

UNIFIED RKHS METHODOLOGY AND ANALYSIS FOR FUNCTIONAL
LINEAR AND SINGLE-INDEX MODELS.

KRISHNA BALASUBRAMANIAN

DEPARTMENT OF STATISTICS, UC DAVIS

ACKNOWLEDGEMENTS

- ▷ Joint work with:
 - ▷ Hans-Georg Muller, UC Davis.
 - ▷ Bharath K Sriperumbudur, Penn State.
- ▷ Paper available in arXiv soon

INTRODUCTION: PROBLEM SETUP

▷ Given:

▷ Domain $S = [0, 1]$.

▷ Input/predictor process $X(t)$, $t \in S$.

▷ Output/response $Y \in \mathbb{R}$.

▷ Functional linear model:

$$Y = \int_S X(t)\beta^*(t) dt + \epsilon = \langle X, \beta^* \rangle_{L^2(S)} + \epsilon,$$

▷ Here, ϵ is an exogenous additive noise such that

$$\mathbb{E}[\epsilon|X] = 0 \text{ and } \mathbb{E}[\epsilon^2] = \sigma^2.$$

- ▷ Functional single-index model:

$$Y = g \left(\int_S X(t) \beta^*(t) dt \right) + \epsilon = g \left(\langle X, \beta^* \rangle_{L^2(S)} \right) + \epsilon, \quad (1)$$

for some function $g : \mathbb{R} \rightarrow \mathbb{R}$.

- ▷ The parameter β^* is called the index and the function g the link function.
- ▷ When $g(a) = a$, the single-index model in (1) becomes the functional linear model.

INTRODUCTION: PROBLEM SETUP

- ▷ Given n observations $\{(X_i, Y_i)\}_{1 \leq i \leq n}$ that are independent and identically distributed copies of (X, Y) , we study how to estimate the index parameter β^* in (1).
- ▷ Estimation procedure is agnostic to the specification of the link function - interaction between the allowed class of link functions and the distribution of the covariate X becomes crucial.
- ▷ Throughout, we assume that X is a zero-mean Gaussian process.

INTRODUCTION: RELATED WORK

- ▷ Yuan and Cai (2010), and Cai and Yuan (2012) considered an RKHS approach for linear setting. They assumed the truth lies inside the RKHS. Avoids restrictive eigen-gap assumptions made in prior FPCA-based works.
- ▷ Muller and Stadtmuller (2005) proposed and analyzed an MLE-based approach for generalized functional linear models (special cases of single-index models) and established consistency results.
- ▷ Shang and Cheng (2015) considered an RKHS approach for the generalized functional linear models and established inferential results when the truth is in RKHS.
- ▷ The above works require knowledge of the link function g to estimate the index parameter.

INTRODUCTION: METHODOLOGICAL CONTRIBUTIONS

- ▷ We provide a unified framework for estimating the index for both the linear and single-index models, for a wide class of *unknown* link functions.
- ▷ Specifically, we illustrate that the standard functional linear RKHS least-squares estimator also provides an efficient estimator of the index parameter in the single-index model under the Gaussian process assumption.
- ▷ Justification based on *infinite-dimensional* analogues of Gaussian Stein's identity.
- ▷ Naturally handles mis-specification with respect to the link function for both the linear and single-index models.

INTRODUCTION: THEORETICAL CONTRIBUTIONS

- ▷ Rates of estimating the index depends on
 - ▷ Integral operator T associated to the RKHS
 - ▷ Covariance operator C of the Gaussian process X .
- ▷ Compared to previous works, we provide results without:
 - ▷ Restrictive commutativity assumptions on T and C .
 - ▷ $\tilde{\beta}^*$ being inside the RKHS under consideration.

Methodology

- ▷ Infinite-dimensional extensions of Gaussian Stein's identity:

For a zero-mean Gaussian random element X in a separable Hilbert space with covariance operator C , and for smooth enough real-valued functions f , we have

$$\mathbb{E}[Xf(X)] = C\mathbb{E}[\nabla f(X)],$$

where ∇ is the Fréchet derivative.

- ▷ In our context, by leveraging the version of Stein's identity for Hilbert-valued random vectors, we have

$$\mathbb{E}[YX] = \mathbb{E}[\nabla g(\langle \beta, X \rangle)] = \vartheta_{g, \beta^*} C \beta^*,$$

where ∇ is the Fréchet derivative.

