

HYBRID PARTIAL LEAST SQUARES REGRESSION WITH MULTIPLE FUNCTIONAL AND SCALAR PREDICTORS

BY JONGMIN MUN^{1,a} AND JEONG HOON JANG^{2,b}

¹*Data Sciences and Operations Department, Marshall School of Business, University of Southern California,*
^ajongmin.mun@marshall.usc.edu

²*Department of Biostatistics and Data Science, University of Texas Medical Branch,* ^bjejang@utmb.edu

Regression of a scalar response on mixed functional- and scalar-valued predictors, such as medical imaging with auxiliary patient information, introduces the new challenge of handling cross-modality correlations. To address this, we propose a hybrid partial least squares (PLS) regression framework that integrates functional and scalar predictors within a unified Hilbert space. We then extend the classical nonlinear iterative PLS (NIPALS) algorithm to this hybrid Hilbert space by iteratively maximizing the empirical cross-covariance between the hybrid predictor and the response. As a result, our method identifies low-dimensional representations that capture both within- and between-modality variation, as well as the response-predictor correlation. The procedure is computationally efficient, requiring only the solution of linear systems at each step. We provide theoretical properties to justify our algorithm and demonstrate its effectiveness through simulations and an application to clinical outcome prediction using renal imaging and scalar covariates from the Emory University renal study.

1. Introduction. Modern biomedical studies frequently collect diverse data types from each subject. As an illustrative example, the Emory University renal study (Chang et al., 2020; Jang, 2021) records both multiple renogram curves (functional data) and multiple renogram variables (scalar data) for each kidney. To effectively analyze such distinct yet related physiological signals, we construct a joint linear regression model incorporating both functional and scalar-valued i.i.d. covariates:

$$(1) \quad Y_i = \beta^\top \mathbf{Z}_i + \sum_{k=1}^K \int \beta_k(t) X_{ik}(t) dt + \epsilon_i, \quad i = 1, \dots, n,$$

where Y_i is a scalar response, $\mathbf{Z}_i = (Z_{i1}, \dots, Z_{ip})^\top$ is a Euclidean vector covariate, $X_{i1}(t), \dots, X_{iK}(t)$ are functional predictors that belong to $\mathbb{L}^2[0, 1]$, and ϵ_i is observational noise. In other words, this is a scalar-on-hybrid regression where the hybrid covariates belong to $(\mathbb{L}^2[0, 1])^K \times \mathbb{R}^p$. For notational convenience, we assume throughout this paper that the responses and predictors have been centered, allowing the intercept term to be ignored.

This hybrid model poses several challenges. First, estimating the infinite-dimensional slope function $\beta_k(t)$ from a finite number of data is an ill-posed problem. Second, we face the high-dimensionality of the predictors with $K + p$ dimensions. Third, strong correlations between functional and scalar predictors are common but often overlooked by separate modeling approaches. This paper addresses these three issues in a unified way by introducing a hybrid partial least squares (PLS) regression framework defined on a novel hybrid Hilbert space. Before summarizing the key contributions of this framework, we first review previous approaches and their limitations.

Keywords and phrases: dimension reduction, functional data analysis, multiple data modalities, multivariate data analysis, multivariate functional data, partial least square.

1.1. *Previous works and limitations.* To address the ill-posedness stemming from the infinite-dimensionality of the functional components, common remedies include basis expansion, using power-series (Goldsmith et al., 2011), B-splines (Cardot, Ferraty and Sarda, 2003; Cai and Hall, 2006), or wavelets (Zhao, Ogden and Reiss, 2012), and structure-aware regularization, such as roughness penalties (JM:) *comment: let's add citations here*.

For the high-dimensionality, a major solution is using derived inputs. A straightforward approach in this direction is principal component analysis (PCA) regression, which applies PCA separately to the functional- (by, for example, Happ and Greven 2018) and scalar predictors, and then performs a classical multivariate regression on the combined scores. PCA regression is well-studied for linear regression with functional predictors alone (Hall and Horowitz, 2007; Reiss and Ogden, 2007; Febrero-Bande, Galeano and González-Manteiga, 2017). However, since the PCA-derived inputs are not informed by the response variable, they are not guaranteed to capture the core regression relationship with a small number of derived inputs. In terms of predictive power, partial least squares (PLS) regression is a powerful alternative. It iteratively constructs a set of orthogonal latent components from the predictors that have the maximal covariance with the response variable, and use the resulting scores for regression. PLS for linear regression in the context of functional data was introduced by Preda and Saporta (2005) for the case of a single predictor, studied from a stochastic process perspective. Motivated by regression on chemometric spectra, Reiss and Ogden (2007), Aguilera et al. (2010), and Aguilera, Aguilera-Morillo and Preda (2016) extended the framework by incorporating basis approximations and roughness penalties to promote smoothness. Delaigle and Hall (2012) provided the first thorough theoretical analysis, while Saricam et al. (2022) proposed a computationally efficient procedure based on Golub–Kahan bidiagonalization. For a comprehensive review, see Febrero-Bande, Galeano and González-Manteiga (2017). More recently, Beyaztas and Lin Shang (2022) extended the framework to accommodate multiple functional predictors.

However, these existing approaches overlook the potentially strong correlations between functional and scalar components, which can lead to multicollinearity and suboptimal predictive performance. Multimodal correlations have mostly been addressed within an unsupervised learning framework. For instance, Kolar, Liu and Xing (2014) studied the estimation of joint undirected graphical models for functional and vector data, while Geng, Kolar and Koyejo (2020) considered joint precision matrix estimation for brain measurements and confounding scalar covariate. Jang (2021) proposed a joint PCA method that accounts for correlations between functional and scalar data. However, the resulting components are still not informed by the response and may fail to capture correlation with the outcome.

