

# Couple of transmission models and deep learning techniques: brief introduction

Pengfei Song

Xi'an Jiaotong University



October 21, 2023

Joint work with Prof Jianhong Wu (York University), Prof Yanni Xiao (XJTU) and Shuangshuang Yin (XJTU)

## 1 Background

## 2 Universal Differential Equations

- Application One: Estimating time vary reproduction number
- Application Two: Learning Unknown Mechanisms
- Application Three: Optimal Control
- Application Four: Misinformation Project

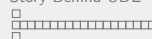
## 3 Story Behind UDE

- Neural Networks
- PINN

## 4 Discussion







# Data driven methods to learn unknown mechanism

- Mechanisms can be characterized by ‘functions’, ‘operators’, ‘Distributions’, ‘Stochastic Processes’. For example,  $f(x) = \exp(x)$  describes exponential growth.
- How to learn mechanisms? Find the function or surrogates of the function from data. Approximation!  $\exp(x) \approx 1 + x + x^2/2 + \dots$  Think about Taylor expansion, Fourier expansion.

**Functions:** Neural Networks, Random feature model (Reservoir Computing, ELM), GPs, Kernel Methods (e.g., SVM), Polynomials, Decision Trees

**Operators:** DeepOnets, Neural Operator,

**Distributions:** GAN, VAE, Auto-regressive models, Normalizing Flows, Diffusion Models, Energy Based Models, Consistency models

**Stochastic Processes:** infinitely deep bayesian neural network as neural SDE



# Background

- Applications of "AI For Science" on Mathematical Epidemiology
- AI4S is a direction combining fundamental researches (Mathematics, Physics, Chemistry, ...) with machine learning techniques.
- New revolutionary scientific research paradigm (Data and Mechanism driven Model)
- More on AI4S: 2022 International Congress of Mathematician 60 Minutes Talk:  
Weinan E: A Mathematical Perspective on Machine Learning
- Scientific Machine Learning (SciML); Scientific Artificial Intelligence (SciAI)

# Outline

## 1 Background

## 2 Universal Differential Equations

- Application One: Estimating time vary reproduction number
- Application Two: Learning Unknown Mechanisms
- Application Three: Optimal Control
- Application Four: Misinformation Project

## 3 Story Behind UDE

- Neural Networks
- PINN

## 4 Discussion

# Universal Differential Equations

"Universal" means "universal approximators" (neural networks, GPs, SVM, random feature models, ...)

UDEs (proposed by Prof. Christopher Rackauckas, MIT, 2020(Christopher Rackauckas et.al. 2020 arxiv) are initial value problems with the following forms:

$$\frac{du}{dt} = f_{\theta_2}(u, t, NN_{\theta_1}(u, t)),$$

where  $f$  is a known mechanistic model and  $NN$  denotes the missing or unknown terms,  $\theta_1$  and  $\theta_2$  are parameters of known mechanisms and neural networks, respectively, which can be estimated simultaneously.

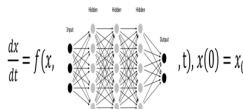
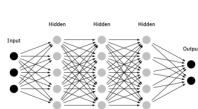
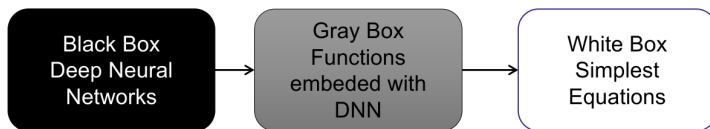




# Black-box, Gray-box, White-box

How to recover the simplest function from deep neural networks?

Equation search methods.