METHODOLOGY

- ▷ $\vartheta_{\mathbf{g},\beta^*}$ is a constant depending on the link function g and the index β^* .
- ▷ The exact form of the constant is irrelevant for our purpose as we focus on estimating the direction of the index parameter.
- ▷ We assume that g is such that $\vartheta_{\mathbf{g},\beta^*} \neq 0$ throughout the rest of the paper.
- ▷ In particular, when g is the identity function, it is easy to see that we have $\vartheta_{\mathbf{g},\beta^*} = 1$.
- ▷ We define $\tilde{\beta}^* := \vartheta_{\mathbf{g},\beta^*} \beta^*$, to handle the single-index and linear model in a unified manner.

- ▷ Based on this, note that we have

$$\tilde{\beta}^* := \arg \min_{\beta \in L^2(S)} \mathbb{E} [Y - \langle X, \beta \rangle]^2.$$

- ▷ Given $(X_1, Y_1), \dots, (X_n, Y_n)$ be n i.i.d. copies of random variables (X, Y) . For some $\lambda > 0$, our estimator based on minimizing the penalized least-squares criterion over the RKHS \mathcal{H} is given by:

$$\hat{\beta}_{n,\lambda} = \arg \min_{\beta \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n [Y_i - \langle \beta, X_i \rangle]^2 + \lambda \|\beta\|_{\mathcal{H}}^2. \quad (2)$$

METHODOLOGY

- ▷ Let \mathcal{H} be an RKHS with the associated kernel $k : S \times S \rightarrow R$.
- ▷ Define $\mathfrak{J} : \mathcal{H} \rightarrow L^2(S)$, $f \mapsto f$, to be the inclusion operator mapping functions in the RKHS \mathcal{H} to $L^2(S)$.
- ▷ We use $\mathfrak{J}^* : L^2(S) \rightarrow \mathcal{H}$ to refer to the adjoint of \mathfrak{J} .
- ▷ We also define the following two important operators that arise in our analysis:

$$T := \mathfrak{J}\mathfrak{J}^* : L^2(S) \rightarrow L^2(S),$$

$$C := \mathbb{E}[X \otimes X] : L^2(S) \rightarrow L^2(S),$$

where \otimes represents the $L^2(S)$ tensor product.

- ▷ Note that the solution of the above optimization problem is given by

$$\hat{\beta}_{n,\lambda} = \left[\mathfrak{J}^* \left(\frac{1}{n} \sum_{i=1}^n X_i \otimes X_i \right) \mathfrak{J} + \lambda I \right]^{-1} \mathfrak{J}^* \left[\frac{1}{n} \sum_{i=1}^n Y_i X_i \right].$$

METHODOLOGY

- ▷ By applying the representer theorem it follows that

$$\hat{\beta} \in \text{span} \left\{ \int_S k(\cdot, t) X_i(t) dt : i = 1, \dots, n \right\},$$

i.e., $\exists \alpha := (\alpha_1, \dots, \alpha_n)^\top \in \mathbb{R}^n$ such that

$$\hat{\beta} = \sum_{i=1}^n \alpha_i \int_S k(\cdot, t) X_i(t) dt.$$

- ▷ Solving for α yields

$$\alpha = (\mathbf{K} + n\lambda I)^{-1} \mathbf{y},$$

where

$$\mathbf{K} \in \mathbb{R}^{n \times n} \text{ with } [\mathbf{K}]_{ij} := \int_S \int_S k(t, s) X_i(t) X_j(t) dt ds$$

and $\mathbf{y} = (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n$.

- ▷ Therefore, $\hat{\beta}$ can be computed by solving a finite dimensional linear system of size n , which is not obvious from the previous expression.

Theory

- ▷ Let $\|T^{-\alpha}\tilde{\beta}^*\| < \infty$, i.e., $\tilde{\beta}^* \in \mathcal{R}(T^\alpha)$ for $\alpha \in (0, 1/2]$.
- ▷ Define

$$\varkappa := \mathbb{E} \left[\left(g(\langle X, \tilde{\beta}^* \rangle) - \langle X, \tilde{\beta}^* \rangle \right)^4 \right]. \quad (3)$$

- ▷ Suppose one of the following conditions hold:
 - (a) $Tr(C^{1/2}) < \infty$ and $\varkappa \in (0, \infty)$,
 - (b) $\varkappa = 0$ and $Tr(C) < \infty$.

