Recent Advances in Learning from High-Dim Dynamical Systems

Ming Zhong

Department of Mathematics University of Houston

April 14, 2025

High Dim Data Analaysis Group

Ming Zhong (UH)

High Dim Data

Table of Contents





Important Applications



Important Applications



Important Applications



cardiac events earlier and save more lives. By Nicole Foster

First some freme i transmission for Hypotham have been enging to get term have to be to dega a dega term have attack. Traditionally, be been relief on tradition and the transmission of the dega and the tradition of the problem of the dega and the tradition of the problem of the dega and the dega and the problem of the dega and the dega and the problem of the dega and the dega and the problem of the dega and the dega and the problem of the dega and the dega and the problem of the dega and the dega and the dega and dega and the dega and the dega and where a dega and the dega and th

Dr. Stephen Weng is an assistant y sor at University of Nortingham (

"

02

UH 2025

Prediction

Given data in the form

$$\{\mathbf{z}_{t_1}, \mathbf{z}_{t_2}, \cdots, \mathbf{z}_{t_L}\}, \quad \mathbf{z}_{t_l} \in \mathbb{R}^D (D \gg 1),$$

with $0 = t_1 < t_2 < \cdots < t_L = T$, can we predict

$$\mathbf{z}_{t_{L+1}}, \mathbf{z}_{t_{L+2}}, \cdots, \quad t_{L+1} = T + \delta t$$

Traditional Methods

 Autoregressive (AR), moving average (MA), ARMA and ARIMA models

Machine Learning methods

• Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), Transformers

Ming Zhong (UH)

ī.

Time Series Data

Dynamical Data

Basic assumption (Markov Chain)

$$\Pr(\mathbf{z}_{L+1} = z_{L+1} | \mathbf{z}_1 = z_1, \mathbf{z}_2 = z_2, \cdots, \mathbf{z}_L = z_L)$$

=
$$\Pr(\mathbf{z}_{L+1} = z_{L+1} | \mathbf{z}_L = z_L)$$

Even further (dynamic structure)

$$\mathbf{z}_{L+1} = \mathbf{z}_L + \mathbf{h}(\mathbf{z}_L)\Delta t + \boldsymbol{\sigma}(\mathbf{z}_L)\Delta \mathbf{w}_L,$$

where

•
$$\mathbf{h}: \mathbb{R}^D o \mathbb{R}^D; \, \boldsymbol{\sigma}: \mathbb{R}^D o \mathbb{R}^{D imes D}.$$

• $\Delta \mathbf{w}_L \in \mathbb{R}^D$ is a Brownian motion.

Dynamics Assumption

Continuously

$$\mathrm{d}\mathbf{z}_t = \mathbf{h}(\mathbf{z}_t)\,\mathrm{d}t + \boldsymbol{\sigma}(\mathbf{z}_t)\,\mathrm{d}\mathbf{w}_t,$$

Learning and Prediction

Given data $\{\mathbf{z}_{t_1}, \mathbf{z}_{t_2}, \cdots, \mathbf{z}_{t_L}\}$, can we predict \mathbf{z}_{L+1} , assuming

$$d\mathbf{z}_t = \mathbf{h}(\mathbf{z}_t) dt + \boldsymbol{\sigma}(\mathbf{z}_t) d\mathbf{w}_t,$$

with h and σ being unknown.

Established approaches

• SINDy, Neural ODE, PINN/PIGP.

Dynamics Assumption

These methods use the normal regression setup

$$\mathcal{E}(\tilde{\mathbf{h}}) = \mathbb{E}\left[\frac{1}{T}\int_{0}^{T}\left|\frac{\mathrm{d}\mathbf{z}_{t}}{\mathrm{d}t} - \tilde{\mathbf{h}}(\mathbf{z}_{t})\right|_{\mathbb{R}^{D}}^{2}\mathrm{d}t.
ight].$$

Remark:

- $\frac{dz_t}{dt}$ is understood loossely.
- $\hat{h} = \operatorname{argmin}_{\tilde{h} \in \mathcal{H}} \mathcal{E}(\tilde{h})$ for some space \mathcal{H} .

Since $D \gg 1$,

- SINDy assumes $\mathbf{h}(\mathbf{z}) \approx \sum_{\eta=1}^{n} \psi_{\eta}(\mathbf{z})$ with sparse coefficients, and pre-defined dictionary ψ_{η} .
- Neural ODE/PINN use deep neural network with over-parametrization and stochastic gradient descent.

Ming Zhong (UH)

High Dim Data

Dynamics Assumption

However

- For SIDNy: the dicionary is chosen to be large enough to contain enough assumptions (non-linearity, derivatives).
 - Sparsity is enforced to reduce the search time.
 - Derivatives can be weaken using weak-SINDy.
 - But "large" is very subjective.
- Deep learning
 - Estimators might have good prediction capability.
 - Estimators do not have clear physical structures (i.e. Gravity is $1/r^2$).
 - Training can be an issue; Setup of the hyperparameters are bit engineering-like.