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

$$t), x(0) = x_0$$

$$\begin{cases} \frac{dS}{dt} = -S\beta \exp(-\alpha I) I^k, \\ \frac{dI}{dt} = S\beta \exp(-\alpha I) I^k - \gamma I, \end{cases}$$

<sup>1</sup>Song, Pengfei and Xiao, Yanni and Wu, Jianhong. (2023) Discovering first-principle behavior change transmission models by deep learning methods. One Chapter of Springer Book. Accepted



# Equation Search Methods: Sparse identification of nonlinear dynamic systems

SINDy applies a set of **candidate functions**  $\Theta(\mathbf{U})$  that would characterize the right-hand side of the governing equations,  $\mathbf{u}' = \mathbf{f}(\mathbf{u}) \approx \Theta(\mathbf{u})\Xi$ , and estimate  $\Xi$  by **sparse regression**.

$$\begin{cases} \frac{dS}{dt} = -\beta SI, \\ \frac{dI}{dt} = \beta SI - \gamma I. \end{cases} \quad (3)$$

We choose the basic functions as

$$\Theta([S, I]) = [S, I, SI, S^2I, S^2I^2, S^2I^2]$$

and now we use SINDy to discover the true equations.

<sup>1</sup>Brunton S L, Proctor J L, Kutz J N. Discovering governing equations from data by sparse identification of nonlinear dynamical systems[J]. Proceedings of the national academy of sciences, 2016, 113(15):3932-3937.



# Application One: Estimating time vary reproduction number

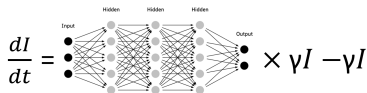
Represent effective reproduction number  $\mathcal{R}_t$  as

$$\mathcal{R}_t = \text{NeuralNetwork}_{\theta}(t, I), \quad (4)$$

and transmission model

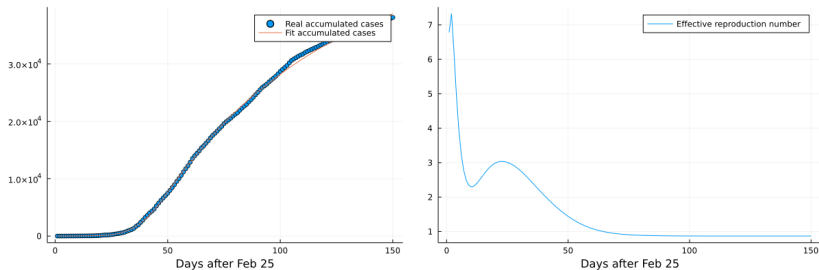
$$\begin{cases} I' = \gamma \text{NeuralNetwork}_{\theta}(t, I)I - \gamma I, \\ H' = \gamma \text{NeuralNetwork}_{\theta}(t, I)I, \end{cases} \quad (5)$$

where  $I(t)$  and  $H(t)$  denote the number of infected individuals and accumulated confirmed cases at time  $t$ .

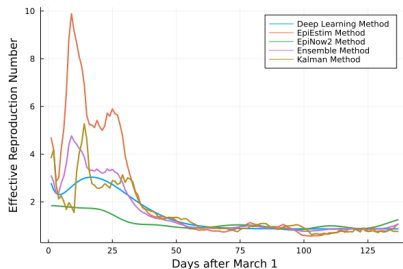


<sup>1</sup>Song, Pengfei and Xiao, Yanni. (2022) Estimating time-varying reproduction number by deep learning techniques (Dedicated to Prof Jibin Li on his 80th birthday). JAAC

**Figure:** Left: Ontario's first wave COVID-19 case data. Right: effective reproduction number estimation by deep learning method.



**Figure:** Left: Ontario's first wave COVID-19 case data fit by different methods. Right: effective reproduction number estimation by different methods.





# Application One: Estimating time vary reproduction number

Methods	Data source	Smooth	Speed	Accuracy of $\mathcal{R}_t$
Deep Learning	Case data, infection period	strong	slow (3682s)	strong
State Space	Case data, infection period	weak	quick (< 1s)	weak
EpiEstim	Case data, serial interval	weak	quick (< 1s)	weak
EpiNow2	Case data, generation time, incubation period, delay distribution	normal	slow (2578s)	strong

**Table:** Comparison of different estimation methods: deep learning, state space, EpiEstim, EpiNow2. Methods. Smooth measures the data fitting abilities.

