

DARTR: Data Adaptive RKHS Tikhonov Regularization for learning kernels in operators

Fei Lu

Department of Mathematics, Johns Hopkins University

Joint with **Quanjun Lang**, **Qingci An**@JHU, **Yue Yu**@Lehigh

MiC Seminars, ORNL; August, 2022



JOHNS HOPKINS
UNIVERSITY



- 1 Learning kernels
- 2 Regression and regularization
- 3 Identifiability and DARTR
- 4 Numerical examples

Learning kernels in operators

Learn the **kernel** ϕ :

$$R_\phi[u] + \epsilon = f$$

from data:

$$\mathcal{D} = \{(u_k, f_k)\}_{k=1}^N, \quad (u_k, f_k) \in \mathbb{X} \times \mathbb{Y}$$

Learning kernels in operators

Learn the **kernel** ϕ :

$$R_\phi[u] + \epsilon = f$$

from data:

$$\mathcal{D} = \{(u_k, f_k)\}_{k=1}^N, \quad (u_k, f_k) \in \mathbb{X} \times \mathbb{Y}$$

- Operator R_ϕ : **linear or nonlinear in u , but linear in ϕ**

- ▶ nonlocal interaction (interacting particles, mean-field)

$$R_\phi[u] = \nabla \cdot [u(K_\phi * u)] = \partial_t u - \sigma \Delta u, \quad K_\phi(x) = \phi(|x|) \frac{x}{|x|} \in \mathbb{R}^d$$

$$R_\phi[\mathbf{X}_t] = \left(-\frac{1}{n} \sum_{j=1}^n K_\phi(\mathbf{X}_t^i - \mathbf{X}_t^j) \right)_i = \dot{\mathbf{X}}_t + \dot{\mathbf{W}}_t, \quad \mathbb{R}^{nd}$$

Learning kernels in operators

Learn the **kernel** ϕ :

$$R_\phi[u] + \epsilon = f$$

from data:

$$\mathcal{D} = \{(u_k, f_k)\}_{k=1}^N, \quad (u_k, f_k) \in \mathbb{X} \times \mathbb{Y}$$

- Operator R_ϕ : **linear or nonlinear in u , but linear in ϕ**

- ▶ nonlocal interaction (interacting particles, mean-field)

$$R_\phi[u] = \nabla \cdot [u(K_\phi * u)] = \partial_t u - \sigma \Delta u, \quad K_\phi(x) = \phi(|x|) \frac{x}{|x|} \in \mathbb{R}^d$$

$$R_\phi[\mathbf{X}_t] = \left(-\frac{1}{n} \sum_{j=1}^n K_\phi(\mathbf{X}_t^i - \mathbf{X}_t^j) \right)_i = \dot{\mathbf{X}}_t + \dot{\mathbf{W}}_t, \quad \mathbb{R}^{nd}$$

- ▶ nonlocal PDE/ fractional diffusion: $R_\phi[u] = \partial_{tt} u - g$

$$R_\phi[u](x) = \int_{\Omega} \phi(x, y)[u(y) - u(x)] dy = \partial_{tt} u - g, \quad \forall x \in \Omega.$$

Learning kernels in operators

Learn the **kernel** ϕ :

$$R_\phi[u] + \epsilon = f$$

from data:

$$\mathcal{D} = \{(u_k, f_k)\}_{k=1}^N, \quad (u_k, f_k) \in \mathbb{X} \times \mathbb{Y}$$

- Operator R_ϕ : **linear or nonlinear in u , but linear in ϕ**

- ▶ nonlocal interaction (interacting particles, mean-field)

$$R_\phi[u] = \nabla \cdot [u(K_\phi * u)] = \partial_t u - \sigma \Delta u, \quad K_\phi(x) = \phi(|x|) \frac{x}{|x|} \in \mathbb{R}^d$$

$$R_\phi[\mathbf{X}_t] = \left(-\frac{1}{n} \sum_{j=1}^n K_\phi(\mathbf{X}_t^i - \mathbf{X}_t^j) \right)_i = \dot{\mathbf{X}}_t + \dot{\mathbf{W}}_t, \quad \mathbb{R}^{nd}$$