Table of Contents





Dynmical Times Series Data Our Approach

Recall $d\mathbf{z}_t = \mathbf{h}(\mathbf{z}_t) dt + \boldsymbol{\sigma}(\mathbf{z}_t) d\mathbf{w}_t$, we consider

- $\sigma = 0 * I_{D \times D}$, determinstic system, **h** is a high-dim function.
- $\sigma = 0 * I_{D \times D}$, determinstic system, **h** is a differential operator.
- $\sigma = \sigma(z)$ is a SPD matrix, **h** is a high-dim function or a differential operator.
- $\sigma = \sigma(\mathsf{z})$ is a singular matrix, h is a high-dim function.

Learning and Prediction

We build different estimating algorithms when both **h** and σ are unknown based on the scenarios, especially on how to handle the curse of dimensionality.

Interacting agent systems for $\mathbf{x}_i \in \mathbb{R}^d$,

$$\frac{\mathsf{d}\mathbf{x}_i}{\mathsf{d}t} = \frac{1}{N}\sum_{j=1, j\neq i}^N \phi(|\mathbf{x}_j - \mathbf{x}_i|)(\mathbf{x}_j - \mathbf{x}_i), \quad i = 1, \cdots, N,$$

With the setup

$$\mathbf{z} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_N \end{bmatrix}, \ \mathbf{h}_{\phi}(\mathbf{z}) = \begin{bmatrix} \vdots \\ \frac{1}{N} \sum_{j=1, j \neq i}^N \phi(|\mathbf{x}_j - \mathbf{x}_i|)(\mathbf{x}_j - \mathbf{x}_i) \\ \vdots \end{bmatrix}$$

Then $\dot{\mathbf{z}} = \mathbf{h}_{\phi}(\mathbf{z})$ with $\mathbf{z} \in \mathbb{R}^{D}$ with D = Nd.

11/26

٠

Instead of learning \mathbf{h}_{ϕ} , we learn ϕ instead from

$$\mathcal{E}(\varphi) = \mathbb{E} ig[rac{1}{T} \int_0^T ig| \dot{\mathbf{z}}(t) - \mathbf{h}_{\phi}(\mathbf{z}(t)) ig|_N^2 dt ig].$$

•
$$\hat{\phi} = \operatorname{argmin}_{\varphi \in \mathcal{H}} \mathcal{E}(\varphi) \approx \phi.$$

- 1D leanring rate.
- \mathbf{h}_{ϕ} has a structure and phisycal meaning (interaction).
- When ${\cal H}$ is finite dimensional, the learning problem becomes solving a linear system.
- Review: Learning Collective Behaviors from Observation, Explorations in the Mathematics of Data Science, 2024.

Ming Zhong (UH)

High Dim Data

Next we consider

$$\dot{z} = h(z) = \mathcal{P}(z), \quad \mathcal{P} \text{ is a differential operator.}$$

For example, compressible Euler system

$$\frac{\partial}{\partial t} \begin{bmatrix} \rho \\ \mathbf{u} \\ e \end{bmatrix} + \begin{bmatrix} \mathbf{u} \cdot \nabla \rho \\ \mathbf{u} \cdot \nabla \mathbf{u} \\ \mathbf{u} \cdot \nabla e \end{bmatrix} + \begin{bmatrix} \rho \nabla \cdot \mathbf{u} \\ \nabla p / \rho \\ p / \rho \nabla \cdot \mathbf{u} \end{bmatrix} = \begin{bmatrix} \mathbf{0} \\ \mathbf{G} \\ \mathbf{0} \end{bmatrix}$$

ρ(t, x) (density), u(t, x) (fluid velocity), e(t, x) (specific internal energy), p(t, x) (presure), G(t, x) (gravitational acceleration).

Consider for $(t, \mathbf{x}) \in \Omega$

$$\begin{split} \dot{\mathbf{z}} &= \mathbf{h}(\mathbf{z}) = \mathcal{P}(\mathbf{z}) \\ \text{subject to } \mathbf{z}(0, \mathbf{x}) = \mathbf{z}_0(\mathbf{x}) \text{ and } \mathcal{B}(\mathbf{z}) = \mathbf{g}(t, \mathbf{x}) \text{ on } \partial\Omega. \end{split}$$

Learn $\mathbf{z}(t, \mathbf{x})$ for $(t, \mathbf{x}) \in \Omega$ through

 $Loss = Data Loss + \lambda * Physics Loss.$

where

Physics Loss = PDE/ODE/Integral Loss + IC Loss + BC Loss.