## Application Two: Learning Unknown Human Behaviour Change Mechanisms

we will use the data of Ontario to fit the following neural differential equation model:

$$\frac{dS}{dt} = -\text{abs}(NN(I, R))S/N,$$

$$\frac{dI}{dt} = -\text{abs}(NN(I, R))S/N - \gamma I,$$

$$\frac{dR}{dt} = \gamma I,$$

$$\frac{dH}{dt} = \text{abs}(NN(I, R))S/N,$$

where  $NN(I, R)$  denotes neural network to learn the human behaviour change, and  $H$  denotes accumulated cases.

<sup>1</sup>Song, Pengfei and Xiao, Yanni and Wu, Jianhong. (2023) Discovering first-principle behavior change transmission models by deep learning methods. One Chapter of Springer Book. Accepted











# Application Three: Optimal Control

Methods	Direct method	Indirect method	HJB method	Deep learning
Transcriptions	Nonlinear programming problem (NLP)	Two point boundary value problem (TPBVP)	Dynamic programming	Parameter optimization
Trajectory or Parameter Optimization	Trajectory	Trajectory	Trajectory	Parameter
$u(x, t)$ or $x(u, t)$	-	$x(u, t)$	$u(x, t)$	$x(u, t)$
OtD or DtO	DtO	OtD	DtO	DtO and OtD
Using frequency	Most often	Often	Seldom	Seldom
Advantages	Mature Optimizers, easy to post, easy to solve		Accurate	Flexible, extendable Bless of dimensionality
Disadvantages	Less accurate	Hard to post, hard to solve, initial guess	Curse of dimensionality	Theoretically under exploring,



## Application Four: Misinformation Project

Question: the impact of misinformation on disease spread and vaccination decision?

- Information Data: classification of information(区分正确和错误信息): NLP(自然语言处理) such as Bert, ChatGPT.
- Evolution of correct and misinformation(未知的信息演化规律): Neural ODE
- Impact of misinformation on disease spread and vaccination decision(信息对传染病传播和疫苗接种未知的影响): UDE

## Application Four: Misinformation Project

## Disease Model

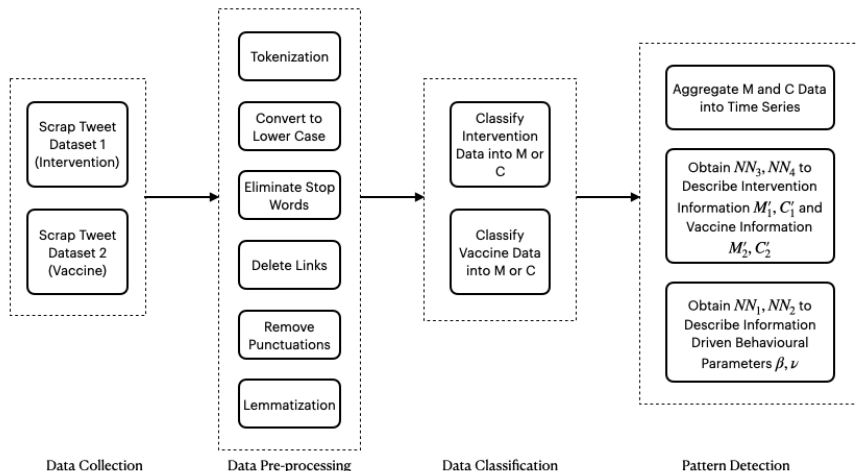
$$\begin{cases} S' &= -\beta_1 S \frac{I}{N} - \nu(t) S \\ V' &= \nu(t) S - \beta_2 V \frac{I}{N} \\ E' &= \beta_1 S \frac{I}{N} + \beta_2 V \frac{I}{N} - \sigma E \\ I' &= \sigma E - \gamma I \\ R' &= \gamma I \end{cases} \quad \text{with} \quad \begin{cases} \beta_1(t) &= NN_1(t, M_1, C_1, M_2, C_2) \\ \beta_2(t) &= (1 - \epsilon) \beta_1 \\ \nu(t) &= NN_2(t, M_1, C_1, M_2, C_2) \end{cases}$$

## Information Model

$$\begin{cases} M'_1 &= NN_3(t, M_1, C_1)[1] \\ C'_1 &= NN_3(t, M_1, C_1)[2] \end{cases} \quad \begin{cases} M'_2 &= NN_4(t, M_2, C_2)[1] \\ C'_2 &= NN_4(t, M_2, C_2)[2] \end{cases}$$

We will investigate the relationship between the two sets of neural networks,  $NN_1$  and  $NN_3$  (for intervention), as well as  $NN_2$  and  $NN_4$  (for vaccine ).