1.2. *Our contributions.* To address the gaps mentioned above, we unify the remedies proposed for these three issues: basis expansion, partial least squares, and accounting for correlations between the functional and scalar components. We propose a hybrid PLS regression framework that integrates functional and vector predictors in a principled and coherent manner. To extract predictive structure from these jointly observed and potentially correlated data types, we define a Hilbert space that treats the tuple of functional and vector components as a single hybrid object, equipped with a suitable inner product. The hybrid PLS direction is then obtained by iteratively maximizing the empirical covariance with the response, subject to a unit-norm constraint in this Hilbert space. Our framework is readily applicable to dense and irregular functional data and supports regularization techniques to prevent overfitting and reduce variance. We also provide the mathematical properties that justify our algorithm.

2. Background on partial least squares and its extension to hybrid predictors. (JM:) *comment: I replaced the formal presentation of the naive algorithm with an the scalar version and a brief outline of the pointwise hybrid extension strategy. The main reason for this*

revision is that the hybrid inner product is introduced specifically to address the extension challenge, so it is more natural to present it after listing the challenges. Previously, the naive algorithm section relied on the hybrid inner product for normalization. For intuition, let us return to the high-dimensional Euclidean predictor setting $Y_i = \beta^\top \mathbf{Z}_i + \epsilon_i$. A common way to address ill-posedness and correlation is to approximate the high-dimensional vector \mathbf{Z}_i using a low-dimensional vector $(\hat{\rho}_1^{[1]}, \dots, \hat{\rho}_1^{[L]})^\top \in \mathbb{R}^L$. To retain the regression relationship, the l -th PLS direction $\hat{\xi}_l$ solves:

$$\max_{\alpha} \widehat{\text{Cov}}^2(\{\langle \alpha, \mathbf{Z}_i \rangle, Y_i\}_{i=1}^n) \text{ s.t. } \|\alpha\|_2 = 1, \alpha^\top \widehat{\text{Cov}}^2(\{\mathbf{Z}_i\}_{i=1}^n) \hat{\xi}_j = 0, \quad j = 1, \dots, l-1,$$

where the two $\widehat{\text{Cov}}^2$ denote sample cross-covariance and sample covariance, respectively. A standard algorithm for solving this problem, called nonlinear iterative partial least squares (NIPALS) is presented in Algorithm 1.

Algorithm 1 Scalar partial least squares regression

- 1: Standardize each $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ so that each feature have mean zero and variance one. Standardize Y_1, \dots, Y_n .
 - 2: **for** $l = 1, 2, \dots, L$ **do**
 - 3: **PLS direction and score estimation:**
 - 4: $\hat{\xi}^{[l]} \leftarrow \arg \max_{\alpha} \widehat{\text{Cov}}^2(\{\langle \alpha, \mathbf{Z}_i^{[l]} \rangle, Y_i^{[l]}\}_{i=1}^n) \text{ s.t. } \|\alpha\|_2 = 1$ ▷ PLS direction
 - 5: $\hat{\rho}_i^{[l]} \leftarrow \langle \hat{\xi}^{[l]}, \mathbf{Z}_i^{[l]} \rangle, i = 1, \dots, n$ ▷ PLS score
 - 6: **Residualization:**
 - 7: $\nu^{[l]} \leftarrow \frac{\sum_{i=1}^n Y_i^{[l]} \hat{\rho}_i^{[l]}}{\sum_{i=1}^n \hat{\rho}_i^{[l]2}}$ ▷ Least squares estimate
 - 8: $Y_i^{[l+1]} \leftarrow Y_i^{[l]} - \nu^{[l]} \hat{\rho}_i^{[l]}, i = 1, \dots, n$
 - 9: $\hat{\delta}^{[l]} \leftarrow \frac{1}{\sum_{i=1}^n \hat{\rho}_i^{[l]2}} \sum_{i=1}^n \hat{\rho}_i^{[l]} \mathbf{Z}_i^{[l]}$ ▷ Least squares estimate
 - 10: $\mathbf{Z}_i^{[l+1]} \leftarrow \mathbf{Z}_i^{[l]} - \hat{\rho}_i^{[l]} \hat{\delta}^{[l]}$
 - 11: **Regression coefficient estimation:**
 - 12: $\hat{\mathbf{c}}^{[1]} \leftarrow \hat{\xi}^{[1]}$
 - 13: **for** $l = 2, \dots, L$ **do**
 - 13: $\hat{\mathbf{c}}^{[l]} \leftarrow \hat{\xi}^{[l]} - \sum_{u=1}^{l-1} \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l]} \rangle \hat{\mathbf{c}}^{[u]}$
 - 13: $\hat{\beta} \leftarrow \sum_{l=1}^L \hat{\rho}_i^{[l]} \hat{\mathbf{c}}^{[l]}$ ▷ regression coefficient estimate
 - 14: **Output:** the regression coefficient estimate
-

To extend Algorithm 1 to the hybrid regression model (1), we must compute hybrid objects in both the PLS direction estimation and predictor residualization steps (lines 4 and 15). These computations cannot be separated into functional and scalar parts, as their correlation must be addressed. A naive approach is to add loops over the observed evaluation points t to lines 4 and 15, treating

$$(X_{i1}(t), \dots, X_{iK}(t), Z_{i1}, \dots, Z_{ip}) \in \mathbb{R}^{K+P}$$

as a Euclidean predictor, applying steps 4 and 15, and aggregating the results. This pointwise method, however, is computationally prohibitive for densely observed functional predictors, infeasible for irregular data, and prone to variability across the domain, resulting in overfitting and unstable predictions. The main challenge of extension to the hybrid setting lies in the following points. First, independent variables consist of multiple highly structured images and scalar predictors. Second, our sample size is small compared to the dimension and number of functional and scalar predictors. Third, existing partial least squares (PLS) methods

can only accommodate (i) univariate or multivariate functional predictors without any scalar predictors (Preda and Saporta, 2005; Delaigle and Hall, 2012; Febrero-Bande, Galeano and González-Manteiga, 2017; Beyaztas and Shang, 2020); or (ii) a univariate functional predictor with other scalar predictors (Wang, 2018). The next section proposes a new extension framework that addresses these issues.

