

Optimal Control and Deep Learning – Turnpikes, Dissipativity and Early Exits

Timm Faulwasser

Institute of Control Systems, Hamburg University of Technology

Joint work with

Jens Püttschneider (TUHH), Simon Heilig & Asja Fischer (RU Bochum),

Stefan Streif (TU Chemnitz) and Jens-Arne Hempel (BA Sachsen)

Workshop on Structured Learning @ ALU Freiburg

tim.faulwasser@ieee.org

Machine learning → Control ?

Machine learning

- Regression
- Support vector machines
- Gaussian Processes & kernel methods
- Neural networks
- ...



Systems & control

- System identification
- Controller approximation
- Data-driven and adaptive control
- Fault detection
- State estimation
- ...

IEEE TRANSACTIONS ON NEURAL NETWORKS, VOL. 1, NO. 1, MARCH 1990

Identification and Control of Dynamical Systems Using Neural Networks

KUMPATI S. NARENDRA FELLOW, IEEE, AND KANNAN PARTHASARATHY



Pergamon

0005-1098(95)00044-5

Automatica, Vol. 31, No. 10, pp. 1443-1451, 1995
Copyright © 1995 Elsevier Science Ltd
Printed in Great Britain. All rights reserved
0005-1098/95 \$9.50 + 0.00

Brief Paper

A Receding-horizon Regulator for Nonlinear Systems and a Neural Approximation*

T. PARISINI[†] and R. ZOPPOLI[†]





Automatica

Volume 172, February 2025, 112006



Meta-learning for model-reference data- driven control ☆

Riccardo Busetto ^a  , Valentina Breschi ^b, Simone Formentin ^a

Control → Machine learning ?

Complete Controllability of
Continuous-Time Recurrent Neural Networks*

Eduardo Sontag
Dept. of Mathematics, Rutgers University
New Brunswick, NJ 08903
sontag@control.rutgers.edu

Héctor Sussmann
Dept. of Mathematics, Rutgers University
New Brunswick, NJ 08903
sussmann@hamilton.rutgers.edu

Math. Control Signals Systems (1989) 2: 303–314

Systems & Control Letters 22 (1994) 235–244
North-Holland

State observability in recurrent
neural networks*

Francesca Albertini** and Eduardo D. Sontag

Mathematics of Control,
Signals, and Systems

© 1989 Springer-Verlag New York Inc.

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†



IEEE TRANSACTIONS ON AUTOMATIC CONTROL, VOL. 68, NO. 5, MAY 2023

Universal Approximation Power of Deep Residual Neural Networks Through the Lens of Control

Paulo Tabuada , *Fellow, IEEE*, and Bahman Ghahsifard , *Member, IEEE*

Systems & control

- Analysis
 - Stability, controllability, ...
- Synthesis
 - Feedback & feedforward controls
- Computational concepts
- ...

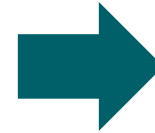
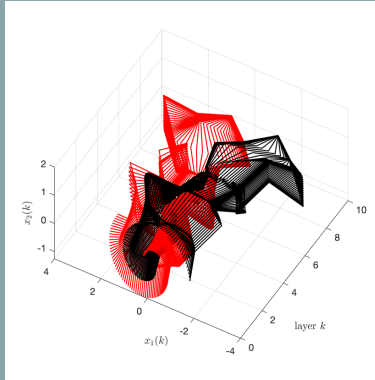


Machine learning



Overview

Optimal Control & Deep Learning?

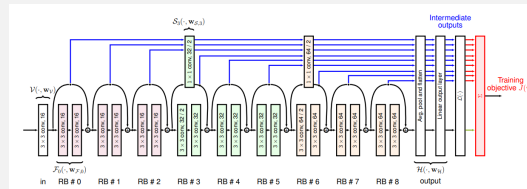
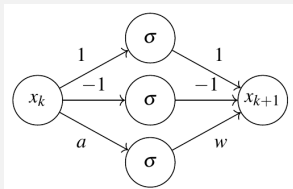


Dissipativity & ResNet Training?

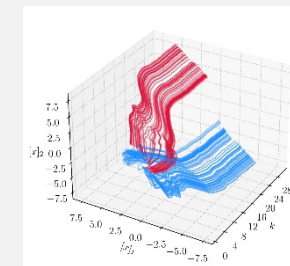
$$S(f(\mathbf{x}, u)) - S(\mathbf{x}) \leq \tilde{\ell}_f(\mathbf{x}, u) + r\|u\|_2^2 - \text{dist}((\mathbf{x}, u), \mathbb{Z}^*)$$



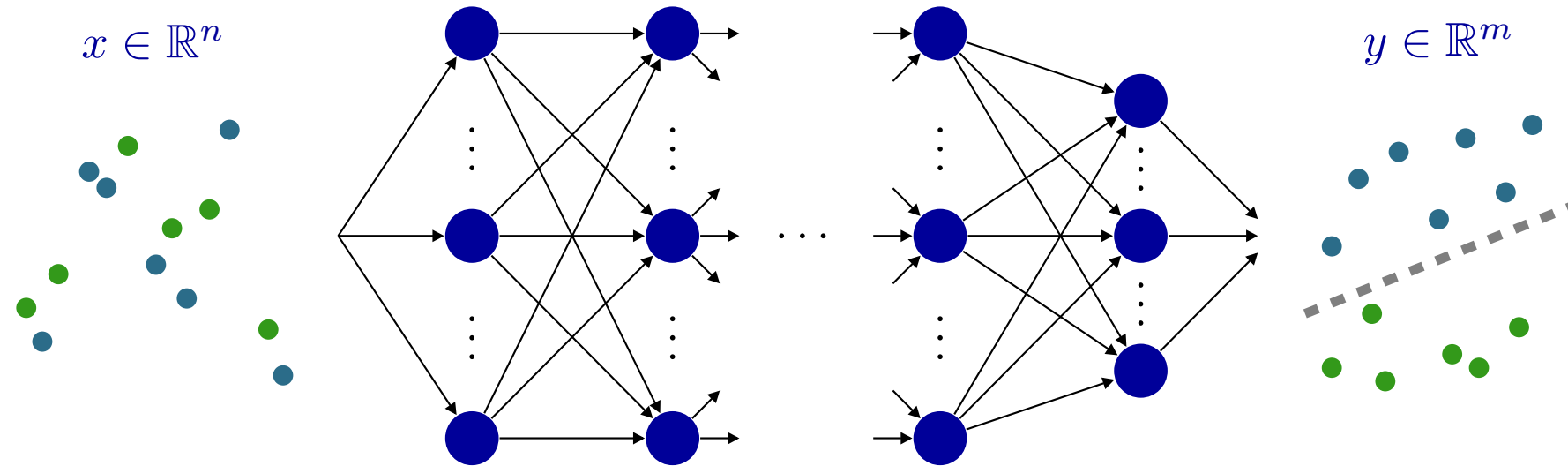
Extension to Other Architectures



Turnpikes in ResNet Training & Experiments



Here – Supervised learning with ResNets




Task – Supervised learning

- Data: $\mathcal{D} = \{(x_i, y_i), i = 1, \dots, M\}$
- Learn a model $f : x \mapsto y$
- **Classification:** $y^i \in \mathbb{N}^m$

Architecture – ResNet

$$v(k+1) = v(k) + \sigma(A(k)v(k) + b(k))$$


$$\frac{dv}{dt} = \sigma(A(t)v(t) + b(t))$$

Optimal control for deep learning?

--

Data propagation:

$$\forall i \in \{1, \dots, M\}, k \in \{0, \dots, N - 1\}$$

$$v_i(k + 1) = v_i(k) + \sigma(A(k)v_i(k) + b(k))$$

$$v_i(0) = x_i$$

$$\hat{y}_i = h(v_i(N; x_i))$$

x_i data points = initial condition y_i labels = targets $A(k)$ weights & $b(k)$ biases = $u(k)$ inputs

Optimal control for deep learning?

$$\min_{A(\cdot), b(\cdot)} \frac{1}{M} \sum_{i=1}^M \ell_f(h(v_i(N; x_i)), y_i)$$

subject to $\forall i \in \{1, \dots, M\}, k \in \{0, \dots, N-1\}$

$$v_i(k+1) = v_i(k) + \sigma(A(k)v_i(k) + b(k))$$

$$v_i(0) = x_i$$

$$\hat{y}_i = h(v_i(N; x_i))$$

} training objective

} ensemble dynamics

x_i data points = initial condition y_i labels = targets $A(k)$ weights & $b(k)$ biases = $u(k)$ inputs

Optimal control for deep learning?

$$\min_{A(\cdot), b(\cdot)} \frac{1}{M} \sum_{i=1}^M \ell_f(h(v_i(N; \mathbf{x}_i)), \mathbf{y}_i)$$

subject to $\forall i \in \{1, \dots, M\}, k \in \{0, \dots, N-1\}$

$$v_i(k+1) = v_i(k) + \sigma(A(k)v_i(k) + b(k))$$
$$v_i(0) = \mathbf{x}_i$$
$$\hat{\mathbf{y}}_i = h(v_i(N; \mathbf{x}_i))$$



ResNet training = optimal control

$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \rho \cdot \|\text{vect}(A(k), b(k))\|^2 + \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

subject to

$$\mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$$
$$\mathbf{x}(0) = [\mathbf{x}_1 \quad \dots \quad \mathbf{x}_M]^\top$$

\mathbf{x}_i data points = initial condition \mathbf{y}_i labels = targets $A(k)$ weights & $b(k)$ biases = $u(k)$ inputs

Optimal control for deep learning?

$$\min_{A(\cdot), b(\cdot)} \frac{1}{M} \sum_{i=1}^M \ell_f(h(v_i(N; \mathbf{x}_i)), \mathbf{y}_i)$$

subject to $\forall i \in \{1, \dots, M\}, k \in \{0, \dots, N-1\}$

$$v_i(k+1) = v_i(k) + \sigma(A(k)v_i(k) + b(k))$$
$$v_i(0) = \mathbf{x}_i$$
$$z_i = h(v_i(N; \mathbf{x}_i))$$



Common practise:
regularizing stage cost / weight decay

ResNet training = optimal control

$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \rho \cdot \|\text{vect}(A(k), b(k))\|^2 + \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

subject to

$$\mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$$

$$\mathbf{x}(0) = [x_1 \quad \dots \quad x_M]^\top$$

x_i data points = initial condition y_i labels = targets $A(k)$ weights & $b(k)$ biases = $u(k)$ inputs

- Training of deep neural networks \Leftrightarrow Optimal control! (LeCun 1988; Li et al. 2017; ...)
- **Now:** Design of training problem? Analysis?

