

Learning differential equations that are easy to solve

Jacob Kelly*
 University of Toronto, Vector Institute
`jkelly@cs.toronto.edu`

Jesse Bettencourt*
 University of Toronto, Vector Institute
`jessebett@cs.toronto.edu`

Matthew James Johnson
 Google Brain
`mattjj@google.com`

David Duvenaud
 University of Toronto, Vector Institute
`duvenaud@cs.toronto.edu`

Abstract

Differential equations parameterized by neural networks become expensive to solve numerically as training progresses. We propose a remedy that encourages learned dynamics to be easier to solve. Specifically, we introduce a differentiable surrogate for the time cost of standard numerical solvers, using higher-order derivatives of solution trajectories. These derivatives are efficient to compute with Taylor-mode automatic differentiation. Optimizing this additional objective trades model performance against the time cost of solving the learned dynamics. We demonstrate our approach by training substantially faster, while nearly as accurate, models in supervised classification, density estimation, and time-series modelling tasks.

1 Introduction

Differential equations describe a system's behavior by specifying its instantaneous dynamics. Historically, differential equations have been derived from theory, such as Newtonian mechanics, Maxwell's equations, or epidemiological models of infectious disease, with parameters inferred from observations. Solutions to these equations usually cannot be expressed in closed-form, requiring numerical approximation.

Recently, ordinary differential equations parameterized by millions of learned parameters, called neural ODEs, have been fit for latent time series models, density models, or as a replacement for very deep neural networks (Rubanova et al., 2019; Grathwohl et al., 2019; Chen et al., 2018). These learned models are not constrained to match a theoretical model, only to optimize an objective on observed data. Learned models with nearly indistinguishable predictions can have substantially different dynamics. This raises the possibility that we can find equivalent models that are easier and faster to solve. Yet standard training methods have no way to penalize the complexity of

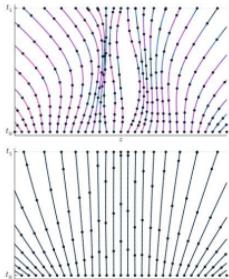


Figure 1: Top: Trajectories of an ODE fit to map $\mathbf{z}(t_1) = \mathbf{z}(t_0) + \mathbf{z}'(t_0)t_1^3$. The learned dynamics are unnecessarily complex and require many evaluations (black dots) to solve.
 Bottom: Regularizing the third total deriva-

Learning Differential Equations that are Easy to Solve

Jacob Kelly, Jesse Bettencourt, Matthew James Johnson, David Duvenaud

Outline

1. Get intuition about: what do we want to affect?
2. Recalls: RK methods, step-size selection
3. =1 + 2: objective from paper

Motivation

$$\dot{z} = z$$

Motivation $\dot{z} = z$

- Consider (some how) we train ODE: $\dot{z} = f_\theta(z, t)$
-
-
-

Motivation $\dot{z} = z$

- Consider (some how) we train ODE: $\dot{z} = f_\theta(z, t)$
- How do we apply it?
By (numerically) solving IVP problem:
$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$

Motivation $\dot{z} = z$

- Consider (some how) we train ODE: $\dot{z} = f_\theta(z, t)$
- How do we apply it?
By (numerically) solving IVP problem:
$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$
- Numerical Solver evaluates: $f_\theta(z(t), t)$
-

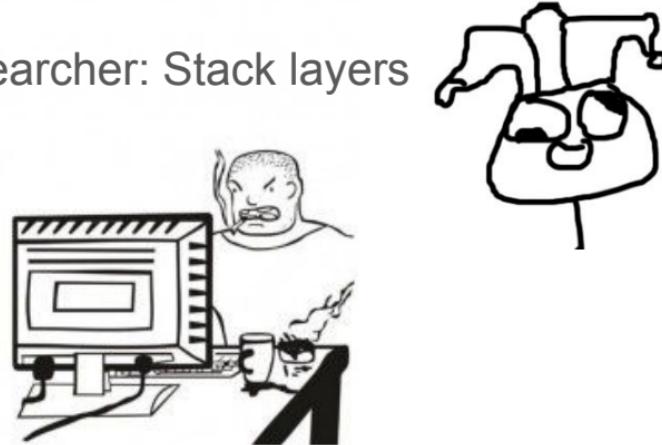
Motivation $\dot{z} = z$

- Consider (some how) we train ODE: $\dot{z} = f_\theta(z, t)$
- How do we apply it?
By (numerically) solving IVP problem:
$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$
- Numerical Solver evaluates: $f_\theta(z(t), t)$
- Deep Learning Researcher: Stack layers



