

# Data-adaptive RKHS regularization for learning kernels in operators

Fei Lu

Department of Mathematics, Johns Hopkins University

CAM Colloquium, U Chicago  
December 5, 2024

Collaborators: Qianjun Lang, Qingci An, Yue Yu, Haibo Li,  
Jinchao Feng, Yvonne Ou



Which norm  $\|x\|_*$  to use in  $\|Ax - b\|^2 + \lambda\|x\|_*^2$   
for an ill-posed inverse problem (e.g.,  $\phi(r) = \sum_{i=1}^n x_i \phi_i(r)$ )?

(A)  $\|x\|_{\mathbb{R}^n}$

(B)  $\|\phi\|_{L^2}$

(C) Total variation

(D) RKHS

- 1 Learning kernels
- 2 Regression and regularization
- 3 Identifiability and DARTR
- 4 Iterative method

# Learning kernels in operators

Learn the **kernel**  $\phi$ :  $R_\phi[u] + \epsilon = f$

from data:

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

$$\text{Operator } R_\phi[u](x) = \int \phi(|x - y|)g[u](x, y)dy$$

# Learning kernels in operators

Learn the **kernel**  $\phi$ :  $R_\phi[u] + \epsilon = f$

from data:

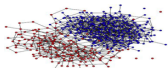
$$\mathcal{D} = \{(u_k, f_k)\}_{k=1}^N, \quad (u_k, f_k) \in \mathbb{X} \times \mathbb{Y}$$

Operator  $R_\phi[u](x) = \int \phi(|x - y|)g[u](x, y)dy$

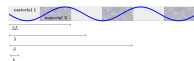
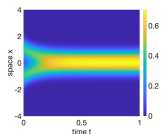
- Interacting particles/agents  $K_\phi(x) = \phi(|x|) \frac{x}{|x|} \in \mathbb{R}^d$

$$R_\phi[\mathbf{X}_t] = \left[ -\frac{1}{n} \sum_{j=1}^n K_\phi(\mathbf{X}_t^i - \mathbf{X}_t^j) \right]_i = \dot{\mathbf{X}}_t + \sqrt{2\nu} \dot{\mathbf{W}}_t,$$

$$R_\phi[u] = \nabla \cdot [u(K_\phi * u)] = \partial_t u - \nu \Delta u,$$



Voter model (wiki)



# Learning kernels in operators

Learn the **kernel**  $\phi$ :  $R_\phi[u] + \epsilon = f$

from data:

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

Operator  $R_\phi[u](x) = \int \phi(|x - y|)g[u](x, y)dy$

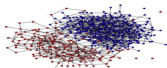
- Interacting particles/agents  $K_\phi(x) = \phi(|x|) \frac{x}{|x|} \in \mathbb{R}^d$

$$R_\phi[\mathbf{X}_t] = \left[ -\frac{1}{n} \sum_{j=1}^n K_\phi(\mathbf{X}_t^i - \mathbf{X}_t^j) \right]_i = \dot{\mathbf{X}}_t + \sqrt{2\nu} \dot{\mathbf{W}}_t,$$

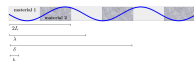
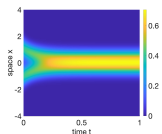
$$R_\phi[u] = \nabla \cdot [u(K_\phi * u)] = \partial_t u - \nu \Delta u,$$

- Nonlocal PDEs:

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



Voter model (wiki)



# Learning kernels in operators

Learn the **kernel**  $\phi$ :  $R_\phi[u] + \epsilon = f$

from data:

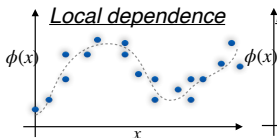
$$\mathcal{D} = \{(u_k, f_k)\}_{k=1}^N, \quad (u_k, f_k) \in \mathbb{X} \times \mathbb{Y}$$

- Operator  $R_\phi[u](x) = \int \phi(|x - y|)g[u](x, y)dy$ :  
**linear or nonlinear in  $u$ , but linear in  $\phi$**
- Statistical learning  $\cap$  inverse problem
  - ▶ random  $\{(u_k, f_k)\}$ : **statistical learning**
  - ▶ deterministic (e.g.,  $N$  small): **inverse problem**

# Learning kernels in operators

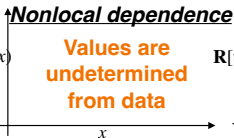
## Classical learning

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



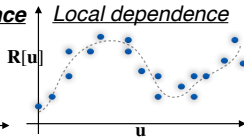
## Learning kernels

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



## Operator learning

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



Inverse problem: well-posed

**ill-posed**

well-posed

This talk:  $\Rightarrow$  introduce a data-adaptive RKHS regularization

- Convergent estimator as mesh refines

$$\mathcal{D} = \{(u_k(x_j), f_k(x_j))\}_{k=1}^N, \quad \Delta x = |x_{j+1} - x_j| \rightarrow 0$$



## Part 2: Regression and regularization

Learn the kernel  $\phi$ :  $R_\phi[u] + \epsilon = f$

from data:

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

Operator  $R_\phi[u](x) = \int \phi(|x - y|)g[u](x, y)dy$

# Nonparametric regression

- Loss functional:  $\mathcal{E}(\phi) = \frac{1}{N} \sum_{k=1}^N \|R_\phi[u_k] - f_k\|_{\mathbb{Y}}^2$ 
  - ▶ Crucial!
  - ▶ Derivative-free Monte Carlo suitable [Lang+Lu22SISC]
- Hypothesis space  $\mathcal{H}_n = \text{span}\{\phi_i\}_{i=1}^n: \phi = \sum_{i=1}^n c_i \phi_i$ ,

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

**Goal:**  $\hat{\phi}_{\mathcal{H}_n}$  converges as data mesh  $\Delta x$  refines

## Challenges

- Choice of  $\mathcal{H}_n: \{\phi_i\}_{i=1}^n$  and  $n = n(\Delta x)$
- $\bar{A}_n^{-1}$ : ill-conditioned/singular

# Regularization

Regularization is necessary:

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

Tikhonov/ridge Regularization:

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

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

# Regularization

Regularization is necessary:

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

Tikhonov/ridge Regularization:

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

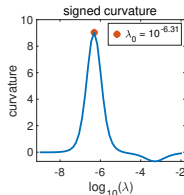
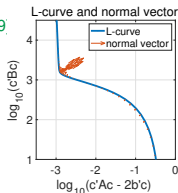
$$\hat{\phi}_{\mathcal{H}_n}^\lambda = \sum_i \hat{c}_{\lambda, i} \phi_i, \quad \text{where } \hat{c}_\lambda = (\bar{A}_n + \lambda \mathbf{B}_*)^{-1} \bar{b}_n,$$

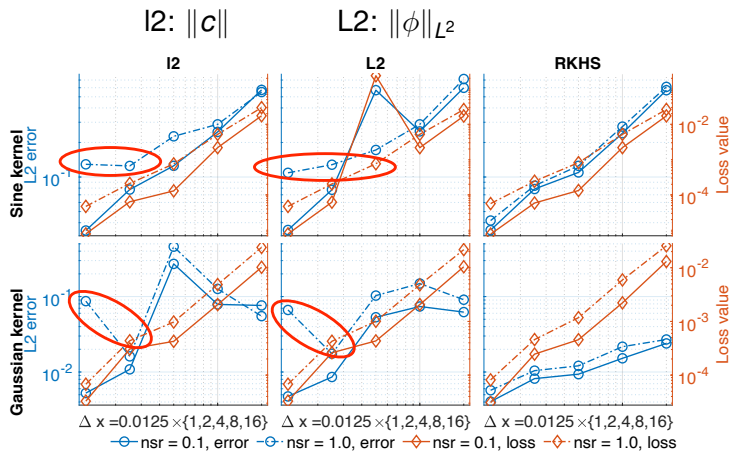
- $\lambda$ : L-curve [Hansen00], GCV [Golub+Heath+Wahba79]

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

$\lambda_*$  = maximal curvature

- Which norm  $\|\cdot\|_*$  to use?  $\|\mathbf{c}\|$ ,  $\|\phi\|$ ?





Convergence of Estimators, nsr = 0.1 & 1

- Risk of blowing up in the small noise limit [Chada-Wang-Lang-Lu22]

Principle: [Stuart2010]

Avoid **discretization** until the last possible moment



Avoid **basis selection** until the last possible moment

---

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

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

Function space of  $\phi$ ? Identifiability?

## Part 3: Identifiability & regularization

DARTR: Data adaptive RKHS Tikhonov regularization

# 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

$$\mathcal{E}(\phi) = \frac{1}{N} \sum_{k=1}^N \|R_\phi[u_k] - f_k\|_{L^2}^2 = \langle \mathcal{L}_{\bar{G}}\phi, \phi \rangle_{L^2_\rho} - 2\langle \phi^D, \phi \rangle_{L^2_\rho} + C$$

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

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

# 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

$$\mathcal{E}(\phi) = \frac{1}{N} \sum_{k=1}^N \|R_\phi[u_k] - f_k\|_{L^2}^2 = \langle \mathcal{L}_{\overline{G}}\phi, \phi \rangle_{L^2_\rho} - 2\langle \phi^D, \phi \rangle_{L^2_\rho} + C$$

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

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

- Function space of identifiability (FSOI):

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

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

## DARTR: Data Adaptive RKHS Tikhonov Regularization

A new task for Regularization:

**ensure that the learning takes place in the FSOI**

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

## DARTR: Data Adaptive RKHS Tikhonov Regularization

A new task for Regularization:

**ensure that the learning takes place in the FSOI**

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

- $\overline{G} \Rightarrow \text{RKHS}: H_G = \mathcal{L}_{\overline{G}}^{-1/2}(L^2_\rho)$
- $\|\phi\|_{H_G}^2 = \langle \mathcal{L}_{\overline{G}}^{-1} \phi, \phi \rangle_{L^2_\rho}$

## DARTR: Data Adaptive RKHS Tikhonov Regularization

A new task for Regularization:

**ensure that the learning takes place in the FSOI**

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

- $\overline{G} \Rightarrow \text{RKHS}: H_G = \mathcal{L}_{\overline{G}}^{-1/2}(L^2_\rho)$

- $\|\phi\|_{H_G}^2 = \langle \mathcal{L}_{\overline{G}}^{-1} \phi, \phi \rangle_{L^2_\rho}$

$\Rightarrow$  Regularization norm:  $\|\phi\|_{H_G}^2$  [Lu+Lang+An22MSML]

$$\mathcal{E}_\lambda(\phi) = \mathcal{E}(\phi) + \lambda \|\phi\|_{H_G}^2 = \langle (\mathcal{L}_{\overline{G}} + \lambda \mathcal{L}_{\overline{G}}^{-1}) \phi, \phi \rangle_{L^2_\rho} - 2 \langle \phi^D, \phi \rangle_{L^2_\rho}$$

$$\hat{\phi}_\lambda = (\mathcal{L}_{\overline{G}} + \lambda \mathcal{L}_{\overline{G}}^{-1})^{-1} \phi^D = (\mathcal{L}_{\overline{G}}^2 + \lambda I)^{-1} \mathcal{L}_{\overline{G}} \phi^D$$

## What DARTR has done:

remove error outside FSOI + regularize in FSOI

- No regularization:

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

- DARTR:  $\mathcal{L}_{\bar{G}} \phi_{H^\perp}^{error} = \mathbf{0}$  or  $\|\phi_{H^\perp}^{error}\|_{H_G}^2 = \infty$

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

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

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

# DARTR: computation

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

**Input:**  $A_n$ ,  $b_n$  and  $B_n = (\langle \phi_i, \phi_j \rangle_{L_p^2})_{i,j}$ .

**Output:** regularized estimator

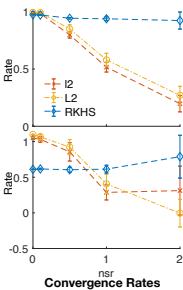
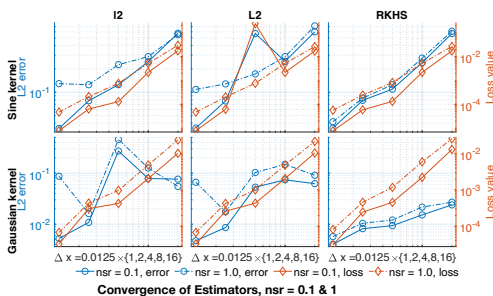
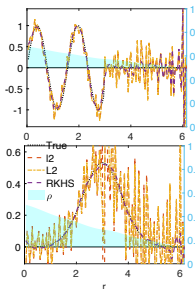
$$\hat{\mathbf{c}}_\lambda = (\mathbf{A}_n + \lambda_* \mathbf{B}_{rkhs})^{-1} \mathbf{b}_n$$

- Generalized eigenvalue problem  $(A_n, B_n) \leftrightarrow \mathcal{L}_{\overline{G}}$   
 $A_n V = B_n V \Lambda$  and  $V^\top B_n V = I_n$   
 $B_{rkhs} = (V \Lambda V^\top)^\dagger$ ; ( $B_{rkhs} = A_n^\dagger$  if  $B_n = I_n$ )
- L-curve to select  $\lambda_*$

# Interaction kernel in a nonlinear operator

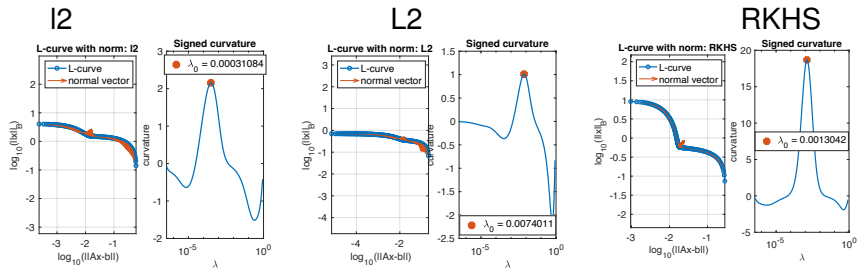
$$R_\phi[u] = \nabla \cdot [u(K_\phi * u)] = f, \quad K_\phi = \phi(|x|) \frac{x}{|x|}$$

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





## More robust L-curve





## Part 4: Iterative method

Large scale  $Ax = b$ ,  $A \in \mathbb{R}^{m \times n}$  ill-conditioned,  $n \gg 1$   
 $b$ : noisy

## Direct method: DARTR for $Ax = b$

$$A_n = A^\top A, b_n = A^\top b: \Rightarrow A_n x = b_n$$

$$\hat{x}_\lambda = (A_n + \lambda_* B_{rkhs})^{-1} b_n$$

- $\rho \propto \sum_j |A_{ij}|$ : measure of  $A$  exploring  $x$
- $B_n = \text{diag}(\rho)$ : pre-conditioning
- Generalized eigenvalue problem  $(A_n, B_n)$   
 $A_n V = B_n V \Lambda$  and  $V^\top B_n V = I_n \Rightarrow B_{rkhs} = (V \Lambda V^\top)^\dagger$   
 $B_{rkhs} = A_n^\dagger$  when  $B_n = I_n$
- L-curve to select  $\lambda_*$

## Direct method: DARTR for $Ax = b$

$$A_n = A^\top A, b_n = A^\top b: \Rightarrow A_n x = b_n$$

$$\hat{x}_\lambda = (A_n + \lambda_* B_{rkhs})^{-1} b_n$$

- $\rho \propto \sum_j |A_{ij}|$ : measure of  $A$  exploring  $x$
- $B_n = \text{diag}(\rho)$ : pre-conditioning
- Generalized eigenvalue problem  $(A_n, B_n)$   
 $A_n V = B_n V \Lambda$  and  $V^\top B_n V = I_n \Rightarrow B_{rkhs} = (V \Lambda V^\top)^\dagger$   
 $B_{rkhs} = A_n^\dagger$  when  $B_n = I_n$
- L-curve to select  $\lambda_*$

---

Direct method: based on **costly** matrix decomposition,  $O(n^3)$ .

Iterative method: without computing  $B_{rkhs}$ ?

# Iterative Data Adaptive RKHS regularization

Solve:  $x_k = \arg \min_{x \in \mathcal{X}_k} \|x\|_{B_{rkhs}}, \mathcal{X}_k = \{x : \min_{x \in \mathcal{S}_k} \|Ax - b\|\}$

$$\mathcal{S}_k = \text{span}\{(B_{rkhs}^\dagger A^\top A)^i B_{rkhs}^\dagger A^\top b\}_{i=0}^k$$

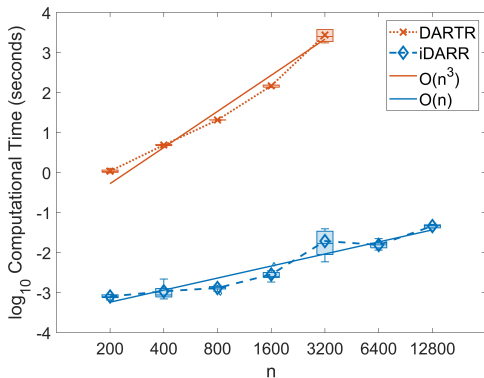
- Use  $B_{rkhs}^\dagger$ , not  $B_{rkhs}$ :  $B_{rkhs}^\dagger = B^{-1} A^\top A B^{-1}$
- generalized Golub-Kahan bidiagonalization (gGKB)  
⇒ construct  $\mathcal{S}_k$  only using matrix-vector product
- $\mathcal{S}_k$  = RKHS-restricted Krylov subspace
- Early stopping: select  $k$

[Li+Feng+Lu, arXiv2401: Scalable iterative data-adaptive RKHS regularization]

# Computational complexity

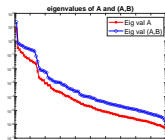
Direct method: DARTR,  $O(n^3)$

Iterative method: iDARR,  $O(3mnk)$

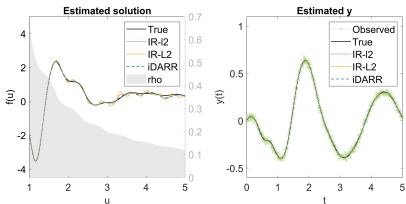


# Fredholm integral equation: 1st kind

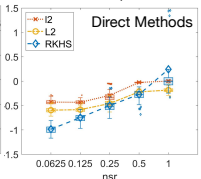
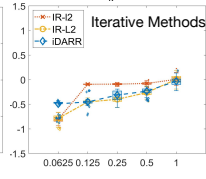
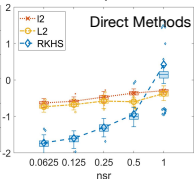
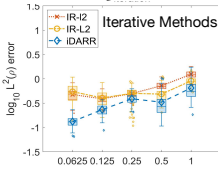
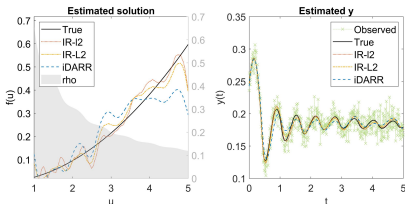
Iterative method  $\approx$  direct method



True function in FSOI



True function outside FSOI

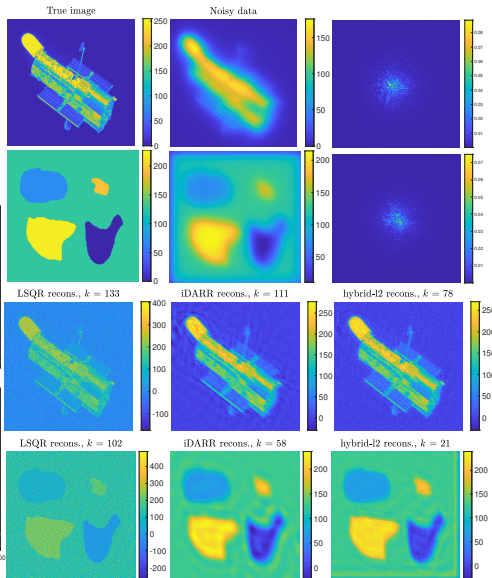
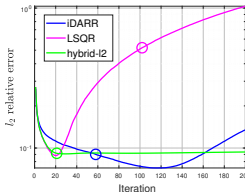
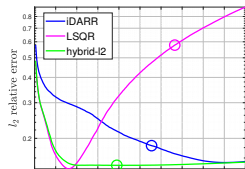




# Image deblurring

## Image deblurring

Gazzola+Hansen+Nagy2019  
256x256; 320x320





# Summary

Learning kernels in operators:

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

## Nonlocal dependence

- Identifiability
- DARTR: data adaptive RKHR Tikhonov-Reg
  - ▶ Synthetic data: convergent, robust to noise
  - ▶ Homogenization: resolution-invariant
- Iterative method: iDARR

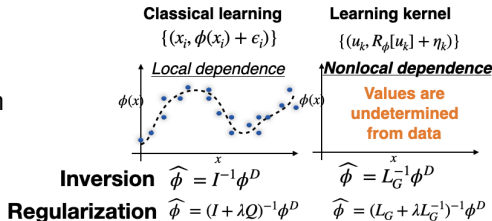
Regularization:  $A_n x_n = b_n \Rightarrow "x_{\lambda,n} = (A_n + \lambda A_n^{-1}) b_n "$

“ What doesn't kill you makes you stronger.”

## Future directions

### Learning with nonlocal dependence

- Minimax rate:  $\Delta x, N$
- Automatic kernel regression
- Regularization for ML:  
 $\|\phi_\theta\|_{rkhs}^2$ , not  $\|\theta\|$



Thank you for your attention!