Training from the dynamic systems perspective

Recursion over network layers

$$v(k+1) = v(k) + \sigma(A(k)v(k) + b(k))$$

Euler-forward integration

$h = 1$ and $t = h \cdot k$

$$\frac{dv}{dt} = \sigma(A(t)v(t) + b(t))$$

NN topology & numerical integration?

Continuous-time training formulation

$$V_T(\mathbf{x}) = \min_{A(\cdot), b(\cdot)} \sum_{i=1}^M \ell_f(h(v_i(T; \mathbf{x}_i)), \mathbf{y}_i)$$

subject to $\forall i \in \{1, \dots, M\}$

$$\frac{dv_i}{dt} = \sigma(A(t)v_i(t) + b(t)), \quad t \in [0, T]$$

$$v_i(0) = \mathbf{x}_i$$

$$y = h(v_i(T; \mathbf{x}_i))$$

- Advantages?
- Insights?

Training from the dynamic systems perspective

Continuous-time training formulation

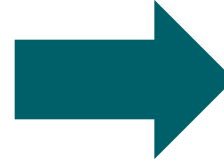
$$V_T(\mathbf{x}) = \min_{A(\cdot), b(\cdot)} \sum_{i=1}^M \ell_f(h(v_i(T; \mathbf{x}_i)), \mathbf{y}_i)$$

subject to $\forall i \in \{1, \dots, M\}$

$$\frac{dv_i}{dt} = \sigma(A(t)v_i(t) + b(t)), \quad t \in [0, T]$$

$$v_i(0) = \mathbf{x}_i$$

$$y = h(v_i(T; \mathbf{x}_i))$$



Optimality conditions

$$\frac{dv_i}{dt} = \sigma(A(t)v_i(t) + b(t)),$$

$$v_i(0) = \mathbf{x}_i, \quad \forall i \in \{1, \dots, M\}$$

$$\frac{d\lambda_i}{dt} = - \left(\frac{\partial}{\partial v_i} \sigma(A(t)v_i(t) + b(t)) \right)^\top \lambda_i,$$

$$\lambda_i(T) = \frac{\partial}{\partial v_i} (\ell_f(h(v_i(T; \mathbf{x}_i))), \mathbf{y}_i)$$

$$0 = \nabla_{A,b} (\lambda^\top \sigma(Av + b))$$

Gradient of optimal value function

$$\nabla_{\mathbf{x}} V_T(\mathbf{x}) = \lambda(0) = [\lambda_i(0)] \quad \frac{\partial V_T}{\partial b} = \frac{\partial V_T}{\partial x} \frac{\partial x}{\partial b}$$

Compute by backward integration

$$\frac{d\lambda_i}{dt} = - \left(\frac{\partial}{\partial v_i} \sigma(A(t)v_i(t) + b(t)) \right)^\top \lambda_i, \quad \lambda_i(T) = \frac{\partial}{\partial v_i} (\ell_f(h(v_i(T; \mathbf{x}_i))), \mathbf{y}_i)$$

back propagation



adjoint sensitivity computation

Training from the dynamic systems perspective

Continuous-time training formulation

$$V_T(\mathbf{x}) = \min_{A(\cdot), b(\cdot)} \sum_{i=1}^M \ell_f(h(v_i(T; \mathbf{x}_i)), \mathbf{y}_i)$$

subject to $\forall i \in \{1, \dots, M\}$

$$\frac{dv_i}{dt} = \sigma(A(t)v_i(t) + b(t)), \quad t \in [0, T]$$

$$v_i(0) = \mathbf{x}_i$$

$$y = h(v_i(T; \mathbf{x}_i))$$



Optimality conditions

$$\frac{dv_i}{dt} = \sigma(A(t)v_i(t) + b(t)),$$

$$v_i(0) = \mathbf{x}_i, \quad \forall i \in \{1, \dots, M\}$$

$$\frac{d\lambda_i}{dt} = - \left(\frac{\partial}{\partial v_i} \sigma(A(t)v_i(t) + b(t)) \right)^\top \lambda_i,$$

$$\lambda_i(T) = \frac{\partial}{\partial v_i} (\ell_f(h(v_i(T; \mathbf{x}_i))), \mathbf{y}_i)$$

$$0 = \nabla_{A,b} (\lambda^\top \sigma(Av + b))$$



Machine Learning

- Regularization with $\|(A, b)\|^2$?
- Back propagation
- Approximation error?
- Network topology?

Optimal Control

- Avoid singular arcs via $\|(A, b)\|^2$
- Adjoint equations
- Reachability
- Numerical integration

Optimal control for deep learning?

$$\min_{A(\cdot), b(\cdot)} \frac{1}{M} \sum_{i=1}^M \ell_f(h(v_i(N; \mathbf{x}_i)), \mathbf{y}_i)$$

subject to $\forall i \in \{1, \dots, M\}, k \in \{0, \dots, N-1\}$

$$v_i(k+1) = v_i(k) + \sigma(A(k)v_i(k) + b(k))$$
$$v_i(0) = \mathbf{x}_i$$
$$z_i = h(v_i(N; \mathbf{x}_i))$$



ResNet training = optimal control

$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \rho \cdot \|\text{vect}(A(k), b(k))\|^2 + \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

subject to

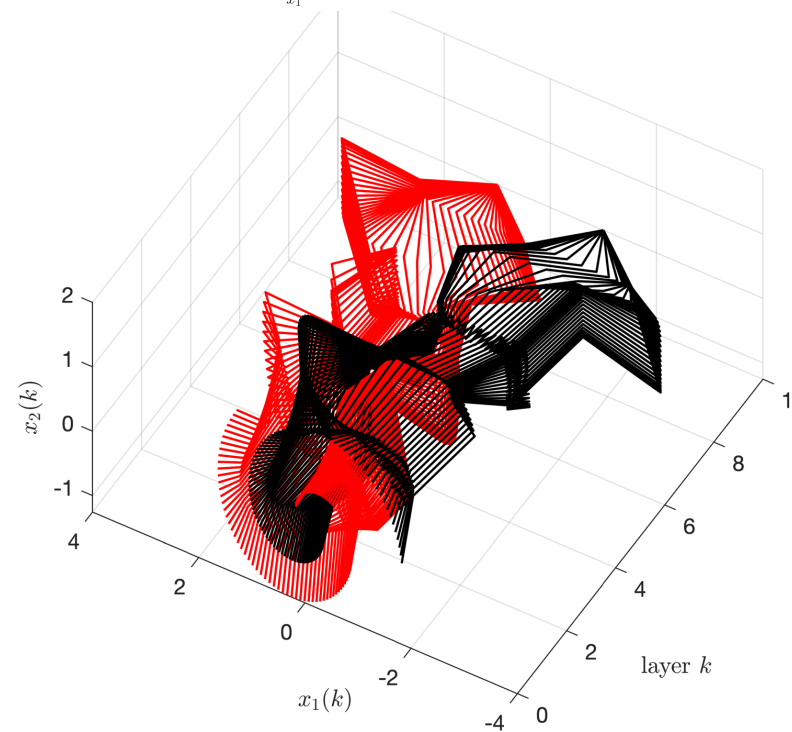
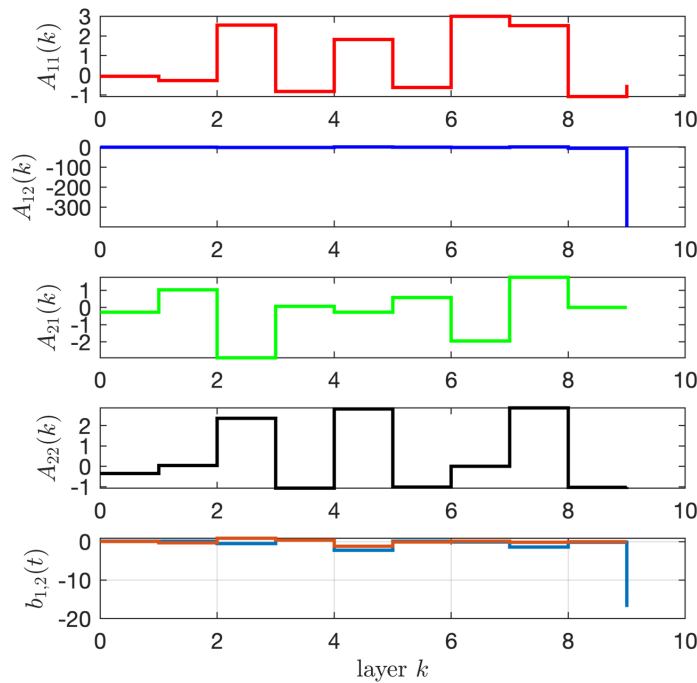
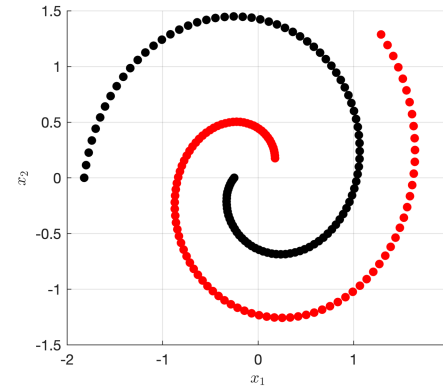
$$\mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$$
$$\mathbf{x}(0) = [x_1 \quad \dots \quad x_M]^\top$$

x_i data points = initial condition y_i labels = targets $A(k)$ weights & $b(k)$ biases = $u(k)$ inputs

- Training of deep neural networks \Leftrightarrow Optimal control! (LeCun 1988; Li et al. 2017; ...)
- **Now:** Design of training problem? Analysis?

Let's start with a toy problem ...