Motivation $\dot{z} = z$

- Consider (some how) we train ODE: $\dot{z} = f_\theta(z, t)$
- How do we apply it?
By (numerically) solving IVP problem:
$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$
- Numerical Solver evaluates: $f_\theta(z(t), t)$
- Deep Learning Researcher: Stack layers



Contradiction

Idea



- ODE should have “simple dynamics”

Idea

- ODE should have “simple dynamics”



Simple vector field

Idea

- ODE should have “simple dynamics”



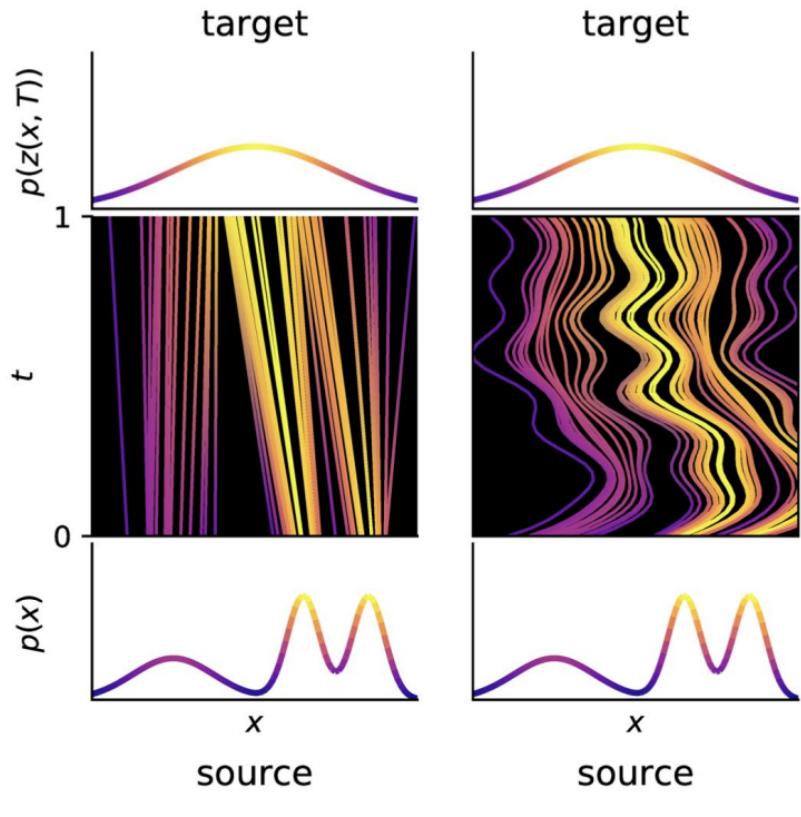
Simple vector field



Not so simple

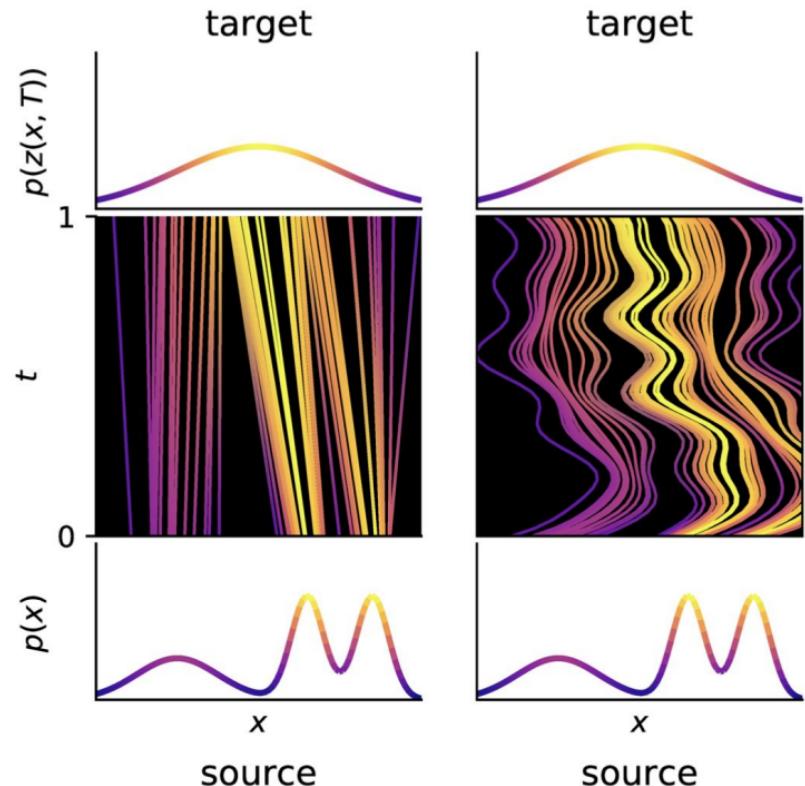
Idea

- ODE should have “simple dynamics”
