

Gaussian process learning of nonlinear dynamics

Scientific Machine Learning Workshop

Dongwei Ye, Mengwu Guo

d.ye-1@utwente.nl

Mathematics of Imaging & AI, Department of Applied Mathematics, University of Twente

6th December, 2023



Mathematics of Imaging & AI

**UNIVERSITY
OF TWENTE.**

- Data-driven learning of dynamical systems from time series is an important component in scientific machine learning, as it bridges the gap between data-driven approximation and physics-based modeling

$$\begin{cases} \dot{x}_1(t) = f_1(\mathbf{x}(t); \boldsymbol{\theta}_1) \\ \dot{x}_2(t) = f_2(\mathbf{x}(t); \boldsymbol{\theta}_2) \\ \dots \dots \\ \dot{x}_N(t) = f_N(\mathbf{x}(t); \boldsymbol{\theta}_N) \end{cases} \quad \text{with} \quad \mathbf{x}(t_0) = \mathbf{x}_0, \quad t \geq t_0.$$

- **Identification and estimation:**

Sparse identification for Nonlinear Dynamics (SINDy)*

The diagram illustrates the SINDy model structure. It shows a vector d (derivatives) equal to a matrix G (basis functions) multiplied by a vector θ (parameters). The vector d has three columns labeled \dot{x}_1 , \dot{x}_2 , and \dot{x}_3 . The matrix G has columns labeled 1 , x_1 , x_2 , and x_3^3 , followed by an ellipsis. The vector θ has three columns labeled θ_1 , θ_2 , and θ_3 . Each column is represented by a vertical bar with a color gradient.

$$\theta = \operatorname{argmin}_{\hat{\theta}} \frac{1}{2} \|d - G\theta\|_2^2 + R(\theta)$$

*Brunton SL, Proctor JL, Kutz JN. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. Proceedings of the national academy of sciences. 2016 Apr 12;113(15):3932-7.

- **Identification and estimation**

Operator Inference*

$$\begin{bmatrix} \dot{x}_1 & \dot{x}_2 \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_1 x_2 \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

$$\begin{matrix} d & G & \theta \end{matrix}$$

Lotka-Volterra

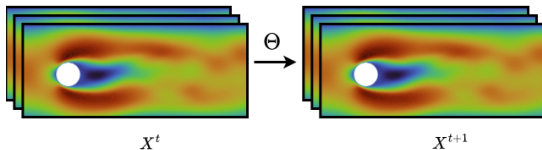
$$\begin{cases} \dot{x}_1 = \alpha x_1 - \beta x_1 x_2 \\ \dot{x}_2 = \delta x_1 x_2 - \gamma x_2 \end{cases}$$

$$\theta = \operatorname{argmin}_{\hat{\theta}} \frac{1}{2} \|d - G\theta\|_2^2 + R(\theta)$$

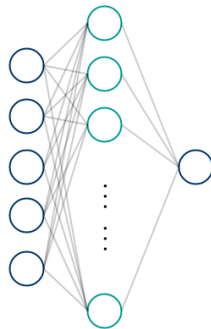
*Peherstorfer B, Willcox K. Data-driven operator inference for nonintrusive projection-based model reduction. Computer Methods in Applied Mechanics and Engineering. 2016 Jul 1;306:196-215.

- **Approximation**

dynamics mode decomposition*



NeuralODE**

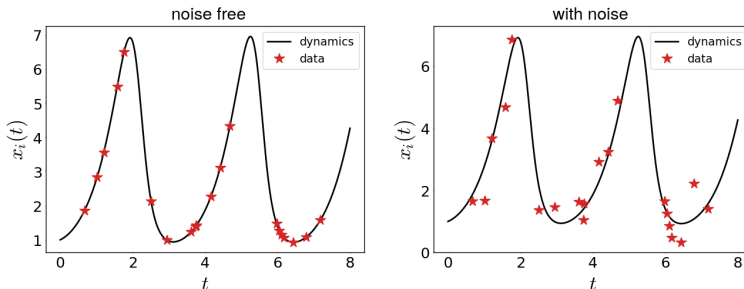


$$\frac{dx}{dt} = f(x(t), \theta)$$

*Schmid PJ. Dynamic mode decomposition of numerical and experimental data. Journal of fluid mechanics. 2010 Aug;656:5-28.