- Spiral data set with 200 samples
- Binary classification of $2D$ data
- NN with depth 10
- Loss $\ell_f(v(N; x_i); y_i) = \|x_i - y_i\|^2$,
 $y_i \in \{(-1, 0), (1, 0)\}$



... and turn it into a teaser

ResNet training

design?

given

$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \underbrace{\mathbf{l}(\mathbf{x}(k)) + \rho \cdot \|\text{vect}(A(k), b(k))\|^2}_{\doteq \mathbf{l}(\mathbf{x}(k), \mathbf{u}(k))} + \gamma \cdot \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

$$\text{subject to } \mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$$

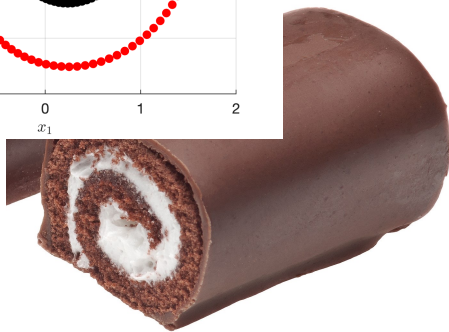
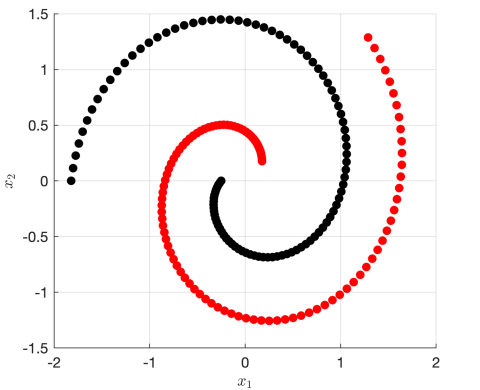
Compute set of optimal *steady states*

$$\bar{\mathbf{x}} \in \mathbb{X}^*(\mathbf{y}) \doteq \underset{\mathbf{x}}{\text{argmin}} \mathbf{l}_f(\mathbf{x}, \mathbf{y}) \neq \emptyset$$

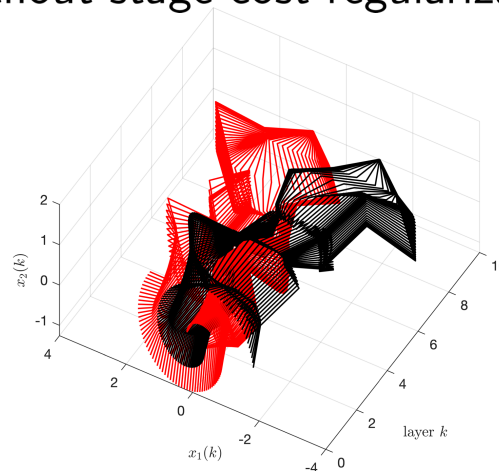
ℓ^2 stage cost

$$\mathbf{l}(\mathbf{x}, \mathbf{u}) = \|\mathbf{x} - \bar{\mathbf{x}}\|_Q^2 + \|\mathbf{u}\|_R^2$$

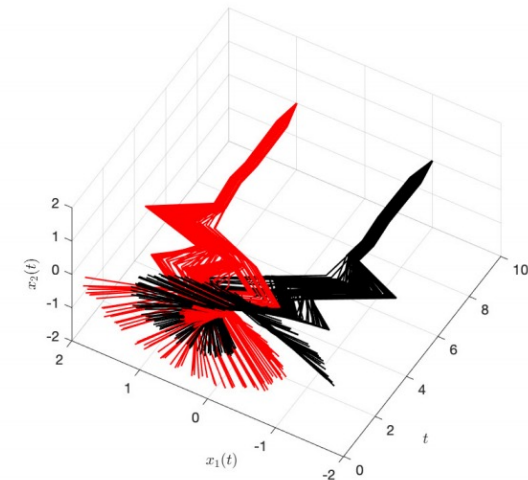
Swiss role



Without stage cost regularization



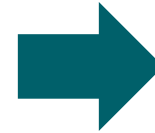
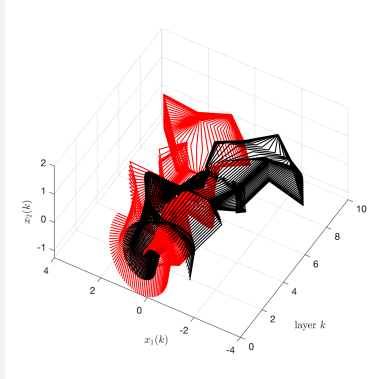
With stage cost regularization



Interplay of stage cost and dynamics on data propagation?

Overview – Dissipativity and Turnpike Properties in NN Training

Optimal Control & Deep Learning?

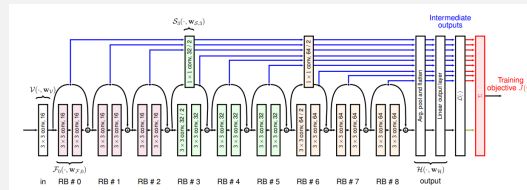
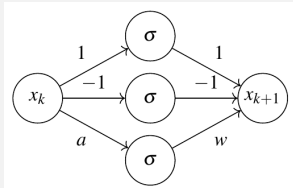


Dissipativity & ResNet Training?

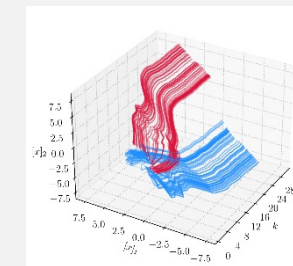
$$S(f(\mathbf{x}, u)) - S(\mathbf{x}) \leq \tilde{\ell}_f(\mathbf{x}, u) + r\|u\|_2^2 - \text{dist}((\mathbf{x}, u), \mathbb{Z}^*)$$



Extension to Other Architectures



Turnpikes in ResNet Training & Experiments



A short detour – Lyapunov-like functions for open systems?

Lyapunov (1892) – Analysis of closed systems

$$\forall x \in \mathbb{R}^n : \quad \frac{dV}{dt} = \nabla V^\top f(x) \leq -\alpha(\|x - \bar{x}\|)$$

change of “energy”
along trajectories decay

$$\Sigma : \quad \begin{aligned} \dot{x} &= f(x) \\ x(0) &= x_0 \end{aligned}$$

Willems (1972, 2007):

“A generalization of Lyapunov functions to open systems, to systems with inputs and outputs”

A short detour – Dissipativity

Lyapunov (1892) – Analysis of closed systems

$$\forall x \in \mathbb{R}^n : \quad \nabla V^\top f(x) \leq -\alpha(\|x - \bar{x}\|)$$

Willems (1972) – Control and analysis of open systems

$$\frac{dS}{dt} = \nabla S^\top f(x, u) \leq w(y, u)$$

change of “energy”
along trajectories

“power”
supply

dissipation inequality

$$\Sigma : \quad \begin{aligned} \dot{x} &= f(x) \\ x(0) &= x_0 \end{aligned}$$

$$\Sigma : \quad \begin{cases} \dot{x} &= f(x, u) \\ y &= h(x, u) \end{cases}$$

← **u**
y →

A short detour – Dissipativity

Lyapunov (1892) – Analysis of closed systems

$$\forall x \in \mathbb{R}^n : \quad \nabla V^\top f(x) \leq -\alpha(\|x - \bar{x}\|)$$

$$\Sigma : \begin{cases} \dot{x} = f(x) \\ x(0) = x_0 \end{cases}$$

Willems (1972) – Dissipation inequality for open systems

$$\frac{dS}{dt} = \nabla S^\top f(x, u) \leq w(y, u)$$

$$\Sigma : \begin{cases} \dot{x} = f(x, u) \\ y = h(x, u) \end{cases} \begin{array}{l} \leftarrow \mathbf{u} \\ \rightarrow \mathbf{y} \end{array}$$

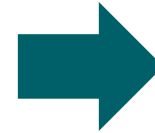
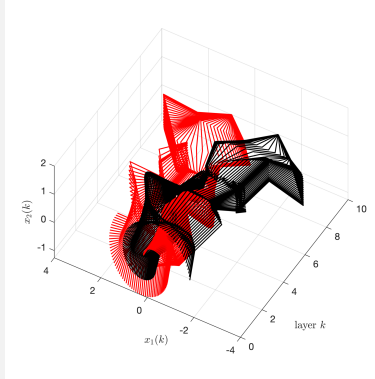
Discrete-time dissipation inequality

$$S(f(x, u)) - S(x) \leq w(x, u) := w(h(x, u), u)$$

$$\Sigma : \begin{cases} x^+ = f(x, u) \\ y = h(x, u) \end{cases} \begin{array}{l} \leftarrow \mathbf{u} \\ \rightarrow \mathbf{y} \end{array}$$

Overview

Optimal Control & Deep Learning?

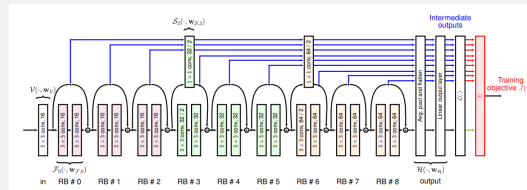
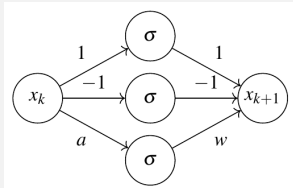


Dissipativity & ResNet Training?

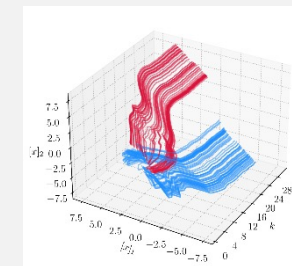
$$S(f(\mathbf{x}, u)) - S(\mathbf{x}) \leq \tilde{\ell}_f(\mathbf{x}, u) + r\|u\|_2^2 - \text{dist}((\mathbf{x}, u), \mathbb{Z}^*)$$



Extension to Other Architectures



Turnpikes in ResNet Training & Experiments



Dissipativity for ResNet training?

ResNet training with stage cost loss

$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \mathbf{l}(\mathbf{x}(k), \mathbf{y}) + \rho \cdot \|\text{vect}(A(k), b(k))\|^2 + \gamma \cdot \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

$$\text{subject to } \mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k)), \mathbf{x}(0) = \mathbf{x}_0 \in \mathbb{R}^{nM}$$

$$\text{Let } \mathbb{X}^*(\mathbf{y}) \doteq \underset{\mathbf{x}}{\text{argmin}} \mathbf{l}(\mathbf{x}, \mathbf{y})$$

Definition: If there exists a bounded storage function $S : \mathbb{R}^{nM} \rightarrow \mathbb{R}_0^+$ such that $\|u\|^2$

$$\underbrace{S(\mathbf{x}(k+1)) - S(\mathbf{x}(k))}_{\text{change of storage}} \leq \underbrace{\mathbf{l}(\mathbf{x}(k), \mathbf{y}) + \rho \cdot \|\text{vect}(A(k), b(k))\|^2}_{\text{supply rate}}$$

then the ResNet dynamics are

w.r.t. the loss function \mathbf{l} .

