

On the intersection of control and machine learning

Borjan¹ Geshkovski
borjang.github.io

GdT Contrôle, LJLL

25 Novembre, 2022



**Massachusetts
Institute of
Technology**

¹Bor-Yann

Supervised learning

Goal: Approximate an unknown function $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ given data

$$\mathcal{D} := \left\{ x^{(i)}, f(x^{(i)}) \right\}_{i \in [n]} \subset \mathbb{R}^d \times \mathbb{R}^m.$$

Supervised learning

Goal: Approximate an unknown function $f : \mathbb{R}^d \rightarrow \mathbb{R}^m$ given data

$$\mathcal{D} := \left\{ x^{(i)}, f(x^{(i)}) \right\}_{i \in [n]} \subset \mathbb{R}^d \times \mathbb{R}^m.$$

We distinguish:

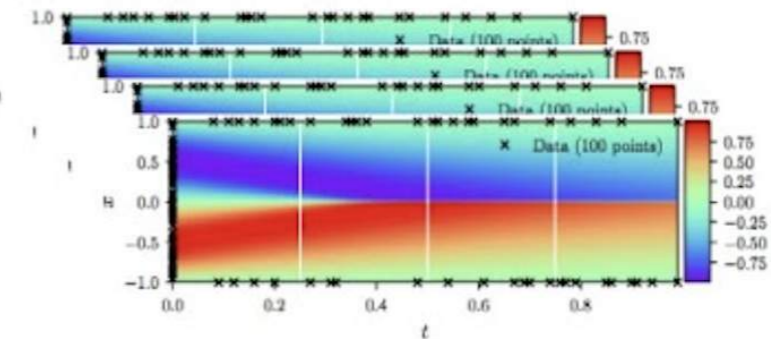
- **Classification:** $\text{ran}(f)$ is $\{e_j\}_{j \in [m]} \subset \mathbb{R}^m$.
Image, audio: $d \geq 10^3, d \gg m$.

- **Regression:** $\text{ran}(f)$ is \mathbb{R}^m .
PDE: fixed initial/boundary data, then learn $(t, x) \mapsto u(t, x)$

$d = 784$
 $m = 10$



$x^{(1)}$ $x^{(2)}$ $x^{(3)}$ $x^{(4)}$



(Feed-forward) neural networks

$$f_{\text{approx}}(x) := P_{\mathbf{X}}^{[n_T]}(x)$$

for $x \in \mathbb{R}^d$.

(Feed-forward) neural networks

$$\boxed{f_{\text{approx}}(x) := P_{\mathbf{X}^{[n_T]}}(x)} \quad \text{for } x \in \mathbb{R}^d.$$

where

- ▶ $P \in \mathbb{R}^{m \times d_{n_T}}$ (suppose given)

(Feed-forward) neural networks

$$f_{\text{approx}}(x) := P_{\mathbf{X}^{[n_T]}}(x) \quad \text{for } x \in \mathbb{R}^d.$$

where

- ▶ $P \in \mathbb{R}^{m \times d_{n_T}}$ (suppose given)
- ▶ $\mathbf{x}^{[n_T]}(x) \in \mathbb{R}^{d_{n_T}}$ is output of neural net with $n_t \geq 1$ *layers*:

$$\begin{cases} \mathbf{x}^{[k+1]} = c^{[k]} \sigma(a^{[k]} \cdot \mathbf{x}^{[k]} + b^{[k]}) & k \in \{0, \dots, n_T - 1\} \\ \mathbf{x}^{[0]} = x, \end{cases}$$

state $\mathbf{x}^{[k+1]} \in \mathbb{R}^{d_{k+1}}$ and **weights** $c^{[k]} \in \mathbb{R}^{d_{k+1}}$, $a^{[k]} \in \mathbb{R}^{d_k}$, $b^{[k]} \in \mathbb{R}$

- ▶ $\sigma \in C^{0,1}(\mathbb{R})$, typically $\sigma(x) = (x)_+$ or $\sigma(x) = \tanh(x)$
- ▶ *widths* d_k given

Residual neural networks (ResNets)

Let $d_k = d$ for all k .

Consider

$$\begin{cases} \mathbf{x}^{[k+1]} = \mathbf{x}^{[k]} + \Delta t c^{[k]} (a^{[k]} \cdot \mathbf{x}^{[k]} + b^{[k]})_+ & k \in \{0, \dots, n_T - 1\} \\ \mathbf{x}^{[0]} = \mathbf{x}, \end{cases}$$

with $c^{[k]}, a^{[k]} \in \mathbb{R}^d$ and $b^{[k]} \in \mathbb{R}$.

Deep residual learning for image recognition

[K He](#), [X Zhang](#), [S Ren](#), [J Sun](#) - Proceedings of the IEEE ..., 2016 - openaccess.thecvf.com

... as learning **residual** functions with ... **residual** networks are easier to optimize, and can gain accuracy from considerably increased depth. On the ImageNet dataset we evaluate **residual** ...

☆ Save [Cite](#) [Cited by 142319](#) [Related articles](#) [All 72 versions](#) [↔](#)

Neural ODEs

Natural to consider [E '17]:

$$\begin{cases} \dot{\mathbf{x}}(t) = c(t)(a(t) \cdot \mathbf{x}(t) + b(t))_+ & t \in (0, T), \\ \mathbf{x}(0) = x \end{cases}$$

And now $f_{\text{approx}}(x) := P_{\mathbf{X}}(T)$.

Useful in practice:

- ▶ Beyond Euler schemes [Chen et al. '18]
- ▶ Structure preserving schemes [Schönlieb et al. '20, '22]
- ▶ Beyond supervised learning (wait until the end of the talk)
- ▶ More compact form for analysis

Approximation theory

Before using the data, a "well-posedness" question can be asked.

Problem (Universal approximation)

Given $f \in \mathcal{H}$ and $\epsilon > 0$, find $(a_\epsilon, b_\epsilon, c_\epsilon) \in L^\infty((0, T); \mathbb{R}^{2d+1})$ such that

$$\|f_{\text{approx}, \epsilon} - f\|_{\mathcal{H}} \leq \epsilon$$

Approximation theory

Before using the data, a "well-posedness" question can be asked.

