

Couple of transmission models and deep learning techniques: brief introduction

Pengfei Song

Xi'an Jiaotong University



February 4, 2024

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

1 Background

2 Universal Differential Equations(微分方程中嵌入了神经网络)

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

3 Story Behind UDE

- Why UDE?
- How to Train UDE?
- Any Theoretical Guarantees?

Background

Epidemic models have proved to be a very powerful tool in guiding public health measures, learning from the past and preparing for the future. Nonetheless, modeling and controlling the emerging infectious disease such as COVID-19 remains a challenge due to the **unknown mechanisms** in transmission dynamics, for example,

- nonstandard incidence rate
- changing human mobility pattern
- wastewater early warning
- impact of misinformation on disease spreading (NLP, UDE, Epi Model)
- ...

生命科学中建模难



Data driven methods to learn unknown mechanism

- Mechanisms can be characterized by 'functions', 'operators', 'Distributions', 'Stochastic Processes', 'Manifolds'. 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 Biology
- 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)

¹Weinan E: A Mathematical Perspective on Machine Learning. 2022 International Congress of Mathematician 60 Minutes Talk

Outline

1 Background

2 Universal Differential Equations(微分方程中嵌入了神经网络)

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

3 Story Behind UDE

- Why UDE?
- How to Train UDE?
- Any Theoretical Guarantees?

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:

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

where f is a known mechanism and UniversalApproximator denotes the missing or unknown terms, θ_1 and θ_2 are parameters of known mechanisms and UAs, respectively, which can be estimated simultaneously.

¹Rackauckas C, Ma Y, Martensen J, et al. Universal differential equations for scientific machine learning[J]. arXiv preprint arXiv:2001.04385, 2020.

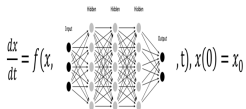
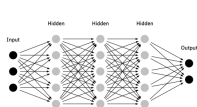
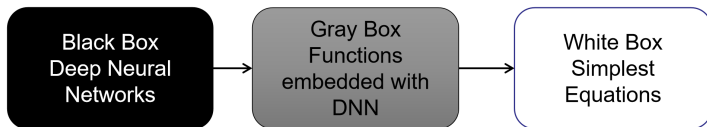
²Yin S, Wu J, Song P*. Optimal control by deep learning techniques and its applications on epidemic models[J]. Journal of Mathematical Biology, 2023, 86(3): 36.

³Song P, Xiao Y. Estimating time-varying reproduction number by deep learning techniques[J]. J Appl Anal Comput, 2022, 12(3): 1077-1089. (Dedicated to Prof Jibin Li on his 80th birthday).

Black-box, Gray-box, White-box

How to recover the simplest function from deep neural networks?

Equation search methods.



$$\frac{dx}{dt} = f(x, t), x(0) = x_0$$

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

¹Song, 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: Symbolic Regression

SR uses **binary-tree** to represent a function, and **no particular formula is provided as a starting point**. SR uses genetic programming, bayesian methods and deep learning methods to discover the simplest equations.

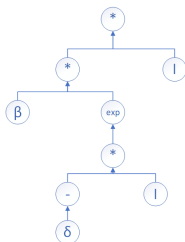


Figure: Expression tree of the function $f(l) = \beta \exp(-\delta l)l$

¹Žegklitz J, Pošík P. Benchmarking state-of-the-art symbolic regression algorithms[J]. Genetic programming and evolvable machines. 2021. 22: 5-33.

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, $u' = f(u) \approx \Theta(u)\Xi$, and estimate Ξ 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

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

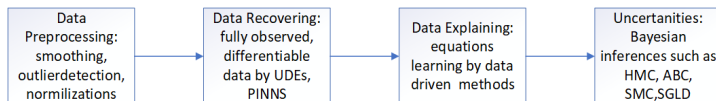
and now we use SINDy to discover the true equations.

¹Brunton 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.

Two-step Learning-Explaining Methods

UDE: representing and learning the unknown mechanisms by neural networks

Equation Search: recover the simplest function from neural networks



Remark: recover the simplest function from neural networks is favored by biologists and mathematician, but NOT in line with the philosophy of deep learning. Improve the generalization ability is.

¹Song, 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 One: Estimating time vary reproduction number

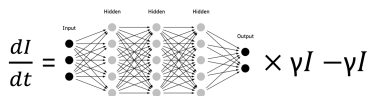
Represent effective reproduction number \mathcal{R}_t as

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

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 .



¹Song P, Xiao Y. Estimating time-varying reproduction number by deep learning techniques[J]. J Appl Anal Comput, 2022, 12(3): 1077-1089. (Dedicated to Prof Jibin Li on his 80th birthday).