Regularized Leanring (Supervised (Data) + Unsupervised (Physics)).

Ming Zhong (UH)

Continue on

PDE Loss =
$$\frac{1}{|\Omega|} \int_{\Omega} \left| \frac{\partial \mathbf{z}_{NN}}{\partial t} - \mathcal{P}(\mathbf{z}_{NN}) \right|_{L^{2}(\Omega)}^{2} \mathrm{d}\Omega.$$

and

$$\mathsf{IC}\;\mathsf{Loss} = \frac{1}{|\Omega_x|} \int_{\Omega_x} \left| \boldsymbol{\mathsf{z}}_{\textit{NN}} - \boldsymbol{\mathsf{z}}_0 \right|^2_{\textit{L}^2(\Omega_x)} \,\mathsf{d}\Omega_x.$$

lastly

$$\mathsf{BC}\;\mathsf{Loss} = \frac{1}{|\partial\Omega|}\int_{\partial\Omega} \big|\mathcal{B}(\mathbf{z}_{\mathit{NN}}) - \mathbf{g}\big|^2_{\mathit{L}^2(\partial\Omega)}\;\mathsf{d}\partial\Omega.$$

The addiitional information from Physics improves prediction capability.



Remarks:

- With additional data in Ω, one can also infer P (or some parametric structure of P).
- IC loss can also considered as a type of data loss (initial data).
- BC loss can be considered as data loss if BC type is Dirichlet.
- Many types of losses lead to multi-objective optimization.
 - PDE loss has larger (in magnitude) gradient, forcing the minimizer to minimize PDE first, ignoring IC and BC, difficult at solving stiff PDEs.
 - In the case of Hyperbolic PDEs (capturing shocks), how to capture the physically meaningful weak solution?
 - Classic formulation and L²-norm are used, leadning to other problems.
 - Meshless loss leads to loss of time flow.

Ming Zhong (UH)

High Dim Data

Some remedies

- Adding data-driven artificial viscosity map to Hyperbolic PDEs: Physics-informed neural networks with adaptive localized artificial viscosity, JCP, 2023.
- Adding hard constrain or building BC into the architecture: Structure Preserving PINN for Solving Time Dependent PDEs with Periodic Boundary, arXiv, 2024.
- Mixture of Experts (MOE), label propagation: Label Propagation Training Schemes for Physics-Informed Neural Networks and Gaussian Processes, arXiv, 2024.

Physics Informed Transformer with Michael Holland on 4/28.

Back to the stochastic case

 $d\mathbf{z}_t = \mathbf{h}(\mathbf{z}_t) dt + \boldsymbol{\sigma}(\mathbf{z}_t) d\mathbf{w}_t$, **h** and $\boldsymbol{\sigma}$ are unknown.

But we assume $\sigma : \mathbb{R}^D \to \mathbb{R}^{D \times D}$ is Symmetric Positive Definite for all z.

• It is possible to learn **h** in the regression setting, since $d\mathbf{z}_t \approx \mathbf{h}(\mathbf{z}_t) dt$, we can

$$\mathbb{E}\big[|\,\mathrm{d}\mathbf{z}_t-\mathbf{h}(\mathbf{z}_t)\,\mathrm{d}t|_{\ell_2}\big].$$

However it does not use any of the information from the noise matrix, and misses the interaction of the noise and the drift.

Hence, we design the loss function as

$$egin{aligned} \mathcal{E}(ilde{\mathbf{h}}) &= \mathbb{E}igg[rac{1}{2}(\int_0^T < ilde{\mathbf{h}}(\mathbf{z}_t), \Sigma^\dagger ilde{\mathbf{h}}(\mathbf{z}_t) > \mathsf{d}t \ &-2 < ilde{\mathbf{h}}(\mathbf{z}_t), \Sigma^\dagger \,\mathsf{d}\mathbf{z}_t >)igg]. \end{aligned}$$

where

- $\Sigma = \sigma \sigma^{\top}$ and Σ^{\dagger} is the pseudo-inverse of Σ (in this case, just the normal inverse).
- $\hat{\mathbf{h}} = \operatorname{argmin}_{\tilde{\mathbf{h}} \in \mathcal{H}} \mathcal{E}(\tilde{\mathbf{h}}) \approx \mathbf{h}$.
- One can think of this is almost

$$|\,\mathrm{d}\mathbf{z}_t - \mathbf{h}(\mathbf{z}_t)\,\mathrm{d}t|_{\sigma^\dagger}$$

Taking into consideration of the effect of the noise.