Problem (Universal approximation)

Given $f \in \mathcal{H}$ and $\epsilon > 0$, find $(a_\epsilon, b_\epsilon, c_\epsilon) \in L^\infty((0, T); \mathbb{R}^{2d+1})$ such that

$$\|f_{\text{approx}, \epsilon} - f\|_{\mathcal{H}} \leq \epsilon$$

- ▶ Feed-forward nets: [Cybenko '89], [Barron '93] ($n_t = 1$ and $\mathcal{H} = C^0([0, 1]^d)$), Pinkus '99 ($n_t \geq 1$)
- ▶ Neural ODEs: [Li, Lin, Shen '22], [Ruiz-Balet, Zuazua '22] ($\mathcal{H} = L^2((0, 1)^d; \mathbb{R}^m)$). ResNets are corollary as controls are piecewise constant
- ▶ Strategies generally **non-algorithmic** and suffer from **curse of dimensionality**

Learning is control

- ▶ To learn, we **can only use the data** \mathcal{D} .

Learning is control

- ▶ To learn, we **can only use the data** \mathcal{D} .
- ▶ So, we consider

Problem

Find controls $(a, b, c) \in L^\infty((0, T); \mathbb{R}^{2d+1})$ such that

$$P_{\mathbf{x}_i}(T) = f(x^{(i)}) \quad \forall i \in [n]$$

where

$$\begin{cases} \dot{\mathbf{x}}_i(t) = c(t)(a(t) \cdot \mathbf{x}_i(t) + b(t))_+ & t \in (0, T), \\ \mathbf{x}_i(0) = x^{(i)}, \end{cases} \quad (1)$$

and hope predictions of $f(x)$ are good if we take initial data points x outside \mathcal{D} (*generalization*).

Learning is control

- ▶ To learn, we **can only use the data** \mathcal{D} .
- ▶ So, we consider

Problem

Find controls $(a, b, c) \in L^\infty((0, T); \mathbb{R}^{2d+1})$ such that

$$P_{\mathbf{x}_i}(T) = f(x^{(i)}) \quad \forall i \in [n]$$

where

$$\begin{cases} \dot{\mathbf{x}}_i(t) = c(t)(a(t) \cdot \mathbf{x}_i(t) + b(t))_+ & t \in (0, T), \\ \mathbf{x}_i(0) = x^{(i)}, \end{cases} \quad (1)$$

and hope predictions of $f(x)$ are good if we take initial data points x outside \mathcal{D} (*generalization*).

- ▶ It's a simultaneous/ensemble control(lability) problem!
- ▶ Nonlinear control-state interaction is necessary

What is done in practice

Least squares, with penalty $\lambda > 0$:

$$\min_{\substack{\theta=(a,b,c) \\ (a,b) \in H^1((0,T);\mathbb{R}^{d+1}) \\ c \in L^2((0,T);\mathbb{R}^d) \\ x_i \text{ solves (1)}}} \underbrace{\frac{1}{n} \sum_{i \in [n]} \left| P_{x_i}(T) - f(x^{(i)}) \right|^2}_{:= \mathbf{E}(\mathbf{X}(t))} + \lambda \|\theta\|_{H^1 \times L^2}^2 \quad (2)$$

- ▶ **Empirical risk minimization:** $\mathbf{E}(\cdot)$ is the empirical risk.
- ▶ H^1 suffices for compactness – if $p \in [1, \infty)$, $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ is s.t. $\varphi \circ u_n \rightharpoonup \varphi \circ u^*$ in $L^p(0, 1)$ for any $u_n \rightharpoonup u^*$ in $L^p(0, 1)$, then φ is affine!
- ▶ Can go way beyond squared Euclidean distance (even non-distances such as cross-entropy for classification)

github.com/borjanG/dynamical.systems

$$T = 5, n_T = 16, n = 3000, \lambda = 0.01$$

Optimal control over long time

- ▶ In practice n_T can be large (*deep learning*)

Optimal control over long time

- ▶ In practice n_T can be large (*deep learning*)
- ▶ But $\Delta t = \frac{T}{n_T}$, so n_T large means T large

Optimal control over long time

- ▶ In practice n_T can be large (*deep learning*)
- ▶ But $\Delta t = \frac{T}{n_T}$, so n_T large means T large

Question

For global minimizer θ_T for (2) and $\mathbf{X}_T \in C^0([0, T]; \mathbb{R}^{d \times n})$ matrix with unique solutions \mathbf{x}_i to (1) as columns, **what happens when $T \rightarrow \infty$?**

Optimal control over long time

- ▶ In practice n_T can be large (*deep learning*)
- ▶ But $\Delta t = \frac{T}{n_T}$, so n_T large means T large

Question

For global minimizer θ_T for (2) and $\mathbf{X}_T \in C^0([0, T]; \mathbb{R}^{d \times n})$ matrix with unique solutions \mathbf{x}_i to (1) as columns, **what happens when $T \rightarrow \infty$?**

Definition (Interpolation)

(1) interpolates \mathcal{D} if $\exists \theta \in H^1((0, 1); \mathbb{R}^{d+1}) \times L^2((0, 1); \mathbb{S}^{d-1})$ such that

$$P_{\mathbf{x}_i}(1) = f(x^{(i)}) \quad \forall i \in [n]$$

where $\mathbf{x}_i \in C^0([0, 1]; \mathbb{R}^d)$ solves (1) with control θ . (I.e., $\mathbf{E}(\mathbf{X}(1)) = 0$.)

Theorem [Esteve-Yagüe, G., Pighin, Zuazua, '22b]

Fix $\lambda > 0$; c in (2) minimized over $L^2((0, T); \mathbb{S}^{d-1})$; (1) interpolates \mathcal{D} . For $T \geq 1$, any global minimizer θ_T for (2) and $\mathbf{X}_T \in C^0([0, T]; \mathbb{R}^{d \times n})$ matrix with columns solutions to (1) satisfy:

Theorem [Esteve-Yagüe, G., Pighin, Zuazua, '22b]

Fix $\lambda > 0$; c in (2) minimized over $L^2((0, T); \mathbb{S}^{d-1})$; (1) interpolates \mathcal{D} . For $T \geq 1$, any global minimizer θ_T for (2) and $\mathbf{X}_T \in C^0([0, T]; \mathbb{R}^{d \times n})$ matrix with columns solutions to (1) satisfy:

1. $\exists C(\mathcal{D}, \lambda) > 0$,

$$\mathbf{E}(\mathbf{X}_T(T)) = \frac{1}{n} \sum_{i \in [n]} \left| P_{\mathbf{x}_i}(T) - f(x^{(i)}) \right|^2 \leq \frac{C}{T}$$