- [“How to Train Your Neural ODE: the World of Jacobian and Kinetic Regularization”](#)



Idea

- ODE should have “simple dynamics”
- [“How to Train Your Neural ODE: the World of Jacobian and Kinetic Regularization”](#)
- Straight lines are good
- “Optimal transport”
formalization of “simple”



(a) Optimal transport map

(b) generic flow

Idea

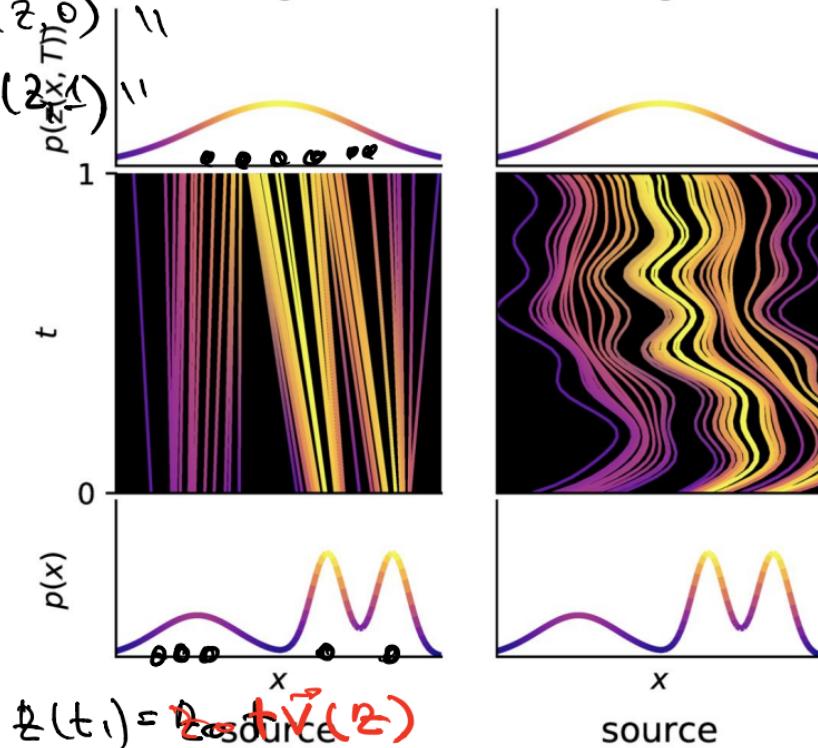
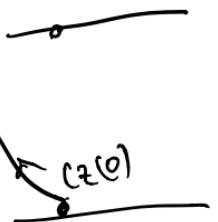
$$\min_{\theta} C(\mathbf{z}(t)) \quad \text{s.t. } 1) \quad \dot{\mathbf{z}}_t = \underbrace{\text{div}(\mathbf{f}(\mathbf{z}(t), t) \mathbf{P})}_{\text{target}}$$
$$2) \quad \mathbf{P}(\mathbf{z}, 0) = \underbrace{\mathbf{P}(\mathbf{z}, T)}_{\text{target}}$$