**Chen RT, Rubanova Y, Bettencourt J, Duvenaud DK. Neural ordinary differential equations. Advances in neural information processing systems. 2018;31.

- Methods such as SINDy or Operator Inference require time derivative of the state data, which is generally not available.
- However, the predictive performance of these dynamics learning techniques may be compromised when data are scarce and/or corrupted by noise.



- Could we also enable uncertainty quantification?

- GP: a stochastic process (a collection of random variables), any finite number of which have a joint Gaussian distribution, e.g.:

$$x_i(t) \sim \mathcal{GP}(0, \kappa_i(t, t')) .$$

- Observations $\{t_j, x_i(t_j)\}_{j=1}^T = \{\mathcal{T}, X_i(\mathcal{T})\}$
- Interpolation/regression

$$\begin{bmatrix} x_i(\mathcal{T}) \\ x_i(t^*) \end{bmatrix} \sim \mathcal{GP} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \kappa_i(\mathcal{T}, \mathcal{T}) & \kappa_i(\mathcal{T}, t^*) \\ \kappa_i(t^*, \mathcal{T}) & \kappa_i(t^*, t^*) \end{bmatrix} \right)$$

$$\Rightarrow x_i(t^*) | t^*, \mathcal{T}, X_i(\mathcal{T}) \sim \mathcal{N}(\kappa_i(t^*, \mathcal{T}) \kappa_i(\mathcal{T}, \mathcal{T})^{-1} X_i(\mathcal{T}),$$

$$\kappa_i(t^*, t^*) - \kappa_i(t^*, \mathcal{T}) \kappa_i(\mathcal{T}, \mathcal{T})^{-1} \kappa_i(\mathcal{T}, t^*))$$

$$\begin{cases} \dot{x}_1(t) = f_1(\mathbf{x}(t); \boldsymbol{\theta}_1) \\ \dot{x}_2(t) = f_2(\mathbf{x}(t); \boldsymbol{\theta}_2) \\ \dots \dots \\ \dot{x}_N(t) = f_N(\mathbf{x}(t); \boldsymbol{\theta}_N) \end{cases} \quad \text{with } \mathbf{x}(t_0) = \mathbf{x}_0, \quad t \geq t_0,$$

- Collect the observed data of $\{x_i(t_k)\}_{k=0}^{K-1}$ at the time-instances \mathcal{T} . Let $\mathbf{U} = [\mathbf{u}_1 \ \dots \ \mathbf{u}_N]^\top = [\mathbf{x}(t_0) \ \dots \ \mathbf{x}(t_{K-1})] \in \mathbb{R}^{N \times K}$ be the full-state data.
- $\mathbf{d}_i \in \mathbb{R}^K$ denote the time derivatives of x_i over \mathcal{T} .

Any linear transformation of a GP such as differentiation and integration is still a GP.

A vector-valued GP for likelihood:

$$\begin{bmatrix} x_i(t) \\ \dot{x}_i(t) \end{bmatrix} \sim \mathcal{GP} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \kappa_i(t, t') & \partial_{t'} \kappa_i(t, t') \\ \partial_t \kappa_i(t, t') & \partial_t \partial_{t'} \kappa_i(t, t') \end{bmatrix} \right).$$

$$\Rightarrow p(\mathbf{d}_i = f_i(\mathbf{U}; \boldsymbol{\theta}_i), \mathbf{u}_i | \boldsymbol{\theta}_i) \propto \exp \left(-\frac{1}{2} \begin{bmatrix} f_i(\mathbf{U}; \boldsymbol{\theta}_i) \\ \mathbf{u}_i \end{bmatrix}^\top \begin{bmatrix} \mathbf{K}_i^{dd} & \mathbf{K}_i^{du} \\ \mathbf{K}_i^{ud} & \mathbf{K}_i^{uu} \end{bmatrix}^{-1} \begin{bmatrix} f_i(\mathbf{U}; \boldsymbol{\theta}_i) \\ \mathbf{u}_i \end{bmatrix} \right),$$

$$\mathbf{K}_i^{dd} = \partial_t \partial_{t'} \kappa_i(\mathcal{T}, \mathcal{T}) + \chi_i^d \mathbf{I}_K \in \mathbb{R}^{K \times K},$$

$$\mathbf{K}_i^{uu} = \kappa_i(\mathcal{T}, \mathcal{T}) + \chi_i^u \mathbf{I}_K \in \mathbb{R}^{K \times K},$$

$$\mathbf{K}_i^{du} = \partial_t \kappa_i(\mathcal{T}, \mathcal{T}) = (\mathbf{K}_i^{ud}) \in \mathbb{R}^{K \times K}.$$