Ming Zhong (UH)

High Dim Data

When $\Sigma = \sigma * \mathbf{I}_{D \times D}$, the loss simplifies to

$$\begin{split} \mathcal{E}(\tilde{\mathbf{h}}) &= \mathbb{E}\big[\frac{1}{2\sigma^2}(\int_0^T < \tilde{\mathbf{h}}(\mathbf{z}_t), \tilde{\mathbf{h}}(\mathbf{z}_t) > \mathsf{d}t \\ &-2 < \tilde{\mathbf{h}}(\mathbf{z}_t), \mathsf{d}\mathbf{z}_t >)\big]. \end{split}$$

Then it is very similar to regression.

• When
$$\Sigma = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_D \end{bmatrix}$$
, the components of **h** is comparable.

• Our method is more effective when σ is not-diagonal.

Ming Zhong (UH)

We have two papers

- Learning Stochastic Dynamics from Data, Guo, Cialenco, Zhong, ICLR, 2024.
- Noise Guided Structural Learning from Observing Stochastic Dynamics, Guo, Cialenco, Zhong, revision, 2025.
- Noise is learned through a modified Qudratic Variation method.

Henry Guo will discuss more about this method on 4/28 on noise interaction and Stochastic Heat Equation.

Lately, we have

$$d\mathbf{z}_t = \mathbf{h}(\mathbf{z}_t) dt + \boldsymbol{\sigma}(\mathbf{z}_t) d\mathbf{w}_t, \quad \boldsymbol{\sigma} \text{ is singular.}$$

WLOG, we assume

$$oldsymbol{\sigma} = egin{bmatrix} oldsymbol{0} & oldsymbol{0} \ oldsymbol{0} & oldsymbol{\sigma}_y \end{bmatrix}$$
 ,

Remark: if not, we juse do SVD, and rotate z_t accordingly. Hence

$$egin{aligned} & \mathsf{d}\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) \, \mathsf{d}t, \ & \mathsf{d}\mathbf{y}_t = \mathbf{g}(\mathbf{x}_t, \mathbf{y}_t) \, \mathsf{d}t + \boldsymbol{\sigma}_y(\mathbf{y}_t) \, \mathsf{d}\mathbf{w}_t^y. \end{aligned}$$

Here σ_y is SPD.

When we set

$$egin{aligned} \mathbf{z} &= egin{bmatrix} \mathbf{x} \ \mathbf{y} \end{bmatrix}, \quad \mathbf{h}(\mathbf{z}) &= egin{bmatrix} \mathbf{f}(\mathbf{x},\mathbf{y}) \ \mathbf{g}(\mathbf{x},\mathbf{y}) \end{bmatrix} \ && \sigma &= egin{bmatrix} \mathbf{0} & \mathbf{0} \ \mathbf{0} & \sigma_{_{Y}} \end{bmatrix}, \end{aligned}$$

and

It goes back to the original high-dim SDE but σ is singular if the normal loss

$$\begin{split} \mathcal{E}(\tilde{\mathbf{h}}) &= \mathbb{E}\big[\frac{1}{2\sigma^2}(\int_0^T < \tilde{\mathbf{h}}(\mathbf{z}_t), \tilde{\mathbf{h}}(\mathbf{z}_t) > \mathrm{d}t \\ &-2 < \tilde{\mathbf{h}}(\mathbf{z}_t), \mathrm{d}\mathbf{z}_t >)\big]. \end{split}$$

is used, the learning will collapse down to only learning \mathbf{g} .

Hence, we treat them separately

$$\mathcal{E}_{f}(\tilde{\mathbf{f}}) = \mathbb{E}\Big[\frac{1}{T}\int_{0}^{T}\Big|\frac{\mathrm{d}\mathbf{x}_{t}}{\mathrm{d}t} - \tilde{\mathbf{f}}(\mathbf{x}_{t},\mathbf{y}_{t})\Big|_{\ell_{2}}^{2}\Big].$$

and

$$egin{aligned} \mathcal{E}_{g}(ilde{\mathbf{g}}) &= \mathbb{E}igg[rac{1}{2}(\int_{0}^{t} < ilde{\mathbf{g}}(\mathbf{x}_{t},\mathbf{y}_{t}), \Sigma_{y}^{\dagger} ilde{\mathbf{g}}(\mathbf{x}_{t},\mathbf{y}_{t}) > \mathrm{d}t \ &- 2 < ilde{\mathbf{g}}(\mathbf{x}_{t},\mathbf{y}_{t}), \Sigma_{y}^{\dagger}\,\mathrm{d}\mathbf{y}_{t} >)igg]. \end{aligned}$$

We will learn σ_{y} from \mathbf{y}_{t} .

- A unified algorithm
 - Given {z_t}_{t∈[0,T]}, we use quadratic variation on z_t to figure out σ.
 - If $\sigma = 0$, we learn it deterministically.
 - If σ is non-singular, we learn it stochasticaly
 - If σ is singular, we perform SVD and find out the x and y direction, then we learn it mixed.

Henry will also discuss that on 4/28.