"How to Train Your Neural ODE: the World of Jacobian and Kinetic Regularization"

- Straight lines are good
- Optimal transport
formalization of "simple"
- Heavy math boils down to
regularization:

$$\left[\|\mathbf{f}(\mathbf{z}(t), t, \theta)\|_2^2 + \|\nabla_{\mathbf{z}} \mathbf{f}(\mathbf{z}(t), t, \theta)\|_F^2 \right] + \text{loss}$$

$\dot{\mathbf{z}} = \mathbf{f}(\mathbf{z}, t)$ 2) optimal b w/ore?



(a) Optimal transport map

(b) generic flow

Motivation

$$\|f(\mathbf{z}(t), t, \theta)\|_2^2$$

$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$
$$f_\theta(z(t), t)$$

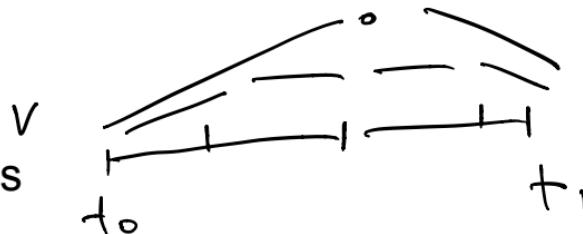
- Less evaluations = **better**
- What Influence on #Evaluations?
(Number of Function Evaluation)

Motivation

$$\|f(\mathbf{z}(t), t, \theta)\|_2^2$$

- Less evaluations = **better**
- What Influence on #Evaluations?
(Number of Function Evaluation)
- Step size: bigger step, less number of steps

$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$
$$f_\theta(z(t), t)$$



Motivation

$$\|f(\mathbf{z}(t), t, \theta)\|_2^2$$

- Less evaluations = better
- What Influence on #Evaluations?
(Number of Function Evaluation)
- Step size: bigger step, less number of steps
- But local error should be bounded, so:
 - Higher order of method: more evaluations
 - Adaptive step size:
 - Rejection of trajectory leads for more evaluations

$$z(t_1) = z(t_0) + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$

$f_\theta(z(t), t)$

Motivation

$$\|f(\mathbf{z}(t), t, \theta)\|_2^2$$

- Less evaluations = better
- What Influence on #Evaluations?
(Number of Function Evaluation)

$$z(t_1) = \boxed{z(t_0)} + \int_{t_0}^{t_1} f_\theta(z(t), t) dt$$
$$\underbrace{f_\theta(z(t), t)}_{f_\theta, (\text{go, } (z(t_0), t_0))}$$

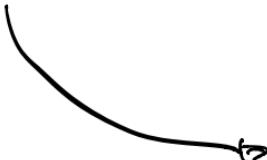
- Step size: bigger step, less number of steps
- But local error should be bounded, so:
 - Higher order of method: more evaluations
 - Adaptive step size:
 - Rejection of trajectory leads for more evaluations

Regularize properties of $f_\theta(z(t), t)$ that adaptive step-size takes in account

Subset of Runge-Kutta methods

- One step [bei δ -yugur zabitat of x_n]
- Multi-stage
- Explicit

$$x(t_{n+1}) = X_{n+1} = X_n + \Delta t \sum_{i=1}^m c_i k_i$$

- 
- $k_1 = f(t_n, x_n)$ slopes
 - $k_2 = f(t_n + \Delta t, x_n + \Delta t \beta_{21} k_1)$ next
 - \vdots
 - $k_m = f(t_n + \Delta m \Delta t, x_n + \Delta t \sum_{j=1}^{m-1} \beta_{mj} k_j)$

