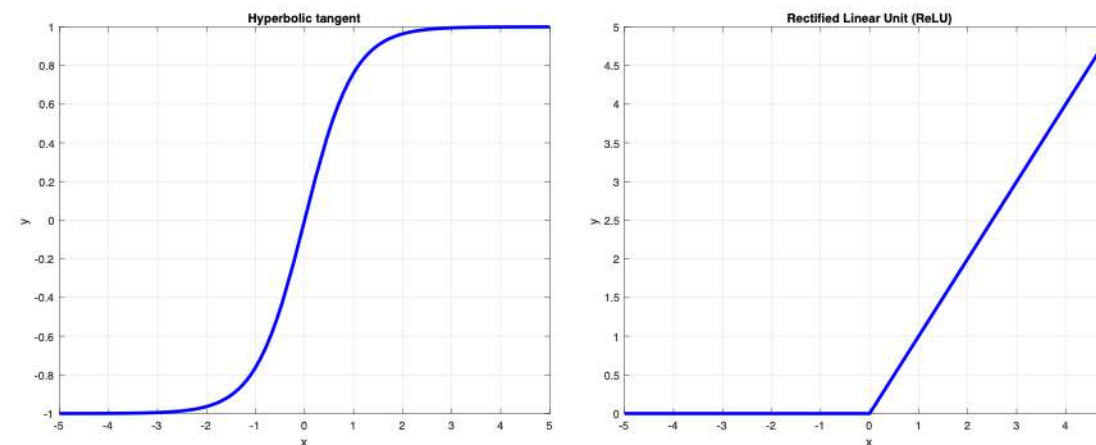


Figure: Sigmoid: $\tanh(x)$ and ReLU: $\max\{x, 0\}$

- ML jargon: multilayer perceptron / fully-connected.

"Training" a NN

Training \iff Optimization:

$$\inf_{\{w^k, b^k\}_{k=0}^{N_{\text{layers}}}} \underbrace{\frac{1}{N} \sum_{i=1}^N \text{loss}(w^{N_{\text{layers}}} \mathbf{x}_i^{N_{\text{layers}}}, \vec{y}_i)}_{\text{training error}} + \lambda \underbrace{\left\| \{w^k, b^k\}_k \right\|_{\ell^p}^p}_{\text{regularization}}$$

- $\text{loss}(x, y) = \|x - y\|_{\ell^p}^p$ for $p = 1, 2$ or $\text{loss}(x, y) = \log(1 + e^{x \cdot y})$;
- $\lambda > 0$ fixed.

Comments:

- $x \mapsto w^{N_{\text{layers}}} x^{N_{\text{layers}}}$ is the candidate approximation of unknown f
- $w^{N_{\text{layers}}} \mathbf{x}_i^{N_{\text{layers}}}$ can be replaced by $\varphi(w^{N_{\text{layers}}} \mathbf{x}_i^{N_{\text{layers}}})$ where $\varphi : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is nonlinear (softmax).
- In practice solved using stochastic gradient descent + backpropagation.

Residual neural networks

ResNets: fix $d_k \equiv d$; for any $i \leq N$

$$\begin{cases} \mathbf{x}_i^{k+1} = \mathbf{x}_i^k + h\sigma(w^k \mathbf{x}_i^k + b^k) & \text{for } k \in \{0, \dots, N_{\text{layers}} - 1\} \\ \mathbf{x}_i^0 = \vec{x}_i \end{cases} \quad (\text{ResNet})$$

where $h = 1$.

layer = timestep¹; $h = \frac{T}{N_{\text{layers}}}$ for given $T > 0$:

$$\begin{cases} \dot{\mathbf{x}}_i(t) = \sigma(w(t)\mathbf{x}_i(t) + b(t)) & \text{for } t \in (0, T) \\ \mathbf{x}_i(0) = \vec{x}_i. \end{cases} \quad (\text{nODE})$$

For (nODE), we shall henceforth assume $\sigma(\lambda x) = \lambda \sigma(x)$ for $\lambda > 0$ (positive homogeneity).

Residual neural networks

- In addition to (NN₁), one can also consider variants:

$$\begin{cases} \mathbf{x}_i^{k+1} = w^k \sigma(\mathbf{x}_i^k) + b^k & \text{for } k \in \{0, \dots, N_{\text{layers}} - 1\} \\ \mathbf{x}_i^0 = \vec{x}_i. \end{cases} \quad (\text{NN}_2)$$

→ motivates considering

$$\begin{cases} \dot{\mathbf{x}}_i(t) = w(t) \sigma(\mathbf{x}_i(t)) + b(t) & \text{for } t \in (0, T) \\ \mathbf{x}_i(0) = \vec{x}_i. \end{cases} \quad (\text{nODE}_2)$$

Training is optimal control

Henceforth fix $P : \mathbb{R}^d \rightarrow \mathbb{R}^m$ s.t. $P^{-1}(\{\vec{y}_i\}) \neq \emptyset$ for all $i \leq N$.

Given $T, \lambda > 0$:

$$\inf_{[w,b]^\top \in H^k(0,T;\mathbb{R}^{d_u})} \frac{1}{N} \sum_{i=1}^N \text{loss}(P\mathbf{x}_i(T), \vec{y}_i) + \lambda \left\| [w, b]^\top \right\|_{H^k(0,T;\mathbb{R}^{d_u})}^2$$

- $k = 0$ for (nODE₂), $k = 1$ for (nODE) (L^2 -regularization **may not be enough** for compactness \longrightarrow enhance to Sobolev regularization)
- $d_u := d^2 + d$

Henceforth denote the training error by

$$\phi(\mathbf{x}(T)) := \frac{1}{N} \sum_{i=1}^N \text{loss}(P\mathbf{x}_i(T), \vec{y}_i)$$

Why ODEs?

ODE formulation has been used to great effect:

Neural ordinary differential equations

[\[PDF\] nips.cc](#)

[RTQ Chen, Y Rubanova, J Bettencourt...](#) - Advances in **neural** ..., 2018 - [papers.nips.cc](#)

... at 3 Replacing residual networks with **ODEs** for supervised learning In this section, we experimentally investigate the training of **neural ODEs** for supervised learning. Software ...