- ▶ nonlocal PDE/ fractional diffusion: $R_\phi[u] = \partial_{tt} u - g$

$$R_\phi[u](x) = \int_{\Omega} \phi(x, y)[u(y) - u(x)] dy = \partial_{tt} u - g, \quad \forall x \in \Omega.$$

- ▶ Integral operators, deconvolution, Toeplitz/Hankel matrix ...
Toeplitz matrix: $R_\phi u = f$, $R_\phi(i, j) = \phi(i - j)$

Learning kernels in operators

Learn the **kernel** ϕ :

$$R_\phi[u] + \epsilon = f$$

from data:

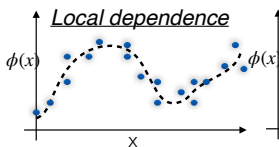
$$\mathcal{D} = \{(u_k, f_k)\}_{k=1}^N, \quad (u_k, f_k) \in \mathbb{X} \times \mathbb{Y}$$

- Operator R_ϕ : **linear or nonlinear in u , but linear in ϕ**
 - ▶ nonlocal interaction (interacting particles, mean-field)
 - ▶ nonlocal PDE/ fractional diffusion
 - ▶ Integral operators, deconvolution, Toeplitz/Hankel matrix ...
- Data: discrete/noisy, **Nonlocal dependence**
 - ▶ random $(u_k, f_k) \sim \mu \otimes \nu$: **statistical learning**
 - ▶ deterministic (e.g., $N=1$): **inverse problem**

Comparison with classical learning

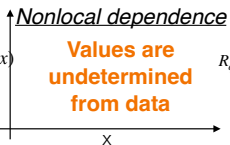
Classical learning

$$\{(x_i, \phi(x_i) + \epsilon_i)\}$$



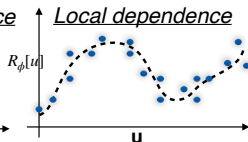
Learning kernel

$$\{(u_k, R_\phi[u_k] + \eta_k)\}$$



Operator learning

$$\{(u_k, R_\phi[u_k] + \eta_k)\}$$



- nonlocal dependence
under-determined (no longer “interpolation”)
- v.s. operator learning: low-dimensional structure
- methods: regression/ML/DL?

This talk: deterministic inverse problems

$$\mathcal{D} = \{u_k, f_k\}_{k=1}^N = \{u_k(x_j, t_l), f_k(x_j, t_l) : j = 1, \dots, \mathcal{J}\}_{i=1}^N,$$

- 1 Learning kernels
- 2 Regression and regularization
- 3 Identifiability and DARTR
- 4 Numerical examples

Nonparametric regression

Loss functional: $\mathcal{E}(\phi) = \frac{1}{N} \sum_{i=1}^N \|R_\phi[u_i] - f_i\|_{L^2}^2.$

► Neural network when ϕ is high-D

Hypothesis space: $\phi = \sum_{i=1}^n \mathbf{c}_i \phi_i \in \mathcal{H}_n = \text{span}\{\phi_i\}_{i=1}^n:$

$$\mathcal{E}(\phi) = \mathbf{c}^\top \bar{\mathbf{A}}_n \mathbf{c} - 2\mathbf{c}^\top \bar{\mathbf{b}}_n + \mathbf{C}_N^f, \Rightarrow \hat{\phi}_{\mathcal{H}_n} = \sum_i \hat{\mathbf{c}}_i \phi_i, \text{ where } \hat{\mathbf{c}} = \bar{\mathbf{A}}_n^{-1} \bar{\mathbf{b}}_n,$$

Three issues

- $\bar{\mathbf{A}}^{-1}$: well-posedness, Identifiability, function space
- Choice of \mathcal{H}_n : $\{\phi_i\}_{i=1}^n$ and n
- Convergence when data mesh refines $\Delta x \rightarrow 0$