3. Proposed PLS Algorithm. Our new framework builds on the fact that all computations in Algorithm 1 are arithmetic operations in an inner product space. We therefore formally define the Hilbert space of hybrid random predictors along with its corresponding complete inner product, and leverage them to extend Algorithm 1. *(JM:) comment: I split the definition into three parts, adding completeness, separability, and the sigma-field.*

DEFINITION 3.1 (Hybrid space). Let \mathbb{H} be a product space defined as the Cartesian product of K copies of the space of square-integrable functions on $[0, 1]$ and the p -dimensional Euclidean space:

$$\mathbb{H} := (\mathbb{L}^2[0, 1])^K \times \mathbb{R}^p.$$

An element $h \in \mathbb{H}$ is an ordered tuple $h = (f_1, \dots, f_K, \mathbf{u})$, where $f_k \in \mathbb{L}^2[0, 1]$ for $k = 1, \dots, K$ and $\mathbf{u} \in \mathbb{R}^p$. This space is equipped with element-wise vector addition and scalar multiplication. We define an inner product on \mathbb{H} for any two elements $h_1 = (f_1, \dots, f_K, \mathbf{u})$ and $h_2 = (g_1, \dots, g_K, \mathbf{v})$ as:

$$(2) \quad \langle h_1, h_2 \rangle_{\mathbb{H}} := \sum_{k=1}^K \int_0^1 f_k(t)g_k(t) dt + \omega \mathbf{u}^\top \mathbf{v},$$

where ω is a positive constant ($\omega > 0$). The inner product induces a norm $\|\cdot\|_{\mathbb{H}}$ on the space, defined as $\|h\|_{\mathbb{H}} := \langle h, h \rangle_{\mathbb{H}}^{1/2}$, and a corresponding metric $d(h_1, h_2) = \|h_1 - h_2\|_{\mathbb{H}}$.

In (2), ω is a positive weight that needs to be pre-specified or estimated. It is mainly used to take into account heterogeneity between functional and scalar parts in terms of measurement scale and/or amount of variation (see Section 3.2). Without loss of generality and for the clarity of illustration, all the following theoretical results will be derived for $\omega = 1$. The results remain valid for any positive weights.

LEMMA 3.2 (Hybrid Hilbert space). *The hybrid space \mathbb{H} is a separable Hilbert space.*

Proof of Lemma 3.2 is provided in Appendix B.

DEFINITION 3.3 (Hybrid Predictor). For the Hilbert space \mathbb{H} defined in Definition 3.1, The Borel σ -field on \mathbb{H} , denoted $\mathfrak{B}(\mathbb{H})$, is the smallest σ -field containing the class \mathfrak{M} of all sets of the form $\{q \in \mathbb{H} \mid \langle q, h \rangle \in O\}$, for any $h \in \mathbb{H}$ and any open subset $O \subseteq \mathbb{R}$ (details can be found in Theorem 7.1.1 of Hsing and Eubank 2015). A hybrid predictor $W_i = (X_{i1}(t), \dots, X_{iK}(t), \mathbf{Z}_i)$ is a measurable mapping from a probability space $(\Omega, \mathfrak{F}, P)$ into $(\mathbb{H}, \mathfrak{B}(\mathbb{H}))$.

Then the joint regression model (1) can be concisely written as

$$(3) \quad Y_i = \langle \beta, W_i \rangle_{\mathbb{H}} + \epsilon_i, \text{ where } \beta := (\beta_1(t), \dots, \beta_K(t), \boldsymbol{\beta}) \in \mathbb{H}.$$

REMARK 1. Our method is applicable to settings where each functional predictor belongs to a different Hilbert space, possibly defined over a distinct compact domain in \mathbb{R}^d , for arbitrary d , and observed at different time points. However, for notational simplicity, our discussion assumes a common Hilbert space over domain $[0, 1]$ for all functional predictors.

Our approach provides an efficient and robust means of producing PLS components and scores in the presence of multiple dense and/or irregular functional predictors and scalar predictors. It also incorporates a regularization scheme that enables the algorithm to borrow strength and exploit structural relationships within and between the functions to avoid overfitting of the PLS components and to improve the generalizability and interpretability of the predictive model. Each iteration of our approach consists of two subroutines: regularized estimation of smoothed PLS components and orthogonalization, detailed in Sections 3.3.1 and 3.3.2, respectively. After a suitable number of iterations, the hybrid regression coefficient is estimated, as described in Section 3.4. For notational simplicity, we omit the iteration index l in the following discussion, with the understanding that the subroutines apply to any iteration. The complete algorithm is summarized in Algorithm 2.

3.1. *Preliminary step 1: finite-basis approximation.* Let $b_m(t)$ be a twice-differentiable basis of $\mathbb{L}^2([0, 1])$ whose second derivatives are also linearly independent, for example, cubic B-splines, the Fourier basis, or an orthonormal polynomial basis of degree greater than three. Using this basis, the j th functional predictor, regression coefficient, PLS component direction, and orthogonalization regression coefficient (with iteration indices suppressed) are represented as follows:

$$X_{ij}(t) = \sum_{m=1}^{\infty} \theta_{ijm} b_m(t), \quad \beta_j(t) = \sum_{m=1}^{\infty} \eta_{jm} b_m(t), \quad \xi_j(t) = \sum_{m=1}^{\infty} \gamma_{jm} b_m(t), \quad \delta_j(t) = \sum_{m=1}^{\infty} \pi_{jm} b_m(t)$$