- ▷ The assumption $\tilde{\beta}^* \in \mathcal{R}(T^\alpha)$ imposes certain smoothness condition on $\tilde{\beta}^*$. It is well-known that $\tilde{\beta}^* \in \mathcal{H}$ when $\alpha = \frac{1}{2}$, which we refer to as the *well-specified setting*. This assumption is equivalent to the condition that $\tilde{\beta}^*$ lies in an interpolation space between $L^2(S)$ and \mathcal{H} with α being the interpolating index.
- ▷ While $\text{Tr}(C) < \infty$ is guaranteed by the well-definedness of the Gaussian process. The following Theorem requires a slightly stronger condition given as $\text{Tr}(C^{1/2}) < \infty$, when $\varkappa \neq 0$.
- ▷ The parameter \varkappa captures the degree of non-linearity of the model. Indeed, $\varkappa = 0$ implies $g(\langle X, \tilde{\beta}^* \rangle) = \langle X, \tilde{\beta}^* \rangle$ with probability 1. Conversely, when the model is linear, $\varkappa = 0$.

▷ Define

$$\begin{aligned}\Theta &:= T^\alpha (CT + \lambda I)^{-1} C (TC + \lambda I)^{-1} T^\alpha, \\ d(\lambda) &:= \frac{\text{Tr}(\Theta)}{\|\Theta\|}, \\ \Xi &:= T(T^{1/2}CT^{1/2} + \lambda I)^{-2} T, \\ N(\lambda) &:= \text{Tr} \left[(T^{1/2}CT^{1/2} + \lambda I)^{-1} T^{1/2}CT^{1/2} \right].\end{aligned}$$

▷ Let $\delta \in (0, 1/e]$, $n \gtrsim (d(\lambda) \vee \log(1/\delta))$ and let

$$\frac{\text{Tr}(T^{1/2}CT^{1/2})}{n} \lesssim \lambda \lesssim \|T^{1/2}CT^{1/2}\|. \quad (4)$$

▷ **Theorem:** With probability at least $1 - 3\delta$, we have

$$\|\hat{\beta} - \tilde{\beta}^*\| \lesssim \text{bias}(\lambda) + \|\Xi\|^{\frac{1}{4}} \sqrt{\frac{(\sigma^2 + \sqrt{\varkappa})N(\lambda)}{n\delta}} + \lambda \|\Xi\|^{\frac{1}{4}} \left(\left\| T^{1/2} C T^{1/2} \right\|^{1/2} + \sqrt{\lambda} \right) \|T\|^{\frac{1}{2}-\alpha} \|T^{-\alpha} \tilde{\beta}^*\| \sqrt{\frac{\|\Theta\| \text{Tr}(\Theta)}{n}},$$

where $\text{BIAS}(\lambda) := \|T(C T + \lambda I)^{-1} C \tilde{\beta}^* - \tilde{\beta}^*\|$.

Commutative Setting

- ▷ Let $\|T^{-\alpha}\tilde{\beta}^*\| < \infty$ for $\alpha \in (0, 1/2]$. Suppose the operators T and C commute and have simple eigenvalues (i.e., of multiplicity one) denoted by μ_i and ξ_i for $i \in \mathbb{N}$, such that,

$$i^{-t} \lesssim \mu_i \lesssim i^{-t} \quad \text{and} \quad i^{-c} \lesssim \xi_i \lesssim i^{-c}, \quad (5)$$

where $t > 1$. Suppose one of the following conditions hold:
(a) $\varkappa \in (0, \infty)$ and $c > 2$, (b) $\varkappa = 0$ and $c > 1$.

▷ **Theorem:** We have that

$$\|\hat{\beta} - \tilde{\beta}^*\| \lesssim_p n^{-\frac{\alpha t}{1+c+2t(1-\alpha)}} \quad (6)$$

for

$$\lambda = n^{-\frac{t+c}{1+c+2t(1-\alpha)}}. \quad (7)$$

- ▷ When $\alpha = 1/2$, i.e., $\tilde{\beta}^* \in \mathcal{H}$ (well-specified case), we obtain

$$\|\hat{\beta} - \tilde{\beta}^*\| \lesssim_{\rho} n^{-\frac{t}{2(1+t+c)}},$$

which matches with the minimax optimal rate obtained in Yuan and Cai (2010), when the model is linear.