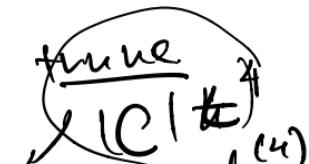
Example

- Let's derive some method as it is easy:
 - Taylor Series of True Trajectory and match its order
 - * Hand derivations here *

$$\text{Blank } \dot{x} = f(x, t) \quad \left\{ \begin{array}{l} \dot{x} \\ f \end{array} \right\} = \int f(\xi, t+3) \quad \dot{x} = f(x)$$

$x(t+\Delta t) = \Phi(x(t), \Delta t) \rightarrow \text{flow map, numeric}$

$x_{n+1} = \Psi(x_n, \Delta t) \rightarrow \text{numeric, not exact.}$

Local Error = $|x_{n+1} - x(t+\Delta t)| = \epsilon_i$, Global = $\sum \epsilon_i$ 

$$\therefore \Phi(x(t), \Delta t) =: x(t) = x(0) + t \dot{x} + \frac{t^2}{2} \ddot{x} + \frac{t^3}{6} \dddot{x} + \dots \approx f^{(4)}$$

$$\dot{x} = f(x), \quad \ddot{x} = \frac{1}{\Delta t} f(x(t)) = f_x \dot{x} = f_x f. \quad f_{j,k} = \frac{\partial f_j}{\partial x_k}$$

$$\ddot{x}_j = \frac{1}{\Delta t} \ddot{x}_j = \frac{1}{\Delta t} \left[\frac{\partial f_j}{\partial x_k} f_k \right] = \frac{1}{\Delta t} f_{j,k} f_k =$$

$$= f_{j,k} f_k f_k + f_{j,k} f_k f_k.$$

$$x_j(t) = x_j + t f_j + \frac{1}{2} t^2 f_{j,k} f_k + \frac{1}{6} t^3 \underbrace{\{$$

3

$$\begin{aligned}
 \text{Blank } \dot{x} = f(t, x) \quad x_{k+1} &= x_k + \Delta t \cancel{\phi} = \\
 &= x_k + \cancel{\Delta t} \left[A f(t_k, x_k) + B f \underbrace{t_k + P \Delta t,}_{x_k + Q f(t_k, x_k) \Delta t} \right] \xrightarrow{\text{vergleich verrechnet}}
 \end{aligned}$$

$$x_{k+1} = x_k + \Delta t \cdot f(t_k, x_k) \cdot (A+B)$$

$$+ \frac{\Delta t^2}{2} B (P f_t + Q f_{xt}) + O(\Delta t^3)$$

④ այս են առջևությունը

$$A+B = \underline{S}$$

$$BP = \frac{1}{2}, BA = \frac{1}{2}$$

R445 $\frac{1}{2}$, $\frac{1}{3}$

R6 0.111738

Example

- Let's derive some method as it is easy:
 - Taylor Series of True Trajectory and match its order
 - * Hand derivations here *
- How to get higher orders in approximation?
 - Higher order -> we need to make $t \cdot t$ some where ✓
 - Use different points to measure slope and average them ✓
 - * Hand derivations here *

Blank

Blank

Банное непрерывн
1. ϵ = allowed error
(tolerance)

Adaptive step-size selection

Одновременно
бесшовный переход
между методами

(Halver)

$K\Delta t$

непрерывный метод

$$e(\Delta t) \leq P\epsilon,$$

$P < 1.$

acc ✓

$$\Delta t' = \underline{\underline{2}} \cdot \Delta t$$

reject ✗

$$\Delta t' = \frac{1}{2} \Delta t.$$

Adaptive step-size selection

$$\tilde{x}_{n+1} = x_{true} + C h^4$$

$$\tilde{x}_{n+1} = x_{true} + C h^5$$

$$|\tilde{x}_{n+1} - x_n| \sim (\Delta t)^5 \rightarrow \text{unstable}$$

$$\tilde{\epsilon} = \frac{(\Delta t)^5}{\Delta t} \rightarrow \text{domino}$$

"Bootstrap"

$$C \leq \sup \left| \frac{d^{k+1}}{dt^{k+1}} \right| \times 1$$

RK5

6 5 толкак
нестабил
надобно го стабилак
која ф.