Simplifying assumption: $\mathbf{l}(\mathbf{x}, \mathbf{y}) = 0$ on $\mathbb{X}^*(\mathbf{y})$

Observation – Dissipativity is closely related to training objective

$$S(\mathbf{x}(1)) - S(\mathbf{x}(0)) \leq -\text{dist}(\mathbf{x}(0), \mathbb{X}^*(\mathbf{y})) + \mathbf{l}(\mathbf{x}(0), \mathbf{y}) + \rho \cdot \|\text{vect}(A(0), b(0))\|^2$$

$$S(\mathbf{x}(2)) - S(\mathbf{x}(1)) \leq -\text{dist}(\mathbf{x}(1), \mathbb{X}^*(\mathbf{y})) + \mathbf{l}(\mathbf{x}(1), \mathbf{y}) + \rho \cdot \|\text{vect}(A(1), b(1))\|^2$$

⋮

$$S(\mathbf{x}(N-1)) - S(\mathbf{x}(N-2)) \leq -\text{dist}(\mathbf{x}(N-2), \mathbb{X}^*(\mathbf{y})) + \dots$$

$$+ \quad S(\mathbf{x}(N)) - S(\mathbf{x}(N-1)) \leq -\text{dist}(\mathbf{x}(N-1), \mathbb{X}^*(\mathbf{y})) + \dots$$

$$S(\mathbf{x}(N)) - S(\mathbf{x}_0) \leq \sum_{k=0}^{N-1} -\text{dist}(\mathbf{x}(k), \mathbb{X}^*(\mathbf{y})) + \sum_{k=0}^{N-1} \mathbf{l}(\mathbf{x}(k), \mathbf{y}) + \rho \cdot \|\text{vect}(A(k), b(k))\|^2$$

Use dissipativity to analyze ResNet training?

Reachability neural ODE \rightarrow Reachability of ResNets

Reachability (Neural ODEs with ReLU activation (Ruiz-Balet))

For any distinct initial point x^i and target z^i , there exists piecewise constant control functions $w(t), A(t), b(t)$ such that all datasamples

$$\dot{x}^i(t) = w(t)\sigma(A(t)x^i(t) + b(t)), \quad x^i(0) = x^i$$

reach $x^i(T) = z^i$ for every $T > 0$



Reachability of ResNet

$$x(k+1) = x(k) + \sigma(A(k)x(k) + b(k))$$

Dissipativity – A very useful system property!

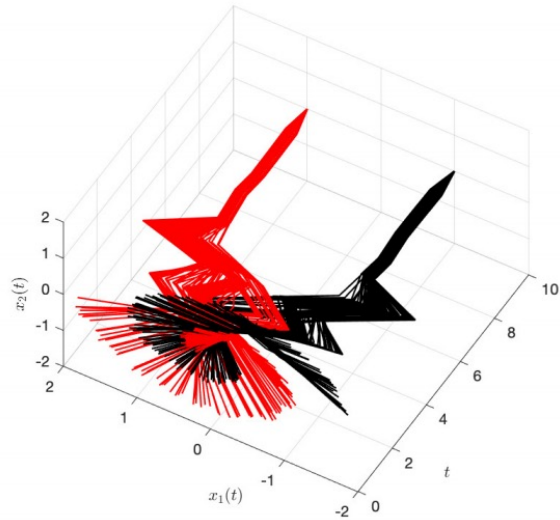
Reachability of $\mathbb{X}^*(\mathbf{y})$: Given \mathbf{x}_0 there exist $A(\cdot), b(\cdot)$ defined on $\{0, \dots, \tilde{N} - 1\}$ and some depth \tilde{N} such that

$$\text{dist}(\mathbb{X}^*(\mathbf{y}), \mathbf{x}(\tilde{N})) = 0$$

$$S(\mathbf{x}^*(N)) - S(\mathbf{x}_0) \leq \sum_{k=0}^{N-1} -\text{dist}(\mathbf{x}^*(k), \mathbb{X}^*(\mathbf{y})) + \mathbf{l}(\mathbf{x}^*(k), \mathbf{y}) + \rho \cdot \|\text{vect}(A^*(k), b^*(k))\|^2 \leq \tilde{C}(\mathbf{x}_0)$$

What is a turnpike? How to exploit it? Does it occur often? Can this be enforced? ...

Turnpike properties in optimal control?



Dorfman et al: “It is exactly like a **turnpike paralleled by a network of minor roads.** [...]

if the origin and destination are far enough apart, it will always pay **to get on to the turnpike** and **cover distance at the best rate of travel**, even if this means adding a little mileage at either end.”

Ramsey (1928); von Neumann (1938); Dorfman et al. (1954); McKenzie (1976); ... Wilde & Kokotovic (1972); Anderson & Kokotovic (1987); Carlson et al. (1991); Rawlings & Amrit (2009); Grüne (2013); F. et al (2014); Trelat & Zuazua (2014); Damm et al. (2014); ...; F. & Grüne (2022)

Turnpikes in optimal control?

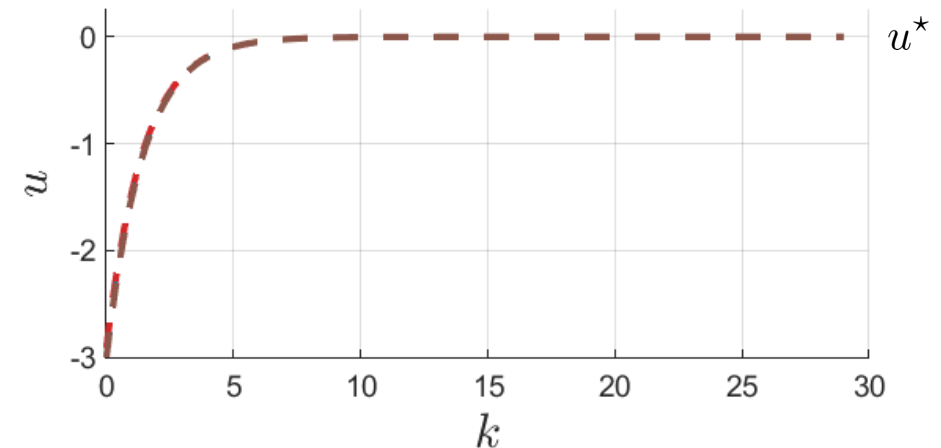
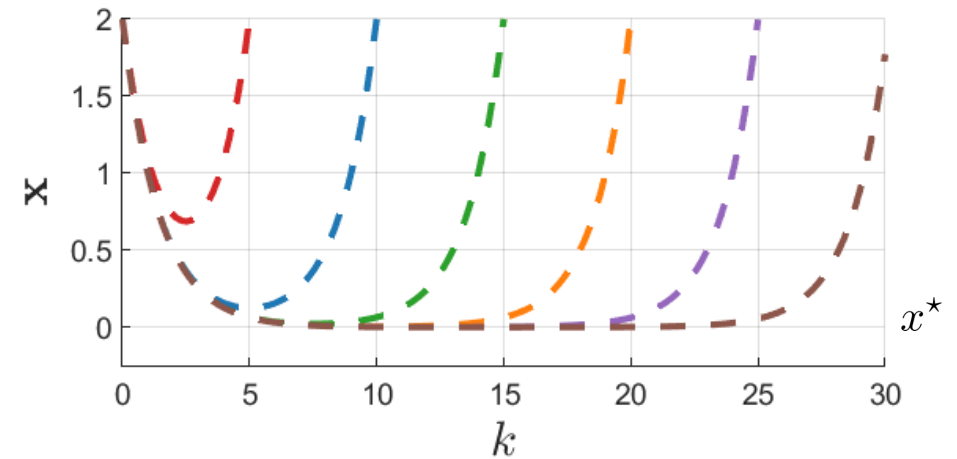
$$\begin{aligned} \min_{u_0, \dots, N-1} \quad & \sum_{k=0}^{N-1} u_k^2 \\ \text{s.t.} \quad & \forall k = 0, \dots, N-1 \\ & \mathbf{x}_{k+1} = 2\mathbf{x}_k + u_k, \quad \mathbf{x}_0 = 2 \\ & \mathbf{x}_k \in [-2, 2], \quad \mathbf{u}_k \in [-3, 3] \end{aligned}$$

Optimal steady state pair

$$\begin{aligned} (\bar{\mathbf{x}}^*, \bar{u}^*) &= \arg \min_{\mathbf{x}, u} \ell(\mathbf{x}, u) \\ \text{s.t.} \quad & \mathbf{x} = \mathbf{f}(\mathbf{x}, u) \end{aligned}$$

Example

$$(\bar{\mathbf{x}}^*, \bar{u}^*) = (0, 0)$$



Optimal steady state pair = **turnpike** = attractor of infinite-horizon optimal solutions

The turnpike phenomenon occurs often in optimal control

J. Math. Biol. (2015) 70:289–327
DOI 10.1007/s00285-014-0768-9

Mathematical Biology

Optimal Control Problems over Large Time

Acta Numerica (2022), pp. 135–263

doi:10.1017/S0962492922000046

Printed in the United Kingdom

in optimal shape design

ce^a, Emmanuel Trélat^a, Enrique Zuazua^b

ersité, CNRS, Université de Paris, Inria, Laboratoire Jacques-Louis Lions (LJLL),
d Analysis, Alexander von Humboldt-Professorship, Department of Mathematics
any;
tational Mathematics, Fundación Deusto Av. de las Universidades 24, 48007 Bilbo,
e Matemáticas, Universidad Autónoma de Madrid, 28049 Madrid, Spain.

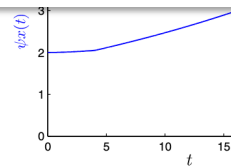
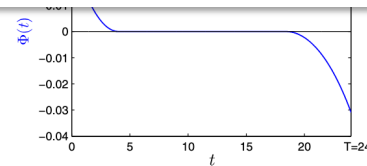
Turnpike in optimal control of PDEs, ResNets, and beyond

Borjan Geshkovski*

Enrique Zuazua

NONLINEAR SYSTEMS

MATHEMATICS OF OPERATIONS RESEARCH
Vol. 17, No. 4, November 1992
Printed in U.S.A.