- ▷ However, the interesting point is that even in the single-index model setting, we obtain the same rate as obtained for the linear model (when $c > 2$) as long as $\varkappa < \infty$.
- ▷ For the linear model setting, the above result extends the results of Cai and Yuan (2010), to the misspecified setting, i.e., $\tilde{\beta}^* \in L^2(S) \setminus \mathcal{H}$
- ▷ The requirement of $c > 2$ ensures that $\text{Tr}(C^{1/2}) < \infty$.

Non-commutative Setting

- ▷ Let $(\zeta_i)_{i \in \mathbb{N}}$ denote the eigenvalues of $T^{1/2}CT^{1/2}$ with $i^{-b} \lesssim \zeta_i \lesssim i^{-b}$, for some $b > 1$.
- ▷ Suppose $\tilde{\beta}^* \in \mathcal{R}(T^{1/2}(T^{1/2}CT^{1/2})^\nu)$ for $\nu \in (0, 1]$ and $\kappa < \infty$.
- ▷ **Theorem:** For

$$\lambda = n^{-\frac{b}{1+b+2b\nu}},$$

we have

$$\|\hat{\beta} - \tilde{\beta}^*\| \lesssim_p n^{-\frac{b\nu}{1+b+2b\nu}}.$$

- ▷ Unlike in the commutative case, the results are presented in terms of the eigen decay behavior of $T^{1/2}CT^{1/2}$. When T and C commute, we obtain $b = t + c$.
- ▷ We would like to highlight that to the best of our knowledge, no result is known in the literature for the estimation error, i.e., $\|\hat{\beta} - \tilde{\beta}^*\|$, in the non-commutative setting, even for linear models.

- ▷ The assumption $\tilde{\beta}^* \in \mathcal{R}(T^{1/2}(T^{1/2}CT^{1/2})^\nu)$, implies $\exists h \in L^2(S)$ such that

$$T^{1/2}(T^{1/2}CT^{1/2})^\nu h = \tilde{\beta}^*,$$

which implies $\tilde{\beta}^* \in \mathcal{R}(T^{1/2}) = \mathcal{H}$.

- ▷ Therefore, the assumption $\tilde{\beta}^* \in \mathcal{R}(T^{1/2}(T^{1/2}CT^{1/2})^\nu)$ is stronger than assuming $\tilde{\beta}^* \in \mathcal{R}(T^{1/2})$.
- ▷ The key reason to make this strong assumption is to control $\text{BIAS}(\lambda)$ in a finer manner and obtain meaningful convergence rates. Indeed, by simply assuming $\tilde{\beta}^* \in \mathcal{R}(T^{1/2})$ ensures $\text{BIAS}(\lambda) \rightarrow 0$ as $\lambda \rightarrow 0$, using which consistency of $\hat{\beta}$ can be established, but with no handle on the convergence rate.

- ▷ Let (ζ_i, ϕ_i) and (μ_i, ψ_i) for $i \in \mathbb{N}$, denote the eigensystems of $T^{1/2}CT^{1/2}$ and T respectively. Suppose

$$i^{-b} \lesssim \zeta_i \lesssim i^{-b} \text{ and } i^{-t} \lesssim \mu_i \lesssim i^{-t}$$

for some $b, t > 1$. Let the eigenfunctions of $T^{1/2}CT^{1/2}$ and T satisfy

$$\sup_{i,l} \frac{1}{\mu_i \mu_l} \left| \sum_j \mu_j \langle \phi_i, \psi_j \rangle \langle \phi_l, \psi_j \rangle \right|^2 < \infty. \quad (8)$$

- ▷ **Theorem:** Assuming $\varkappa < \infty$ and $\tilde{\beta}^* \in \mathcal{R}(T^{1/2}(T^{1/2}CT^{1/2})^\nu)$ for some $\nu \in (0, \frac{1}{2} - \frac{t}{2b}]$, we have

$$\|\hat{\beta} - \tilde{\beta}^*\| \lesssim_p n^{-\frac{b\nu+(t-1)/2}{t+b+2b\nu}}$$

for

$$\lambda = n^{-\frac{b}{t+b+2b\nu}}. \quad (9)$$

THEORY

- ▷ For $\nu \in (0, \frac{1}{2} - \frac{t}{2b}]$, the rate in latter Theorem is clearly faster than that in former Theorem.