Regularization

Regularization is necessary:

- \bar{A}_n ill-conditioned/singular
- \bar{b}_n : noise or numerical error

Tikhonov/ridge Regularization:

$$\mathcal{E}_\lambda(\phi) = \mathcal{E}(\phi) + \lambda \|\phi\|_*^2 \Rightarrow \mathbf{c}^\top \bar{A}_n \mathbf{c} - 2 \bar{b}_n^\top \mathbf{c} + \lambda \|\mathbf{c}\|_{B_*}^2$$

$$\hat{\phi}_{\mathcal{H}_n}^\lambda = \sum_i \hat{\mathbf{c}}_i^\lambda \phi_i, \quad \text{where } \hat{\mathbf{c}} = (\bar{A}_n + \lambda B_*)^{-1} \bar{b}_n,$$

Regularization

Regularization is necessary:

- \bar{A}_n ill-conditioned/singular
- \bar{b}_n : noise or numerical error

Tikhonov/ridge Regularization:

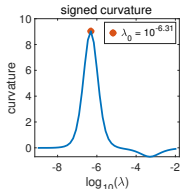
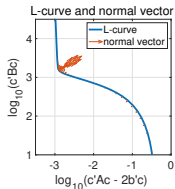
$$\mathcal{E}_\lambda(\phi) = \mathcal{E}(\phi) + \lambda \|\phi\|_*^2 \Rightarrow \mathbf{c}^\top \bar{A}_n \mathbf{c} - 2\bar{b}_n^\top \mathbf{c} + \lambda \|\mathbf{c}\|_{B_*}^2$$

$$\hat{\phi}_{\mathcal{H}_n}^\lambda = \sum_i \hat{c}_i^\lambda \phi_i, \quad \text{where } \hat{\mathbf{c}} = (\bar{A}_n + \lambda B_*)^{-1} \bar{b}_n,$$

- λ by the L-curve method [Hansen00]

$$l(\lambda) = (x(\lambda), y(\lambda)) := (\log(\mathcal{E}(\hat{\mathbf{c}}_\lambda)), \log(\|\hat{\mathbf{c}}_\lambda\|_*^2)),$$

$$\lambda_0 = \arg \max_{\lambda_{\min} \leq \lambda \leq \lambda_{\max}} \frac{x' y'' - x'' y'}{(x'^2 + y'^2)^{3/2}},$$



Regularization

Regularization is necessary:

- \bar{A}_n ill-conditioned/singular
- \bar{b}_n : noise or numerical error

Tikhonov/ridge Regularization:

$$\mathcal{E}_\lambda(\phi) = \mathcal{E}(\phi) + \lambda \|\phi\|_*^2 \Rightarrow \mathbf{c}^\top \bar{A}_n \mathbf{c} - 2\bar{b}_n^\top \mathbf{c} + \lambda \|\mathbf{c}\|_{B_*}^2$$

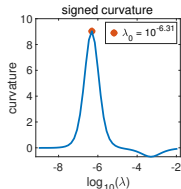
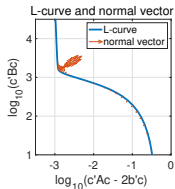
$$\hat{\phi}_{\mathcal{H}_n}^\lambda = \sum_i \hat{c}_i^\lambda \phi_i, \quad \text{where } \hat{\mathbf{c}} = (\bar{A}_n + \lambda B_*)^{-1} \bar{b}_n,$$

- λ by the L-curve method [Hansen00]

$$l(\lambda) = (x(\lambda), y(\lambda)) := (\log(\mathcal{E}(\hat{c}_\lambda)), \log(\|\hat{c}_\lambda\|_*^2)),$$

$$\lambda_0 = \arg \max_{\lambda_{\min} \leq \lambda \leq \lambda_{\max}} \frac{x' y'' - x'' y'}{(x'^2 + y'^2)^{3/2}},$$

- Which norm $\|\cdot\|_*$ to use?