IFAC Journal of Systems and Control

Volume 30, December 2024, 100290



TURNPIKE SETS AND THEIR ANALYSIS IN STOCHASTIC PRODUCTION PLANNING PROBLEMS*

S. SETHI, H. M. SONER, Q. ZHANG AND J. JIANG

This paper considers optimal infinite horizon stochastic production planning problems with capacity and demand to be finite state Markov chains. The existence of the optimal feedback

Full length article

On the turnpike to design of deep neural networks: Explicit depth bounds ☆

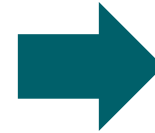
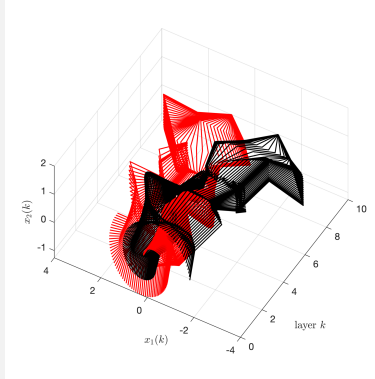
Timm Faulwasser^a ✉, Arne-Jens Hempel^b ✉, Stefan Streif^c ✉

Abstr
over
class
solut
two
trans

Con
con
con

Overview – Dissipativity and Turnpike Properties in NN Training

Optimal Control & Deep Learning?

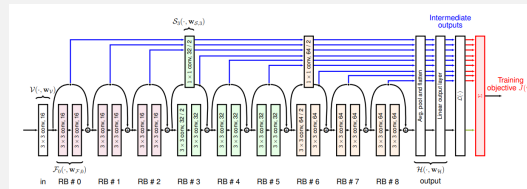
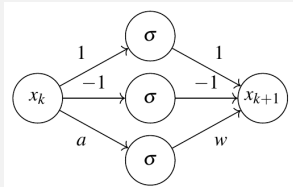


Dissipativity & ResNet Training?

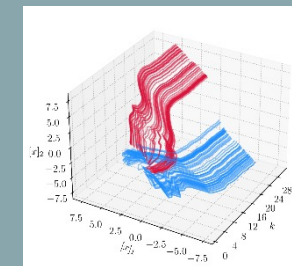
$$S(f(\mathbf{x}, u)) - S(\mathbf{x}) \leq \tilde{\ell}_f(\mathbf{x}, u) + r\|u\|_2^2 - \text{dist}((\mathbf{x}, u), \mathbb{Z}^*)$$



Extension to Other Architectures



Turnpikes in ResNet Training & Experiments



How to exploit turnpikes for deep learning?

Suppose the training OCP has the turnpike property:

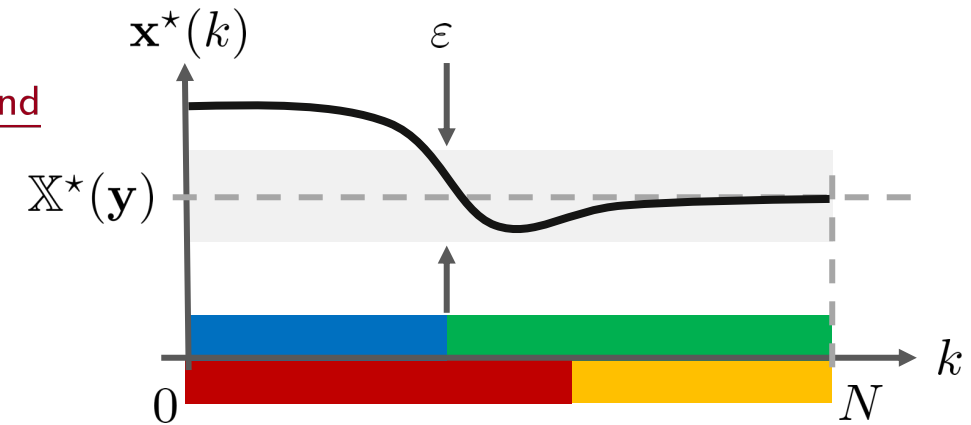
$$\#\{k \in \{0, N - 1\} \mid \mathbf{x}^*(k) \notin \mathcal{B}_\varepsilon(\mathbb{X}^*(\mathbf{y}))\} \leq \frac{\text{Perf+Storage Bnd}}{\varepsilon}$$

Split depth horizon $\{0, \dots, N - 1\}$:

$$\#\{k \mid \mathbf{x}^*(k) \notin \mathcal{B}_\varepsilon\} \cup \#\{k \mid \mathbf{x}^*(k) \in \mathcal{B}_\varepsilon\}$$

Fix $\varepsilon > 0$, **pick** N such that $0 < N - \frac{\text{Perf+Storage Bnd}}{\varepsilon}$

$$\mathbb{X}^*(\mathbf{y}) \doteq \underset{\mathbf{x}}{\operatorname{argmin}} \mathbf{l}_f(\mathbf{x}, \mathbf{y})$$



Hence, **design** training OCP with **turnpike** \Leftarrow **strict dissipativity w.r.t. \mathbf{l}_f !**

Implications of dissipative design of ResNet training

- Strict dissipativity + reachability $\Rightarrow N > \frac{\text{Perf+Storage Bnd}}{\varepsilon} \Rightarrow |\mathbf{l}_f(\mathbf{x}^*(N), \mathbf{y}) - \mathbf{l}^*| \leq C \cdot \varepsilon$
- Problems where ε -close to $\mathbb{X}^*(\mathbf{y})$ is good enough \Rightarrow empirical risk = 0
(ML: performance on data)

Dissipative OCP design for deep learning?

ResNet training

design

given

$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \underbrace{\mathbf{l}(\mathbf{x}(k)) + \rho \cdot \|\text{vect}(A(k), b(k))\|^2}_{\doteq \mathbf{l}(\mathbf{x}(k), \mathbf{u}(k))} + \gamma \cdot \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

subject to $\mathbf{x}(k+1) = \mathbf{x}(k) + \sigma ((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$

Compute set of optimal *steady states*

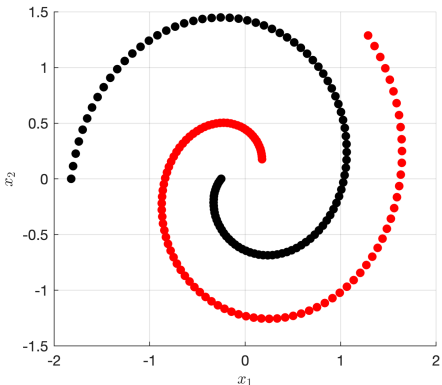
$$\bar{\mathbf{x}} \in \mathbb{X}^*(\mathbf{y}) \doteq \underset{\mathbf{x}}{\text{argmin}} \mathbf{l}_f(\mathbf{x}, \mathbf{y}) \neq \emptyset$$

Stage cost

$$\mathbf{l}(\mathbf{x}, \mathbf{u}) = \|\mathbf{x} - \bar{\mathbf{x}}\|_Q^2 + \|\mathbf{u}\|_R^2$$

$Q, R \succ 0 \Rightarrow$ strict dissipativity & $S(x) = 0$

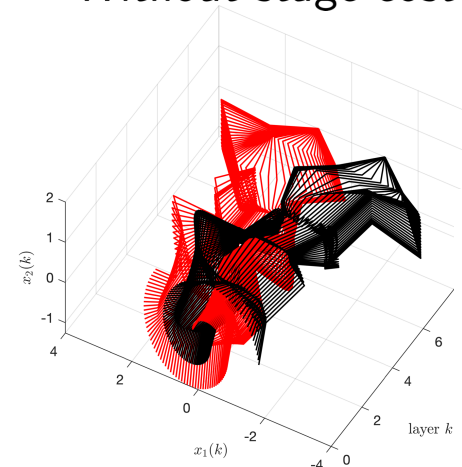
Swiss role



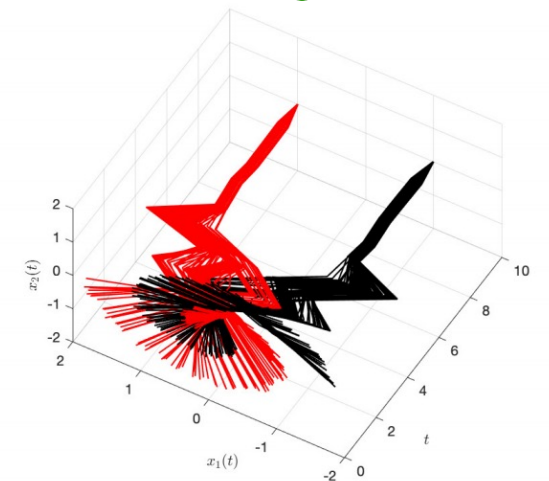
Results with label noise

N	M	β	ρ	$\hat{N}_2(\beta, \rho)$	\hat{N}_2	\hat{N}_∞
5	20	0.87	0.83	$5.24 \cdot 10^2$	$1.45 \cdot 10^2$	7.42
5	50	1.48	0.83	$8.88 \cdot 10^2$	$2.88 \cdot 10^2$	7.20
5	100	3.67	0.83	$2.20 \cdot 10^3$	$6.27 \cdot 10^2$	7.15
5	250	8.83	0.83	$5.30 \cdot 10^3$	$1.58 \cdot 10^3$	8.57
5	500	15.9	0.83	$9.53 \cdot 10^3$	$3.56 \cdot 10^3$	8.40
10	20	1.01	0.75	$4.02 \cdot 10^2$	$1.82 \cdot 10^2$	11.50
10	50	1.65	0.68	$5.07 \cdot 10^2$	$2.35 \cdot 10^2$	8.02
10	100	2.86	0.79	$1.34 \cdot 10^3$	$5.89 \cdot 10^2$	11.38
10	250	7.91	0.72	$2.80 \cdot 10^3$	$1.11 \cdot 10^3$	8.50
10	500	13.8	0.83	$7.94 \cdot 10^3$	$3.30 \cdot 10^3$	18.75

