

Introduction to Nonparametric Learning of Kernels in Operators

Fei Lu

Department of Mathematics, Johns Hopkins University

Plan:

Lecture 1. Overview and a review of classical learning theory

Lecture 2. Learning interaction kernels in interacting particle systems

Lecture 3. Coercivity condition and minimax rate of convergence

Lecture 4. Learning interaction kernels in mean-field equations

Lecture 5. Data adaptive RKHS Tikhonov regularization

Lecture 6. Small noise analysis of RKHS regularizations

Lecture 5. DARTR: data adaptive RKHS Tikhonov regularization

Learning interaction kernel $K_\phi(x - y) = \phi(|x - y|) \frac{x-y}{|x-y|}$

$$dX_t^i = \frac{1}{N} \sum_{j=1}^N K_\phi(X_t^j - X_t^i) dt + \sqrt{2\nu} dB_t^i \quad \Leftrightarrow R_\phi(\mathbf{X}_t) = \dot{\mathbf{X}}_t - \sqrt{2\nu} \dot{\mathbf{B}}_t$$

$$\partial_t u = \nu \Delta u + \nabla \cdot [u(K_\phi * u)] \quad \Leftrightarrow R_\phi[u(\cdot, t)] = f(\cdot, t)$$

- ▶ N small: M trajectories $\{\mathbf{X}_{t_1:t_L}^{(m)}\}_{m=1}^M \Rightarrow$ statistical learning
- ▶ $N = \infty$: $\{u_N(x, t_l) = \frac{1}{N} \sum_i \delta_{X_{t_l}^i}(x)\}$ or $\{u(x_m, t_l)\}_{m,l} \Rightarrow$ inverse problem

Nonparametric regression: $\phi = \sum_{k=1}^n c_k \phi_k$

$$A_n c = b_n$$

Ill-conditioned A_n \Rightarrow **Regularization**

$$\|A_n c - b_n\|^2 + \lambda \|c\|_*^2$$

How to choose $\|c\|_*^2$?

1. Learning kernels
2. Regression and regularization
3. Identifiability and DARTR
4. Adaptive prior for Bayesian inverse

- LLA22: Lu, Lang, and An. MSML22.

- LAY22: Lu, An and Yu. arXiv2205

- CLLW22: Chada, Lang, L. & Wang. arXiv2212

Outline

1. Learning kernels
2. Regression and regularization
3. Identifiability and DARTR
4. Adaptive prior for Bayesian inverse

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

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 $R_\phi[u](x) = \int \phi(|x - y|)g[u](x, y)dy$
 - interacting particles/agents

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

$$R_\phi[\mathbf{X}_t] = \left[-\frac{1}{N} \sum_{j=1}^N K_\phi(X_t^i - X_t^j) \right]_i = \dot{\mathbf{X}}_t + \dot{\mathbf{W}}_t, \quad u = \frac{1}{N} \sum_{i=1}^N \delta(X_t^i - x)$$

- nonlocal PDEs: $R_\phi[u] = \partial_{tt}u - v$

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

- 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 $R_\phi[u]$: **linear in ϕ**
- ▶ Data: discrete/noisy, **Nonlocal dependence**
 - random $(u_k, f_k) \sim \mu \otimes \nu$: **statistical learning**
 - deterministic (e.g., N small): **inverse problem**

Outline

1. Learning kernels
2. Regression and regularization
3. Identifiability and DARTR
4. Adaptive prior for Bayesian inverse

Nonparametric regression

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

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

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

Nonparametric regression

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

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

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

Three issues

- ▶ \bar{A}_n^{-1} : ill-conditioned/singular
- ▶ 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
- ▶ \bar{b}_n : noisy or with numerical error

Tikhonov/ridge Regularization: ($\|c\|_{B_*}^2 = c^\top B_* c$)

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

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

Regularization

Regularization is necessary:

- ▶ \bar{A}_n ill-conditioned
- ▶ \bar{b}_n : noisy or with numerical error

Tikhonov/ridge Regularization: ($\|c\|_{B_*}^2 = c^\top B_* c$)