Theorem [Esteve-Yagüe, G., Pighin, Zuazua, '22b]

Fix $\lambda > 0$; c in (2) minimized over $L^2((0, T); \mathbb{S}^{d-1})$; (1) interpolates \mathcal{D} . For $T \geq 1$, any global minimizer θ_T for (2) and $\mathbf{X}_T \in C^0([0, T]; \mathbb{R}^{d \times n})$ matrix with columns solutions to (1) satisfy:

1. $\exists C(\mathcal{D}, \lambda) > 0$,

$$\mathbf{E}(\mathbf{X}_T(T)) = \frac{1}{n} \sum_{i \in [n]} \left| P_{\mathbf{X}_i}(T) - f(x^{(i)}) \right|^2 \leq \frac{C}{T}$$

2. $\mathbf{X}_{T_k}(T_k) \rightarrow \mathbf{X}^*$ for some subsequence $T_k > 0$, $T_k \rightarrow \infty$ ($k \rightarrow \infty$) and $\mathbf{X}^* \in \mathbb{R}^{d \times n}$ with $\mathbf{E}(\mathbf{X}^*) = 0$

Theorem [Esteve-Yagüe, G., Pighin, Zuazua, '22b]

Fix $\lambda > 0$; c in (2) minimized over $L^2((0, T); \mathbb{S}^{d-1})$; (1) interpolates \mathcal{D} . For $T \geq 1$, any global minimizer θ_T for (2) and $\mathbf{X}_T \in C^0([0, T]; \mathbb{R}^{d \times n})$ matrix with columns solutions to (1) satisfy:

1. $\exists C(\mathcal{D}, \lambda) > 0$,

$$\mathbf{E}(\mathbf{X}_T(T)) = \frac{1}{n} \sum_{i \in [n]} \left| P_{\mathbf{X}_i}(T) - f(x^{(i)}) \right|^2 \leq \frac{C}{T}$$

2. $\mathbf{X}_{T_k}(T_k) \rightarrow \mathbf{X}^*$ for some subsequence $T_k > 0$, $T_k \rightarrow \infty$ ($k \rightarrow \infty$) and $\mathbf{X}^* \in \mathbb{R}^{d \times n}$ with $\mathbf{E}(\mathbf{X}^*) = 0$
3. Set $\theta_k(t) := (T_k a_{T_k}(tT_k), T_k b_{T_k}(tT_k), c_{T_k}(tT_k))$ for $t \in [0, 1]$. Then $\theta_k \rightarrow \theta^*$ strongly in $H^1 \times L^2$ where θ^* is some solution to

$$\inf_{\substack{\theta=(a,b,c) \\ (a,b) \in H^1((0,1); \mathbb{R}^{d+1}) \\ c \in L^2((0,1); \mathbb{S}^{d-1}) \\ \mathbf{E}(\mathbf{X}(1))=0}} \|\theta\|_{H^1 \times L^2}^2$$

1. $\theta_1 = (a_1, b_1, c_1) \in H^1((0, 1); \mathbb{R}^{d+1}) \times L^2((0, T); \mathbb{S}^{d-1})$ yields x_i^1 solution to (1) on $[0, 1]$. Then

$$\theta_T(\cdot) = \left(\frac{1}{T} a_1 \left(\frac{\cdot}{T} \right), \frac{1}{T} b_1 \left(\frac{\cdot}{T} \right), c_1 \left(\frac{\cdot}{T} \right) \right)$$

defined on $[0, T]$, yields solution $x_i^T(\cdot) \equiv x_i^1\left(\frac{\cdot}{T}\right)$ to (1)

1. $\theta_1 = (a_1, b_1, c_1) \in H^1((0, 1); \mathbb{R}^{d+1}) \times L^2((0, T); \mathbb{S}^{d-1})$ yields x_i^1 solution to (1) on $[0, 1]$. Then

$$\theta_T(\cdot) = \left(\frac{1}{T} a_1 \left(\frac{\cdot}{T} \right), \frac{1}{T} b_1 \left(\frac{\cdot}{T} \right), c_1 \left(\frac{\cdot}{T} \right) \right)$$

defined on $[0, T]$, yields solution $x_i^T(\cdot) \equiv x_i^1(\frac{\cdot}{T})$ to (1)

2. In turn,

$$\begin{aligned} \mathbf{E}(\mathbf{X}_T(T)) + \lambda \int_0^T |\theta_T(t)|^2 dt \\ = \mathbf{E}(\mathbf{X}_1(1)) + \frac{\lambda}{T} \int_0^1 |(a_1(s), b_1(s))|^2 ds + \lambda. \end{aligned}$$

1. $\theta_1 = (a_1, b_1, c_1) \in H^1((0, 1); \mathbb{R}^{d+1}) \times L^2((0, T); \mathbb{S}^{d-1})$ yields x_i^1 solution to (1) on $[0, 1]$. Then

$$\theta_T(\cdot) = \left(\frac{1}{T} a_1 \left(\frac{\cdot}{T} \right), \frac{1}{T} b_1 \left(\frac{\cdot}{T} \right), c_1 \left(\frac{\cdot}{T} \right) \right)$$

defined on $[0, T]$, yields solution $x_i^T(\cdot) \equiv x_i^1(\frac{\cdot}{T})$ to (1)

2. In turn,

$$\begin{aligned} \mathbf{E}(\mathbf{X}_T(T)) + \lambda \int_0^T |\theta_T(t)|^2 dt \\ = \mathbf{E}(\mathbf{X}_1(1)) + \frac{\lambda}{T} \int_0^1 |(a_1(s), b_1(s))|^2 ds + \lambda. \end{aligned}$$

3. Take interpolation control (on $(0, 1)$), stretch it out to $(0, T)$, and compare with θ_T .

1. $\theta_1 = (a_1, b_1, c_1) \in H^1((0, 1); \mathbb{R}^{d+1}) \times L^2((0, T); \mathbb{S}^{d-1})$ yields x_i^1 solution to (1) on $[0, 1]$. Then

$$\theta_T(\cdot) = \left(\frac{1}{T} a_1 \left(\frac{\cdot}{T} \right), \frac{1}{T} b_1 \left(\frac{\cdot}{T} \right), c_1 \left(\frac{\cdot}{T} \right) \right)$$

defined on $[0, T]$, yields solution $x_i^T(\cdot) \equiv x_i^1(\frac{\cdot}{T})$ to (1)

2. In turn,

$$\begin{aligned} \mathbf{E}(\mathbf{X}_T(T)) + \lambda \int_0^T |\theta_T(t)|^2 dt \\ = \mathbf{E}(\mathbf{X}_1(1)) + \frac{\lambda}{T} \int_0^1 |(a_1(s), b_1(s))|^2 ds + \lambda. \end{aligned}$$

3. Take interpolation control (on $(0, 1)$), stretch it out to $(0, T)$, and compare with θ_T .

Corollary

In this setting, $T \rightarrow \infty$ is equivalent to $\lambda \rightarrow 0$.