## Application Four: Misinformation



# Outline

## 1 Background

## 2 Universal Differential Equations

- Application One: Estimating time vary reproduction number
- Application Two: Learning Unknown Mechanisms
- Application Three: Optimal Control
- Application Four: Misinformation Project

## 3 Story Behind UDE

- Neural Networks
- PINN

## 4 Discussion

# Neural Networks

- input layer:  $\mathcal{N}^0(\mathbf{U}) = \mathbf{U} \in \mathbb{R}^{d_{\text{ta}}}$
- hidden layers:  $\mathcal{N}^l(\mathbf{U}) = \sigma(W^l \mathcal{N}^{l-1}(\mathbf{U}) + b^l) \in \mathbb{R}^{N_l}$  for  $1 \leq l \leq L-1$
- output layer:  $\mathcal{N}^L(\mathbf{U}) = W^L \mathcal{N}^{L-1}(\mathbf{U}) + b^L \in \mathbb{R}^{d_{\text{out}}}$  where  $W^l$  is a Matrix or Tensor.

NN is composition of operators. Any abstract operator can be a layer, such as solution map of PDE, solution of implicit function. See more in NeurIPS 2020 tutorial: Deep Implicit Layers.

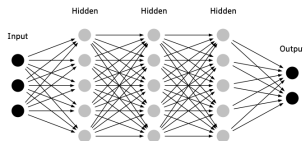


Figure: Scheme of deep neural network.



## Neural Networks: Universal Apprimators

Neural Networks are universal apprimators. It can approximate **unknown nonlinear operators in Banach space**, such as elliptic, nonlocal diffusion.

### Theorem (Lu et.al. 2021 Nat Mach Intell)

*Suppose that  $X$  is a Banach space,  $K_1 \subset X$ ,  $K_2 \subset \mathbb{R}^d$  are two compact sets in  $X$  and  $\mathbb{R}^d$ , respectively,  $V$  is a compact set in  $C(K_1)$ . Assume that  $G: V \rightarrow C(K_2)$  is a nonlinear continuous operator. Then, for any  $\epsilon > 0$ , there exist positive integers  $m, p$ , continuous vector functions  $g: \mathbb{R}^m \rightarrow \mathbb{R}^p$ ,  $f: \mathbb{R}^d \rightarrow \mathbb{R}^p$ , and  $x_1, x_2, \dots, x_m \in K_1$ , such that*

$$|G(u)(y) - \langle g(u(x_1), u(x_2), \dots, u(x_m)), f(y) \rangle| < \epsilon$$

*holds for all  $u \in V$  and  $y \in K_2$ , where  $\langle \cdot, \cdot \rangle$  denotes the dot product in  $\mathbb{R}^p$ . Furthermore, the functions  $g$  and  $f$  can be chosen as diverse classes of neural networks, which satisfy the classical universal approximation theorem of functions, for example, (stacked/unstacked) fully connected neural networks, residual neural networks and convolutional neural networks.*

Consequently, the  $\beta$  values are not significantly different from zero, and the  $\alpha$  values are not significantly different from one. The  $\alpha$  and  $\beta$  values are also not significantly different from each other. The  $\alpha$  and  $\beta$  values are also not significantly different from the  $\alpha$  and  $\beta$  values of the other countries. The  $\alpha$  and  $\beta$  values are also not significantly different from the  $\alpha$  and  $\beta$  values of the other countries.





# Why Neural Networks? Powerful Generalization, Implicit Regularization

- Escaping saddle point
- Sharpness aware
- Local minimum is enough: easily find good solutions or surrogates, and the surrogate **doesn't need to be unique**

<sup>1</sup>Du S S, Jin C, Lee J D, et al. Gradient descent can take exponential time to escape saddle points. Advances in neural information processing systems, 2017, 30.