Embedded
RK4

Regularization

$$\mathcal{R}_K(\theta) = \int_{t_0}^{t_1} \left\| \frac{d^K \mathbf{z}(t)}{dt^K} \right\|_2^2 dt \quad \dot{\mathbf{z}} = \mathbf{f}(\mathbf{z}, t)$$

Assume, $\mathcal{R}_K(\theta) = 0$

- $K=0 \times \int_{t_0}^{t_1} \langle \mathbf{z}_j, \dot{\mathbf{z}}_j \rangle dt = 0, \quad \mathbf{z}_j(t) = 0$

Regularization

$$\mathcal{R}_K(\theta) = \int_{t_0}^{t_1} \left\| \frac{d^K \mathbf{z}(t)}{dt^K} \right\|_2^2 dt$$

Assume, $\mathcal{R}_K(\theta) = 0$

- $K = 0$

$$\langle \dot{\mathbf{z}}, \dot{\mathbf{z}} \rangle = 0 \quad \dot{\mathbf{z}}_j = 0 \Rightarrow z_j = \underset{A_j}{\text{const}}$$

- $K = 1$

Regularization

$$\mathcal{R}_K(\theta) = \int_{t_0}^{t_1} \left\| \frac{d^K \mathbf{z}(t)}{dt^K} \right\|_2^2 dt$$

Assume, $\mathcal{R}_K(\theta) = 0$

\sum

- $K = 0$
- $K = 1$

$$z_j = C_1^j t + C_2^j$$

- $K = 2$

$$z(t) = z(0) + \underbrace{\frac{\dot{z}}{1} t}_{\text{arrow}} + \cancel{\frac{\ddot{z}}{2} t^2} + \dots$$

Experiments

R_{3.}

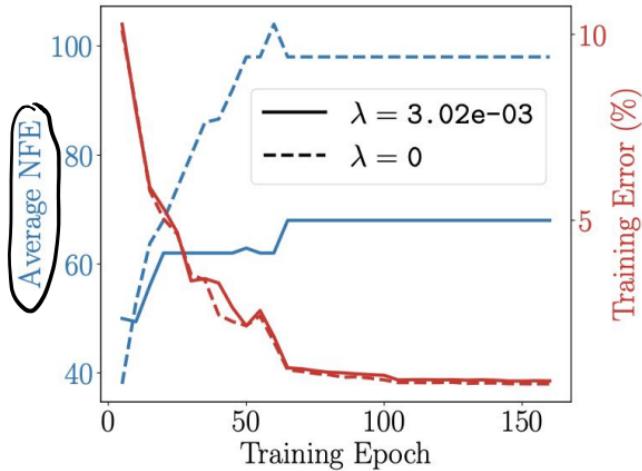


Figure 3: Number of function evaluations (NFE) and training error during training. Speed regularization (solid) decreases the NFE throughout training without substantially changing the training error.

training time ↑

5.1 Supervised Learning

We construct a model for MNIST classification: it takes in as input a flattened MNIST image and integrates it through dynamics given by a simple MLP, then applies a linear