☆ [🔗](#) Cited by 729 [Related articles](#) [All 20 versions](#) [🔗](#)

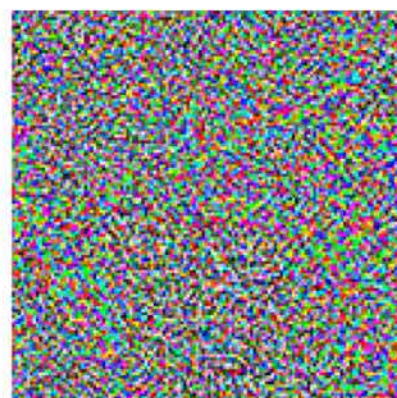
- **adaptive schemes, solvers** (Chen et al. '18, Dupont et al. '19, Benning et al. '19)
- **PMP-based training algos** (E et al. '19)
- **Stability to adversarial perturbations** (Haber, Ruthotto et al. '18)



"panda"

57.7% confidence

+ ϵ



=



"gibbon"

99.3% confidence

Artificial intelligence / Machine learning

A radical new neural network design could overcome big challenges in AI

Researchers borrowed equations from calculus to redesign the core machinery of deep learning so it can model continuous processes like changes in health.

by **Karen Hao**

December 12, 2018

MIT Tech Review, 2018

Why ODEs?

$$T \rightarrow \infty \quad \sim \quad N_{\text{layers}} \rightarrow \infty.$$

- Set $\mathbf{x}^0 = [\vec{x}_1, \dots, \vec{x}_N]^\top$, $u = [w, b]^\top$, and put both (nODE) and (nODE₂) in the form

$$\begin{cases} \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}(t), u(t)) & \text{in } (0, T) \\ \mathbf{x}(0) = \mathbf{x}^0 \in \mathbb{R}^{d_x}. \end{cases} \quad (\text{nODE})$$

- And so

$$\inf_{\substack{u \in H^k(0, T; \mathbb{R}^{d_u}) \\ \text{subject to (nODE)}}} \phi(\mathbf{x}(T)) + \lambda \|u\|_{H^k(0, T; \mathbb{R}^{d_u})}^2 \quad (\text{SL}_1)$$

- $T \gg 1$ would imply that $N_{\text{layers}} \gg 1 \rightarrow \text{deep learning regime}$. So,

Question: What happens to a minimizer u^T solving (SL₁), and corresponding state \mathbf{x}^T to (nODE) when $T \rightarrow \infty$?

Qualitative results

Scaling

Key idea: *Time-Scaling*.

- Assumptions on $\sigma \longrightarrow \mathbf{f}(\mathbf{x}, u)$ is positively homogeneous w.r.t. u , i.e. $\mathbf{f}(\mathbf{x}, \alpha u) = \alpha \mathbf{f}(\mathbf{x}, u)$ for $\alpha > 0$.
- Hence, given $u^T(t)$ and the solution $\mathbf{x}^T(t)$ to

$$\begin{cases} \dot{\mathbf{x}}^T(t) = \mathbf{f}(\mathbf{x}^T(t), u^T(t)) & \text{in } (0, T) \\ \mathbf{x}^T(0) = \mathbf{x}^0, \end{cases} \quad (1)$$

then $u^1(t) := Tu^T(tT)$ is such that $\mathbf{x}^1(t) := \mathbf{x}^T(tT)$ solves (1) for $t \in [0, 1]$.

Then:

$$\begin{aligned} \inf_{u^T} \phi(\mathbf{x}^T(T)) + \lambda \int_0^T \|u^T(t)\|^2 dt &= \inf_{u^T} \phi(\mathbf{x}^T(T)) + \frac{\lambda}{T} \int_0^1 \|Tu^T(sT)\|^2 ds \\ &= \frac{1}{T} \inf_{u^1} T\phi(\mathbf{x}^T(T)) + \lambda \int_0^1 \|Tu^T(sT)\|^2 ds \\ &= \frac{1}{T} \inf_{u^1} \phi(\mathbf{x}^1(1)) + \lambda \int_0^1 \|u^1(s)\|^2 ds. \end{aligned}$$

Theorem (Esteve, G., Pighin, Zuazua, '20): Fix $\lambda > 0$ and suppose $\{\phi = 0\} \neq \emptyset$. For any $T > 0$, let u^T be minimizer in (SL_1) , \mathbf{x}^T associated solution to (nODE) . Assume that (nODE) is controllable. Then

$$\phi(\mathbf{x}^T(T)) \longrightarrow 0 \quad \text{as } T \rightarrow +\infty.$$

Moreover, $\exists \{T_n\}_{n=1}^{+\infty}$ positive times and $\exists \mathbf{x}^\dagger \in \mathbb{R}^{d_x}$, $\phi(\mathbf{x}^\dagger) = 0$, such that

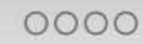
$$\|\mathbf{x}^{T_n}(T_n) - \mathbf{x}^\dagger\| \longrightarrow 0 \quad \text{as } n \rightarrow +\infty.$$

Moreover

$$\left\| \frac{1}{T_n} u^{T_n} \left(\frac{\cdot}{T_n} \right) - u^* \right\|_{H^k(0,1;\mathbb{R}^{d_u})} \longrightarrow 0 \quad \text{as } n \rightarrow +\infty$$

where u^* solves

$$\inf_{\substack{u \in H^k(0,1;\mathbb{R}^{d_u}) \\ \text{subject to } (\text{nODE}) \text{ with } T=1 \\ \text{and} \\ \phi(\mathbf{x}(1))=0}} \lambda \|u\|_{H^k(0,1;\mathbb{R}^{d_u})}^2.$$



T=16



T=36



T=81



Figure: Here $N_{\text{layers}} = \left\lfloor T^{\frac{3}{2}} \right\rfloor$ and thus $h = \frac{1}{\sqrt{T}}$, and we consider $\alpha = 1$.

$$T \rightarrow \infty \iff \lambda \rightarrow 0$$

Back to

$$\begin{aligned} \phi(\mathbf{x}^T(T)) + \lambda \int_0^T \|u^T(t)\|^2 dt &= \phi(\mathbf{x}^T(T)) + \frac{\lambda}{T} \int_0^1 \|Tu^T(sT)\|^2 ds \\ &= \phi(\mathbf{x}^T(T)) + \frac{\lambda}{T} \int_0^1 \|u^1(s)\|^2 ds \\ &= \phi(\mathbf{x}^1(1)) + \frac{\lambda}{T} \int_0^1 \|u^1(s)\|^2 ds. \end{aligned}$$

Corollary: Fix $T > 0$. Under the assumptions of the Theorem,

$$\phi(\mathbf{x}^T(T)) \longrightarrow 0 \quad \text{as } \lambda \rightarrow 0.$$

Moreover

$$\|u^T\|_{H^k(0,T;\mathbb{R}^{d_u})} \longrightarrow \|u^*\|_{H^k(0,1;\mathbb{R}^{d_u})} \quad \text{as } \lambda \rightarrow 0$$

where u^* solves

$$\begin{aligned} &\inf_{u \in H^k(0,1;\mathbb{R}^{d_u})} \|u\|_{H^k(0,1;\mathbb{R}^{d_u})}^2 \\ &\text{subject to (nODE) with } T=1 \\ &\quad \text{and } \phi(\mathbf{x}(1))=0 \end{aligned}$$

Discussion

- Back to

$$\min_{w \in \mathbb{R}^d} \sum_{i=1}^N \text{loss} \left(w^\top \text{sign}(\vec{x}_i), \vec{y}_i \right) + \lambda \|w\|_{\ell^p}^p$$

where e.g. $\text{loss}(x, y) = \log(1 + e^{x \cdot y})$; shown² that

$$\lim_{\lambda \rightarrow 0} \hat{w}^\lambda / \|\hat{w}^\lambda\| = w^*$$

where w^* is maximum margin separator:

$$w^* = \operatorname{argmax}_{\|w\|_p=1} \min_i \vec{y}_i w^\top \text{sign}(\vec{x}_i).$$

- Compared to other convergence results of generalization nature: **implicit bias property of gradient descent**³:

"In the overparametrized regime, after training a neural network with gradient-based methods until zero training error, with $\lambda = 0$, among the many classifiers which overfit on the training dataset, the algorithm selects the one which performs best on the test dataset."