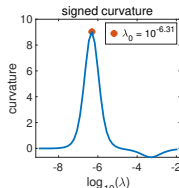
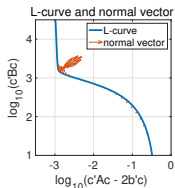
$$\mathcal{E}_\lambda(\phi) = \mathcal{E}(\phi) + \lambda \|\phi\|_*^2 \Rightarrow c^\top \bar{A}_n c - 2\bar{b}_n^\top c + \lambda \|c\|_{B_*}^2$$

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

- ▶ hyper-parameter λ : CV, truncated SVD, ...
- ▶ **The L-curve method** [Hansen00]

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

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



Regularization

Regularization is necessary:

- ▶ \bar{A}_n ill-conditioned
- ▶ \bar{b}_n : noisy or with numerical error

Tikhonov/ridge Regularization: ($\|c\|_{B_*}^2 = c^\top B_* c$)

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

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

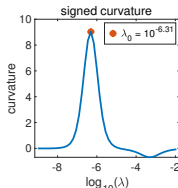
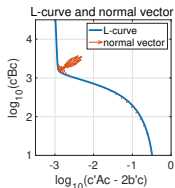
- ▶ hyper-parameter λ : CV, truncated SVD, ...
- ▶ The L-curve method [Hansen00]

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

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

- ▶ Which norm $\|\cdot\|_*$?

2 Regression and regularization



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 ϕ ? Identifiability?

Outline

1. Learning kernels
2. Regression and regularization
- 3. Identifiability and DARTR**
4. Adaptive prior for Bayesian inverse

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

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

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

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

- ▶ Function space of identifiability (FSOI):

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

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

DARTR: Data Adaptive RKHS Tikhonov Regularization

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

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

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

A new task for Regularization:

ensure that the learning takes place in the FSOI

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

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

▶ 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: $\|\phi\|_{H_G}^2$

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

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 $C = \sum_i \phi_i \otimes \phi_i$ 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

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

Output: reguarized estimator

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

- ▶ Generalized eigenvalue problem (A_n, B_n) : $A_n V = B V \Lambda$
- ▶ $B_{rkhs} = (V \Lambda V^\top)^{-1}$: avoid inverse ↓
Set $D = B_n^{-1} A_n^{1/2}$, then, $\hat{c}_\lambda = D(D^\top A_n D + \lambda I)^{-1} D^\top b_n$
- ▶ L-curve to select λ_*

Connecting computation with theory

$$\hat{c}_\lambda = (A_n + \lambda_* B_{rkhs})^{-1} b_n \quad \leftrightarrow \quad \hat{\phi} = \mathcal{L}_{\bar{G}}^{-1} \phi^D$$

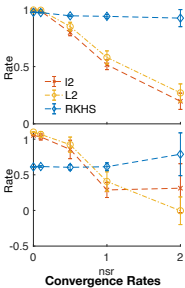
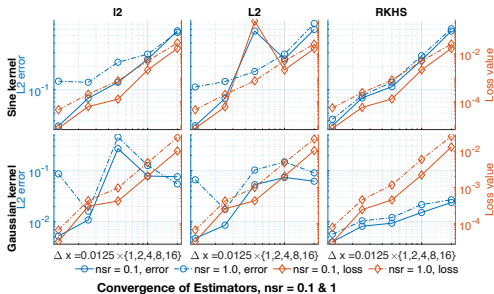
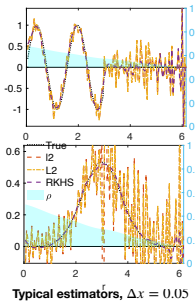
Theorem: If $\mathcal{L}_{\bar{G}}(L^2(\rho)) \subset \mathcal{H}$, then, $\mathcal{L}_{\bar{G}}$ eigenvalues \leftrightarrow GEP (A_n, B_n) .

