Supervised Learning In Banach Space Approximation Theory For Operators

Andrew Stuart

Computing and Mathematical Sciences California Institute of Technology

AFOSR, ARL, NSF, ONR

* * *

IPAM Closing Workshop

June 8th 2020

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

Collaborators

- K Bhattacharya, B Hosseini, NB Kovachki and AMS Model Reduction And Neural Networks For Parametric PDEs arXiv:2005.03180
- NH Nelsen and AMS The Random Feature Model For Input-Output Maps Between Banach Spaces arXiv:2005.10224

 Z Li, NB Kovachki, K Azizzadenesheli, B Liu, K Bhattacharya, AMS and A Anandkumar Neural Operator: Graph Kernel Network for Partial Differential Equations arXiv:2003.03485

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Overview

Supervised Learning

The Talk In A Couple Of Slides

Model Reduction Approach

Conclusions and References

Random Features Approach

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Graph Kernel Approach

The Setting

▲□▶ ▲□▶ ▲ 三▶ ▲ 三 ● ● ●

The Problem

Setting

Input-Output Map: $\Psi^{\dagger} : \mathcal{X} \mapsto \mathcal{Y}$. Regression: $\mathcal{X} = \mathbb{R}^{q}, \mathcal{Y} = \mathbb{R}^{r}$. Classification: $\mathcal{X} = \mathbb{R}^{q}, \mathcal{Y} = \{1, \dots, K\}$.

Parametric Approximation Of Ψ^{\dagger}

$$\begin{array}{ll} \mathsf{Data:} & \{x_n, y_n\}_{n=1}^N, \quad y_n = \Psi^{\dagger}(x_n).\\ \mathsf{Parameter Space:} & \Theta \subseteq \mathbb{R}^p.\\ \mathsf{Approximation Class:} & \Psi : \mathcal{X} \times \Theta \mapsto \mathcal{Y}.\\ & \mathsf{Goal:} & \theta^* \in \Theta : \Psi(\cdot; \theta^*) \approx \Psi^{\dagger}. \end{array}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

[1] Y LeCun, Y Bengio, G Hinton. (deep learning and applications)

Training

μ a probability measure on \mathcal{X} ; $\ell: \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}^+$ a distance.

Perfect Setting

$$\begin{array}{ll} \text{Objective Function:} & \mathsf{J}(\theta) = \mathbb{E}^{x \sim \mu} \ell \Big(\Psi(x; \theta), \Psi^{\dagger}(x) \Big). \\ \\ \text{Minimizer:} & \theta^{\star} = \operatorname{argmin}_{\theta \in \Theta} \mathsf{J}(\theta). \end{array}$$

Let
$$x_n \overset{i.i.d.}{\sim} \mu$$
. and $y_n = \Psi^{\dagger}(x_n)$.

Data-Driven Setting

Objective Function:
$$J^{N}(\theta) = \frac{1}{N} \sum_{n=1}^{N} \ell \Big(\Psi(x_{n}; \theta), y_{n} \Big),$$

Minimizer: $\theta^{\star, N} = \operatorname{argmin}_{\theta \in \Theta} J^{N}(\theta).$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Testing

Perfect Setting

$${ t Error}$$
 Measure: $e_{\scriptscriptstyle N} = \mathbb{E}^{ imes \sim \mu} \ell \Big(\Psi(x; heta^{ imes, N}), \Psi^{\dagger}(x) \Big).$

Let
$$x'_m \stackrel{i.i.d.}{\sim} \mu'$$
. and $y'_m = \Psi^{\dagger}(x'_m)$. $\{x'_m\}_{m=1}^M \perp \{x_n\}_{n=1}^N$.

Data-Driven Setting

Error Measure:
$$e_{N,M} = \frac{1}{M} \sum_{m=1}^{M} \ell \left(\Psi(\mathbf{x}'_m; \theta^{\star, N}), \mathbf{y}'_m \right).$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

This talk: $\mu' = \mu$; but $\mu \neq \mu'$ also relevant.

Neural Networks

Deep And/Or Wide

$$egin{aligned} & v_0 = Ax + a, \ & v_{k+1} = \sigma(W_k v_k + \omega_k), \quad k = 0, 1, \cdots, K-1, \ & \Psi(x; heta) = B v_K + b. \end{aligned}$$

Notation

$$\sigma(z) = \max(0, z) \qquad (\mathsf{Nemytskii}),$$
$$\theta = (A, a, B, b) \bigcup_{k=1}^{K} \{W_k, \omega_k\},$$

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

[2] I Goodfellow, Y Bengio, A Courville. (deep learning text book)

The Talk In A Couple Of Slides

Previous work at PDE/neural networks interface:

[3] M Raissi, P Perdikaris, GE Karniadakis;

[4] Weinan E and B Yu.