- We clearly **exhibit the explicit L^2 -regularization** of the control parameters.

²Rosset, Hu, Hastie '04

³Zhang et al. '16, Soudry et al. '18, Gunasekar et al. '18, Chizat & Bach '20

Proof of Theorem

For simplicity, suppose $k = 0$.

Part 1). We first show that

$$\phi(\mathbf{x}^T(T)) \longrightarrow 0 \quad \text{as } T \rightarrow \infty.$$

Recall that $\min_{\mathbb{R}^{d_x}} \phi = 0$.

1 By controllability, $\exists u^1 \in L^2(0, 1)$ such that $\phi(\mathbf{x}^1(1)) = 0$.

2 Since u^T is a minimizer,

$$\begin{aligned} \phi(\mathbf{x}^T(T)) + \lambda \|u^T\|_{L^2(0, T)}^2 &\leq \phi(\mathbf{x}^1(1)) + \lambda \left\| \frac{\cdot}{T} u^1 \left(\frac{\cdot}{T} \right) \right\|_{L^2(0, T)}^2 \\ &= \phi(\mathbf{x}^1(1)) + \frac{\lambda}{T} \|u^1\|_{L^2(0, 1)}^2 \end{aligned}$$

3 Since $\phi(\mathbf{x}^1(1)) = 0$,

$$0 \leq \phi(\mathbf{x}^T(T)) \leq \frac{\lambda}{T} \|u^1\|_{L^2(0, 1)}^2.$$

Part 2). We now show that $\exists \{T_n\}_{n=1}^{+\infty}$ of positive times and $\exists \mathbf{x}^\dagger \in \mathbb{R}^{d_x}$, $\phi(\mathbf{x}^\dagger) = 0$ such that

$$\left\| \mathbf{x}^{T_n}(T_n) - \mathbf{x}^\dagger \right\| \longrightarrow 0 \quad \text{as } n \rightarrow +\infty.$$

1 Integral formulation of ODE + Grönwall + Cauchy Schwarz + scaling:

$$\begin{aligned} \left\| \mathbf{x}^T(T) - \mathbf{x}^0 \right\| &\lesssim_{N,\sigma} \sqrt{T} \left\| u^T \right\|_{L^2(0,T)} \exp \left(\sqrt{T} \left\| u^T \right\|_{L^2(0,T)} \right) \\ &\lesssim_{N,\sigma} \left\| u^1 \right\|_{L^2(0,1)} \exp \left(\left\| u^1 \right\|_{L^2(0,1)} \right) \end{aligned}$$

Thus $\{\mathbf{x}^T(T)\}_{T>0}$ is bounded (subset of \mathbb{R}^{d_x});

2 $\longrightarrow \exists \{T_n\}_{n=1}^{+\infty}$ of positive times and $\exists \mathbf{x}^\dagger \in \mathbb{R}^{d_x}$ such that

$$\left\| \mathbf{x}^{T_n}(T_n) - \mathbf{x}^\dagger \right\| \longrightarrow 0 \quad \text{as } n \rightarrow +\infty.$$

3 By **Part 1)**, $\phi(\mathbf{x}^{T_n}(T_n)) \longrightarrow 0$. By lower semicontinuity of ϕ ,

$$\phi(\mathbf{x}^\dagger) \leq \liminf_{n \rightarrow +\infty} \phi(\mathbf{x}^{T_n}(T_n)) = 0.$$

Part 3). We finally show that $u_n(t) := \frac{1}{T_n} u^{T_n}(\frac{t}{T_n})$ for $t \in [0, T_n]$ satisfies

$$\|u_n - u^*\|_{L^2(0,1)} \longrightarrow 0 \quad \text{as } n \rightarrow +\infty$$

where u^* solves

$$\inf_{u \in L^2(0,1)} \lambda \|u\|_{L^2(0,1)}^2 \cdot$$

subject to (nODE) with $T=1$
and
 $\phi(\mathbf{x}(1))=0$

- 1 Let $u^0 \in L^2(0,1)$ solution to above minimization problem;
- 2 By contradiction: $\|u_n\|_{L^2(0,1)} \leq \|u^0\|_{L^2(0,1)}$ for every $n \geq 1$;
- 3 Banach-Alaoglu: $\exists u^* \in L^2(0,1)$ such that

$$u_n \rightharpoonup u^* \quad \text{weakly in } L^2(0,1)$$

\mathbf{x}^* trajectory associated to u^* ; compactness properties of ODE:

$$\mathbf{x}_n \longrightarrow \mathbf{x}^* \quad \text{strongly in } C^0[0,1]$$

- 4 But $\mathbf{x}^{T_n}(T_n) = \mathbf{x}_n(1)$ thus $\mathbf{x}^*(1) = \mathbf{x}^\dagger$ by **Part 1)**, so $\phi(\mathbf{x}^*(1)) = 0$.
- 5 Weak lower semicontinuity of L^2 -norm:

$$\|u^0\|_{L^2(0,1)}^2 \leq \|u^*\|_{L^2(0,1)}^2 \leq \liminf_{n \rightarrow \infty} \|u_n\|_{L^2(0,1)}^2 \leq \limsup_{n \rightarrow \infty} \|u_n\|_{L^2(0,1)}^2 \leq \|u^0\|_{L^2(0,1)}^2$$

so strong L^2 -convergence and u^* solves the desired problem.

Quantitative results

Setting

Question: Quantitative estimates for the time T required to reach the zero training error regime $\phi(\mathbf{x}(T)) = 0$?