In practice, the full set of coefficients can not be obtained with finite sample size, as functional data are measured on a finite grid. Thus, we truncate the expansion at M terms. We choose a moderately large M (e.g., 15 or 20) to capture functional variation without fine-tuning, as smoothness is handled via penalization (see Section 3.3.1). The truncated expansions of the predictor and coefficient are denoted as

$$\tilde{X}_{ij}(t) := \sum_{m=1}^M \theta_{ijm} b_m(t), \quad \tilde{\beta}_j(t) := \sum_{m=1}^M \eta_{jm} b_m(t),$$

and our suggest method restricts each PLS component direction and orthogonalization regression coefficient to admit the following expansion:

$$(4) \quad \xi_j(t) = \sum_{m=1}^M \gamma_{jm} b_m(t), \quad \delta_j(t) = \sum_{m=1}^M \pi_{jm} b_m(t)$$

This implies that all computations in this paper are carried out entirely within the subspace

$$(5) \quad \tilde{\mathbb{H}} := \text{span}(b_1(t), \dots, b_M(t))^K \times \mathbb{R}^p \subset \mathbb{H}.$$

The i th hybrid predictor, projected on $\tilde{\mathbb{H}}$, is represented by the tuple

$$\tilde{W}_i := (\tilde{X}_{i1}, \dots, \tilde{X}_{iK}, \mathbf{Z}_i).$$

Let $\boldsymbol{\theta}_{ij}$, $\boldsymbol{\eta}_j$, $\boldsymbol{\gamma}_j$, and $\boldsymbol{\pi}_j$ denote the M -dimensional vectors of coefficients:

$$(6) \quad \boldsymbol{\theta}_{ij} := (\theta_{ij1}, \dots, \theta_{ijM})^\top, \quad \boldsymbol{\eta}_j := (\eta_{j1}, \dots, \eta_{jM})^\top, \quad \boldsymbol{\gamma}_j := (\gamma_{j1}, \dots, \gamma_{jM})^\top, \quad \boldsymbol{\pi}_j := (\pi_{j1}, \dots, \pi_{jM})^\top$$

For the predictors, we stack the coefficient vectors across observations into the matrix

$$(7) \quad \Theta_j := (\boldsymbol{\theta}_{1j}, \dots, \boldsymbol{\theta}_{nj})^\top \in \mathbb{R}^{n \times M},$$

and construct the full design matrix

$$(8) \quad \Theta := (\Theta_1, \dots, \Theta_K, \mathbf{Z}) \in \mathbb{R}^{n \times (MK+p)}.$$

Let us denote the response vector as $\mathbf{y} := (y_1, \dots, y_n)^\top$. Let $B, B'' \in \mathbb{R}^{M \times M}$ denote the Gram matrices of the basis functions and their second derivatives, with entries

$$(9) \quad B_{m,m'} := \int_0^1 b_m(t) b_{m'}(t) dt, \quad B''_{m,m'} := \int_0^1 b''_m(t) b''_{m'}(t) dt,$$

for $m, m' = 1, \dots, M$. We then define the block-diagonal matrices

$$(10) \quad \mathbb{B} := \text{blkdiag}(B, \dots, B, I_p), \quad \mathbb{B}'' := \text{blkdiag}(B'', \dots, B'', I_p),$$

Then the full data for the hybrid PLS problem at the l -th iteration can be represented by the tuple

$$(11) \quad (\mathbb{B}, \mathbb{B}'', \Theta, \mathbf{y}) \in \mathbb{R}^{(MK+p) \times (MK+p)} \times \mathbb{R}^{(MK+p) \times (MK+p)} \times \mathbb{R}^{n \times (MK+p)} \times \mathbb{R}^n,$$

with the index l omitted for brevity.

REMARK 2. While different bases could be used for each functional predictor, we adopt a common basis for simplicity. The definitions of \mathbb{B} and \mathbb{B}'' remain general enough to accommodate distinct bases if needed.

3.2. Preliminary step 2 : data preprocessing. Functional and scalar elements of the hybrid predictors often have incompatible units and/or exhibit different amounts of variation. This can be problematic for our PLS framework which is not scale invariant as: i) each predictor has different chance of contributing to the predictor/response structure; and ii) a predictor with high correlation to Y but relatively low variance may be overlooked.

To obtain PLS components that have a meaningful interpretation, we standardize the predictor data via the following steps. The first step is to account for discrepancies *within* respective functional and scalar parts, if needed. If the functional parts $\tilde{X}_{i1}(t), \dots, \tilde{X}_{iK}(t)$ are measured in different units or have quite different domains, one can standardize them to have mean zero and integrated variance of one. If multivariate scalar predictors $\mathbf{Z}_i = (Z_{i1}, \dots, Z_{ip})^\top$ exhibit different amounts of variation, one can standardize them to have mean zero and unit variance. The second step is to eliminate the discrepancies *between* functional and scalar parts. To accomplish this aim, we choose an appropriate weight ω in the hybrid inner product (2) that ensures functional and vector parts have comparable variance. A sensible data-driven approach to choosing an appropriate weight is to set

$$\omega = \frac{\sum_{i=1}^n \sum_{k=1}^K \|\tilde{X}_i\|_{\mathbb{L}^2[0,1]}^2}{\sum_{i=1}^n \|\mathbf{Z}_i\|_2^2},$$

In practice, instead of using inner product weighted by $\omega^{1/2}$, one can implement this weighting scheme by formulating the hybrid object as $\tilde{W}_i = (\tilde{X}_{i1}(t), \dots, \tilde{X}_{iK}(t), \omega^{1/2} \mathbf{Z}_i)$, whose vector part has been scaled by a factor of $\omega^{1/2}$.

3.3. Iterative steps. The iterative process presented here yields an orthonormal hybrid basis that effectively captures the predictor-response relationships. It proceeds through two intermediate steps: the estimation of the PLS component direction (Section 3.3.1) and residualization (Section 3.3.2). The proofs are deferred to Appendix C. The properties of the resulting estimates are introduced in Section 4.