Interpolation, Controllability

1. *Combinatorics*. For (1): [Li, Lin, Shen '22], [Ruiz-Balet, Zuazua '22]. Distinct targets if $d = m$.

2. *Lie algebra*. For

$$\dot{x}_i(t) = \theta(t)\sigma(x_i(t)), \quad (3)$$

with $\sigma \in C^{0,1} \cap C^1(\mathbb{R})$ element-wise: [Agrachev, Sarychev '22]

Interpolation, Controllability

1. *Combinatorics*. For (1): [Li, Lin, Shen '22], [Ruiz-Balet, Zuazua '22]. Distinct targets if $d = m$.

2. *Lie algebra*. For

$$\dot{\mathbf{x}}_i(t) = \theta(t)\sigma(\mathbf{x}_i(t)), \quad (3)$$

with $\sigma \in C^{0,1} \cap C^1(\mathbb{R})$ element-wise: [Agrachev, Sarychev '22]

A digression - little homotopy method inspired by [Coron-Trélat '04]:

Proposition [Esteve-Yagüe, G., Pighin, Zuazua '22b]

Suppose $d \geq n$. Fix $\mathbf{X}^1 \in \mathbb{R}^{d \times n}$ with

$$\text{span}\{\sigma(\mathbf{x}_1^1), \dots, \sigma(\mathbf{x}_n^1)\} = \mathbb{R}^d$$

Then $\exists r, C > 0$ such that $\forall \mathbf{X}^0 \in B_r(\mathbf{X}^1)$, $\exists \theta \in L^\infty((0, 1); \mathbb{R}^{d \times d})$ for which the solutions \mathbf{x}_i to (3) with $\mathbf{x}_i(0) = \mathbf{x}_i^0$ satisfy $\mathbf{x}_i(1) = \mathbf{x}_i^1 \forall i \in [n]$.

Moreover

$$\|\theta\|_{L^\infty} \leq \frac{C}{T} |\mathbf{X}^1 - \mathbf{X}^0|.$$

Generalization: a statistical approach

Focus on dynamics (3).

Generalization: a statistical approach

Focus on dynamics (3).

- ▶ Look at $\{x^{(i)}, y^{(i)}\}_{i \in [n]} \subset \mathbb{R}^d \times \mathbb{R}^m$ as i.i.d. samples from unknown joint law $\mu \in \mathcal{P}_c(\mathbb{R}^d \times \mathbb{R}^m)$. Then $f(x) := \mathbb{E}(y|x)$ which minimizes $\mathbb{E}_{(x,y) \sim \mu} |f(x) - y|^2$ over all functions f .

Generalization: a statistical approach

Focus on dynamics (3).

- ▶ Look at $\{x^{(i)}, y^{(i)}\}_{i \in [n]} \subset \mathbb{R}^d \times \mathbb{R}^m$ as i.i.d. samples from unknown joint law $\mu \in \mathcal{P}_c(\mathbb{R}^d \times \mathbb{R}^m)$. Then $f(x) := \mathbb{E}(y|x)$ which minimizes $\mathbb{E}_{(x,y) \sim \mu} |f(x) - y|^2$ over all functions f .
- ▶ Associated to (2): **population risk minimization**

$$\min_{\substack{\theta \in L^2((0,T); \mathbb{R}^{d \times d}) \\ \mathbf{x}^\theta \text{ solves (3)} \\ \mathbf{x}^\theta(0) = x}} \mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{x}^\theta}(T) - y|^2 + \lambda \int_0^T |\theta(t)|^2 dt \quad (4)$$

Generalization: a statistical approach

Focus on dynamics (3).

- ▶ Look at $\{x^{(i)}, y^{(i)}\}_{i \in [n]} \subset \mathbb{R}^d \times \mathbb{R}^m$ as i.i.d. samples from unknown joint law $\mu \in \mathcal{P}_c(\mathbb{R}^d \times \mathbb{R}^m)$. Then $f(x) := \mathbb{E}(y|x)$ which minimizes $\mathbb{E}_{(x,y) \sim \mu} |f(x) - y|^2$ over all functions f .
- ▶ Associated to (2): **population risk minimization**

$$\min_{\substack{\theta \in L^2((0,T); \mathbb{R}^{d \times d}) \\ \mathbf{x}^\theta \text{ solves (3)} \\ \mathbf{x}^\theta(0) = x}} \mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{X}^\theta}(T) - y|^2 + \lambda \int_0^T |\theta(t)|^2 dt \quad (4)$$

- ▶ **Generalization:** θ_n minimizer of J_n in (2), and θ^* of J in (4), then $\exists \alpha > 0$:

$$\mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{X}^{\theta_n}}(T) - y|^2 - \mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{X}^{\theta^*}}(T) - y|^2 = \mathcal{O}\left(\frac{1}{n^\alpha}\right)$$

What is known

1. [E, Han, Li '19]:

- ▶ Pontryagin Maximum Principle for both (2) and (4);
- ▶ Hamiltonian $\theta \mapsto H(x, p, \theta) = p \cdot \theta \sigma(x) + \lambda |\theta|^2$ strongly concave for any (x, p)
- ▶ given θ^* , with high probability $\exists \theta_n$ critical point of Hamiltonian for (2) such that

$$\mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{X}}^{\theta_n}(T) - y|^2 - \mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{X}}^{\theta^*}(T) - y|^2 \leq \frac{C(d)}{n^{\frac{1}{2} - \epsilon}}$$

with high probability, for any $\epsilon > 0$.

- ▶ Ensuring that θ_n is global minimizer: true when $T \ll 1$, so $\lambda \gg 1$!