- We might need to change the cost function!
- Consider $\text{loss}(x, y) = \|x - y\|^2$ and so we recall

$$\phi(\mathbf{x}(T)) := \frac{1}{N} \sum_{i=1}^N \|P\mathbf{x}_i(T) - \vec{y}_i\|^2$$

- Modified supervised learning problem:

$$\min_{\substack{u \in H^k(0, T; \mathbb{R}^{d_u}) \\ \text{subject to (nODE)}}} \left[J_T(u) := \int_0^T \phi(\mathbf{x}(t)) dt + \lambda \|u\|_{H^k(0, T; \mathbb{R}^{d_u})}^2 \right] \quad (\text{SL}^*)$$

Exponential stabilization

Theorem (Esteve, G., Pighin, Zuazua '20): Fix $\lambda > 0$ and suppose that (nODE) is controllable. There exist $T^* > 0$ such that for any $T \geq T^*$, any solution (u^T, \mathbf{x}^T) to (SL*)–(nODE) satisfies

$$\phi(\mathbf{x}^T(t)) \leq C_1 e^{-\mu t} \quad \forall t \in [0, T]$$

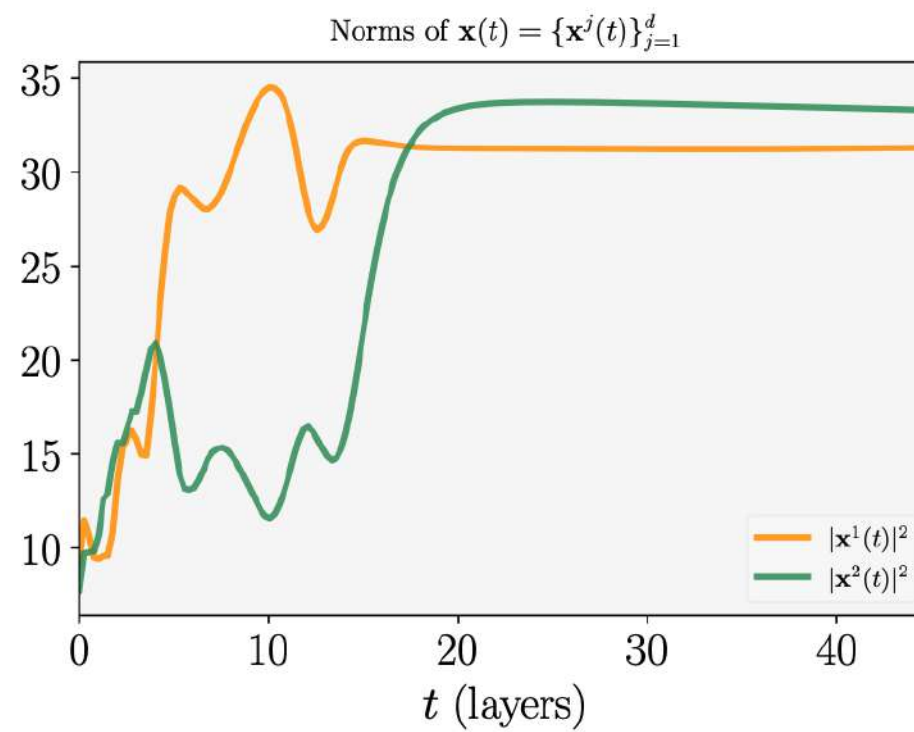
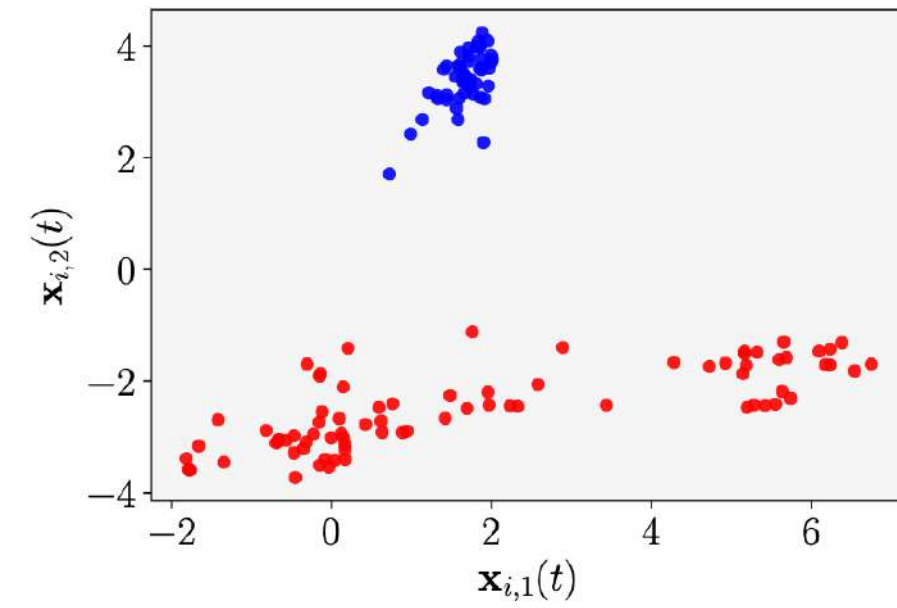
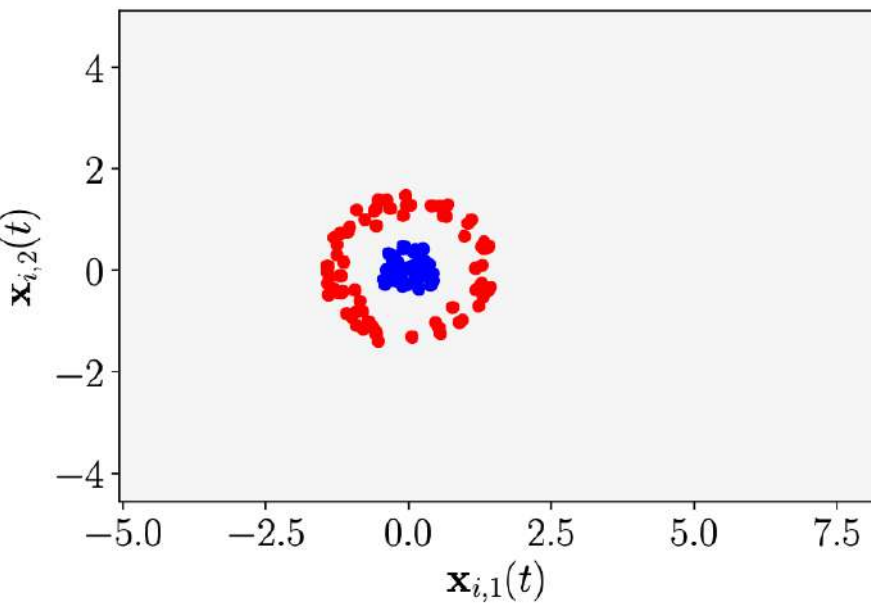
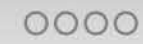
and