3.3.1. Iterative step 1: regularized estimation of PLS component direction. We begin by formally introducing the core optimization problem pertinent to the PLS direction estimation, which is formulated as a generalized Rayleigh quotient (Proposition 3.4). Building upon this foundational concept, we present our regularized PLS component direction estimation step that promotes smoothness (Proposition 3.5). Furthermore, we detail an efficient computational scheme (Proposition 3.6).

Core optimization problem. We present the core optimization problem that directly estimates the hybrid PLS component direction. It fully leverages the continuous nature of the functional components and the function-scalar hybrid structure. We describe the strategy at the l -th iteration. The PLS component direction is estimated by the unit-norm direction $\xi^{[l]} \in \tilde{\mathbb{H}}$ that maximizes the squared empirical covariance, which quantifies the linear dependence between the PLS scores $\langle \widetilde{W}_1, \xi \rangle_{\mathbb{H}}, \dots, \langle \widetilde{W}_n, \xi \rangle_{\mathbb{H}}$ and the responses y_1, \dots, y_n , defined as

$$(12) \quad \widehat{\text{Cov}}(\langle \widetilde{W}, \xi \rangle_{\mathbb{H}}, Y) := \frac{1}{n} \sum_{i=1}^n y_i \langle \widetilde{W}_i, \xi \rangle_{\mathbb{H}}.$$

We denote the estimated PLS component direction as

$$(13) \quad \hat{\xi} := \arg \max_{\xi \in \tilde{\mathbb{H}}} \widehat{\text{Cov}}^2(\langle \widetilde{W}, \xi \rangle_{\mathbb{H}}, Y) \text{ s.t. } \|\xi\|_{\mathbb{H}} = 1.$$

Here, $\hat{\xi} \in \tilde{\mathbb{H}}$ is an ordered pair expanded as:

$$\hat{\xi} = (\hat{\xi}_1(t), \dots, \hat{\xi}_K(t), \hat{\zeta}) = \left(\sum_{m=1}^M \hat{\gamma}_{1m} b_m(t), \dots, \sum_{m=1}^M \hat{\gamma}_{Km} b_m(t), \hat{\zeta} \right),$$

where $\hat{\zeta}$ is the scalar part. Obtaining these coefficients is equivalent to solving the maximization problem (13). The following proposition formulates this coefficients obtaining procedure as a generalized Rayleigh quotient:

PROPOSITION 3.4. *Let $(\mathbb{B}, \mathbb{B}'', \Theta, \mathbf{y})$ denote the observed data defined in (11). At the l -th iteration of the PLS algorithm, the coefficients of the squared covariance maximizer defined in (13), is obtained as*

$$(14) \quad \left(\hat{\gamma}_{11}, \dots, \hat{\gamma}_{1M}, \dots, \hat{\gamma}_{K1}, \dots, \hat{\gamma}_{KM}, \hat{\zeta}^\top \right)^\top = \arg \max_{\xi \in \mathbb{R}^{MK+p}} \xi^\top V \xi \quad \text{subject to} \quad \xi^\top \mathbb{B} \xi = 1.$$

where

$$(15) \quad V := \frac{1}{n^2} (\mathbb{B} \Theta^\top \mathbf{y}) (\mathbb{B} \Theta^\top \mathbf{y})^\top \in \mathbb{R}^{(MK+p) \times (MK+p)}.$$

The proof of Proposition 3.4 is provided in Appendix C.1.

Proposed regularized estimation procedure. Although Proposition 3.4 offers an efficient way to estimate the PLS component direction $\hat{\xi}$, its functional components, $\hat{\xi}_1, \dots, \hat{\xi}_K$, may not be smooth. (JM:) added: Since the final regression coefficient $\hat{\beta}$ is constructed as a linear combination of these directions (Section 3.4), this can lead to a non-smooth $\hat{\beta}$ which complicates interpretation and can lead to overfitting and unstable predictions. To address

this, we propose a regularized extension that balances predictive performance with smoothness. Specifically, we penalize the roughness of each ξ_j using its integrated squared second derivative

$$(16) \quad \text{PEN}(\xi_j) := \int_0^1 \{\xi_j''(t)\}^2 dt.$$

Instead of solely maximizing the squared empirical covariance $\widehat{\text{Cov}}^2(\langle \widetilde{W}, \xi \rangle_{\mathbb{H}}, Y)$, we incorporate this roughness penalty to simultaneously control the complexity of the estimated functional components. One possible approach is to extend the smoothed functional PCA framework of [Rice and Silverman \(1991\)](#) by modifying the objective in (13) to

$$\widehat{\text{Cov}}^2(\langle \widetilde{W}, \xi \rangle_{\mathbb{H}}, Y) - \sum_{j=1}^K \lambda_j \text{PEN}(\xi_j),$$

where the smoothing parameters $\{\lambda_j\}_{j=1}^K$ control the trade-off between maximizing covariance and penalizing roughness. However, this approach [Rice and Silverman \(1991\)](#) assumes that the functional predictors admit an orthogonal expansion in the \mathbb{L}^2 sense.