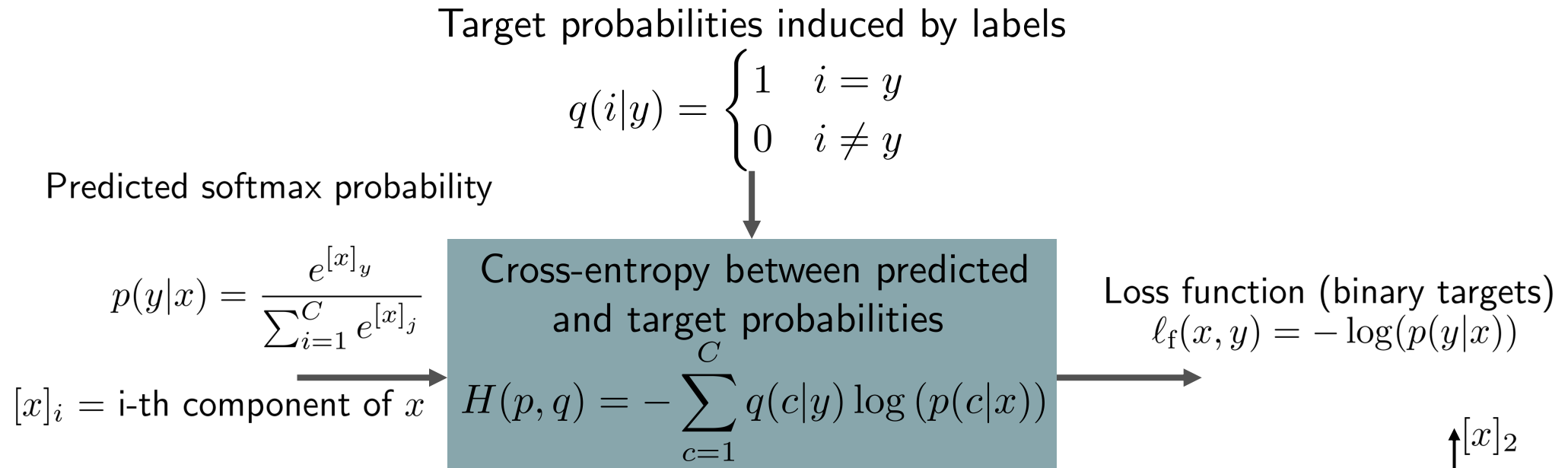
Without stage cost



With stage cost

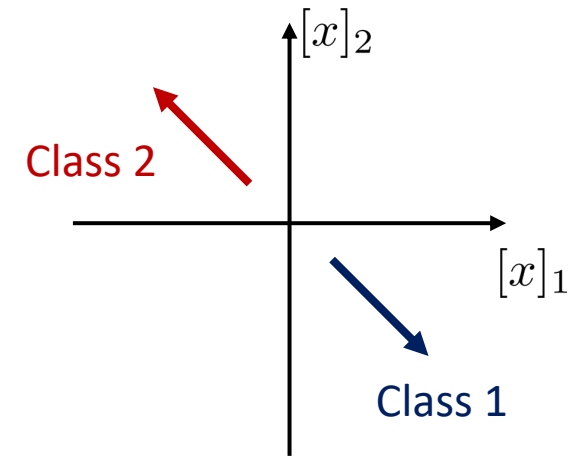


But ... most ML problems are much more complex – Cross entropy?



Minimization of the cross entropy?

$$\min_{x \in \mathbb{R}^2} \ell_f(x, 1)$$



No finite minimizer in \mathbb{R}^C !

How to get finite loss minimizers?

Label smoothing target probabilities
(Szegedy et al. 2016)

$$\tilde{q}(c|y) = \begin{cases} p_d & c = y \\ \frac{1-p_d}{C-1} & c \neq y \end{cases}$$

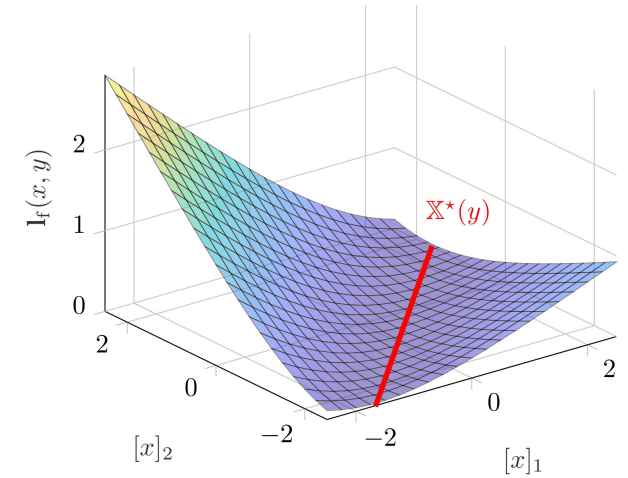
Predicted softmax probability

$$p(y|x) = \frac{e^{[x]_y}}{\sum_{i=1}^C e^{[x]_i}}$$

Soft cross entropy

$$H(p, q) = - \sum_{c=1}^C q(c|y) \log(p(c|x))$$

Soft cross entropy (label smoothing)



$$l_f(x, y) \geq \alpha \circ \text{dist}(x, \mathbb{X}^*(y))$$

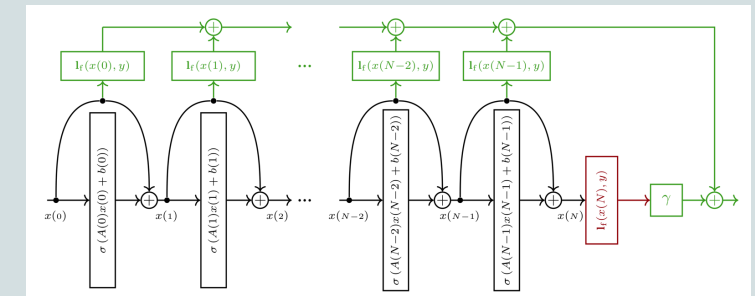
$$l_f(\mathbf{x}, \mathbf{y}) \geq \alpha \circ \text{dist}(\mathbf{x}, \mathbb{X}^*(\mathbf{y}))$$

$\alpha(0) = 0$ & strictly mono. increasing

Soft cross entropy in **stage cost** & **end penalty** \Rightarrow **strictly dissipative training OCP**

$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} l_f(\mathbf{x}(k), \mathbf{y}) + \rho \cdot \|\text{vect}(A(k), b(k))\|^2 + \gamma \cdot l_f(\mathbf{x}(N), \mathbf{y})$$

subject to $\mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$



Evaluation on MNIST



Standard training

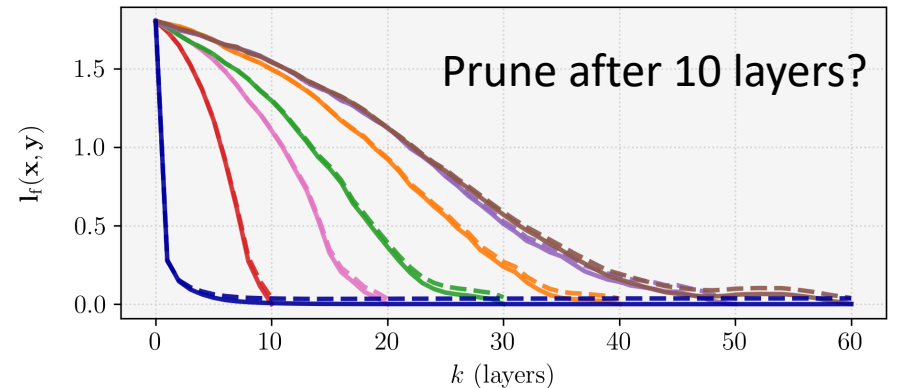
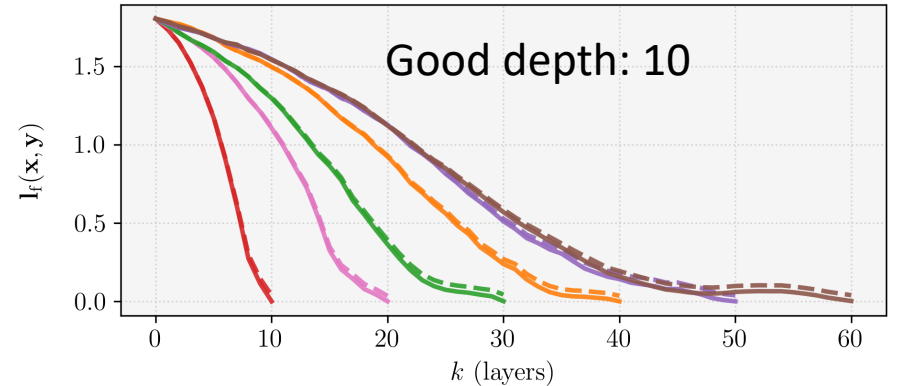
$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \rho \cdot \|\text{vect}(A(k), b(k))\|^2 + \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