Example What Goes Wrong If You Discretize And Then Apply Standard Neural Network Algorithms



[5a] Y Zhu and N Zabaras; see also [5b] Y Khoo, J Lu and L Ying

The Problem (Revisited)

Key Idea

Design Architecture On Banach Space Then Discretize

Setting

Input-Output Map: $\Psi^{\dagger} : \mathcal{X} \mapsto \mathcal{Y}$, Separable Banach Spaces Data: $\{x_n, y_n\}_{n=1}^N, y_n = \Psi^{\dagger}(x_n), x_n \stackrel{i.i.d.}{\sim} \mu$.

Goal

 $\begin{array}{ll} \mbox{Parameter Space} & \Theta \subseteq \mathbb{R}^{\rho} \\ & \mbox{Operator Class:} & \Psi: \mathcal{X} \times \Theta \mapsto \mathcal{Y} \\ \mbox{Operator Approximation:} & \Psi(\cdot; \theta^{\star}) \approx \Psi^{\dagger} \end{array}$

Model Reduction Approach

K Bhattacharya, B Hosseini, NB Kovachki and AMS Model Reduction And Neural Networks For Parametric PDEs arXiv:2005.03180

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬる

The Idea

In A Picture



In Equations

$$\begin{split} G_{\mathcal{X}} \circ F_{\mathcal{X}} &\approx I_{\mathcal{X}}, \\ G_{\mathcal{Y}} \circ F_{\mathcal{Y}} &\approx I_{\mathcal{Y}}, \\ G_{\mathcal{Y}} \circ \varphi \circ F_{\mathcal{X}} &\approx \Psi^{\dagger}. \end{split}$$

▲□▶ ▲圖▶ ▲国▶ ▲国▶ - 国 - のへで

Principal Components Analysis

Required Information

Hilbert Space \mathcal{H} Probility Measure: π Mean: $m = \mathbb{E}^{h \sim \pi} h$ Covariance: $C = \mathbb{E}^{h \sim \pi} (h - m) \otimes (h - m)$

Projection

$$\begin{aligned} \mathsf{EVP:} \quad C\psi_j &= \lambda_j\psi_j\\ \mathsf{F:}\,\mathcal{H} &\mapsto \mathbb{R}^d \quad Fh = \{v_j := \langle h, \psi_j \rangle\}_{j=1}^d\\ \mathsf{G:}\,\mathbb{R}^d &\mapsto \mathcal{H} \quad Gv = \sum_{j=1}^d v_j\psi_j. \end{aligned}$$

The Approximations

PCA (\mathcal{X}, \mathcal{Y} Hilbert spaces)

$$\begin{array}{ccc} \mathcal{X} & \xrightarrow{F_{\mathcal{X}}} & \mathbb{R}^{d_{\mathcal{X}}} & \xrightarrow{G_{\mathcal{X}}} & \mathcal{X} \\ & & & & & \\ \Psi^{\dagger} & & & & & \\ \mathcal{Y} & \xrightarrow{F_{\mathcal{Y}}} & & \mathbb{R}^{d_{\mathcal{Y}}} & \xrightarrow{G_{\mathcal{Y}}} & \mathcal{Y} \end{array}$$

$$\begin{split} \varphi &:= F_{\mathcal{Y}} \circ \Psi^{\dagger} \circ G_{\mathcal{X}}, \\ \Psi_{PCA} &:= G_{\mathcal{Y}} \circ \varphi \circ F_{\mathcal{X}}. \end{split}$$

Neural Networks

$$\psi \approx \varphi,$$

 $\Psi_{NN} := G_{\mathcal{Y}} \circ \psi \circ F_{\mathcal{X}}.$

Main Theorem

Theorem

For any $\epsilon > 0$, there are dimensions $(d_{\mathcal{X}}, d_{\mathcal{Y}})(\epsilon)$, a requisite amount of data $N = N(d_{\mathcal{X}}, d_{\mathcal{Y}})$, and a zero-extended neural network ψ depending on $\epsilon, d_{\mathcal{X}}, d_{\mathcal{Y}}$ such that

$$\mathbb{E}_{\{x_j\}\sim \mu}\mathbb{E}_{x\sim \mu} \|\Psi_{\scriptscriptstyle NN}(x)-\Psi^{\dagger}(x)\|_{\mathcal{Y}}^2 < \epsilon.$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

[6] L Zwald, O Bousquet, G Blanchard (used in proof)

[7] D Yarotsky (used in proof)

Related work:

[8] G Kutyniok, P Petersen M Raslan and R Schneider ; (finite dimensional nn approximation)

[9] C Schwab, J Zech; (finite dimensional nn approximation)

[10] A Chkifa, A Cohen, R DeVore and C Schwab; (parametric pde approximation)