To avoid the orthogonal basis assumption of [Rice and Silverman \(1991\)](#), we adopt the strategy of [Silverman \(1996\)](#), which replaces the standard orthonormality constraint with a weaker one based on a modified inner product that incorporates roughness. Accordingly, our estimation procedure at the l -th iteration, iteration index omitted and assuming the observations have been residualized in previous steps, solves the following optimization problem:

$$(17) \quad \hat{\xi} := \arg \max_{\xi \in \mathbb{H}} \widehat{\text{Cov}}^2(\langle \widetilde{W}, \xi \rangle_{\mathbb{H}}, Y) \text{ s.t. } \|\xi\|_{\mathbb{H}} + \sum_{j=1}^K \lambda_j \text{PEN}(\xi_K) = 1.$$

Here, $\hat{\xi} \in \widetilde{\mathbb{H}}$ is an ordered pair expanded as:

$$\hat{\xi} = (\hat{\xi}_1(t), \dots, \hat{\xi}_K(t), \hat{\zeta}) = \left(\sum_{m=1}^M \hat{\gamma}_{1m} b_m(t), \dots, \sum_{m=1}^M \hat{\gamma}_{Km} b_m(t), \hat{\zeta} \right),$$

where $\hat{\zeta}$ is the scalar part. This formulation maximizes the squared covariance over a class of smooth functions. Obtaining these coefficients is equivalent to solving the maximization problem (17). The following proposition formulates this coefficients obtaining procedure as a generalized Rayleigh quotient:

PROPOSITION 3.5 (Regularized estimation of PLS component direction). *Let $(\mathbb{B}, \mathbb{B}'', \Theta, \mathbf{y})$ denote the data given at the l -th iteration, as defined in (11). Recall from (15) that $V = n^{-2}(\mathbb{B}\Theta^\top \mathbf{y})(\mathbb{B}\Theta^\top \mathbf{y})^\top$. Let $\Lambda \in \mathbb{R}^{(MK+p) \times (MK+p)}$ be defined as:*

$$(18) \quad \Lambda := \text{blkdiag}(\lambda_1 I_M, \dots, \lambda_K I_M, 0_{p \times p}), \text{ where } \lambda_1, \dots, \lambda_K \geq 0.$$

Here, $0_{p \times p}$ denotes the $p \times p$ zero matrix. The coefficients of the squared covariance maximizer defined in (17), are obtained as

$$(19) \quad \left(\hat{\gamma}_{11}, \dots, \hat{\gamma}_{1M}, \dots, \hat{\gamma}_{K1}, \dots, \hat{\gamma}_{KM}, \hat{\zeta}^\top \right)^\top = \arg \max_{\xi \in \mathbb{R}^{MK+p}} \xi^\top V \xi \text{ s.t. } \xi^\top (\mathbb{B} + \Lambda \mathbb{B}'') \xi = 1.$$

The proof of Proposition 3.5 is provided in Appendix C.2. The constraint $\xi^\top (\mathbb{B} + \Lambda \mathbb{B}'') \xi = 1$ enforces the orthonormality of the estimated PLS component directions with respect to a modified inner product (see Section 4.2 for details). The smoothing parameter λ_k balances

goodness of fit and smoothness in $\hat{\xi}_j$. Smaller λ_k yields components that better fit the data but risks overfitting; setting $\lambda_k = 0$ recovers the unregularized solution in Proposition 3.4. Larger λ_k enforces greater smoothness, and in the limit $\lambda_k \rightarrow \infty$, $\hat{\xi}_j(t)$ approaches a linear form $a + bt$. In practice, both $\{\lambda_k\}$ and the number of components L can be selected via cross-validation using a predictive criterion such as mean squared error.

Computation. The generalized eigenproblem presented in Proposition 3.5 may be computationally unstable in practice. However, by leveraging the rank-one structure of the matrix V Proposition 3.6 derives a closed-form solution that requires only the solution of linear systems.

PROPOSITION 3.6 (Closed-form solution). *Consider the optimization problem described in Proposition 3.5. Define the following quantities, which depend on the observed data but are not decision variables:*

$$\mathbf{u}_j := B\Theta_j^\top \mathbf{y} \in \mathbb{R}^M \quad \text{for } j = 1, \dots, K, \quad \text{and} \quad \mathbf{v} := \mathbf{Z}^\top \mathbf{y} \in \mathbb{R}^p.$$

Let

$$q := \sum_{j=1}^K \mathbf{u}_j^\top (B + \lambda_j B'')^{-1} \mathbf{u}_j + \mathbf{v}^\top \mathbf{v}.$$

Then the unique (up to sign) solution to the regularized maximization problem is given in closed form by

$$\hat{\gamma}_j = \frac{1}{\sqrt{q}} (B + \lambda_j B'')^{-1} \mathbf{u}_j \quad \text{for } j = 1, \dots, K, \quad \text{and} \quad \hat{\zeta} = \frac{1}{\sqrt{q}} \mathbf{v}.$$

The proof of Proposition 3.6 is provided in Appendix C.3. The expressions above involve solving linear systems for the functional and scalar components separately, followed by normalization by a common factor. Although the unnormalized coefficients are obtained independently, the normalization step couples the functional and scalar parts, allowing them to influence one another. This coupling enables the procedure to capture the correlation between the functional and scalar components of the PLS direction.

3.3.2. Iterative step 2: residualization via hybrid-on-scalar regression. The l -th iteration's second step involves residualization of both predictors and responses. We first compute the individual PLS score:

$$(20) \quad \hat{\rho}_i^{[l]} := \langle \widetilde{W}_i^{[l]}, \hat{\xi}^{[l]} \rangle_{\mathbb{H}},$$

using the estimated PLS component direction $\hat{\xi}^{[l]}$ obtained from Propositions 3.5 and 3.6. Since $\widetilde{W}_i^{[l]}$ are assumed to have a sample mean of zero, these PLS scores will also have a sample mean of zero. To obtain the $(l+1)$ -th iteration's responses and hybrid predictors, we regress the (l) -th iteration's responses and hybrid predictors on these PLS scores by least squares and then residualize. Specifically, the $(l+1)$ -th predictor is computed as a residual of hybrid-on-scalar linear regression model:

$$\widetilde{W}_i^{[l]} = \hat{\rho}_i^{[l]} \hat{\delta}^{[l]} + \epsilon_i,$$

where $\hat{\delta}^{[l]} \in \widetilde{\mathbb{H}}$ is the regression coefficient. In the same spirit as the PLS component direction estimation step, rather than treating the hybrid object as a long vector of concatenated function evaluations at time points and scalar vectors, we employ a basis expansion approach to fit

the entire hybrid object in one step. Therefore, our method is computationally efficient, and applicable for dense or irregular functional data. Consequently, $\delta^{[l]}$ is obtained by minimizing a least squares criterion: *(JM:) comment: I removed the penalization because it makes proving the orthonormality of the PLS components in Proposition 4.7 impossible. Smoothness of β seems to solely rely on the smoothness of ξ . See Section 3.4.*

$$(21) \quad \hat{\delta}^{[l]} := \arg \min_{\delta \in \mathbb{H}} \sum_{i=1}^n \|\widetilde{W}_i^{[l]} - \hat{\rho}_i^{[l]} \delta\|_{\mathbb{H}}^2.$$