- 1 Learning kernels
- 2 Regression and regularization
- 3 Identifiability and DARTR
- 4 Numerical examples

Identifiability

- An exploration measure: $\rho(dr) \Rightarrow \phi \in L^2(\rho)$

$$R_\phi[u](x) = \int_{\Omega} \phi(|x-y|)g[u](x,y)dy, \quad \rho(dr) \propto \int \int \delta_{|x-y|}(dr)|g[u](x,y)|dxdy$$

Identifiability

- An exploration measure: $\rho(dr) \Rightarrow \phi \in L^2(\rho)$

$$R_\phi[u](x) = \int_{\Omega} \phi(|x-y|)g[u](x,y)dy, \quad \rho(dr) \propto \int \int \delta_{|x-y|}(dr)|g[u](x,y)|dxdy$$

- An integral operator \Leftarrow the Fréchet derivative of loss functional

$$\mathcal{E}(\psi) = \frac{1}{N} \sum_{i=1}^N \|R_\psi[u_i] - f_i\|_{L^2}^2 = \langle \mathcal{L}_{\bar{G}}\psi, \psi \rangle_{L^2(\rho)} - 2\langle \phi^D, \psi \rangle_{L^2(\rho)}$$

$$\nabla \mathcal{E}(\psi) = 2\mathcal{L}_{\bar{G}}\psi - 2\phi^D = 0 \Rightarrow \hat{\phi} = \mathcal{L}_{\bar{G}}^{-1}\phi^D$$

- ▶ $\mathcal{L}_{\bar{G}}$ is a nonnegative compact operator: $\{(\lambda_i, \psi_i)\}$, $\lambda_i \downarrow 0$
- ▶ $\phi^D = \mathcal{L}_{\bar{G}}\phi_{true} + \phi^{error}$

Identifiability

- An exploration measure: $\rho(dr) \Rightarrow \phi \in L^2(\rho)$

$$R_\phi[u](x) = \int_\Omega \phi(|x-y|)g[u](x,y)dy, \quad \rho(dr) \propto \int \int \delta_{|x-y|}(dr)|g[u](x,y)|dxdy$$

- An integral operator \Leftarrow the Fréchet derivative of loss functional

$$\mathcal{E}(\psi) = \frac{1}{N} \sum_{i=1}^N \|R_\psi[u_i] - f_i\|_{L^2}^2 = \langle \mathcal{L}_{\overline{G}}\psi, \psi \rangle_{L^2(\rho)} - 2\langle \phi^D, \psi \rangle_{L^2(\rho)}$$

$$\nabla \mathcal{E}(\psi) = 2\mathcal{L}_{\overline{G}}\psi - 2\phi^D = 0 \Rightarrow \hat{\phi} = \mathcal{L}_{\overline{G}}^{-1}\phi^D$$

- ▶ $\mathcal{L}_{\overline{G}}$ is a nonnegative compact operator: $\{(\lambda_i, \psi_i)\}$, $\lambda_i \downarrow 0$
- ▶ $\phi^D = \mathcal{L}_{\overline{G}}\phi_{true} + \phi^{error}$

- Function space of identifiability (FSOI):

$$\hat{\phi} = \mathcal{L}_{\overline{G}}^{-1}(\mathcal{L}_{\overline{G}}\phi_{true} + \phi^{error}) \Rightarrow H = \text{span}\{\psi_i\}_{i:\lambda_i>0}$$

- ▶ ill-defined beyond H ; ill-posed in H

DARTR: Data Adaptive RKHS Tikhonov Regularization

$$\hat{\phi} = \mathcal{L}_{\bar{G}}^{-1} \phi^f = \mathcal{L}_{\bar{G}}^{-1} (\mathcal{L}_{\bar{G}} \phi^{true} + \phi^{error})$$