Interpreting Range Space Conditions on $\tilde{\beta}^*$

RANGE SPACE CONDITIONS

- ▷ **Proposition:** For $x, y \in [0, 1]$, suppose that the reproducing kernel k and the covariance function c are given respectively by

$$k(x, y) = \sum_{i \geq 1} a_i \phi_i(x) \phi_i(y), \quad c(x, y) = \sum_{m \geq 1} b_m \psi_m(x) \psi_m(x),$$

where $a_i \geq 0$ for all i , $b_m \geq 0$ for all m , $\sum_{i \geq 1} a_i \leq \infty$, $\sum_{m \geq 1} b_m \leq \infty$ and $(\phi_i)_i$ and $(\psi_m)_m$ form an orthonormal basis of $L^2([0, 1])$. Define $\tau_j := \sum_i a_i \eta_{ij}^2$ where $\eta_{ij} := \sum_{m \geq 1} b_m \theta_{mi} \theta_{mj}$ and $\theta_{mj} := \langle \psi_m, \phi_j \rangle$, and assume $\sup_j \tau_j < \infty$.

RANGE SPACE CONDITIONS

▷ Then the following hold:

(i) The RKHS induced by the kernel k is given by

$$\mathcal{H} = \left\{ f(x) = \sum_{i \geq 1} f_i \phi_i(x), x \in [0, 1] : \sum_i \frac{f_i^2}{a_i} < \infty \right\},$$

with the associated inner product defined by

$$\langle f, g \rangle_{\mathcal{H}} = \sum_i a_i^{-1} f_i g_i.$$

RANGE SPACE CONDITIONS

(ii) The space $\mathcal{R}(T^{1/2}(T^{1/2}CT^{1/2}))$ satisfies the inclusion

$$\mathcal{R}(T^{1/2}(T^{1/2}CT^{1/2})) \subset \tilde{\mathcal{H}} \subset \mathcal{H},$$

where

$$\tilde{\mathcal{H}} = \left\{ f(x) = \sum_i f_i \phi_i(x), x \in [0, 1] : \sum_i \frac{f_i^2}{a_i \tau_i} < \infty \right\},$$

is an RKHS induced by the kernel

$\tilde{k}(x, y) = \sum_{i \geq 1} a_i \tau_i \phi_i(x) \phi_i(y)$ with inner product

$$\langle f, g \rangle_{\tilde{\mathcal{H}}} = \sum_{i \geq 1} f_i g_i (\tau_i a_i)^{-1}.$$

A CONCRETE EXAMPLE

- ▷ Suppose $\phi_i(x) = \cos(i\pi x)$, $x \in [0, 1]$ and $\psi_m(\cdot) = \cos(\omega_m\pi\cdot)$ where $\omega_m = am + b$ for some $a, b \in \mathbb{R}$ such that $\omega_m \notin \mathbb{Z}$ and $m \in \mathbb{N}$. Let $b_m \lesssim m^{-(1+\delta)}$, for some $\delta > 0$.
- ▷ Then, we have

$$\theta_{mi} = \frac{\pi\omega_m}{\pi^2\omega_m^2 - (i\pi)^2} \sin(\pi\omega_m)(-1)^i.$$

Furthermore,

$$\eta_{ij} \lesssim (ij)^{-\min(1, \frac{\delta+1}{2})}, \quad (10)$$

- ▷ This implies that $\tau_j \lesssim j^{-\min(\delta+1, 2)}$ and $\sup_j |\tau_j| < \infty$. Hence, the inclusion $\mathcal{R}(T^{1/2}(T^{1/2}CT^{1/2})) \subset \tilde{\mathcal{H}} \subset \mathcal{H}$, follows, where $\tilde{\mathcal{H}}$ consists of functions that are $\min(1, \frac{1+\delta}{2})$ more smoother than the functions in \mathcal{H} .

MORE RESULTS

- ▷ In the paper we also provide:
 - ▷ similar results for prediction.
 - ▷ several other examples.

Thank you!