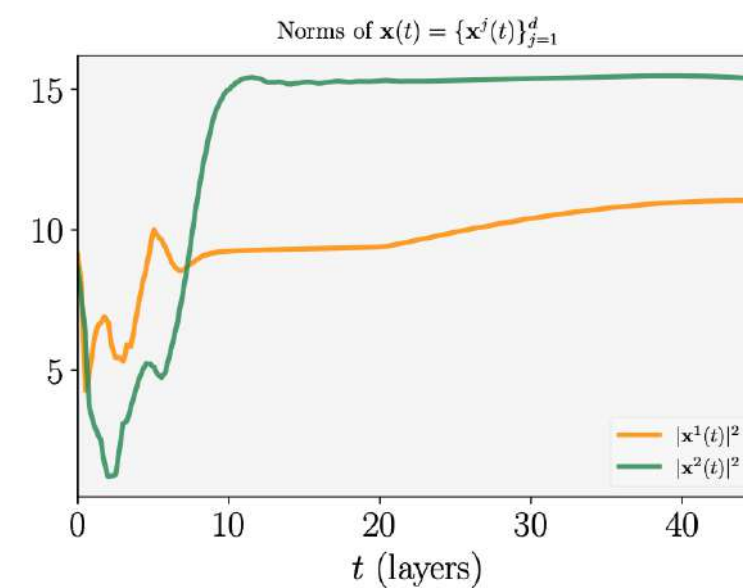
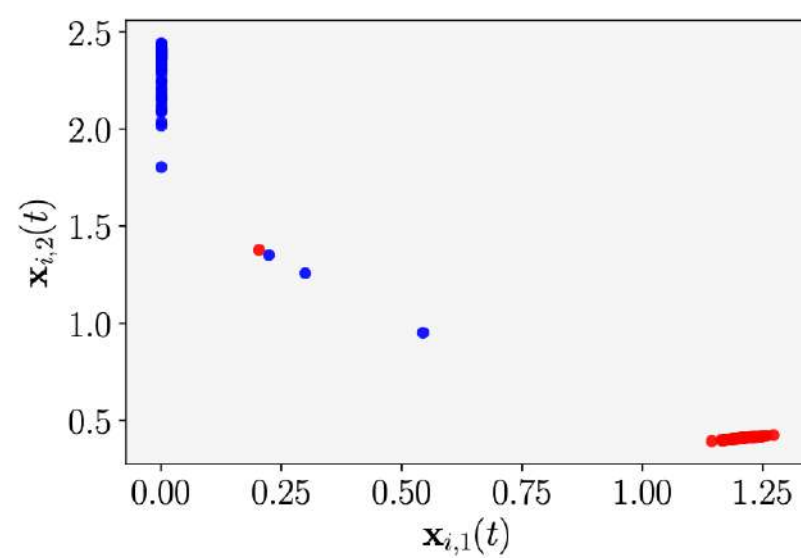
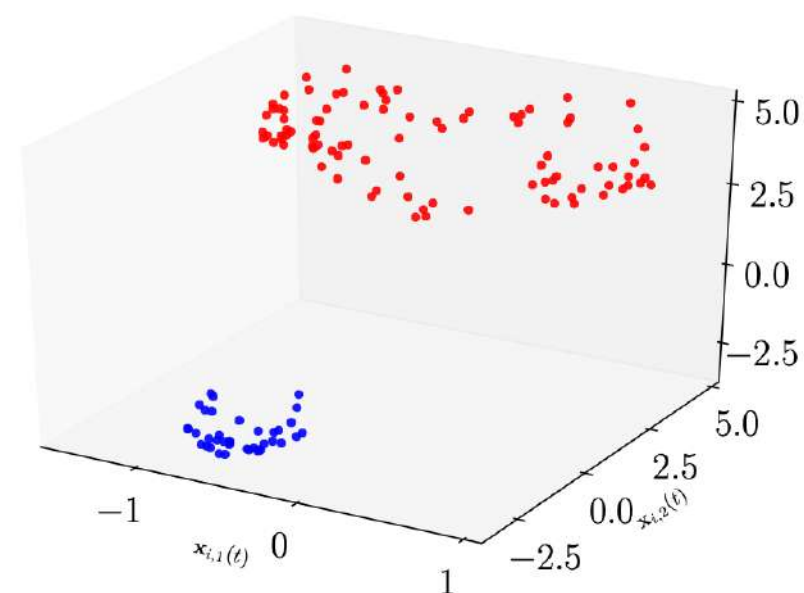
$$\|u^T(t)\| \leq C_2 e^{-\mu t} \quad \text{for a.e. } t \in [0, T]$$

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

Remarks:

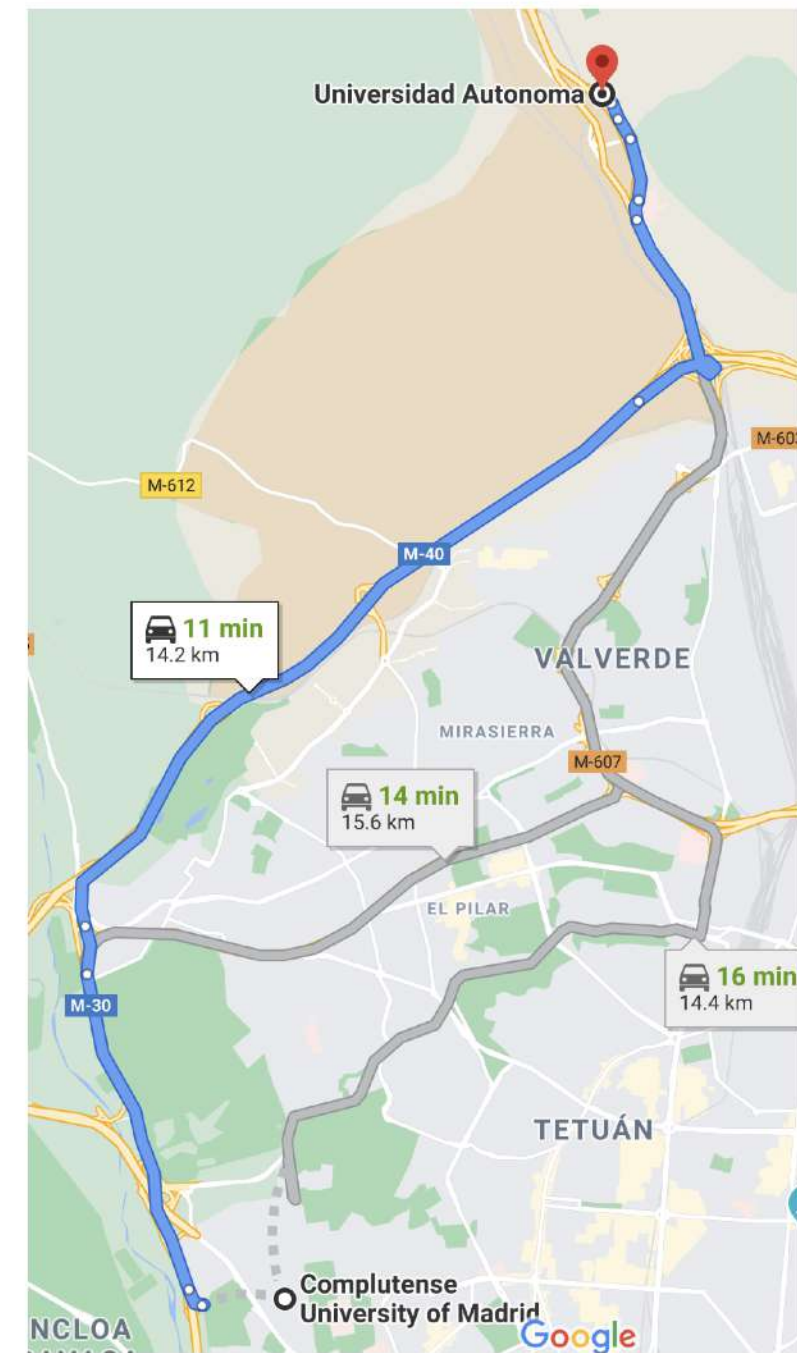
- Akin to *universal approximation*: given tolerance $\varepsilon > 0$, there exists $T_\varepsilon > 0$ (number of layers) and control parameters u^ε such that the neural network output is ε -close to the desired target.
- The difference with universal approximation is that our parameters may be computed explicitly via a training procedure.