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

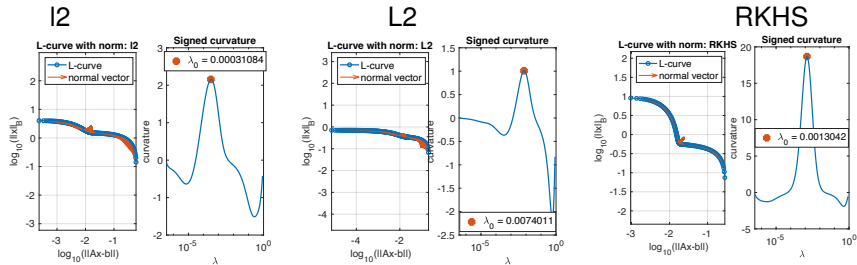
Interaction kernel in a nonlinear operator

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

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

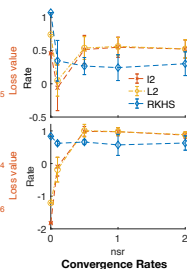
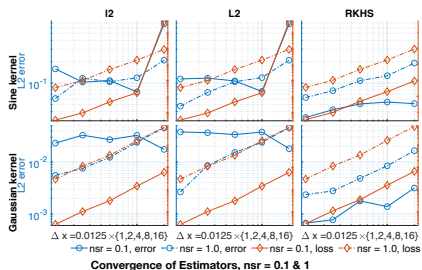
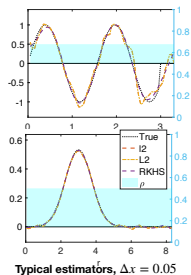


More robust L-curve



Linear integral operator:

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

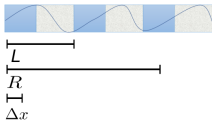


Robust in accuracy, consistent rates

Homogenization of wave propagation in meta-material

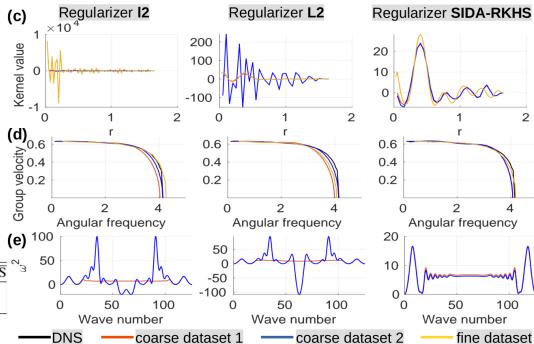
- ▶ heterogeneous bar with microstructure + DNS \Rightarrow Data
- ▶ Homogenization: $R_\phi[u] = \int_\Omega \phi(|y|)[u(x+y) - u(x)]dy = \partial_{tt}u - g$.

(a) Wave propagation in a heterogeneous bar



(b) Displacement error on a cross-validation dataset

Resolution	l^2	L^2	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 L^2 leading to non-physical kernel

Outline

1. Learning kernels
2. Regression and regularization
3. Identifiability and DARTR
4. Adaptive prior for Bayesian inverse

Bayesian inverse problem

Variational inverse problem: **ill-posed**

$$\hat{\phi} = \arg \min_{\phi} \mathcal{E}(\phi) = \langle \mathcal{L}_{\bar{G}} \phi, \phi \rangle_{L^2(\rho)} - 2 \langle \phi^D, \phi \rangle_{L^2(\rho)} + C.$$

Finite dimensional case:

- ▶ **Prior** $\mathcal{N}(0, \mathcal{Q}_0)$ with $\mathcal{Q}_0 > 0$: $\frac{d\pi_0(\phi)}{d\phi} \propto e^{-\frac{1}{2} \langle \phi, \mathcal{Q}_0^{-1} \phi \rangle_{L^2(\rho)}}$
- ▶ **Likelihood** of data: $\pi_L \sim \mathcal{N}(\mathcal{L}_{\bar{G}}^{-1} \phi^D, \sigma_\eta^2 \mathcal{L}_{\bar{G}}^{-1})$

$$\frac{d\pi_L(\phi)}{d\phi} \propto \exp\left(-\frac{1}{2\sigma_\eta^2} \mathcal{E}(\phi)\right)$$