A new task for Regularization:

ensure that the learning takes place in the FSOI

data-dependent $H = \text{span}\{\psi_i\}_{i:\lambda_i > 0}$

DARTR: Data Adaptive RKHS Tikhonov Regularization

$$\hat{\phi} = \mathcal{L}_{\bar{G}}^{-1} \phi^f = \mathcal{L}_{\bar{G}}^{-1} (\mathcal{L}_{\bar{G}} \phi^{true} + \phi^{error})$$

A new task for Regularization:

ensure that the learning takes place in the FSOI

data-dependent $H = \text{span}\{\psi_i\}_{i:\lambda_i>0} = \overline{H_G}^{L^2(\rho)}$

- $\bar{G} \Rightarrow$ RKHS: $H_G = \mathcal{L}_{\bar{G}}^{-1/2}(L^2(\rho))$ System Intrinsic Data Adaptive
- For $\phi = \sum_k c_k \psi_k$, $\|\phi\|_{L^2(\rho)}^2 = \sum_k c_k^2$, $\|\phi\|_{H_G}^2 = \sum_k \lambda_k^{-1} c_k^2$

\Rightarrow Regularization norm: $\|\phi\|_{H_G}^2$

$$\mathcal{E}_\lambda(\phi) = \mathcal{E}(\phi) + \lambda \|\phi\|_{H_G}^2 \Rightarrow \mathbf{c}^\top \bar{\mathbf{A}}_n \mathbf{c} - 2\bar{\mathbf{b}}_n^\top \mathbf{c} + \lambda \|\mathbf{c}\|_{B_{rkhs}}^2$$

Why DARTR is good: FSOI is fundamental:

$$\hat{\phi} = \mathcal{L}_{\bar{G}}^{-1} \phi^D = \mathcal{L}_{\bar{G}}^{-1} (\mathcal{L}_{\bar{G}} \phi_{true} + \phi_H^{error} + \phi_{H^\perp}^{error})$$

- DARTR: $\|\phi_{H^\perp}^{error}\|_{H_G}^2 = \infty$

$$(\mathcal{L}_{\bar{G}} + \lambda \mathcal{L}_{\bar{G}}^{-1})^{-1} \phi^D = (\mathcal{L}_{\bar{G}} + \lambda \mathcal{L}_{\bar{G}}^{-1})^{-1} (\mathcal{L}_{\bar{G}} \phi_{true} + \phi_H^{error})$$

- ℓ^2 or L^2 regularizer: with $C = \sum \phi_i \otimes \phi_j$ or $C = I$

$$(\mathcal{L}_{\bar{G}} + \lambda C)^{-1} \phi^D = (\mathcal{L}_{\bar{G}} + \lambda C)^{-1} (\mathcal{L}_{\bar{G}} \phi_{true} + \phi_H^{error} + \phi_{H^\perp}^{error})$$

Why DARTR is good: FSOI is fundamental:

$$\hat{\phi} = \mathcal{L}_{\bar{G}}^{-1} \phi^D = \mathcal{L}_{\bar{G}}^{-1} (\mathcal{L}_{\bar{G}} \phi_{true} + \phi_H^{error} + \phi_{H^\perp}^{error})$$

- DARTR: $\|\phi_{H^\perp}^{error}\|_{H_G}^2 = \infty$

$$(\mathcal{L}_{\bar{G}} + \lambda \mathcal{L}_{\bar{G}}^{-1})^{-1} \phi^D = (\mathcal{L}_{\bar{G}} + \lambda \mathcal{L}_{\bar{G}}^{-1})^{-1} (\mathcal{L}_{\bar{G}} \phi_{true} + \phi_H^{error})$$

- ℓ^2 or L^2 regularizer: with $C = \sum \phi_i \otimes \phi_j$ or $C = I$

$$(\mathcal{L}_{\bar{G}} + \lambda C)^{-1} \phi^D = (\mathcal{L}_{\bar{G}} + \lambda C)^{-1} (\mathcal{L}_{\bar{G}} \phi_{true} + \phi_H^{error} + \phi_{H^\perp}^{error})$$