What is known

1. [E, Han, Li '19]:

- ▶ Pontryagin Maximum Principle for both (2) and (4);
- ▶ Hamiltonian $\theta \mapsto H(x, p, \theta) = p \cdot \theta \sigma(x) + \lambda |\theta|^2$ strongly concave for any (x, p)
- ▶ given θ^* , with high probability $\exists \theta_n$ critical point of Hamiltonian for (2) such that

$$\mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{X}}^{\theta_n}(T) - y|^2 - \mathbb{E}_{(x,y) \sim \mu} |P_{\mathbf{X}}^{\theta^*}(T) - y|^2 \leq \frac{C(d)}{n^{\frac{1}{2} - \epsilon}}$$

with high probability, for any $\epsilon > 0$.

- ▶ Ensuring that θ_n is global minimizer: true when $T \ll 1$, so $\lambda \gg 1$!

2. [Bonnet et al. '22]:

- ▶ For $\lambda \gg 1$, J_n strongly convex on any L^2 ball
- ▶ Mean-field PMP ... rate $\mathcal{O}\left(\frac{1}{n^{\frac{1}{d}}}\right)$

An observation with C. Letrouit and P. Rigollet

1. Strong convexity on any $B \subset L^2((0, T); \mathbb{R}^{d \times d})$

$$\|\theta_1 - \theta_2\|_{L^2}^2 \lesssim |\nabla J_n(\theta_1) - \nabla J_n(\theta_2)|, \quad \forall \theta_1, \theta_2 \in B.$$

An observation with C. Letrouit and P. Rigollet

1. Strong convexity on any $B \subset L^2((0, T); \mathbb{R}^{d \times d})$

$$\|\theta_1 - \theta_2\|_{L^2}^2 \lesssim |\nabla J_n(\theta_1) - \nabla J_n(\theta_2)|, \quad \forall \theta_1, \theta_2 \in B.$$

2. Then

$$\begin{aligned} \|\theta_n - \theta^*\|_{L^2}^2 &\lesssim |\nabla J_n(\theta_n) - \nabla J_n(\theta^*)| \\ &= |\nabla J(\theta^*) - \nabla J_n(\theta^*)| \end{aligned}$$

$$= \left| \mathbb{E}_{(x,y) \sim \mu} \nabla_{\theta} \ell(\mathbf{x}^{\theta^*}(T), y) - \frac{1}{n} \sum_{i \in [n]} \nabla_{\theta} \ell(\mathbf{x}_i^{\theta^*}(T), y^{(i)}) \right|$$

$$|\mathbb{P}_{\mathbf{x}_i^{\theta^*}(T)} - y^{(i)}|^2$$



An observation with C. Letrouit and P. Rigollet

1. Strong convexity on any $B \subset L^2((0, T); \mathbb{R}^{d \times d})$

$$\|\theta_1 - \theta_2\|_{L^2}^2 \lesssim |\nabla J_n(\theta_1) - \nabla J_n(\theta_2)|, \quad \forall \theta_1, \theta_2 \in B.$$

2. Then

$$\begin{aligned} \|\theta_n - \theta^*\|_{L^2}^2 &\lesssim |\nabla J_n(\theta_n) - \nabla J_n(\theta^*)| \\ &= |\nabla J(\theta^*) - \nabla J_n(\theta^*)| \\ &= \left| \mathbb{E}_{(x,y) \sim \mu} \nabla_{\theta} \ell(\mathbf{x}^{\theta^*}(T), y) - \frac{1}{n} \sum_{i \in [n]} \nabla_{\theta} \ell(\mathbf{x}_i^{\theta^*}(T), y^{(i)}) \right| \end{aligned}$$

3. θ^* is fixed, not random. Sum of bounded, independent random variables, compared with expectation \Rightarrow **concentration of measure** (Hoeffding inequality): $\forall \delta > 0, \exists \kappa > 0$ such that $\forall n \geq 1$,

$$\mathbb{P}(\|\theta_n - \theta^*\|_{L^2}^2 \leq \kappa/\sqrt{n}) \geq 1 - \delta.$$

Bound $J(\theta_n) - J(\theta^*)$ from above and get $\mathcal{O}(\frac{1}{\sqrt{n}})$ rate.

Improving $1/T$

We consider

$$\inf_{\substack{c \in L^2((0,T); \mathbb{R}^d) \\ x_i(\cdot) \text{ solves (1)}}} \int_0^T \frac{1}{n} \sum_{i \in [n]} \left| P_{x_i}(t) - f(x^{(i)}) \right|^2 dt + \int_0^T |c(t)|^2 dt \quad (5)$$

Controls $a \in L^\infty(\mathbb{R}_+; \mathbb{R}^d)$ and $b \in L^\infty(\mathbb{R}_+)$ assumed fixed in (1) (need L^2 -penalties and compactness simultaneously)

Can also consider dynamics as (3), or (a, b) can be optimized over \mathbb{S}^d .

Two assumptions

Recall, for $\mathbf{X} = [x_1 \cdots x_n] \in \mathbb{R}^{d \times n}$:

$$\mathbf{E}(\mathbf{X}) := \frac{1}{n} \sum_{i \in [n]} \left| P_{x_i} - f(x^{(i)}) \right|^2.$$

Assumption 1.

Set $\mathcal{Z} := \{\mathbf{Z} \in \mathbb{R}^{d \times n} : \mathbf{E}(\mathbf{Z}) = 0\}$. Then $P \in \mathbb{R}^{m \times d}$ is such that

$$\kappa_1 \text{dist}(\mathbf{X}, \mathcal{Z})^2 \leq \mathbf{E}(\mathbf{X}) \leq \kappa_2 \text{dist}(\mathbf{X}, \mathcal{Z})^2$$

for some $\kappa_2 \geq \kappa_1 > 0$, and $\forall \mathbf{X} \in \mathbb{R}^{d \times n}$.

- ▶ Lower bound: [global Lojasiewicz inequality](#) for analytic functions

Assumption 2.

Fix $\mathbf{X}^0 = [x_1^0 \cdots x_n^0] \in \mathbb{R}^{d \times n}$. We assume $\exists c \in L^2((0, 1); \mathbb{R}^d)$ such that the matrix $\mathbf{X} \in C^0([0, 1]; \mathbb{R}^{d \times n})$ with columns $x_i(\cdot)$ solutions to (1) with $x_i(0) = x_i^0$, satisfies $\mathbf{X}(1) \in \mathcal{Z}$.

Moreover, $\exists C(n) > 0$,

$$\int_0^1 |c(t)|^2 dt \leq C(n) \text{dist}(\mathbf{X}^0, \mathcal{Z})^2.$$

Assumption 2.