- ▶ **Posterior** = Gaussian $\mathcal{N}(\mu_1, \mathcal{Q}_1)$ **well-posed**

$$\frac{d\pi_1(\phi)}{d\phi} \propto \exp\left(-\frac{1}{2} [\sigma_\eta^{-2} \mathcal{E}(\phi) + \langle \mathcal{Q}_0^{-1} \phi, \phi \rangle_{L^2(\rho)}]\right)$$

$$\mu_1 = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1} \phi^D, \text{ and } \mathcal{Q}_1 = \sigma_\eta^2 (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1}.$$

Infinite-dimensional prior and posterior

$\phi \in L^2_\rho$, infinite-dimensional:

- ▶ **Prior** $\mathcal{N}(0, \mathcal{Q}_0)$, where $\mathcal{Q}_0 > 0$, trace-class operator;
- ▶ **Posterior** $\mathcal{N}(\mu_1, \mathcal{Q}_1)$: Radon–Nikodym derivative w.r.t the prior

$$\frac{d\pi_1}{d\pi_0} \propto \exp\left(-\frac{1}{2}\sigma_\eta^{-2}\mathcal{E}(\phi)\right) = \exp\left(-\frac{1}{2}\sigma_\eta^{-2}[\langle \mathcal{L}_{\bar{G}}\phi, \phi \rangle_{L^2(\rho)} - 2\langle \phi^D, \phi \rangle_{L^2(\rho)} + C]\right),$$

$$\mu_1 = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1} \phi^D, \text{ and } \mathcal{Q}_1 = \sigma_\eta^2 (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1}.$$

Is it a good idea to choose $\mathcal{Q}_0 > 0$?

Risk in a non-degenerate prior

Theorem ([CLLW22])

A non-degenerate prior has the risk of leading to a divergent posterior mean in the small noise limit.

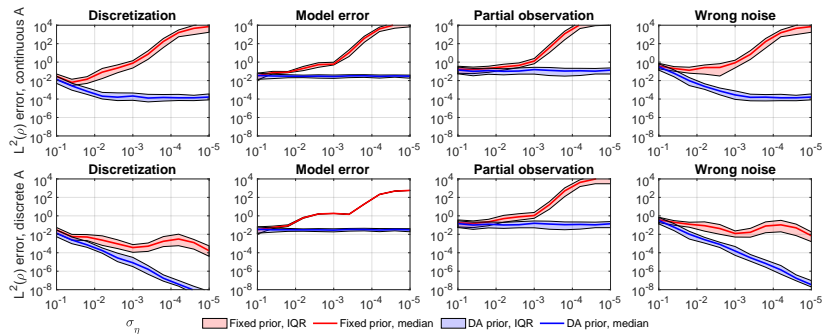
I.e., $\lim_{\sigma_\eta \rightarrow 0} (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1} \phi^D = \infty$ under Assumption (next page).

Risk in a non-degenerate prior

Theorem ([CLLW22])

A non-degenerate prior has the risk of leading to a divergent posterior mean in the small noise limit.

i.e., $\lim_{\sigma_\eta \rightarrow 0} (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1} \phi^D = \infty$ under Assumption (next page).



Assumptions: $\mathcal{L}_{\bar{G}} \sim \{(\lambda_k, \psi_k)\}$

▶ *Finite rank $\mathcal{L}_{\bar{G}}$:*

$\lambda_k > 0$ for $1 \leq k \leq K$, $\lambda_k = 0$ for $k > K$. $\dim(L^2(\rho)) > K$.

As a result, the FSOI $H = \text{span}\{\psi_i\}_{i=1}^K \subsetneq L^2_\rho$.

▶ *Prior covariance commutes with $\mathcal{L}_{\bar{G}}$:*

The prior $\mathcal{N}(0, \mathcal{Q}_0)$ satisfies $\mathcal{Q}_0\psi_i = r_i\psi_i$ with $r_i > 0$ for all i .

▶ *Error outside of FSOI.*

$$\phi^D = \mathcal{L}_{\bar{G}}\phi_{true} + \epsilon^\xi + \epsilon^\eta,$$

The model error $\epsilon^\xi = \sum_i \epsilon_i^\xi \psi_i$ has $\epsilon_{i_0}^\xi \neq 0$ with $i_0 > K$.