Turnpike property

- Theorem is a special manifestation of the well-known **turnpike property** in optimal control and economics.
- *For suitable optimal control problems in a sufficiently large T , any optimal solution (u^T, \mathbf{x}^T) remains, during most of the time, $\mathcal{O}(e^{-t} + e^{-(T-t)})$ -close to the optimal solution of a corresponding “static” optimal control problem.*
Optimal static solution is referred to as the turnpike – the name stems from the idea that a turnpike is the fastest route between two points which are far apart, even if it is not the most direct route.
- Since $\mathbf{f}(\mathbf{x}, 0) = 0$ for all \mathbf{x} , $\bar{\mathbf{y}}_i$ may be seen as the turnpike for $P\mathbf{x}_i$. Since this is a steady state, we do not see an exit from the turnpike and we stabilize.



Proof of $\phi(\mathbf{x}^T(t)) \lesssim e^{-t}$

Suppose $k = 0$, $d = m$ and $N = 1$ for simplicity. Then $\phi(\mathbf{x}^T(t)) = \|\mathbf{x}^T(t) - \vec{y}\|^2$.

Part 1). For $T \geq 1$, we first prove that

$$\left\| \mathbf{x}^T(t) - \vec{y} \right\|^2 + \left\| \mathbf{x}^T - \vec{y} \right\|_{L^2(0,T)}^2 + \left\| u^T \right\|_{L^2(0,T)}^2 \lesssim_{\sigma} \|\vec{x} - \vec{y}\|^2 \quad (2)$$

for all $t \in [0, T]$ uniformly in T .

1 Controllability: $\exists u^1 \in L^2(0, 1)$ such that $\mathbf{x}^1(1) = \vec{y}$ and $\|u^1\|_{L^2} \lesssim \|\vec{x} - \vec{y}\|$.

2 Grönwall: $\|\mathbf{x}^1(t) - \vec{y}\| \lesssim_{\sigma} \|\vec{x} - \vec{y}\|$

3 Set

$$u^{\text{aux}}(t) := \begin{cases} u^1(t) & \text{for } t \in (0, 1) \\ 0 & \text{for } t \in (1, T). \end{cases}$$

Then $\mathbf{x}^{\text{aux}}(t) = \vec{y}$ for $t \in [1, T]$.

4 u^T minimizer, so

$$\begin{aligned} \left\| \mathbf{x}^T - \vec{y} \right\|_{L^2(0,T)}^2 + \left\| u^T \right\|_{L^2(0,T)}^2 &\leq \left\| \mathbf{x}^1 - \vec{y} \right\|_{L^2(0,1)}^2 + \left\| u^1 \right\|_{L^2(0,1)}^2 \\ &\lesssim_{\sigma} \|\vec{x} - \vec{y}\|. \end{aligned}$$

5 Conclude by Grönwall.

Part 2). Fix $\tau > C_\sigma^4 + 1$ ($C_\sigma > 0$ appears in (2)) and let $T \geq 2\tau + 1$.

1 For $t \in [0, \tau + 1]$, desired estimate follows from (2):

$$\|\mathbf{x}^T(t) - \vec{y}\| \lesssim_\sigma \|\vec{x} - \vec{y}\| \lesssim_{\sigma, \tau} e^{-t} \|\vec{x} - \vec{y}\|$$

2 (2) + contradiction argument:

$$\|\mathbf{x}^T(t) - \vec{y}\| \leq \frac{C_\sigma^2}{\sqrt{\tau}} \|\vec{x} - \vec{y}\|$$

for $t \in [\tau, T]$.

3 Bootstrap: for $n \leq \frac{T}{2\tau}$

$$\|\mathbf{x}^T(t) - \vec{y}\| \leq \left(\frac{C_\sigma^2}{\sqrt{\tau}} \right)^n \|\vec{x} - \vec{y}\|$$

for $t \in [n\tau, T]$.

4 Suppose $t \in [\tau + 1, T]$. Set $n(t) = \lfloor \frac{t}{\tau+1} \rfloor$. Then $n(t) \leq \frac{T}{2\tau}$, $t \in [n(t)\tau, T]$ and $n(t) \geq \frac{t}{\tau+1} - 1$, so

$$\begin{aligned} \|\mathbf{x}^T(t) - \vec{y}\| &\leq \exp \left(-n(t) \log \left(\frac{\sqrt{\tau}}{C_\sigma^4} \right) \right) \|\vec{x} - \vec{y}\| \\ &\lesssim_{\tau, \sigma} \exp \left(-\frac{\log \left(\frac{\sqrt{\tau}}{C_\sigma^4} \right)}{\tau + 1} t \right) \|\vec{x} - \vec{y}\|. \end{aligned}$$

Proof of $\|u^T(t)\| \lesssim e^{-t}$

Let $t \in [0, T)$ and $0 < h \ll 1$ s.t. $t + 2h \in [0, T]$.