On the other hand, the $(l+1)$ -th response is computed as a residual of a scalar-on-scalar linear regression model:

$$Y_i^{[l]} = \hat{\nu}^{[l]} \hat{\rho}_i^{[l]} + \epsilon_i.$$

The following proposition demonstrates that this residualization step can be performed simply, analogous to scalar PLS.

LEMMA 3.7 (Closed-form solution). *(JM:) Previously stated as "we can show". Let us denote $\hat{\rho}^{[l]} := (\hat{\rho}_1^{[l]}, \dots, \hat{\rho}_n^{[l]})^\top$. The $(l+1)$ -th iteration's predictors and responses are computed as*

$$\widetilde{W}_i^{[l+1]} := \widetilde{W}_i^{[l]} - \hat{\rho}_i^{[l]} \hat{\delta}^{[l]}, \text{ where } \delta^{[l]} := \frac{1}{\|\hat{\rho}^{[l]}\|_2^2} \sum_{i=1}^n \hat{\rho}_i^{[l]} \widetilde{W}_i^{[l]},$$

and

$$(22) \quad Y_i^{[l+1]} = Y_i^{[l]} - \hat{\nu}^{[l]} \hat{\rho}_i^{[l]}, \text{ where } \hat{\nu}^{[l]} := \frac{\mathbf{y}^{[l]\top} \hat{\rho}^{[l]}}{\|\hat{\rho}^{[l]}\|_2^2}.$$

Proof of Lemma 3.7 is provided in Appendix C.4. *(JM:) proof needs to be revised* Since $\widetilde{W}_i^{[l]}$ and $Y_i^{[l]}$ are assumed to have a sample mean of zero, their respective residuals, $\widetilde{W}_i^{[l+1]}$ and $Y_i^{[l+1]}$, also maintain a zero sample mean.

3.4. *Final step: estimating the hybrid regression coefficient.* The hybrid regression coefficient β in model (3) can be written as a linear combination of PLS directions:

LEMMA 3.8. *Let us define $\hat{\iota}^{[1]} := \hat{\xi}^{[1]}$. For $l \geq 2$, we recursively define:*

$$\hat{\iota}^{[l]} = \hat{\xi}^{[l]} - \sum_{u=1}^{l-1} \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l]} \rangle_{\mathbb{H}} \hat{\iota}^{[u]}.$$

Then we have:

$$\hat{\rho}_i^{[l]} = \langle W_i^{[l]}, \hat{\xi}^{[l]} \rangle_{\mathbb{H}} = \langle W_i, \hat{\iota}^{[l]} \rangle_{\mathbb{H}}.$$

Proof of Lemma 3.8 is provided in Appendix D.

Next, (22) leads to the following model:

$$Y_i = \sum_{l=1}^L \hat{\nu}^{[l]} \hat{\rho}_i^{[l]} + \epsilon_i.$$

This model lets us to express Y_i as:

$$Y_i = \sum_{l=1}^L \hat{\nu}^{[l]} \langle W_i^{[l]}, \hat{\xi}^{[l]} \rangle_{\mathbb{H}} + \epsilon_i = \langle W_i, \sum_{l=1}^L \hat{\nu}^{[l]} \hat{\iota}^{[l]} \rangle_{\mathbb{H}} + \epsilon_i,$$

which, given the uniqueness of β , leads to

$$\hat{\beta} = \sum_{l=1}^L \hat{\nu}^{[l]} \hat{\tau}^{[l]} = \textcolor{teal}{(JM:)} \text{ added: } \sum_{l=1}^L \left(\hat{\nu}^{[l]} - \sum_{k=l+1}^L \hat{\nu}^{[k]} \langle \hat{\delta}^{[l]}, \hat{\xi}^{[k]} \rangle_{\mathbb{H}} \right) \hat{\xi}^{[l]}.$$

(JM:) added: Note that if L is too large, $\hat{\beta}$ can be wiggly even if $\hat{\xi}^{[l]}$'s are smooth. Thus the number of PLS components L to be estimated can be chosen by cross-validation.