Note: $\epsilon^\eta \sim \mathcal{N}(0, \sigma_\eta^2 \mathcal{L}_{\bar{G}})$ comes from noise.

Proof:

$$\phi^D = \sum_i \phi_i^D \psi_i, \text{ with } \phi_i^D = \lambda_i \phi_{true,i} + \sigma_\eta \lambda_i^{1/2} \epsilon_i^\eta + \epsilon_i^\xi.$$

The posterior mean $\mu_1 = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1} \phi^D$ becomes

$$\mu_1 = \sum_i (\lambda_i + \sigma_\eta^2 r_i^{-1})^{-1} \phi_i^D \psi_i = \sum_{i=1}^K (\lambda_i + \sigma_\eta^2 r_i^{-1})^{-1} \phi_i^D \psi_i + \sum_{i>K} \sigma_\eta^{-2} r_i \epsilon_i^\xi \psi_i.$$

The first term $\rightarrow \sum_{1 \leq i \leq K} (\phi_{true,i} + \lambda_i^{-1} \epsilon_i^\xi) \psi_i$

But the 2nd term $\sum_{i>K} \sigma_\eta^{-2} r_i \epsilon_i^\xi \psi_i$ diverges because $\exists \epsilon_{i_0}^\xi \neq 0$.

Data-adaptive prior

Fixed prior $\phi \sim \pi_0 = \mathcal{N}(0, \mathcal{Q}_0)$, $\pi_1 = \mathcal{N}(\mu_1, \mathcal{Q}_1)$

$$\mu_1 = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1} \phi^D$$

Data-adaptive prior $\phi \sim \pi_0 = \mathcal{N}(0, \lambda_*^{-1} \mathcal{L}_{\bar{G}})$, $\pi_1 = \mathcal{N}(\mu_1, \mathcal{Q}_1^D)$

$$\mu_1^D = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \lambda_* \mathcal{L}_{\bar{G}}^{-1})^{-1} \phi^D$$

$$\mathcal{Q}_1^D = \sigma_\eta^2 (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \lambda_* \mathcal{L}_{\bar{G}}^{-1})^{-1}$$

Data-adaptive prior

Fixed prior $\phi \sim \pi_0 = \mathcal{N}(0, \mathcal{Q}_0)$, $\pi_1 = \mathcal{N}(\mu_1, \mathcal{Q}_1)$

$$\mu_1 = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \mathcal{Q}_0^{-1})^{-1} \phi^D$$

Data-adaptive prior $\phi \sim \pi_0 = \mathcal{N}(0, \lambda_*^{-1} \mathcal{L}_{\bar{G}})$, $\pi_1 = \mathcal{N}(\mu_1^D, \mathcal{Q}_1^D)$

$$\mu_1^D = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \lambda_* \mathcal{L}_{\bar{G}}^{-1})^{-1} \phi^D$$

$$\mathcal{Q}_1^D = \sigma_\eta^2 (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \lambda_* \mathcal{L}_{\bar{G}}^{-1})^{-1}$$

Theorem (Small noise limit of the posterior mean [CLLW22])

Fixed prior with $\mathcal{Q}_0 = I_d$: blows up;

Data-adaptive prior: stable and convergent

Proof With $\sigma_\eta^2 \lambda_* > 0$:

$$\begin{aligned} (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \lambda_* \mathcal{L}_{\bar{G}}^{-1})^{-1} \psi_i &= (\mathcal{L}_{\bar{G}}^2 + \sigma_\eta^2 \lambda_* I)^{-1} \mathcal{L}_{\bar{G}} \psi_i \\ &= \frac{\lambda_i}{\lambda_i^2 + \sigma_\eta^2 \lambda_*} \psi_i = 0, \text{ if } i \geq K + 1. \end{aligned}$$

Then, the posterior mean is:

$$\mu_1^D = (\mathcal{L}_{\bar{G}} + \sigma_\eta^2 \lambda_* \mathcal{L}_{\bar{G}}^{-1})^{-1} \phi^D = \sum_{1 \leq i \leq K} (\lambda_i + \sigma_\eta^2 \lambda_* \lambda_i^{-1})^{-1} \phi_i^D \psi_i.$$