Fix $\mathbf{X}^0 = [x_1^0 \cdots x_n^0] \in \mathbb{R}^{d \times n}$. We assume $\exists c \in L^2((0, 1); \mathbb{R}^d)$ such that the matrix $\mathbf{X} \in C^0([0, 1]; \mathbb{R}^{d \times n})$ with columns $x_i(\cdot)$ solutions to (1) with $x_i(0) = x_i^0$, satisfies $\mathbf{X}(1) \in \mathcal{Z}$.

Moreover, $\exists C(n) > 0$,

$$\int_0^1 |c(t)|^2 dt \leq C(n) \text{dist}(\mathbf{X}^0, \mathcal{Z})^2.$$

► When $d > m$, then $P \in \mathbb{R}^{d \times n}$ is generically surjective and

$$\mathcal{Z} = \left\{ [z_1 \cdots z_n] \in \mathbb{R}^{d \times n} : z_i \in P^{-1} \{f(x^{(i)})\} \right\}$$

So $x_i(1) \in P^{-1} \{f(x^{(i)})\}$ for all $i \in [n]$ and

$$\int_0^1 |c(t)|^2 dt \leq C(n) \inf_{\substack{[z_1 \cdots z_n] \in \mathbb{R}^{d \times n} \\ z_i \in P^{-1} \{f(x^{(i)})\}}} \sum_{i \in [n]} |x_i^0 - z_i|^2.$$

Theorem [Esteve-Yagüe, G., Pighin, Zuazua, '22a]

Suppose $m = d$ and $P = \text{Id}$. Then $\exists T_*, \omega > 0, C \geq 1$ such that for $T \geq T_*$, any global minimizer $c_T \in L^2((0, T); \mathbb{R}^d)$ to (5) and x_i^T solution to (1) satisfy

$$\sum_{i \in [n]} \left| x_i^T(t) - f(x^{(i)}) \right|^2 + |c_T(t)|^2 \leq \left(C \sum_{i \in [n]} \left| x^{(i)} - f(x^{(i)}) \right|^2 \right) e^{-\omega t}$$

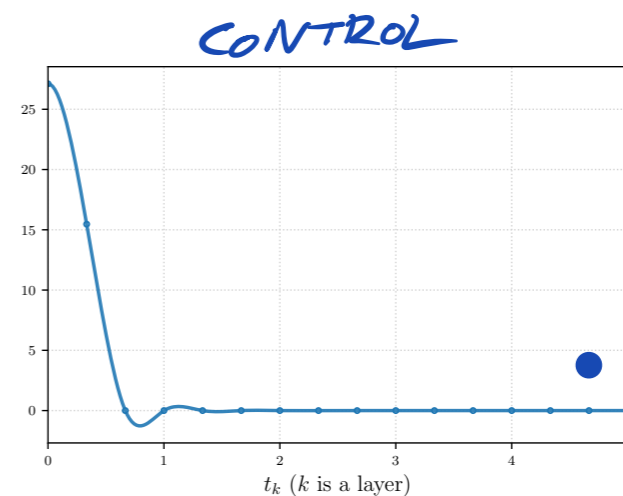
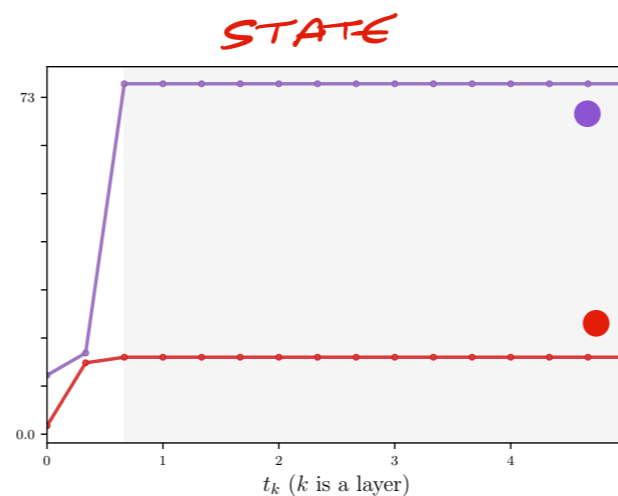
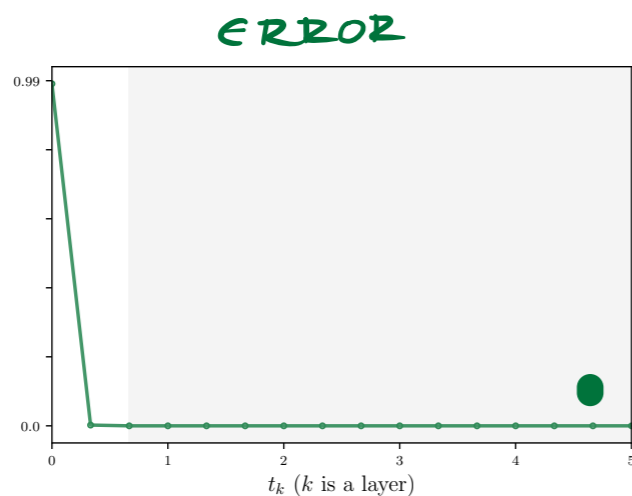
$\forall t \in [0, T]$.

Theorem [Esteve-Yagüe, G., Pighin, Zuazua, '22b]

Replace $(\cdot)_+$ by $\sigma \in L^\infty(\mathbb{R})$ in (1) and let $d > m$. Then $\exists T_*, \omega > 0$, $C \geq 1$ such that for $T \geq T_*$, any global minimizer $c_T \in L^2((0, T); \mathbb{R}^d)$ to (5) and x_i^T solution to (1) satisfy

$$\sum_{i \in [n]} \left| P_{x_i^T}(t) - f(x^{(i)}) \right|^2 + \inf_{\substack{[z_1 \dots z_n] \in \mathbb{R}^{d \times n} \\ z_i \in P^{-1}\{f(x^{(i)})\}}} \sum_{i \in [n]} \left| x_i^T(t) - z_i \right|^2 + |c_T(t)|^2 \leq \left(C \inf_{\substack{[z_1 \dots z_n] \in \mathbb{R}^{d \times n} \\ z_i \in P^{-1}\{f(x^{(i)})\}}} \sum_{i \in [n]} \left| x^{(i)} - z_i \right|^2 \right) e^{-\omega t}$$

$\forall t \in [0, T]$.