<sup>2</sup>Wu L, Su W J. The Implicit Regularization of Dynamical Stability in Stochastic Gradient Descent. ICML, 2023.

<sup>3</sup>Foret P, Kleiner A, Mobahi H, et al. Sharpness-aware minimization for efficiently improving generalization. COLT, 2020.

## Why Not Neural Networks Only?

Feed on "big" data, Difficult to interpret, . Training DNN by "big" data is an effective but unwise way.

- Alphago in 2016 costs about 35 million dollars
- GPT-3 in 2020 has 175 billion parameters.
- Megatron-Turing in 2021 has 530 billion parameters.
- Recent Persia can have 100 Trillion parameters.
- GPT3.5: 200 billion; GPT4 not known.

It is AI but not human intelligence!?

# Why Not Neural Networks Only? Incorporating Knowledge

Feed on "big" data, Difficult to interpret.

- Newton's law, simple, it is science.
- Lorenz system, simple, it is science.
- SIR model, simple, it is science.
- Neural Networks, complex, it is black-box.

Researchers hate and love deep neural networks. The art of a good deep learning model is incorporating

Knowledge.

Join AI and Human Intelligence together. Universal differential equations is one way.

# Universal Differential Equations

"Universal" means "universal approximators"

UDEs (proposed by Prof. Christopher Rackauckas, MIT, 2020(Christopher Rackauckas et.al. 2020 arxiv) are initial value problems with the following forms:

$$u' = f_{\theta_2}(u, t, NN_{\theta_1}(u, t)), \quad (8)$$

where  $f$  is a known mechanistic model and  $NN$  denotes the missing or unknown terms,  $\theta_1$  and  $\theta_2$  are parameters of known mechanisms and neural networks, respectively, which can be estimated simultaneously.

# How to Train Universal Differential Equations?

The essence of training UDE is to solve the following **abstract evolution equations constrained optimization problem or optimal control problem (or inverse problems or bayesian inversion problems)**:

$$\begin{cases} \min_{\theta} J = \mathbb{E} \int_0^T g(X, \text{NeuralNetwork}_{\theta}(t, X), t) dt + \phi(X(T), T) \\ \text{s.t. } F(t, X(t), X(\alpha(t)), \text{UA}(t, X(t), X(\beta(t))), W(t)) = 0, t \in [0, T], \end{cases} \quad (9)$$

$F$  denotes the abstract evolution equations such as stochastic partial functional differential equations.

# How to Train Universal Differential Equations?

For ODE case, its essence is to solve the following optimal control problem:

$$\begin{cases} \min_{\theta} J = \int_0^T g(x, \text{NeuralNetwork}_{\theta}(t, x), t) dt + \phi(x(T), T) \\ \text{s.t.} \quad \frac{dx}{dt} = f(x, \text{NeuralNetwork}_{\theta}(t, x), t), x(0) = x_0, t \in [0, T]. \end{cases} \quad (10)$$

# How to Train Universal Differential Equations?

"backpropagation" for differential equations: Adjoint sensitivity analysis in optimal control theory (Cao et.al. 2003 SIAM J. Sci. Comput).

## Theorem

Let

$$J(\theta) = \int_0^T e(y, t) dt, y' = m(y, \theta, t), y(0) = y_0,$$

where  $\theta \in \mathbb{R}^k$  and the functions  $m : \mathbb{R}^n \times \mathbb{R}^k \times \mathbb{R} \rightarrow \mathbb{R}^n$ ,  $e : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$  are continuously differentiable. Then we have

$$\begin{cases} \frac{dJ}{d\theta} = \int_0^T \lambda(t) m_\theta dt, \\ \lambda'(t) = -e_y - m_y \lambda(t), \lambda(T) = 0, \\ y' = m(y, \theta, t), y(0) = y_0. \end{cases} \quad (11)$$



# Story Behind: Universal Differential Equations

UDEs are **beyond the extension of DNN or NDEs**.

- The success behind UDEs are **high-performance differential programing and scientific computation, especially high-performance adjoint sensitivity analysis**.
- UDEs can be more general and abstract

$$u' = Lu + f_{\theta_2}(u, t, UA_{\theta_1}(u, t)) \quad (12)$$

where  $L$  is some kind of abstract evolution operators such as elliptic, nonlocal diffusion in transmission models.