1 Set

$$u^{\text{aux}}(s) := \begin{cases} u^T(s) & \text{for } s \in (0, t) \\ \frac{1}{2} u^T\left(t + \frac{s-t}{2}\right) & \text{for } s \in (t, t+2h) \\ u^T(s-h) & \text{for } s \in (t+2h, T). \end{cases}$$

2 Since u^T minimizer, by $J_T(u^T) \leq J_T(u^{\text{aux}})$, we will find

$$\frac{1}{2} \int_t^{t+h} \|u^T(s)\|^2 ds \leq \int_t^{t+h} \|x^T(s) - \vec{y}\|^2 ds.$$

3 Combined with $\|x^T(s) - \vec{y}\|^2 \lesssim e^{-t}$,

$$\int_t^{t+h} \|u^T(s)\|^2 ds \lesssim \int_t^{t+h} e^{-s} ds \lesssim h e^{-t}$$

4 Lebesgue differentiation theorem: for a.e. $t \in [0, T]$,

$$\|u^T(t)\|^2 = \lim_{h \rightarrow 0^+} \frac{1}{h} \int_t^{t+h} \|u^T(s)\|^2 ds \lesssim e^{-t}.$$

L^1 -regularization

Theorem (Esteve, G., Pighin, Zuazua, '20): Fix $M > 0$ and assume $\{\phi = 0\} \neq \emptyset$. Suppose (nODE₂) is controllable. Consider

$$\inf_{\substack{u \in L^1(0, T; \mathbb{R}^{d_u}) \\ \text{esssup} \|u\| \leq M \\ \text{subject to (nODE}_2\text{)}}} \int_0^T \phi(\mathbf{x}(t)) dt + \lambda \|u\|_{L^1(0, T; \mathbb{R}^{d_u})}.$$

Then there exists $T_M > 0$ such that for any $T > T_M$, any optimal u^T and corresponding state \mathbf{x}^T , unique solution to (nODE₂), satisfy

$$\phi(\mathbf{x}^T(t)) = 0, \quad \text{for all } t \in [T^*, T]$$

and

$$\begin{aligned} \|u^T(t)\| &= M, & \text{for a.e. } t \in (0, T^*) \\ \|u^T(t)\| &= 0, & \text{for a.e. } t \in (T^*, T). \end{aligned}$$

for some $0 < T^* \leq T_M$.

Proof

1. **Part 1:** Show that for

$$T^* := \min \left\{ t \in [0, T] : \phi(\mathbf{x}^T(t)) = \min_{s \in [0, T]} \phi(\mathbf{x}^T(s)) \right\}$$

it holds

$$\begin{aligned} \phi(\mathbf{x}^T(t)) &> \phi(\mathbf{x}^T(T^*)) \quad \text{for all } t \in [0, T^*), \\ \text{and } \phi(\mathbf{x}^T(t)) &= \phi(\mathbf{x}^T(T^*)) \quad \text{for all } t \in [T^*, T]. \end{aligned}$$

2. **Part 2:** Show that the optimal control parameters u^T satisfy

$$\begin{aligned} \|u^T(t)\| &= M \quad \text{for a.e. } t \in (0, T'), \\ \text{and } \|u^T(t)\| &= 0 \quad \text{for a.e. } t \in (T', T). \end{aligned}$$

This relies on the fact that the L^1 -norm is invariant to the time-scaling and Lebesgue measure theory.

3. **Part 3:** Show that there exists $T_M > 0$ such that if $T > T_M$, then

$$T^* \leq T_M \quad \text{and} \quad \phi(\mathbf{x}^T(T^*)) = 0.$$

Zero training error regime

Recall:

- We are given a training dataset $\{\vec{x}_i, \vec{y}_i\}$;
- the training error $\phi \in C^0(\mathbb{R}^{d \times N}; \mathbb{R}_+)$ is defined by

$$\phi(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \text{loss}(P\mathbf{x}_i, \vec{y}_i)$$

Our results stipulate: **when depth N_{layers} is increased, trained trajectories of neural networks approach the zero training error regime $\phi = 0$.**

Question: Can we reach exactly the zero training error regime, and how big are the control parameters in this case?

A lower bound

Weights need to be at least of a certain size! \rightarrow depend on the way the dataset is "spread out".

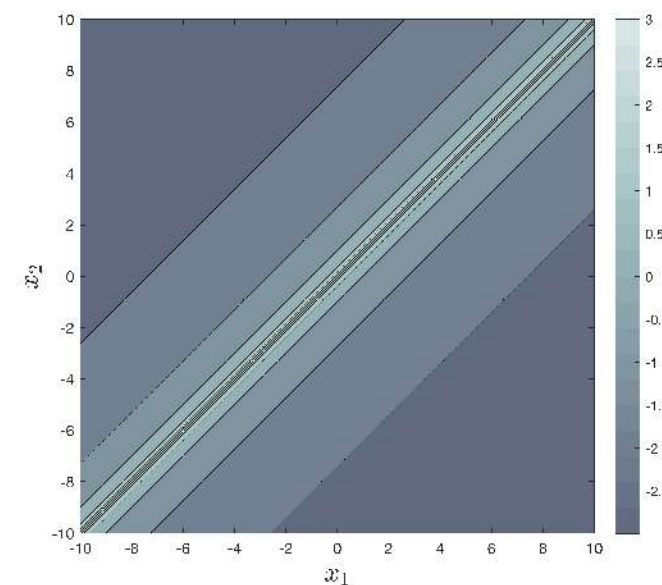
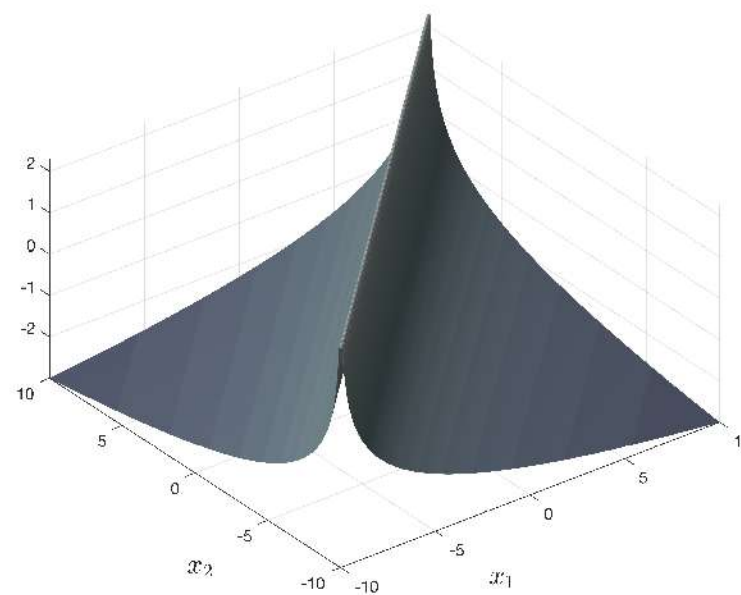
Theorem (Esteve et al. '20): Let $T > 0$. Assume that for some control parameters $u := [w, b]^\top$, the solutions $\mathbf{x}_i(t)$ to either (nODE) or (nODE₂) satisfies

$$P\mathbf{x}_i(T) = \vec{y}_i \quad \text{for all } i \in \{1, \dots, N\}.$$

Then

$$\|w\|_{L^1(0, T; \mathbb{R}^{d_u})} \geq L_\sigma \max_{\substack{(i,j) \in \{1, \dots, N\}^2 \\ i \neq j}} \inf_{\substack{\mathbf{x}_i^1 \in P^{-1}(\{\vec{y}_i\}) \\ \mathbf{x}_j^1 \in P^{-1}(\{\vec{y}_j\})}} \log \left(\frac{\|\mathbf{x}_i^1 - \mathbf{x}_j^1\|}{\|\mathbf{x}_i^0 - \mathbf{x}_j^0\|} \right)$$

where $L_\sigma > 0$ is the Lipschitz constant of σ .