It's faster

Takeaway: when possible, proceed in *model predictive control manner*: start with small T , evaluate error, and proceed by increasing T adaptively.

Important tool

Focus on $d = m$, $P = \text{Id}$.

Lemma

$\exists C_1 > 0$ independent of T , $\forall c \in L^2$ and x_i solution to (1):

$$\begin{aligned} \sup_{t \in [0, T]} \sum_{i \in [n]} \left| x_i(t) - f(x^{(i)}) \right|^2 &\leq C_1 \left(\sum_{i \in [n]} \left| x^{(i)} - f(x^{(i)}) \right|^2 \right. \\ &\quad + \int_0^T \sum_{i \in [n]} \left| x_i(t) - f(x^{(i)}) \right|^2 dt \\ &\quad \left. + \int_0^T |c(t)|^2 dt \right) \end{aligned}$$

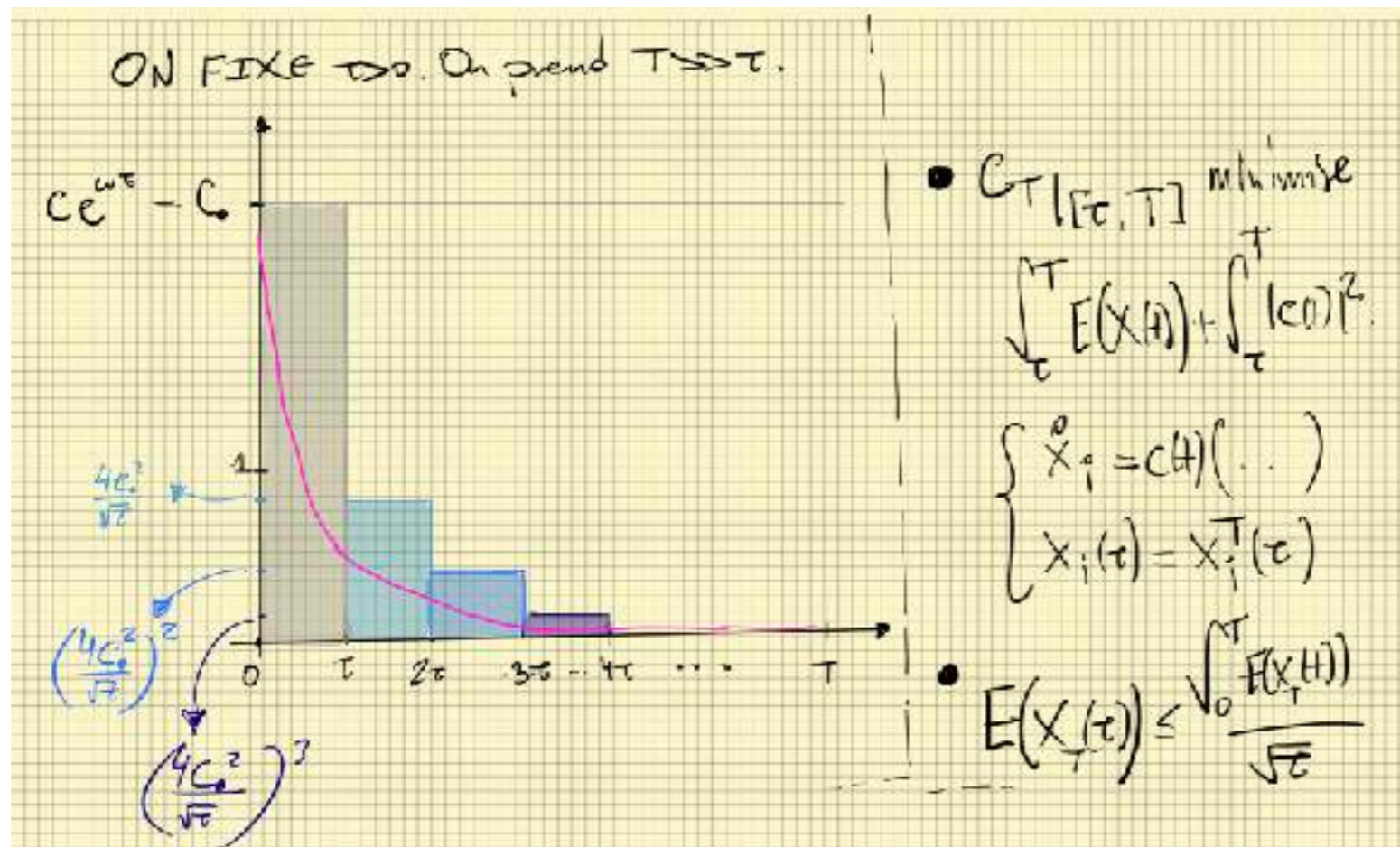
Ingredients

1. $c_{\text{aux}}(t) := c_1(t)1_{[0,1]}(t)$, where c_1 ensures controllability. As c_T is optimal and c_{aux} is not, $\exists C_2 > 0$ independent of T :

$$\int_0^T \sum_{i \in [n]} \left| x_i^T(t) - f(x^{(i)}) \right|^2 dt + \int_0^T |c_T(t)|^2 dt \leq C_2 \sum_{i \in [n]} \left| x^{(i)} - f(x^{(i)}) \right|^2$$

Lemma yields pointwise bound uniform in T .

- 2.



Extensions

- ▶ Using this method, we **can't have exponential decay with BV -penalty** for (a, b) as norm tracks singularities unlike L^2 . We can at most get decay of time-averages of the error and controls.
- ▶ **Results are more general.** Per [Esteve-Yagüe, G., Pighin, Zuazua '22a; G., Zuazua '22]: controllable PDE

$$y_t(t, x) - Ay(t, x) + Bu(t, x) = f(y(t, x))$$

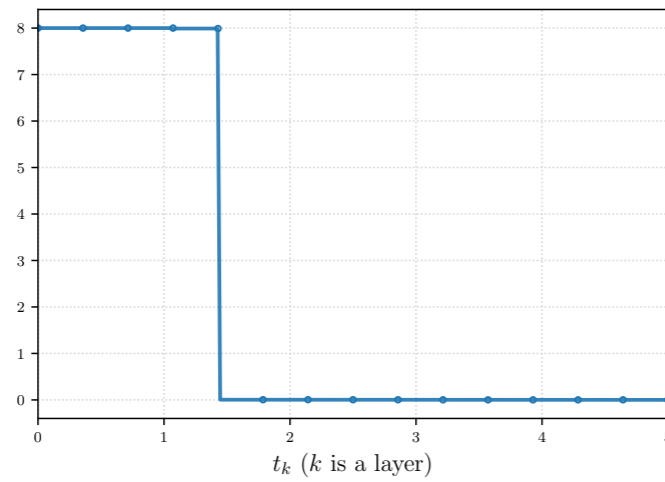
with Lipschitz (possibly non-smooth) nonlinearity f and cost

$$\phi(y(T)) + \int_0^T |y(t) - \bar{y}|_{\mathcal{H}}^2 dt + \int_0^T |u(t)|_{\mathcal{U}}^2 dt$$

where \bar{y} is any steady state \Rightarrow exponential turnpike without smoothness or smallness assumptions!