- Prior of parameters θ_i

$$p(\theta_i) \propto \exp \left(-\frac{\lambda_i}{\nu} \|\theta_i\|_\nu^\nu \right),$$

- Posterior via Bayesian inference

$$\begin{aligned} p(\theta_i | \mathbf{d}_i = f_i(\mathbf{U}; \theta_i), \mathbf{u}_i) \\ \propto p(\mathbf{d}_i = f_i(\mathbf{U}; \theta_i), \mathbf{u}_i | \theta_i) p(\theta_i) \\ \propto \exp \left(-\frac{1}{2} \left(f_i(\mathbf{U}; \theta_i)^\top \mathbf{R}_i^{dd} f_i(\mathbf{U}; \theta_i) + 2 f_i(\mathbf{U}; \theta_i)^\top \mathbf{R}_i^{du} \mathbf{u}_i + \frac{2\lambda_i}{\nu} \|\theta_i\|_\nu^\nu \right) \right). \end{aligned}$$

where

$$\begin{bmatrix} \mathbf{K}_i^{dd} & \mathbf{K}_i^{du} \\ \mathbf{K}_i^{ud} & \mathbf{K}_i^{uu} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{R}_i^{dd} & \mathbf{R}_i^{du} \\ \mathbf{R}_i^{ud} & \mathbf{R}_i^{uu} \end{bmatrix}$$

- Furthermore, we define an estimate of the time derivatives over \mathcal{T} (i.e., $\dot{x}_i(\mathcal{T})$) using a Gaussian process regression with the state observations $(\mathcal{T}, \mathbf{u}_i)$ at the same instances, given as

$$\hat{\mathbf{d}}_i := \mathbf{K}_i^{du} (\mathbf{K}_i^{uu})^{-1} \mathbf{u}_i.$$

- Considering the fact that $\mathbf{R}_i^{dd} \hat{\mathbf{d}}_i = -\mathbf{R}_i^{du} \mathbf{u}_i$

$$\propto \exp \left(-\frac{1}{2} \left(f_i(\mathbf{U}; \boldsymbol{\theta}_i)^\top \mathbf{R}_i^{dd} f_i(\mathbf{U}; \boldsymbol{\theta}_i) + 2 f_i(\mathbf{U}; \boldsymbol{\theta}_i)^\top \mathbf{R}_i^{du} \mathbf{u}_i + \frac{2\lambda_i}{\nu} \|\boldsymbol{\theta}_i\|_\nu^\nu \right) \right).$$

$$\propto \exp \left(-\frac{1}{2} \left(\left\| f_i(\mathbf{U}; \boldsymbol{\theta}_i) - \hat{\mathbf{d}}_i \right\|_{\mathbf{R}_i^{dd}}^2 + \frac{2\lambda_i}{\nu} \|\boldsymbol{\theta}_i\|_\nu^\nu \right) \right),$$

where $\|\mathbf{z}\|_{\mathbf{R}_i^{dd}} := \sqrt{\mathbf{z}^\top \mathbf{R}_i^{dd} \mathbf{z}}$, $\mathbf{z} \in \mathbb{R}^K$.

- A physically meaningful parametrization for the dynamical system is independent of initial conditions, so should the inference of parameters θ_i be:

$$p(\theta_i | \mathbf{d}_i = f_i(\mathbf{U}; \theta_i), \mathbf{u}_i) \propto \exp \left(-\frac{1}{2} \left(\sum_{\text{ICs}} \|f_i(\mathbf{U}; \theta_i) - \hat{\mathbf{d}}_i\|_{\mathbf{R}_i^{dd}}^2 + \frac{2\lambda_i}{\nu} \|\theta_i\|_{\nu}^{\nu} \right) \right).$$