Experiments

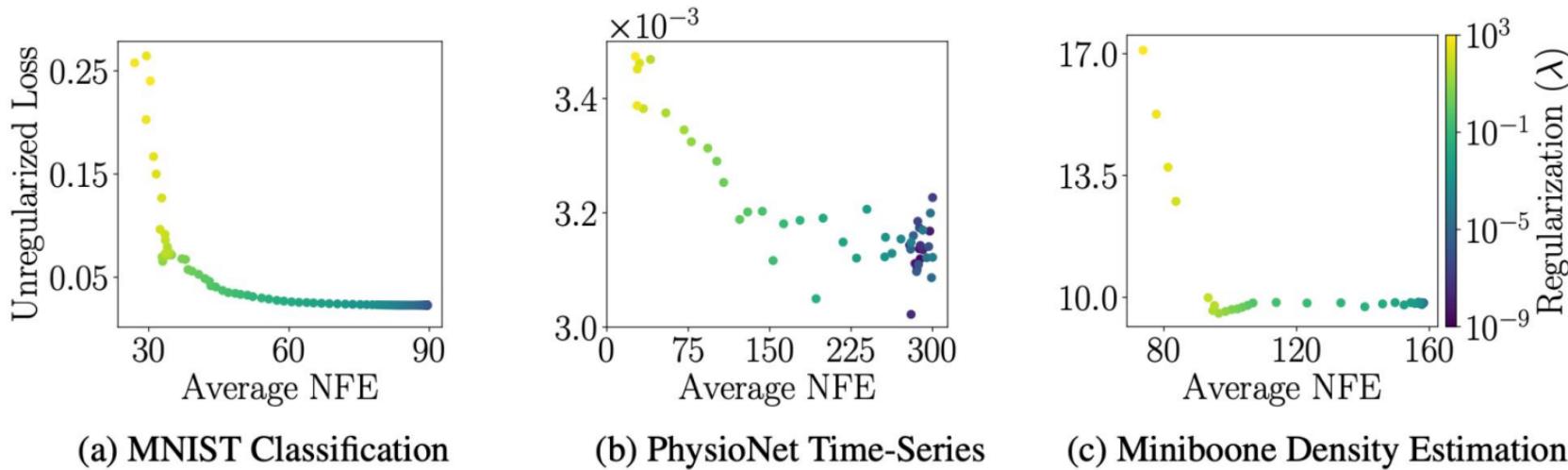


Figure 5: Tuning the regularization of \mathcal{R}_2 trades off between training loss and solver speed in three different applications of neural ODEs. Horizontal axes show average number of function evaluations, and vertical axes show unregularized training loss, both at the end of training.

controlling \mathcal{R}_2

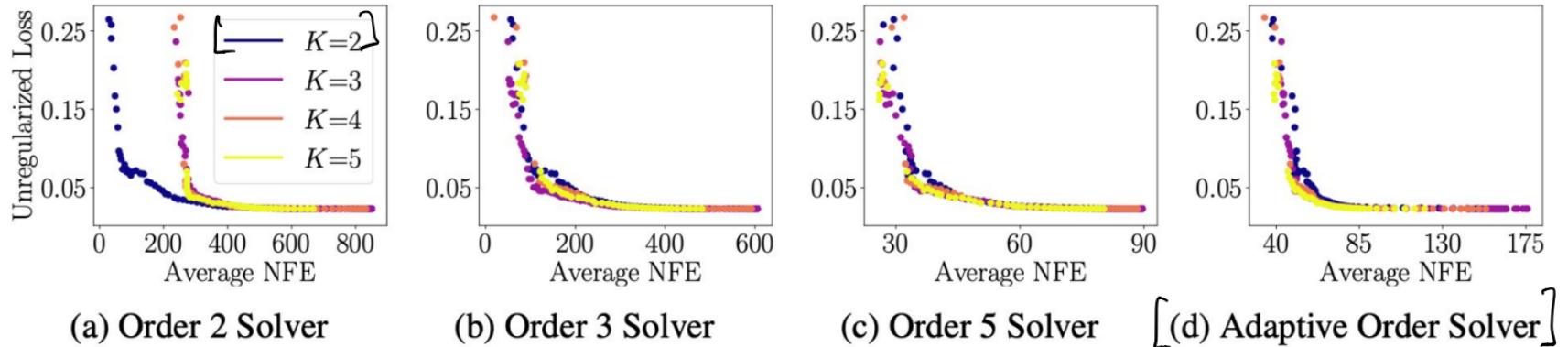


Figure 6: Comparing tradeoff between speed and performance when regularizing different orders. **6a)**: For a 2nd-order solver, regularizing the 2nd total derivative gives the best tradeoff. **6b)**: For a 3rd-order solver, regularizing the 3rd total derivative gives the best tradeoff, but the difference is small. **6c)**: For a 5th-order solver, results are mixed. **6d)**: For an adaptive-order solver, the difference is again small but regularizing higher orders works slightly better.