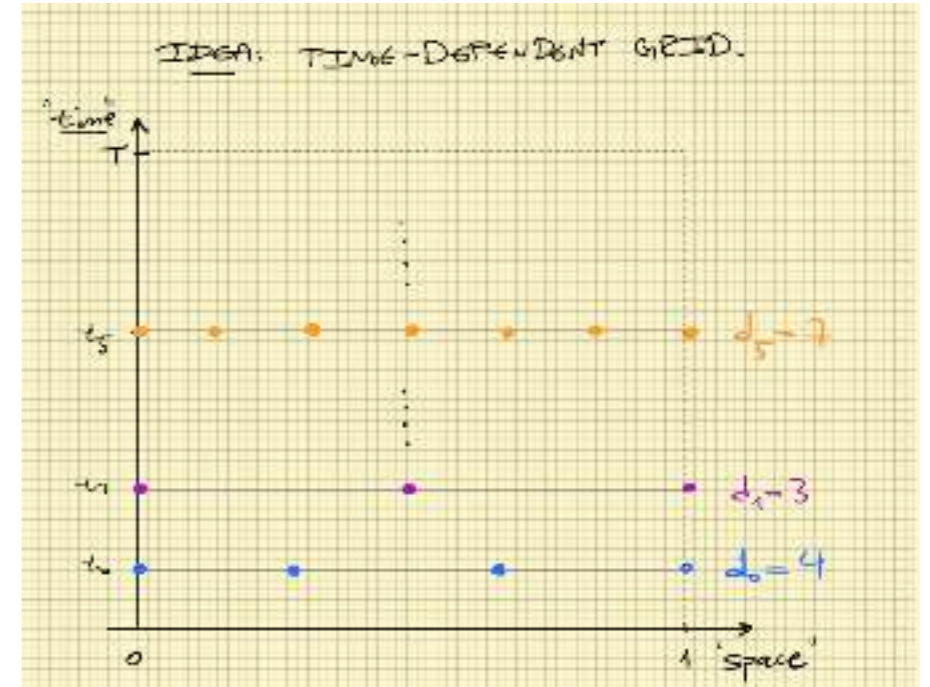
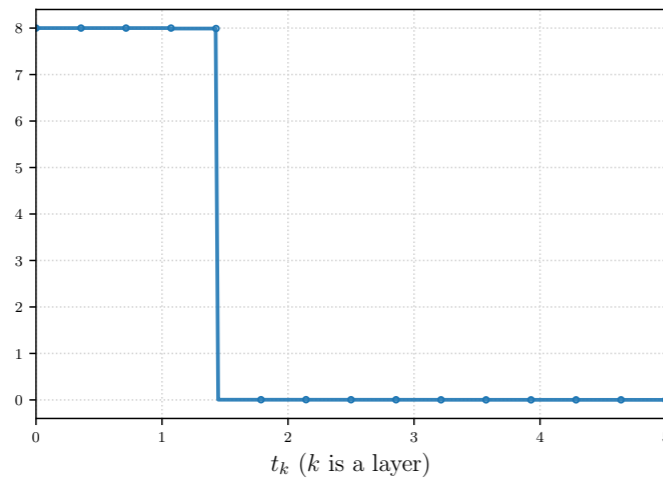
Further comments

1. $L^1(0, T)$ -penalties for (3) [Esteve-Yagüe, G. '22]:



Further comments

1. $L^1(0, T)$ -penalties for (3) [Esteve-Yagüe, G. '22]:



2. Variable-width ResNets:

$$\mathbf{x}^{[k+1]} = \Pi^{[k]} \mathbf{x}^{[k]} + c^{[k]} \sigma(a^{[k]} \mathbf{x}^{[k]})$$

where $\Pi^{[k]} : \mathbb{R}^{d_k} \rightarrow \mathbb{R}^{d_{k+1}}$, $a^{[k]} \in \mathbb{R}^{d_{k+1} \times d_k}$, $c^{[k]} \in \mathbb{R}^{d_{k+1} \times d_{k+1}}$. So,

$$\partial_t \mathbf{x}(t, z) = \int_0^1 c(t, z, \zeta) \sigma(a(t, z, \zeta) \mathbf{x}(t, \zeta)) d\zeta \quad (0, T) \times (0, 1)$$

Helpful for structured controls (convolutional neural networks).

Outlook

1. Control

- ▶ Exponential decay/turnpike with BV -penalty for (a, b) ?
- ▶ Using control: are feedback controls for $n \gg 1$ trajectories possible, useful?
- ▶ Extrapolating to control: can we get robustness using the lens of many data and statistics?

2. Unsupervised learning/**generative modeling** with normalizing flows (E. Vanden-Eijnden et al.).

- ▶ NF: diffeomorphism $\mathcal{T} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ optimized to transport $\{z^{(i)}\}_{i \in [n]} \subset \mathbb{R}^d$ samples from a known law ρ_0 (Gaussian with unit variance) to unknown target law ρ_1 of which we know samples $\{x^{(i)}\}_{i \in [n]} \subset \mathbb{R}^d$.
- ▶ Parametrizing \mathcal{T} by the flow of a neural ODE, and then solving an optimal control problem (KL divergence) is quite practical.

Merci pour votre attention!

Bibliographie:

1. [Esteve-Yagüe, G., Pighin, Zuazua '22a]:
<https://arxiv.org/abs/2011.11091>, Nonlinearity
2. [Esteve-Yagüe, G., Pighin, Zuazua '22b]:
<https://arxiv.org/abs/2008.02491>
3. [G., Zuazua '22]:
<https://arxiv.org/abs/2202.04097>, Acta Numerica
4. [Esteve-Yagüe, G. '22]:
<https://arxiv.org/abs/2102.13566>, submitted