- With shared parameters:

$$\begin{aligned} & p(\theta_i, \theta_j | \mathbf{d}_i = f_i(\mathbf{U}; \theta_i), \mathbf{d}_j = f_j(\mathbf{U}; \theta_j), \mathbf{u}_i, \mathbf{u}_j) \\ & \propto \exp \left(-\frac{1}{2} \left(\|f_i(\mathbf{U}; \theta_i) - \hat{\mathbf{d}}_i\|_{\mathbf{R}_i^{dd}}^2 + \|f_j(\mathbf{U}; \theta_j) - \hat{\mathbf{d}}_j\|_{\mathbf{R}_j^{dd}}^2 + \frac{2\lambda_i}{\nu} \|\theta_i\|_{\nu}^{\nu} + \frac{2\lambda_j}{\nu} \|\theta_j\|_{\nu}^{\nu} \right) \right). \end{aligned}$$

- Bayesian prediction by marginalizing over the posterior of parameters

$$p(\mathbf{x}(t) | \mathbf{d}_1, \mathbf{u}_1, \dots, \mathbf{d}_N, \mathbf{u}_N) = \int p(\mathbf{x}(t) | \theta_1, \dots, \theta_N) \prod_{i=1}^N p(\theta_i | \mathbf{d}_i = f_i(\mathbf{U}; \theta_i), \mathbf{u}_i) d\theta_i.$$

- Consider an affine parametrization of $f_i(\cdot; \theta_i)$ that is linear with respect to θ_i , i.e.,

$$f_i(\mathbf{x}; \theta_i) = \mathbf{g}_i(\mathbf{x})^\top \theta_i.$$

Hence, $f_i(\mathbf{U}; \theta_i)$ is rewritten as $\mathbf{G}_i \theta_i$ with $\mathbf{G}_i := [\mathbf{g}_i(\mathbf{x}(t_0)) \cdots \mathbf{g}_i(\mathbf{x}(t_{K-1}))]^\top \in \mathbb{R}^{K \times p_i}$

- The proposed inference method can be represented as:

$$\begin{aligned} p(\theta_i | \mathbf{d}_i = \mathbf{G}_i \theta_i, \mathbf{u}_i) &\propto \exp \left(-\frac{1}{2} \left(\left\| \mathbf{G}_i \theta_i - \hat{\mathbf{d}}_i \right\|_{\mathbf{R}_i^{dd}}^2 + \lambda_i \|\theta_i\|_2^2 \right) \right) \\ &\propto \exp \left(-\frac{1}{2} (\theta_i - \mu_i)^\top \Sigma_i^{-1} (\theta_i - \mu_i) \right), \end{aligned}$$

with mean vector:

$$\mu_i = (\theta_i)_{\text{MAP}} = (\mathbf{G}_i^\top \mathbf{R}_i^{dd} \mathbf{G}_i + \lambda_i \mathbf{I}_{p_i})^{-1} \mathbf{G}_i^\top \mathbf{R}_i^{dd} \hat{\mathbf{d}}_i,$$

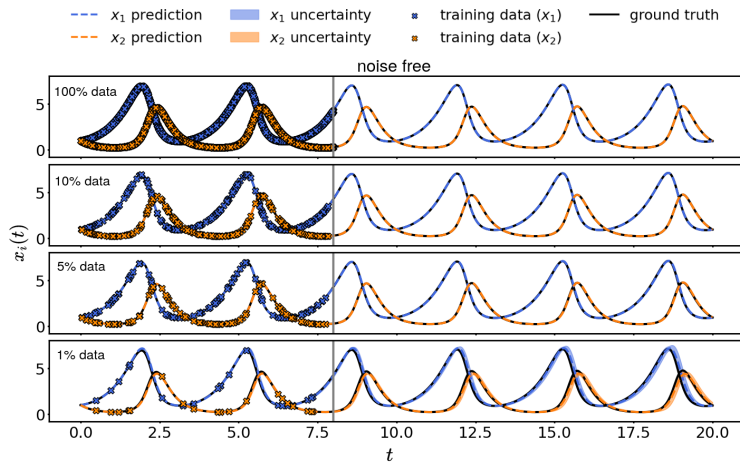
and the posterior covariance:

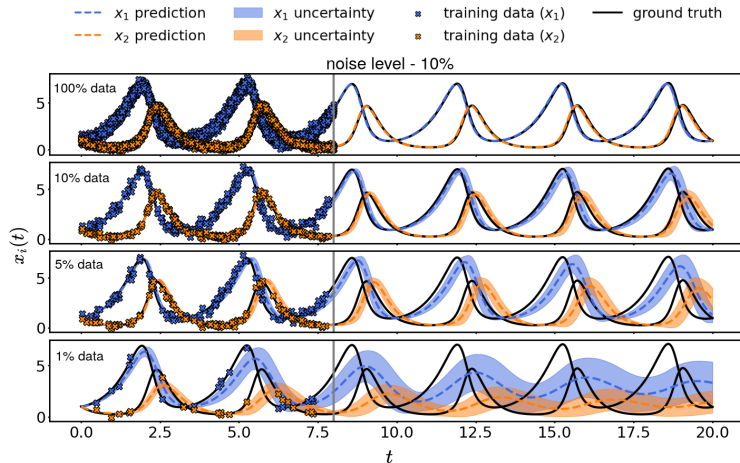
$$\Sigma_i = (\mathbf{G}_i^\top \mathbf{R}_i^{dd} \mathbf{G}_i + \lambda_i \mathbf{I}_{p_i})^{-1}.$$

- Example: Lotka-Volterra model

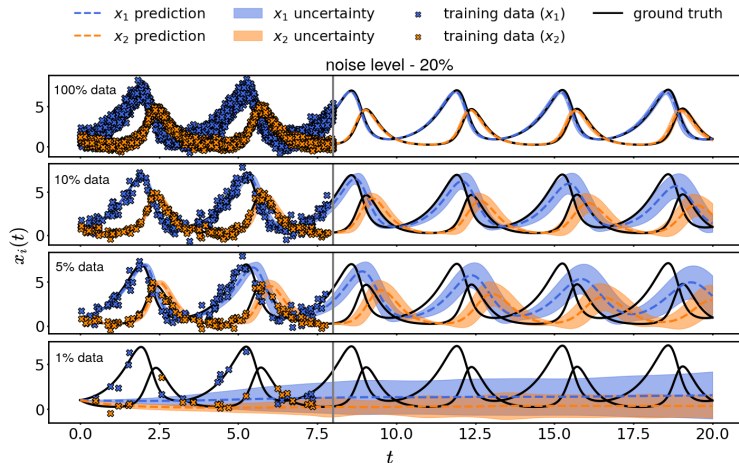
$$\begin{cases} \dot{x}_1 = \alpha x_1 - \beta x_1 x_2, \\ \dot{x}_2 = \delta x_1 x_2 - \gamma x_2, \end{cases}$$

$$\Rightarrow \quad \begin{aligned} \mathbf{g}_1(\mathbf{x}) &= [x_1 \quad x_1 x_2]^\top, & \boldsymbol{\theta}_1 &= [\alpha \quad -\beta]^\top; \\ \mathbf{g}_2(\mathbf{x}) &= [x_1 x_2 \quad x_2]^\top, & \boldsymbol{\theta}_2 &= [\delta \quad -\gamma]^\top. \end{aligned}$$





Case I: Lotka-Volterra model

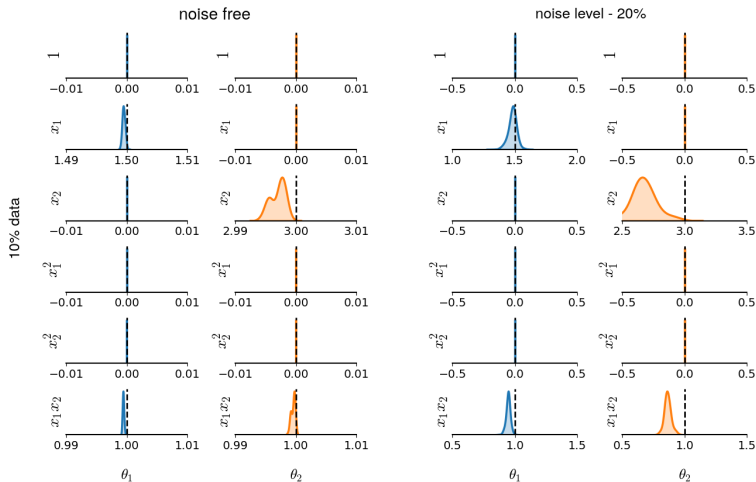


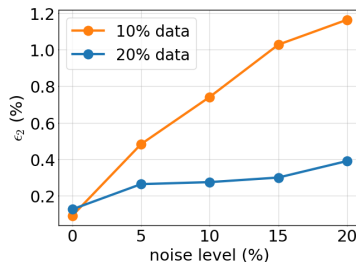
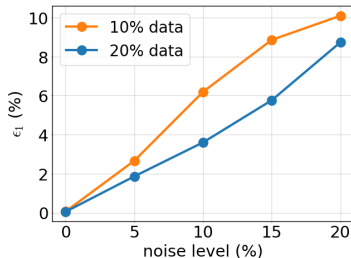
- The proposed inference method can be represented as:

$$p(\boldsymbol{\theta}_i | \mathbf{d}_i = \mathbf{G}_i \boldsymbol{\theta}_i, \mathbf{u}_i) \propto \exp \left(-\frac{1}{2} \left(\left\| \mathbf{G}_i \boldsymbol{\theta}_i - \hat{\mathbf{d}}_i \right\|_{\mathbf{R}_i^{dd}}^2 + 2\lambda_i \|\boldsymbol{\theta}_i\|_1 \right) \right).$$

$$\mathbf{g}_1(\mathbf{x}) = [1 \ x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1 x_2]^\top, \quad \boldsymbol{\theta}_1 = [0 \ \alpha \ 0 \ 0 \ 0 \ -\beta]^\top;$$

$$\mathbf{g}_2(\mathbf{x}) = [1 \ x_1 \ x_2 \ x_1^2 \ x_2^2 \ x_1 x_2]^\top, \quad \boldsymbol{\theta}_2 = [0 \ 0 \ -\gamma \ 0 \ 0 \ \delta]^\top$$





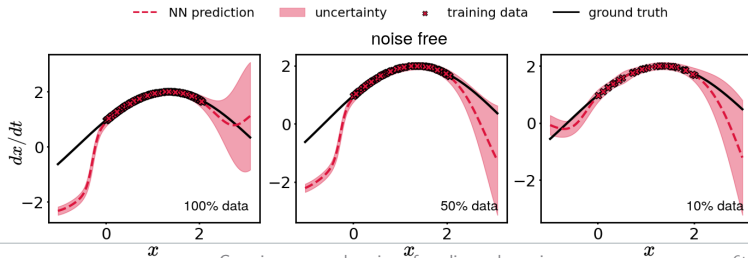
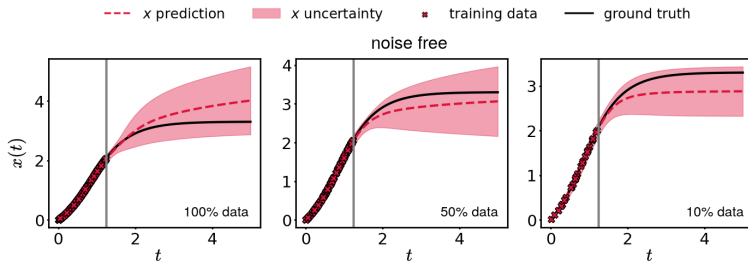
$$\epsilon_1 = 100\% \times \left(\frac{\sum_{i=1}^2 \|\boldsymbol{\theta}_i - \mathbb{E}_{\text{post}}[\boldsymbol{\theta}_i]\|_2^2}{\sum_{i=1}^2 \|\boldsymbol{\theta}_i\|_2^2} \right)^{1/2}, \quad \epsilon_2 = 100\% \times \left(\frac{\sum_{i=1}^2 \text{Var}_{\text{post}}[\boldsymbol{\theta}_i]}{\sum_{i=1}^2 \|\boldsymbol{\theta}_i\|_2^2} \right)^{1/2}$$

Nonlinear parametrization with a shallow neural network

$$f_i(\mathbf{x}; \boldsymbol{\theta}_i) = \sum_{l=1}^L v_{il} \sigma(\mathbf{w}_{il}^T \mathbf{x} + b_{il}) , \quad \text{with } \boldsymbol{\theta}_i := \{v_{il}, b_{il}, \mathbf{w}_{il}\}_{l=1}^L ,$$

- 1D synthetic example

$$\dot{x} = \gamma \sin(\alpha x + \beta)$$



- Does not require a direct finite-difference estimation of time-derivatives from solution data as in Oplnf or SINDy, but instead (implicitly) evaluates the derivatives Gaussian process approximation.
- Improves predictiveness and facilitates uncertainty quantification for dynamics learning with noisy and/or scarce data.
- Key features: simplicity, robustness, generality.

Thank you for your attention!

Questions?