A Bayesian perspective:

- Prior $\mathcal{N}(0, \mathcal{L}_{\bar{G}})$; v.s. $\mathcal{N}(0, C)$: singular or equivalent
- Posterior $\mathcal{N}(\hat{\phi}_*, (\mathcal{L}_{\bar{G}} + \lambda \mathcal{L}_{\bar{G}}^{-1})^{-1})$ v.s. $\mathcal{N}(\hat{\phi}_\dagger, (\mathcal{L}_{\bar{G}} + \lambda C)^{-1})$
- Zellner's g-prior $\mathcal{N}(0, \bar{A}_n^{-1})$ if $\mathcal{H}_n = \text{span}\{\phi_i\}_{i=1}^n$ o.n.b.

DARTR: compute the SIDA-RKHS norm

Let $B_n = (\langle \phi_i, \phi_j \rangle_{L^2(\rho)})_{i,j}$.

Theorem (Generalized eigenvalue problem)

If $\mathcal{L}_{\bar{G}}(L^2(\rho)) \subset \mathcal{H}$, then $\mathcal{L}_{\bar{G}}$ eigenvalues are solved by the generalized eigenvalue problem (\bar{A}_n, B_n) and $B_{rkhs} = (V \Lambda V^\top)^{-1}$.

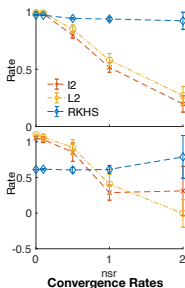
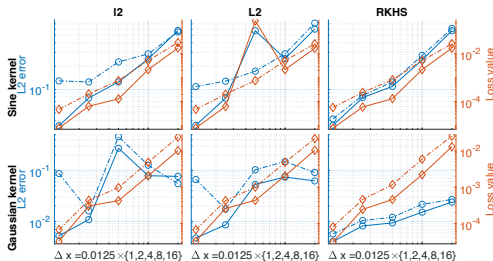
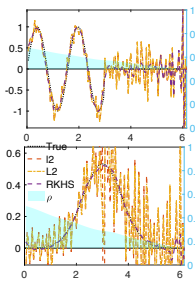
- If $B_n = I_n$: $B_{rkhs} = \bar{A}_n^{-1}$, the Zellner's g-prior;
- If $\phi_i = \psi_i$: $\bar{A}_n = \text{diag}(\lambda_i)$, $B_{rkhs} = \bar{A}_n^{-1}$:
 $\hat{c} = \sum_{i:\lambda_i > 0} (\lambda_i + \lambda)^{-1} (v_i^\top \bar{b}) v_i$
- l^2 or L^2 regularizer: $\hat{c} = [\sum_{i:\lambda_i > 0} + \sum_{i:\lambda_i = 0}] (\lambda_i + \lambda)^{-1} (v_i^\top \bar{b}) v_i$

- 1 Learning kernels
- 2 Regression and regularization
- 3 Identifiability and DARTR
- 4 Numerical examples

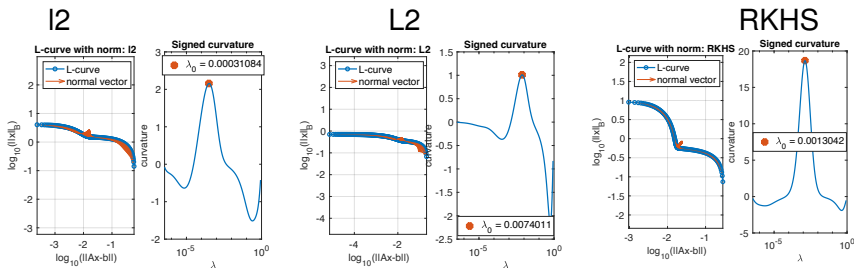
Interaction kernel in a nonlinear operator

$$R_\phi[u] = \nabla \cdot [u(K_\phi * u)] = f, \quad K_\phi = \phi(|x|) \frac{x}{|x|}$$