Controllability

Theorem (Esteve et al. '20): Let $T > 0$ and assume that $N \leq d$. Fix $\mathbf{x}^1 \in \mathbb{R}^{d_x}$ and assume that $\sigma \in C^1(\mathbb{R})$ is such that

$$\left\{ \sigma(\mathbf{x}_1^1), \dots, \sigma(\mathbf{x}_i^1), \dots, \sigma(\mathbf{x}_N^1) \right\}$$

is a system of linearly independent vectors in \mathbb{R}^d .

There exists $r > 0$ such that for any $\mathbf{x}^0 \in \mathbb{R}^{d_x}$ satisfying $\|\mathbf{x}^0 - \mathbf{x}^1\| \leq r$, there exists weights $w \in L^\infty(0, T; \mathbb{R}^{d^2})$ s.t. the solution \mathbf{x} to

$$\begin{cases} \dot{\mathbf{x}}(t) = \text{diag}(w(t))\sigma(\mathbf{x}(t)) & \text{in } (0, T) \\ \mathbf{x}(0) = \mathbf{x}^0, \end{cases}$$

satisfies

$$\mathbf{x}(T) = \mathbf{x}^1,$$

and the estimate

$$\|w\|_{L^\infty(0, T; \mathbb{R}^{d^2})} \leq \frac{C}{T} \|\mathbf{x}^0 - \mathbf{x}^1\|,$$

holds for some $C > 0$ independent of T .

Extensions

Variable width

Variable width ResNets via **integro-differential equation**: for $i \leq N$

$$\partial_t \mathbf{z}_i(t, x) = \sigma \left(\int_{\Omega} w(t, x, \xi) \mathbf{z}_i(t, \xi) d\xi + b(t, x) \right) \quad \text{in } (0, T) \times \Omega.$$

- e.g. $\Omega = \text{image} \times (0, 1) \subset \mathbb{R}^3$;
- All previous asymptotics theorems apply here;
- Variable width ResNets can be obtained by semi-discretizing via time-dependent mesh.

Switched systems: Changing widths over layers as switched systems over time:

$$\dot{x}(t) = f_{\rho(t)}(x(t), u(t))$$

given M vector fields f_1, \dots, f_M and switching signal $\rho : [0, T] \rightarrow \{1, \dots, M\}$;

Quasi-turnpike strategy:

#1 increase the dimension to the "optimal system" f_{j^*} ,

#2 use the stabilization/turnpike for fixed width

The optimal system f_{j^*} ? \longrightarrow optimal with respect to cost.

What are the switching times? How many?

Outlook

- 1 Long-time behavior depends on the cost functional to be minimized.
- 2 Results should be complemented by ML subfields (e.g. CNN design, training algorithms..)

Some additional problems:

- Asymptotics remain to be proven when $P : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is optimizable variable?
- Stabilization for non L^2 –loss with L^2 –regularization?
- Further characterization of limit control u^* in first theorem?
- Extensive bibliography can be found in

LARGE-TIME ASYMPTOTICS IN DEEP LEARNING

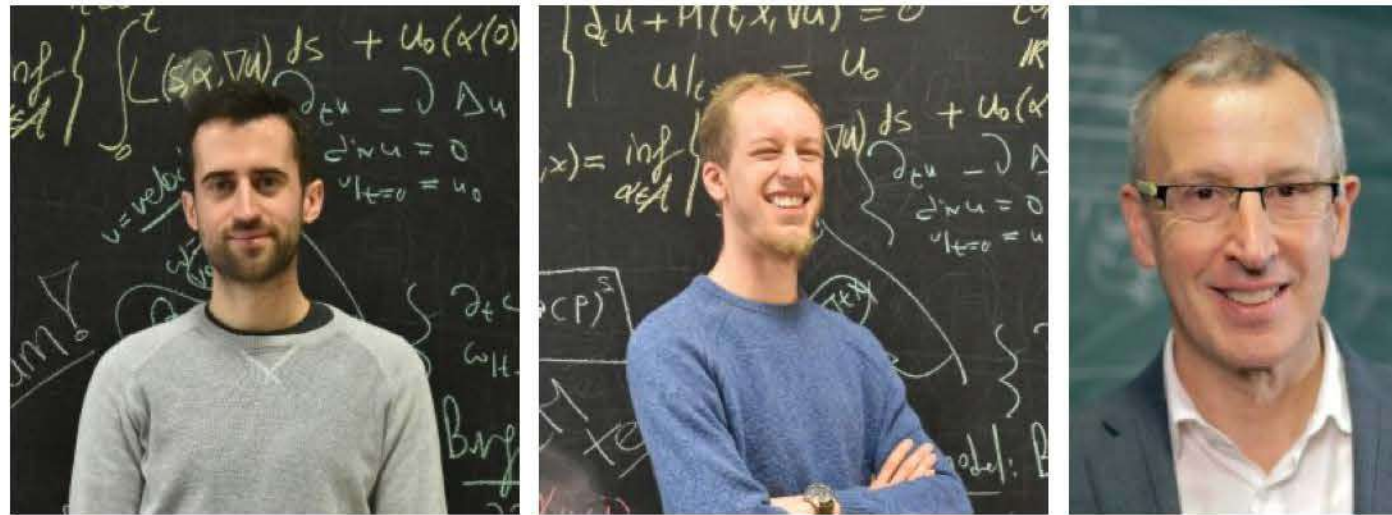
CARLOS ESTEVE, BORJAN GESHKOVSKI, DARIO PIGHIN, AND ENRIQUE ZUAZUA

ABSTRACT. It is by now well-known that practical deep supervised learning may roughly be cast as an optimal control problem for a specific discrete-time, nonlinear dynamical system called an artificial neural network. In this work, we consider the

<https://arxiv.org/abs/2008.02491>

Thank you for your attention!

Collaborators:



- C. Esteve (UAM/Deusto), D. Pighin (PhD @ UAM, 2020), E. Zuazua (FAU/Deusto/UAM).



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



FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG