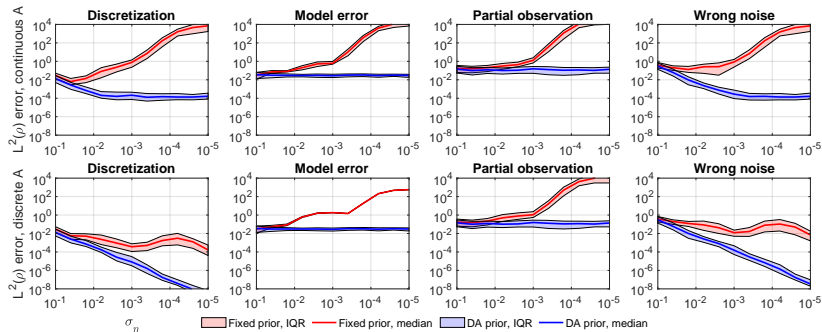
The limit exist as $\sigma_\eta \rightarrow 0$:

(recall that $\phi^D = \sum_i \phi_i^D \psi_i$, with $\phi_i^D = \lambda_i \phi_{true,i} + \sigma_\eta \lambda_i^{1/2} \epsilon_i^\eta + \epsilon_i^\xi$)

$$\lim_{\sigma_\eta \rightarrow 0} \mu_1^D = \sum_{i=1}^K \left(\phi_{true,i} + \lambda_i^{-1} \epsilon_i^\xi \right) \psi_i.$$

Small noise limit of posterior mean

Integral operator, ϕ_{true} inside the FSOI:

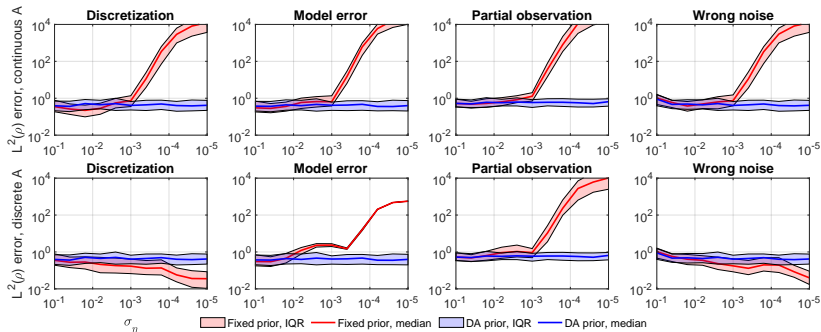


Fixed prior: **blowup**

Data-adaptive prior: **stable and convergent**

Small noise limit of posterior mean

Integral operator, ϕ_{true} outside the FSOI:



Fixed prior: **blowup**

Data-adaptive prior: **stable and convergent**

Summary

Learning kernels in operators:

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

Nonlocal dependence

- ▶ Identifiability: data-dependent FSOI
- ▶ DARTR: data adaptive RKHR Tikhonov-Reg
 - Synthetic data: convergent, robust to noise
 - Homogenization: resolution-independent
- ▶ Small noise analysis in Bayesian

Next:

How do we theoretically compare different regularizations?

Small noise analysis_[LO23,LL23]

Future directions

Inverse problems \leftrightarrow Learning with nonlocal dependence

- ▶ Iterative methods for DARTR
- ▶ Convergence: $\Delta x, N$
- ▶ Regularization for ML: $\|\phi_\theta\|_{rhs}^2$, not $\|\theta\|$

Next:

How do we theoretically compare different regularizations?

Small noise analysis_[LO23,LL23]

Future directions

Inverse problems \leftrightarrow Learning with nonlocal dependence

- ▶ Iterative methods for DARTR
- ▶ Convergence: $\Delta x, N$
- ▶ Regularization for ML: $\|\phi_\theta\|_{rhs}^2$, not $\|\theta\|$

References

(@ <http://www.math.jhu.edu/~feilu>)

- LLA22: Lu, Lang, and An. MSML22. (Matlab code available)
- LAY22: Lu, An and Yu. arXiv2205
- CLLW22: Chada, Lang, Lu, and Wang. arXiv2212