- Recover kernel from **discrete noisy data**
- **Robust in accuracy, consistent rates** as mesh refines

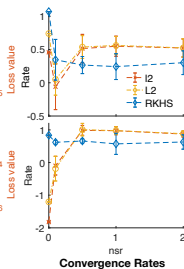
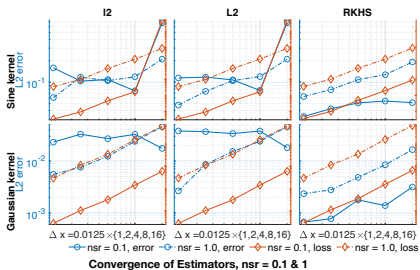
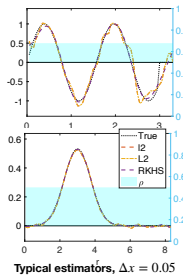


More robust L-curve



Linear integral operator:

$$R_\phi[u](x) = \int_{\Omega} \phi(|y-x|)u(y)dy = f(x).$$



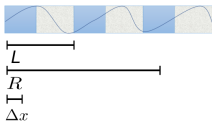
Robust in accuracy, consistent rates

Homogenization of wave propagation in meta-material

- heterogeneous bar with microstructure + DNS \Rightarrow Data
- Homogenization: $R_\phi[u] = \partial_{tt}u - g$.

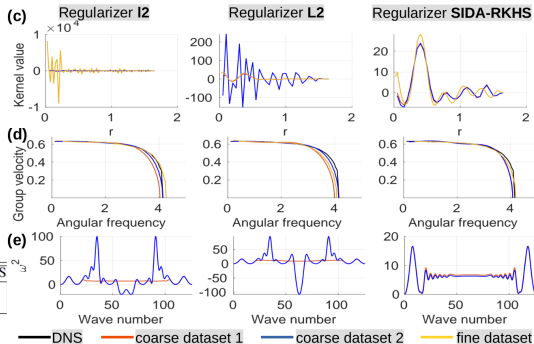
$$R_\phi[u](x) = \int_{\Omega} \phi(|y|)[u(x+y) - u(x)]dy.$$

(a) Wave propagation in a heterogeneous bar



(b) Displacement error on a cross-validation dataset

Resolution	I2	L2	SIDA-RKHS
Coarse ($\Delta x = 0.05$)	23.5%	28.4%	21.8%
Fine ($\Delta x = 0.025$)	INF	23.4%	19.2%



- (c): resolution-invariant
- (e): l^2 and $L2$ leading to non-physical kernel

Summary

Learning kernels in operators:

$$R_\phi[u] = f \quad \Leftarrow \quad \mathcal{D} = \{(u_k, f_k)\}_{k=1}^N$$

Nonlocal dependence

- Identifiability theory: FSOI
- DARTR: data adaptive RKHR Tikhonov-Reg
- Numerical tests: robust accuracy, consistent rates

Future directions

- Convergence of regularized estimator $(\Delta x, N)$
- Inverse problems with nonlocal dependence
- Regularization for neural network: $\|\phi_\theta\|_{rkhs}^2$, not $\|\theta\|$

	Data	Goal: ϕ	Inversion*	FSOI	Regularization
Classical learning	$\{(x_i, y_i)\}$	$Y = \phi(X)$	$\hat{\phi} = I^{-1}\phi^D$	$L^2(\rho)$	no need of FSOI
Learning kernels	$\{(u_i, f_i)\}$	$R_\phi[u] = f$	$\hat{\phi} = \mathcal{L}_G^{-1}\phi^D$	SIDA	FSOI necessary

References

(@ <http://www.math.jhu.edu/~feilu>) QR-code →

- F.Lu, Q .Lang and Q. An. Data adaptive RKHS Tikhonov regularization for learning kernels in operators. MSML22. (Matlab code available)

- F.Lu, Q .An and Y. Yu. Nonparametric learning of kernels in nonlocal operators. arXiv2205