[11] RA DeVore. (parametric pde approximation)



$$-\nabla \cdot (a\nabla u) = f, \quad z \in D$$
$$u = g, \quad z \in \partial D.$$

Operators Of Interest

$$\begin{array}{ll} \mathsf{Linear} \ (\mathsf{Boundary}) & \Psi^{\dagger} : g \in H^{\frac{1}{2}}(\partial D) \mapsto u \in H^{1}(D); \\ \mathsf{Linear} \ (\mathsf{Forcing}) & \Psi^{\dagger} : f \in L^{2}(D) \mapsto u \in H^{1}(D); \\ & \mathsf{Nonlinear} & \Psi^{\dagger} : a \in L^{\infty}(D) \mapsto u \in H^{1}(D). \end{array}$$

Distance

$$\ell(\mathbf{v}, u) = \|\mathbf{v} - u\|_{L^2(D)} / \|u\|_{L^2(D)}.$$

Input-Output

 $\begin{array}{lll} \mbox{Input:} & a \in L^\infty(D) & (Left), \\ \mbox{Output:} & u \in H^1(D). & (Right), \end{array}$







▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへ(で)

Conclusions and References

▲□▶ ▲□▶ ▲国▶ ▲国▶ ▲国 ● のへで

Conclusions

1. Neural networks: empirical success in function approximation.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

- 2. Typically:
 - regression: $\mathbb{R}^m \mapsto \mathbb{R}^n$;
 - classification: $\mathbb{R}^m \mapsto \{1, \ldots, K\}$.
- 3. We consider: $\mathcal{X} \mapsto \mathcal{Y}$, \mathcal{X}, \mathcal{Y} function spaces.
- 4. Bad: Discretize then apply 2.
- 5. Good: Conceive of architecture then discretize.
- 6. Purely data-driven ("equation-free").
- 7. Applications: PDEs
- 8. Applications: Other?

References



[1] Y LeCun, Y Bengio, G Hinton .

Deep Learning Nature 2015.



[2] I Goodfellow, Y Bengio, A Courville Deep Learning MIT Press 2016



[3] M Raissi, P Perdikaris, GE Karniadakis

Physics-informed neural networks... J. Comp. Phys. 2019.



[4] Weinan E and B Yu

The Deep Ritz Method... Communications in Mathematics and Statistics 2018.



[5a] Y Zhu and N Zabaras

Bayesian Deep Convolutional Encoder-Decoder Networks... J. Comp. Phys. 2018.



[5b] Y Khoo, J Lu, L Ying

Solving parametric PDE problems with artificial neural networks arXiv: 1707.03351,



[6] L Zwald, O Bousquet, G Blanchard

Statistical properties of kernel principal component analysis... International Conference on Computational Learning Theory, Springer, 2004.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

References



[7] D Yarotsky

Error bounds for approximations with deep ReLU networks Neural Networks 2017



[8] G Kutyniok, P Petersen M Raslan and R Schneider

A theoretical analysis of deep neural networks and parametric PDEs arXiv:1904.00377



[9] C Schwab, J Zech

Deep learning in high dimension: Neural network expression rates for generalized polynomial chaos expansions in UQ

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Analysis and Applications 2019



Sparse adaptive Taylor approximation algorithms for parametric and stochastic elliptic PDEs ESAIM: Mathematical Modelling and Numerical Analysis 2013



[11] RA DeVore

The Theoretical Foundation of Reduced Basis Methods Model Reduction and Approximation, SIAM, 2014



[12] A Rahimi, B Recht

Uniform approximation of functions with random bases 46th Annual Allerton Conference on Communication. Control. and Computing, 2008

Random Features Approach

NH Nelsen and AMS

The Random Feature Model For Input-Output Maps Between Banach Spaces arXiv:2005.10224

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬる

The Architecture



Random Feature Model

Parameter Space : $\Theta \subseteq \mathbb{R}^m$ Architecture : $\Psi_m(x; \theta) := \frac{1}{m} \sum_{j=1}^m \theta_j \psi(x; \gamma_j), \quad \gamma_j \sim \nu \text{ i. i. d.}$

[12] A Rahimi and B Recht. (random features and approximation theory)

Using ResNet As Random Features

Input Space: $\mathcal{X} = \mathbb{R}^2, A \in \mathbb{R}^{1 \times 2}$ Output Space: $\mathcal{Y} = \mathbb{R}, a \in \mathbb{R}$ Weights: $W \in C([0, 1]; \mathbb{R}^{2 \times 2}), \ \omega \in C([0, 1]; \mathbb{R}^2),$ Random Evolution: $\dot{v}(t) = \sigma(W(t)v(t) + \omega(t)), v(0) = x,$ Random Features: $\psi(x; \gamma) = Av(1) + a, \ \gamma = (W, \omega, A, a).$





Reproducing Kernel Hilbert Space (RKHS)