15 / 42

16 / 42

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:

$$\begin{aligned}\frac{dS}{dt} &= -\text{abs}(NN(I, R))S/N, \\ \frac{dI}{dt} &= -\text{abs}(NN(I, R))S/N - \gamma I, \\ \frac{dR}{dt} &= \gamma I, \\ \frac{dH}{dt} &= \text{abs}(NN(I, R))S/N,\end{aligned}$$

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

¹Song, 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 Two: Learning Unknown Human Behaviour Change Mechanisms

To start with, learn the data by UDEs.

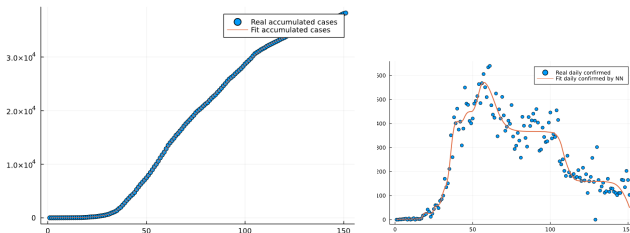


Figure: Learn Ontario first wave data by universal differential equations.

Application Three: Optimal Control

Consider the following optimal control problem in *Bolza* form:

$$\begin{cases} \max_{u(t) \in \Omega(t)} J = \int_0^T g(x, u, t) dt + \phi(x(T), T) \\ \text{s.t.} \quad \frac{dx}{dt} = f(x, u, t), x(0) = x_0, \end{cases} \quad (6)$$

where the functions $f: \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}^n$, $g: \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$ and $\phi: \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$ are assumed to be continuously differentiable. By representing the optimal control function $u(t)$ as a neural network

$$u(t) = \text{NeuralNetwork}_\theta(t, x) \quad (7)$$

receiving t and x as inputs.——

¹Yin S, Wu J, Song P*. Optimal control by deep learning techniques and its applications on epidemic models[J]. Journal of Mathematical Biology, 2023, 86(3): 36.

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	Accurate
Disadvantages	Less accurate	Hard to post, hard to solve, initial guess	Curse of dimensionality	Flexible, extendable Bless 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 or text from Twitter(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

¹Team York University, Team XJTU. A neural differential equation model for cognitive behavioural response to information and disease. Preprint.

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).

Outline

1 Background

2 Universal Differential Equations(微分方程中嵌入了神经网络)

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

3 Story Behind UDE

- Why UDE?
- How to Train UDE?
- Any Theoretical Guarantees?

Neural Networks

- input layer: $\mathcal{N}^0(\mathbf{U}) = \mathbf{U} \in \mathbb{R}^{d_{ta}}$
- 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: $\mathcal{N}^L(\mathbf{U}) = W^L \mathcal{N}^{L-1}(\mathbf{U}) + b^L \in \mathbb{R}^{d_{out}}$ 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.

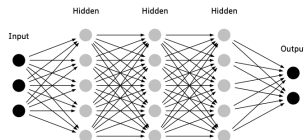


Figure: Scheme of deep neural network.

Neural Networks: Universal Apprimators of Functions

Neural Networks are universal apprimators. It can approximate **unknown mappings and their derivatives**. (Speical "Taylor expansion", **Difference: Projection VS Composition**)

Theorem (Allan Pinkus 1999 Acta Numer)

Let $m_i \in \mathbb{Z}^{d+}, i = 1, \dots, s$, and set $m = \max_{i=1, \dots, s} |m_i|$. Assume $\sigma \in C(\mathbb{R}^d)$ and that σ is not a polynomial. Then a single hidden layer neural network:

$$\mathcal{M}(\sigma) := \text{span} \left\{ \sigma(\mathbf{w} \cdot \mathbf{U} + b) : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$

is dense in

$$\mathcal{C}^{m^1, \dots, m^s}(\mathbb{R}^d) := \cap_{i=1}^s \mathcal{C}^{m^i}(\mathbb{R}^d)$$

Neural Networks: Universal Apprimators of Operators

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.

Neural Networks: Universal Apprimators of Distributions

Generative models (GAN, VAE, Auto-regressive models, Normalizing Flows, Diffusion Models, Energy Based Models, Consistency models) are universal approximators for **distributions**.

Continuous normalizing flow:

$$x' = NN(x, t), x(0) \sim p$$

From an unified perspective of **Optimal Transport**.

¹Grathwohl W, Chen R T Q, Bettencourt J, et al. Ffjord: Free-form continuous dynamics for scalable reversible generative models[J]. arXiv preprint arXiv:1810.01367, 2018.