$$\text{subject to } \mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$$

Proposed dissipative formulation

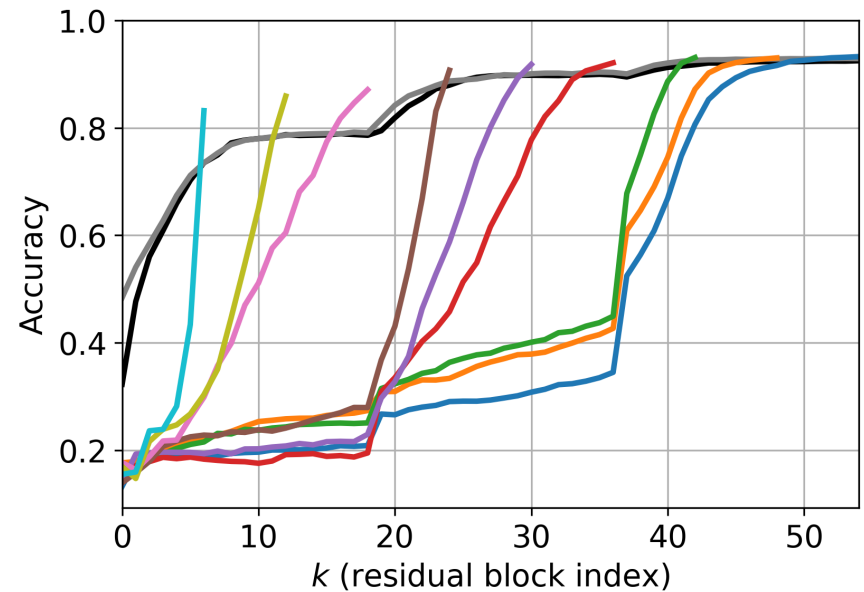
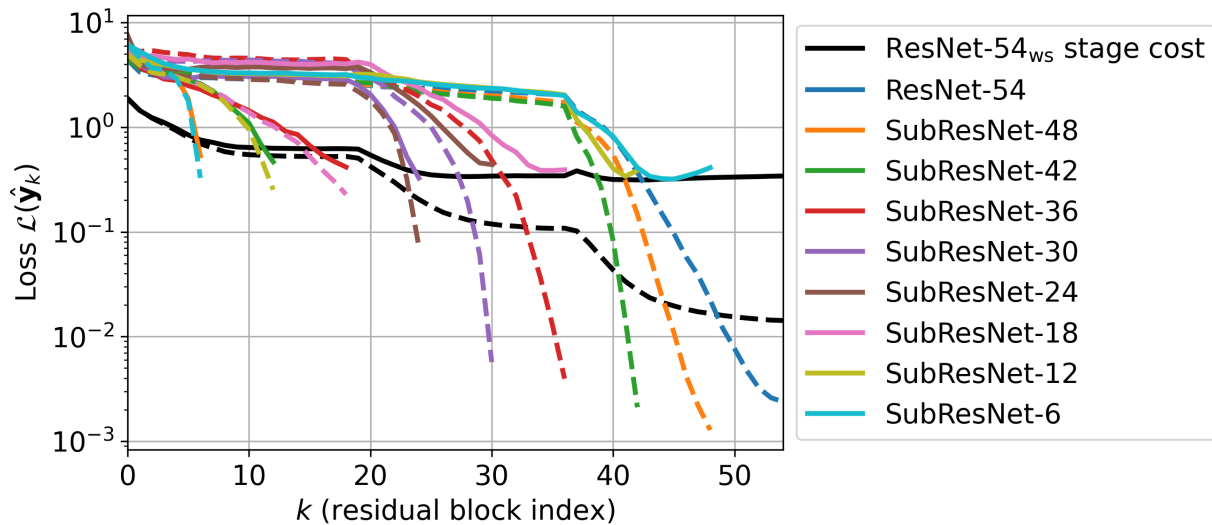
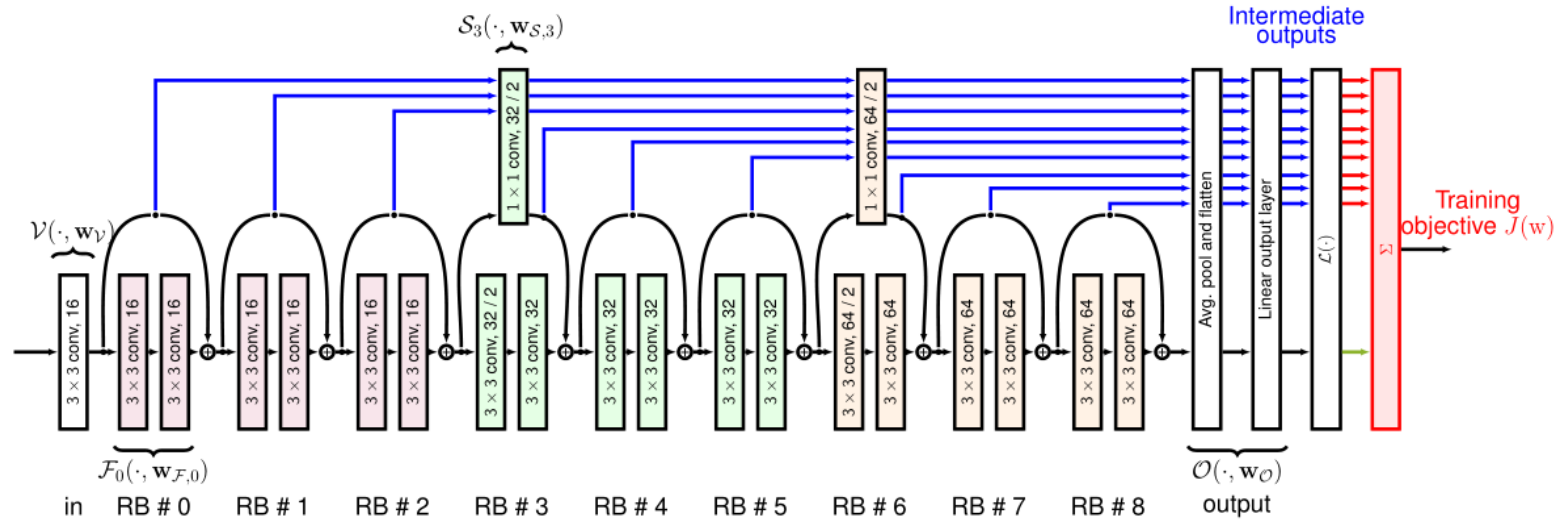
$$\min_{A(\cdot), b(\cdot)} \sum_{k=0}^{N-1} \mathbf{l}_f(\mathbf{x}(k), \mathbf{y}) + \rho \cdot \|\text{vect}(A(k), b(k))\|^2 + \gamma \cdot \mathbf{l}_f(\mathbf{x}(N), \mathbf{y})$$

$$\text{subject to } \mathbf{x}(k+1) = \mathbf{x}(k) + \sigma((I \otimes A(k))\mathbf{x}(k) + \mathbf{1} \otimes b(k))$$



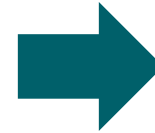
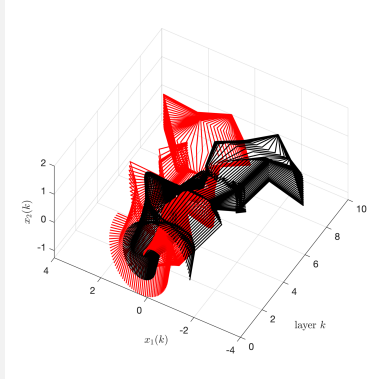
Dissipative training formulation \Rightarrow Reveals early exits & enables formal analysis

Results for more complex ResNets and CIFAR-10



Overview – Dissipativity and Turnpike Properties in NN Training

Optimal Control & Deep Learning?

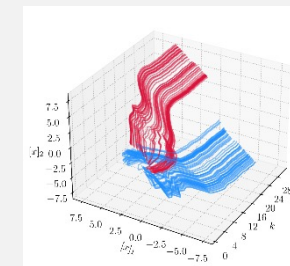


Dissipativity & ResNet Training?

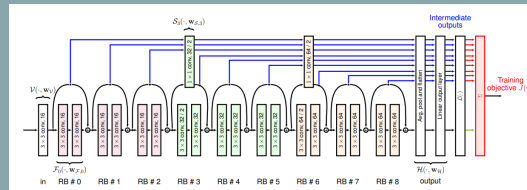
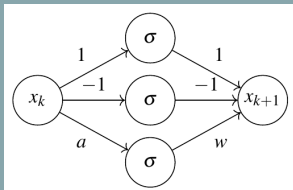
$$S(f(\mathbf{x}, u)) - S(\mathbf{x}) \leq \tilde{\ell}_f(\mathbf{x}, u) + r\|u\|_2^2 - \text{dist}((\mathbf{x}, u), \mathbb{Z}^*)$$



Turnpikes in ResNet Training & Experiments



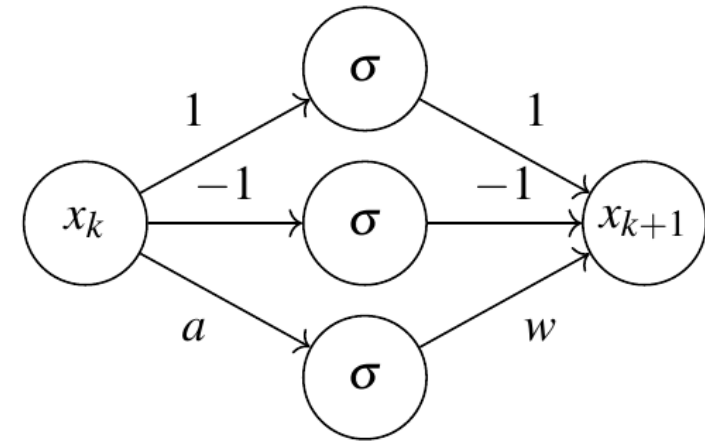
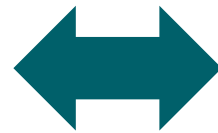
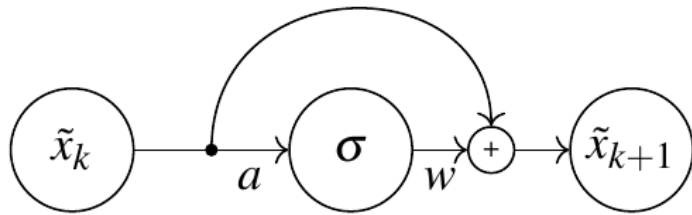
Extension to Other Architectures



ReLU activation and ResNets

Observation – ReLU identity

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



ReLU residual connection neuron

Three equivalent ReLU neurons

From a neuron to a neural network?

Petersen, P., & Voigtlaender, F. (2018). Optimal approximation of piecewise smooth functions using deep ReLU neural networks. *Neural Networks*, 108, 296–330.

Liu, C., Liang, E., & Chen, M. (2024). Characterizing ResNet's universal approximation capability. *Proceedings of Machine Learning Research*, 235, 31477–31515.

Püttchneider, J., Heilig, S., Adilova, L., Fischer, A., & Faulwasser, T. (2026). Turnpikes in deep learning: Beyond ResNets and neural ODEs? In *2026 European Control Conference (ECC)*

Equivalence of ResNets and fully connected networks

ResNet

$$\tilde{\mathbf{x}}_{k+1} = \tilde{\mathbf{x}}_k + (I^D \otimes \tilde{W}_k) \sigma \left((I^D \otimes \tilde{A}_k) \tilde{\mathbf{x}}_k + \mathbf{1}^D \otimes \tilde{b}_k \right)$$

Two linear layers

-State dimension n to hidden dimension \tilde{h}

$$\tilde{A}_k \in \mathbb{R}^{\tilde{h} \times n}, \tilde{b}_k \in \mathbb{R}^{\tilde{h}}$$

-Project back to state dimension

$$\tilde{W}_k \in \mathbb{R}^{n \times \tilde{h}}$$

Fully connected network

$$\mathbf{x}_{k+1} = (I^D \otimes W_k) \sigma \left((I^D \otimes A_k) \mathbf{x}_k + \mathbf{1}^D \otimes b_k \right)$$

Two linear layers

-State dimension n to hidden dimension h

$$A_k \in \mathbb{R}^{h \times n}, b_k \in \mathbb{R}^h$$

-Project back to state dimension n

$$W_k \in \mathbb{R}^{n \times h}$$

Equivalence property for $h = \tilde{h} + 2n$

ResNet parameters

$$\tilde{A}_k \quad \tilde{b}_k \quad \tilde{W}_k$$



Fully connected network parameters

$$A_k = \begin{bmatrix} I^n \\ -I^n \\ \tilde{A}_k \end{bmatrix} \quad b_k = \begin{bmatrix} \mathbf{0}^n \\ \mathbf{0}^n \\ \tilde{b}_k \end{bmatrix} \quad W_k = \begin{bmatrix} I^n & -I^n & \tilde{W}_k \end{bmatrix}$$

Lead to the same state trajectory $\mathbf{x}_k = \tilde{\mathbf{x}}_k$, for all $k = 0, \dots, N - 1$.