Kernel From Features

$$\begin{array}{ll} \mathsf{Kernel:} & k: \mathcal{X} \times \mathcal{X} \to \mathcal{L}(\mathcal{Y}, \mathcal{Y}), \\ \mathsf{Features:} & k(x, z) := \mathbb{E}^{\gamma \sim \nu} \Big(\psi(x; \gamma) \otimes \psi(z; \gamma) \Big) \\ \mathsf{RKHS:} & \mathcal{H}_k \subset L^2_{\mu}(\mathcal{X}; \mathcal{Y}), \end{array}$$

Integral Representation

$$egin{aligned} f(z) &= \int heta(\gamma) \psi(z;\gamma)
u(d\gamma), \ & heta(\gamma) &= \langle f, \psi(\cdot;\gamma)
angle_{\mathcal{H}_k} \end{aligned}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

From Nonparametric to Parametric Regression

Gaussian Process Regression

Penalized LSQ:
$$I(F) := \sum_{n=1}^{N} \frac{1}{2} \|y_n - F(x_n)\|_{\mathcal{Y}}^2 + \frac{\lambda}{2} \|F\|_{\mathcal{H}_k}^2,$$

Regression: $F^* = \operatorname{argmin}_{F \in \mathcal{H}_k} I(F), \qquad F^*(x) := \sum_{\ell=1}^{N} k(x, x_\ell) \beta_\ell.$

Link To Random Features Model

R

Parametric:
$$k(x, z) \approx k^{(m)}(x, z) = \frac{1}{m} \sum_{j=1}^{m} \psi(x; \gamma_j) \otimes \psi(z; \gamma_j),$$

Objective: $J^{\lambda}(\theta) := \sum_{n=1}^{N} \frac{1}{2} \|y_n - \frac{1}{m} \sum_{\ell=1}^{m} \theta_{\ell} \psi(x_n; \gamma_{\ell})\|_{\mathcal{Y}}^2 + \frac{\lambda}{2m} \|\theta\|_2^2,$
Regularized RFM: $\theta^* := \operatorname{argmin}_{\theta \in \mathbb{R}^m} J^{\lambda}(\theta).$

Burgers Equation

$$\partial_t u + \partial_s (u^2/2) = \delta \partial_{ss}^2 u, \quad (s,t) \in \mathbb{T}^1 \times (0,1],$$

 $u|_{t=0} = u_0, \quad s \in \mathbb{T}^1.$

Operator Of Interest

Nonlinear
$$\Psi^{\dagger}: u_0 \in L^2(\mathbb{T}^1) \mapsto u|_{t=1} \in H^r(\mathbb{T}^1).$$

Input-Output





▲□▶ ▲圖▶ ▲園▶ ▲園▶ 三国 - 釣A@

Graph Kernel Approach

Z Li, NB Kovachki, K Azizzadenesheli, B Liu, K Bhattacharya, AMS and A Anandkumar Neural Operator: Graph Kernel Network for Partial Differential Equations arXiv:2003.03485

▲□▶▲□▶▲≡▶▲≡▶ ≡ めぬる

Integral Kernel Neural Networks

Deep And/Or Wide

$$f(s) = (s, x(s), x_{\epsilon}(s), \nabla_{s} x_{\epsilon}(s))^{T},$$

$$v_{0}(s) = Af(s) + a,$$

$$v_{k+1}(s) = \sigma \Big(Wv_{k}(s) + \int k(s, r, x(s), x(r); \vartheta) v_{k}(r) \lambda(dr) \Big),$$

$$\Psi(x; \theta) = Bv_{K}(s) + b.$$

In Abstract Notation

$$k(\cdot,\vartheta) \in \{\varphi : \mathbb{R}^{2d+2} \mapsto \mathbb{R}^{p \times p}\},\$$
$$\theta = (A, a, B, b, W, \vartheta).$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○○



$$-\nabla \cdot (a\nabla u) = f, \quad z \in D$$
$$u = g, \quad z \in \partial D.$$

Operators Of Interest

$$\begin{array}{ll} \mathsf{Linear} \ (\mathsf{Boundary}) & \Psi^{\dagger} : g \in H^{\frac{1}{2}}(\partial D) \mapsto u \in H^{1}(D); \\ \mathsf{Linear} \ (\mathsf{Forcing}) & \Psi^{\dagger} : f \in L^{2}(D) \mapsto u \in H^{1}(D); \\ \mathsf{Nonlinear} & \Psi^{\dagger} : a \in L^{\infty}(D) \mapsto u \in H^{1}(D). \end{array}$$

Distance

$$\ell(v, u) = \|v - u\|_{L^2(D)} / \|u\|_{L^2(D)}.$$

Input-Output Nonlinear







Input-Output Linear (Green's Function D = [0, 1])