Why Neural Networks? Solving Curse of dimensionality

Universal approximators like Polinomial spaces, Fourier expansion , Chebyshev expansion, Decision trees , Gaussian process face curse of dimensionality (COD). Many techniques are needed to handle COD, such as sparsity, parallel computing.

However, NN is believed to some extent share bless of dimensionality (many practical findings and few theoretical results).

¹Weinan E: A Mathematical Perspective on Machine Learning. 2022 International Congress of Mathematician 60 Minutes Talk

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**
- Form the View of Optimal Transport: Not only leaning the measure, but also the manifold structure. (Theoretically under exploring).

¹Du 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.

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

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

Why Not Neural Networks Only?

Question: DNN is so powerfull, why we need UDE?

Feed on "big" and high quality data, Difficult to interpret, . Trianing 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.
- LLMs: start from "billion"

It is AI but not human intelligence!?



Why Not Neural Networks Only? Incorporating Knowledge

Feed on "big" and high quality 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:

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

where f is a known mechanisms and UniversalApproximator denotes the missing or unknown terms, θ_1 and θ_2 are parameters of known mechanisms and neural networks, respectively, which can be estimated simultaneously.

¹Rackauckas C, Ma Y, Martensen J, et al. Universal differential equations for scientific machine learning[J]. arXiv preprint arXiv:2001.04385, 2020.

²Yin S, Wu J, Song P*. Optimal control by deep learning techniques and its applications on epidemic models[J]. Journal of Mathematical Biology, 2023, 86(3): 36.

³Song P, Xiao Y. Estimating time-varying reproduction number by deep learning techniques[J]. J Appl Anal Comput, 2022, 12(3): 1077-1089. (Dedicated to Prof Jibin Li on his 80th birthday).

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)**:

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

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. Back to Pontryagin.

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

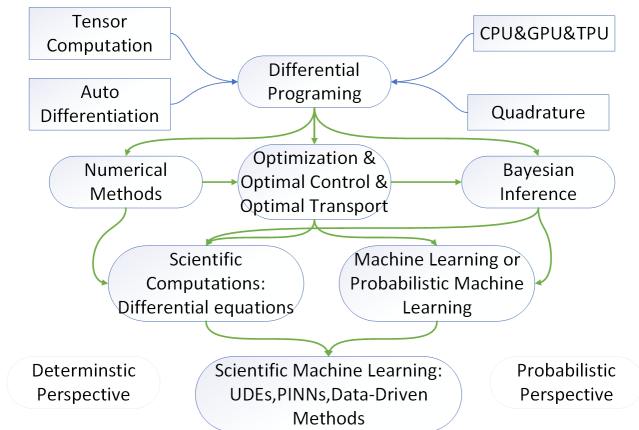
$$\begin{cases} \frac{dJ}{d\theta} = \int_0^T \lambda(t) m_{\theta} dt, \\ \lambda'(t) = -e_y - m_y \lambda(t), \lambda(T) = 0, \\ y' = m(y, \theta, t), y(0) = y_0. \end{cases} \quad (11)$$

¹Cao Y, Li S, Petzold L, et al. Adjoint sensitivity analysis for differential-algebraic equations: The adjoint DAE system and its numerical solution[J]. SIAM journal on scientific computing, 2003, 24(3): 1076-1089.

²Yin S, Wu J, Song P*. Optimal control by deep learning techniques and its applications on epidemic models[J]. Journal of Mathematical Biology, 2023, 86(3): 36.

Practical Engine: Scientific Machine Learning

The success behind machine learning is everything can and should be auto differentiable.



Any Theoretical Guarantees?

Theorem

Total error

$$\begin{aligned}
 & \hat{J}(u_{N,h,\delta}^k) - \hat{J}(u^*) \\
 & \leq \underbrace{(L_{gx}L_S + L_{gu})\text{Aprox}_{u^*}(N)}_{\text{Approximation error}} \\
 & \quad + \underbrace{C_{gx}\text{Num}(h)\|S(u_{N,h,\delta}^k)\| + C_{gx}\text{Num}(h)\|S(u_N^*)\|}_{\text{Numerical error}} \\
 & \quad + \underbrace{4(L_{gx}L_{S,h} + L_{gu})\mathcal{R}_\delta(\mathcal{U}_N)}_{\text{Generalization error}} \\
 & \quad + \underbrace{(L_{gx}L_{S,h} + L_{gu})C_{\text{opt}}\text{Opt}(k)}_{\text{Optimization error}},
 \end{aligned} \tag{12}$$

¹Song Pengfei, Lou Yuan, Yin Shuangshuang. Theoretical guarantees for deep learning based optimal control method for ordinary differential equations. preparing.

Thanks!