- $UA$  can be **other universal approximators** such as decision trees, Chebyshev expansion, or Gaussian Process to handle uncertainty.
- Universal partial differential equations (UPDE), delayed differential equations (UDAE), Filippov systems, impulsive differential equations, stochastic differential equations (USDE).

# Universal DDE

Generate Data from the following SIR model with lag effect of media impact:

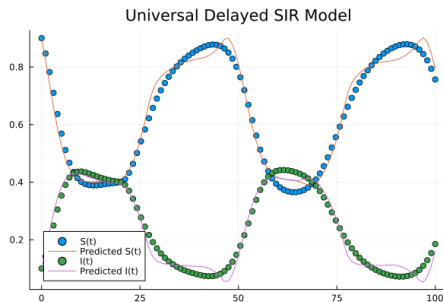
$$\begin{cases} \frac{dS}{dt} = \Lambda - S\beta I \exp(-\alpha I(t - \tau)) - dS, \\ \frac{dI}{dt} = S\beta I \exp(-\alpha I(t - \tau)) - dI - \gamma I, \end{cases} \quad (13)$$

where  $S, I$  denote suspected and infected individuals.

# Universal DDE

Using Generated data to learn

$$\begin{cases} \frac{dS}{dt} = \Lambda - SNN(I(t), I(t - \tau)) - dS, \\ \frac{dI}{dt} = SNN(I(t), I(t - \tau)) - dI - \gamma I, \end{cases} \quad (14)$$



# Universal Filippov System

Generate Data from the following switched SIR model:

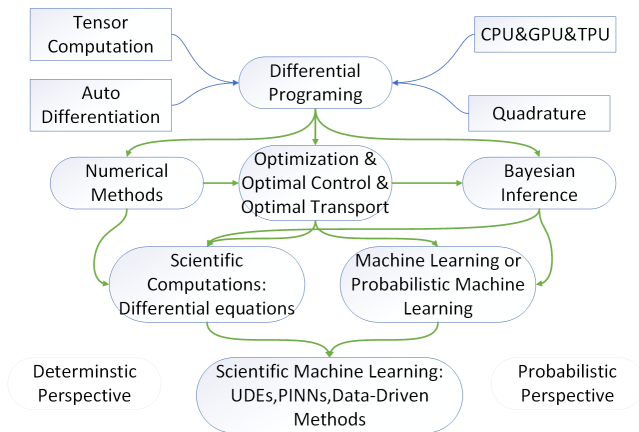
$$\begin{cases} \frac{dS}{dt} = \Lambda - \epsilon\beta SI - dS, \\ \frac{dI}{dt} = \epsilon\beta SI - dI - \gamma I, \\ \epsilon = 1.0(I \leq 0.3), \epsilon = 0.5(I > 0.3), \end{cases} \quad (15)$$

where  $S, I$  denote suspected and infected individuals.



# Practical Engine: Scientific Machine Learning Computation

The success behind machine learning is everything can and should be auto differentiable. (draw by myself, may be not that correct.)









# Outline

## 1 Background

## 2 Universal Differential Equations

- Application One: Estimating time vary reproduction number
- Application Two: Learning Unknown Mechanisms
- Application Three: Optimal Control
- Application Four: Misinformation Project

## 3 Story Behind UDE

- Neural Networks
- PINN

## 4 Discussion

## Publications on this topic

- Song, Pengfei and Xiao, Yanni and Wu, Jianhong. (2023) Methods coupling transmission models and deep learning. Preprint.
- Song, Pengfei and Xiao, Yanni and Wu, Jianhong. (2023) Discovering first-principle behavior change transmission models by deep learning methods. One Chapter of Springer Book. Accepted
- Yin, Shuangshuang and Song, Pengfei and Wu, Jianhong. (2023) Optimal epidemic control by deep learning techniques. JMB.
- Song, Pengfei and Xiao, Yanni. (2022) Estimating time-varying reproduction number by deep learning techniques (Dedicated to Prof Jibin Li on his 80th birthday). JAAC.

Thanks!