Implications?

Optimal steady states of fully connected networks?

Steady state properties

ResNet

$$\tilde{\mathbf{x}}_{k+1} = \tilde{\mathbf{x}}_k + (I^D \otimes \tilde{W}_k) \sigma \left((I^D \otimes \tilde{A}_k) \tilde{\mathbf{x}}_k + \mathbf{1}^D \otimes \tilde{b}_k \right)$$

Any state $\tilde{\mathbf{x}}$ = steady state for $\tilde{u}_k = (\tilde{A}_k, \tilde{b}_k, \tilde{W}_k) = 0$

Fully connected network

$$\mathbf{x}_{k+1} = (I^D \otimes W_k) \sigma \left((I^D \otimes A_k) \mathbf{x}_k + \mathbf{1}^D \otimes b_k \right)$$

Any state \mathbf{x} , there exists a nonempty set of steady-state control inputs $\bar{u} \in \bar{U}(\bar{\mathbf{x}})$, but in general $0 \notin \bar{U}(\bar{\mathbf{x}})$

Optimal steady states

ResNet

$$\ell(\tilde{\mathbf{x}}, \tilde{u}) = \ell_f(\tilde{\mathbf{x}}, \mathbf{y}) + r \|u\|_2^2$$

Set of optimal steady states

$$\mathbb{Z}_{\mathbf{y}}^* = \{(\bar{\mathbf{x}}, \mathbf{0}) \mid \bar{\mathbf{x}} \in \mathbb{X}_{\mathbf{y}}^*\}$$

Fully connected network

$$\ell(\tilde{\mathbf{x}}, \tilde{u}) = \ell_f(\tilde{\mathbf{x}}, \mathbf{y})$$

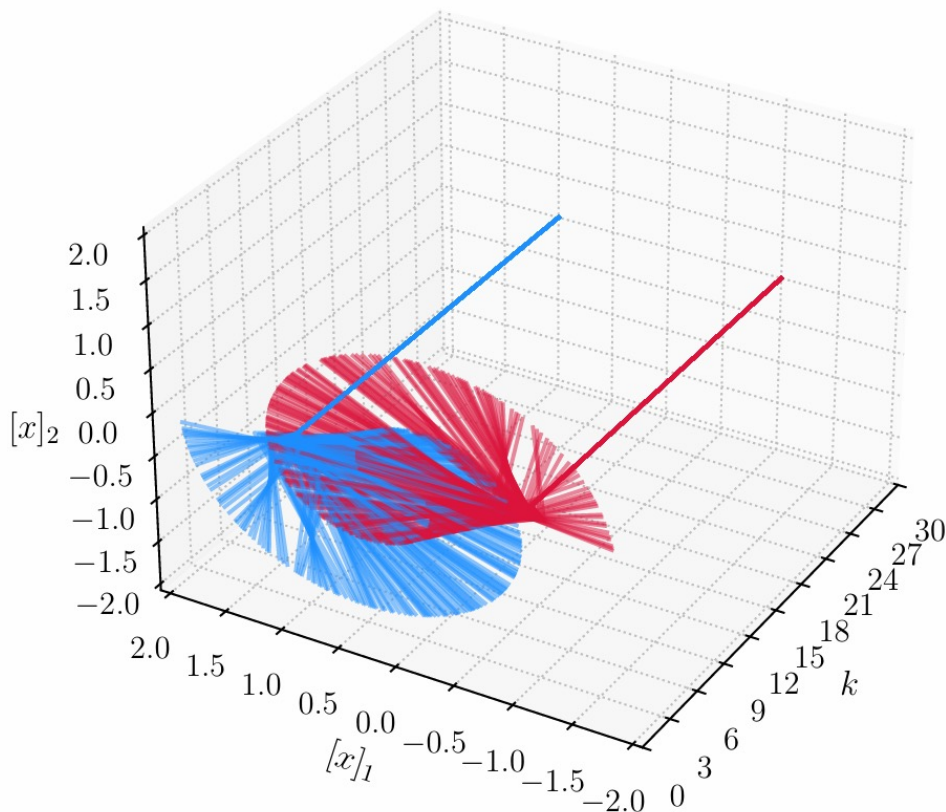
Set of optimal steady states

$$\mathbb{Z}_{\mathbf{y}}^* = \{(\bar{\mathbf{x}}, \bar{u}) \mid \bar{\mathbf{x}} \in \mathbb{X}_{\mathbf{y}}^*, \bar{u} \in \bar{U}(\bar{\mathbf{x}})\}$$

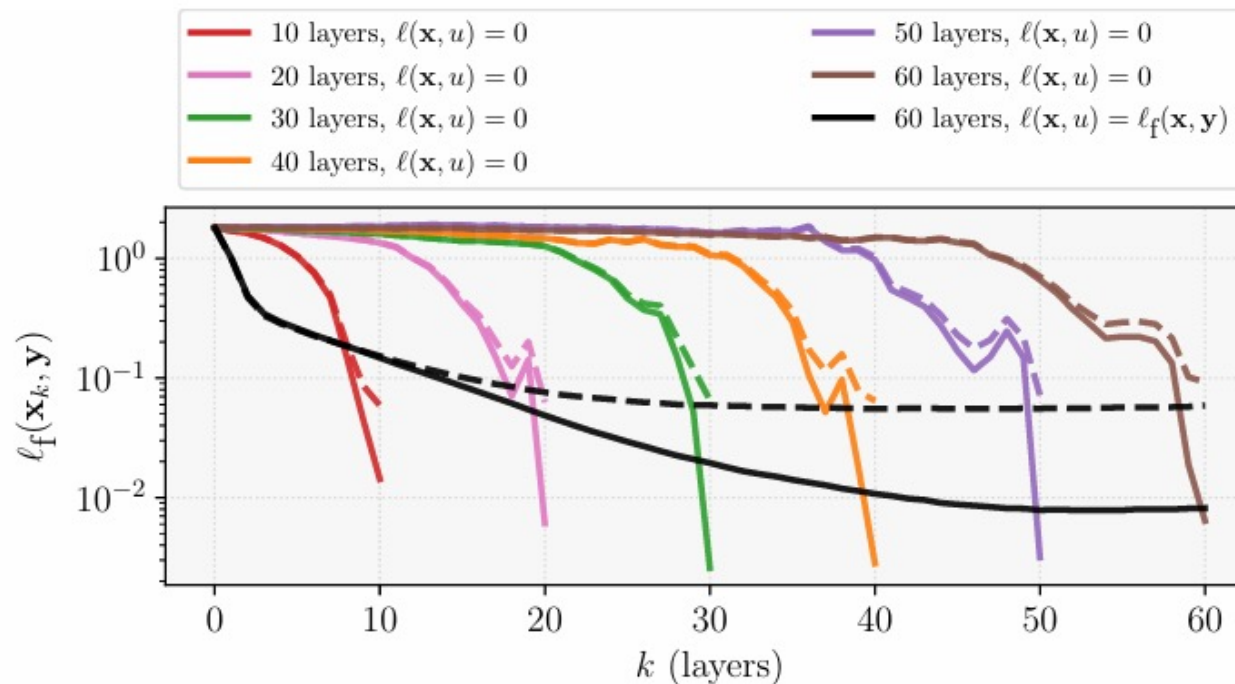
Reachability and analysis results carry over!

Turnpikes with fully connected networks – Experiments

Two spirals + l_2 loss



MNIST + soft cross-entropy loss



Reachability and analysis results carry over!

Summary

Turnpike properties in optimal control

- Very useful tool for non-local analysis: properties of **OCPs parametric in x_0**
- Dissipativity and turnpikes are closely related
- Results for ODEs, PDEs, stochastic systems, ...
- **Problem analysis and OCP synthesis**
- ...

Turnpikes & dissipativity in deep learning

- Known – shape loss landscape:
 - BranchyNet: Teerapittayanon et al. (2015); ...
- **Dissipativity-based formulation of training uses loss as stage cost**
- **Dissipativity + early exits = avenue for analysis of deep networks**
- ***We are just at the beginning ...***

References

- He, Zhang, Ren & Sun (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*
- LeCun (1988) A theoretical framework for back-propagation. In *Proceedings of the 1988 connectionist models summer school*
- Li, Chen, Tai & E (2018). Maximum principle based algorithms for deep learning. *Journal of Machine Learning Research*
- Esteve, Geshkovski, Pighin & Zuazua (2020). Large-time asymptotics in deep learning. *arXiv*
- Ruiz-Balet & Zuazua (2023). Neural ODE control for classification, approximation, and transport. *SIAM Review*
- Tabuada & Gharesifard (2023). Universal approximation power of deep residual neural networks through the lens of control. *IEEE Transactions on Automatic Control*
- ...
- Faulwasser, Hempel & Streif (2024). On the turnpike to design of deep neural networks: Explicit depth bounds. *IFAC Journal of Systems and Control*
- Pan, Stomberg, Engelmann & Faulwasser (2021). First results on turnpike bounds for stabilizing horizons in NMPC. *IFAC-PapersOnLine*
- Püttschneider & Faulwasser (2024). On dissipativity of cross-entropy loss in training ResNets. *Automatica*
- Püttschneider, Heilig, Fischer & Faulwasser (2025). Towards an Optimal Control Perspective of ResNet Training. *arXiv & ICML workshop*
- Püttschneider, J., Heilig, S., Adilova, L., Fischer, A., & Faulwasser, T. (2026). Turnpikes in deep learning: Beyond ResNets and neural ODEs? In *2026 European Control Conference (ECC)*

Summary

Turnpike properties in optimal control

- Very useful tool for non-local analysis: properties of **OCPs parametric in x_0**
- Dissipativity and turnpikes are closely related
- Results for ODEs, PDEs, stochastic systems, ...
- **Problem analysis and OCP synthesis**
- ...

Turnpikes & dissipativity in deep learning

- Known – shape loss landscape:
 - BranchyNet: Teerapittayanon et al. (2015); ...
- **Dissipativity-based formulation of training uses loss as stage cost**
- **Dissipativity + early exits = avenue for analysis of deep networks**
- ***We are just at the beginning ...***

Thank you & and all co-authors...

Funded by



Deutsche
Forschungsgemeinschaft
German Research Foundation



Transregio
391

FOR 5785

ALeSCo

Active Learning for Systems and Control