Algorithm 2 Hybrid partial least squares regression

- 1: **Initialize:** $(\mathbb{B}, \mathbb{B}'', \Theta^{[1]}, \mathbf{y}^{[1]})$ as the data objects after basis expansion, following Section 3.1.
 - 2: $\widetilde{W}_1^{[1]}, \dots, \widetilde{W}_n^{[1]}, Y_1^{[1]}, \dots, Y_n^{[1]} \leftarrow$ standardized versions of $W_1, \dots, W_n, Y_1, \dots, Y_n$, following Section 3.2
 - 3: **for** $l = 1, 2, \dots, L$ **do**
 - 4: **PLS direction and score estimation (Proposition 3.6):**
 - 5: $\mathbf{u}_j^{[l]} \leftarrow B \Theta_j^{[l]\top} \mathbf{y}^{[l]}, j = 1, \dots, K$
 - 6: $\mathbf{v}^{[l]} \leftarrow \mathbf{Z}^{[l]\top} \mathbf{y}^{[l]}$
 - 7: $q^{[l]} \leftarrow \sum_{j=1}^K \mathbf{u}_j^{[l]\top} (B + \lambda_j B'')^{-1} \mathbf{u}_j^{[l]} + \mathbf{v}^{[l]\top} \mathbf{v}^{[l]}$
 - 8: $(\hat{\gamma}_{j1}^{[l]}, \dots, \hat{\gamma}_{jM}^{[l]})^\top \leftarrow \frac{1}{\sqrt{q}} (B + \lambda_j B'')^{-1} \mathbf{u}_j^{[l]}, j = 1, \dots, K$
 - 9: $\hat{\zeta}^{[l]} \leftarrow \frac{1}{\sqrt{q}} \mathbf{v}^{[l]}$
 - 10: $\hat{\xi}^{[l]} \leftarrow \left(\sum_{m=1}^M \hat{\gamma}_{1m}^{[l]} b_m(t), \dots, \sum_{m=1}^M \hat{\gamma}_{Km}^{[l]} b_m(t), \hat{\zeta}^{[l]} \right)$ ▷ PLS direction
 - 11: $\hat{\rho}_i^{[l]} \leftarrow \langle \hat{\xi}^{[l]}, \widetilde{W}_i^{[l]} \rangle, i = 1, \dots, n$ ▷ PLS score
 - 12: **Residualization (Proposition 3.7):**
 - 13: $\nu^{[l]} \leftarrow \frac{\sum_{i=1}^n Y_i^{[l]} \hat{\rho}_i^{[l]}}{\sum_{i=1}^n \hat{\rho}_i^{[l]2}}$ ▷ Least squares estimate
 - 14: $Y_i^{[l+1]} \leftarrow Y_i^{[l]} - \nu^{[l]} \hat{\rho}_i^{[l]}, i = 1, \dots, n$
 - 15: $\hat{\delta}^{[l]} \leftarrow \frac{1}{\sum_{i=1}^n \hat{\rho}_i^{[l]2}} \sum_{i=1}^n \hat{\rho}_i^{[l]} \widetilde{W}_i^{[l]}$ ▷ Least squares estimate
 - 16: $\widetilde{W}_i^{[l+1]} \leftarrow \widetilde{W}_i^{[l]} - \hat{\rho}_i^{[l]} \hat{\delta}^{[l]}$
 - 17: **Regression coefficient estimation (Section 3.4):**
 - 17: $\hat{\tau}^{[1]} \leftarrow \hat{\xi}^{[1]}$
 - 18: **for** $l = 2, \dots, L$ **do**
 - 19: $\hat{\tau}^{[l]} \leftarrow \hat{\xi}^{[l]} - \sum_{u=1}^{l-1} \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l]} \rangle \hat{\tau}^{[u]}$
 - 20: $\hat{\beta} \leftarrow \sum_{l=1}^L \hat{\nu}^{[l]} \hat{\tau}^{[l]}$ ▷ regression coefficient estimate
 - 20: **Output:** the regression coefficient estimate $\hat{\beta}$
-

4. Properties of the hybrid PLS. This section provides the mathematical properties that support the algorithm suggested in Section 3. Section 4.1 shows that the core optimization problem for the Partial Least Squares (PLS) direction estimation step, presented in Proposition 3.4, is well-defined under mild conditions. Section 4.2 demonstrates that our algorithm preserves the core properties of PLS, namely the orthonormality of the derived directions and the orthogonality of the scores.

4.1. *Tucker's Criterion* . The fundamental property of classical scalar PLS is that the PLS component direction corresponds to the first eigenvector of the squared covariance matrix, satisfying the Tucker's Criterion (Tucker, 1938). Under a mild assumption, We derive an analogous result in our scalar-on-hybrid regression model setting (3). This justifies our

estimation procedure introduced in Section 3. We omit l in the notations and first define the cross-covariance term between the response and the functional predictors (observed versions if $l = 1$ and residualized versions if $l \geq 2$):

$$(23) \quad \sigma_{YX} = (\sigma_{YX,1}(t), \dots, \sigma_{YX,K}(t)) := (\mathbb{E}[Y_1 X_{11}], \dots, \mathbb{E}[Y_1 X_{1K}]) \in \mathcal{F}^{\otimes K},$$

and between the response and the scalar predictor:

$$\sigma_{YZ} = (\sigma_{YZ,1}, \dots, \sigma_{YZ,p})^\top := (\mathbb{E}[Y_1 Z_{11}], \dots, \mathbb{E}[Y_1 Z_{1p}])^\top \in \mathbb{R}^p.$$

Based on these definitions, we define the cross-covariance term between the response and the hybrid predictor as

$$\Sigma_{YW} := \mathbb{E}[Y_1 W_1] = (\sigma_{YX,1}, \dots, \sigma_{YX,K}, \sigma_{YZ}) \in \mathbb{H}.$$

Based on these definitions, we introduce two cross-covariance operators and their properties.

LEMMA 4.1. *Define the operator $\mathcal{C}_{YW} = \mathbb{E}(W_1 \otimes_{\mathbb{H}} Y_1) : \mathbb{H} \rightarrow \mathbb{R}$ such that for any $h = (f_1, \dots, f_K, \mathbf{v}) \in \mathbb{H}$, it maps h to a real number as:*

$$\mathcal{C}_{YW}h := \mathbb{E}[\langle W_1, h \rangle_{\mathbb{H}} Y_1] = \langle \Sigma_{YW}, h \rangle_{\mathbb{H}} = \sum_{k=1}^K \int_0^1 \sigma_{YX,k}(t) f_k(t) dt + \sigma_{YZ}^\top \mathbf{v}.$$

(JM:) *instead of showing \mathcal{U} is compact, let's show \mathcal{C}_{YW} is compact. This way, we don't need uniform continuity of σ_{YX} . If there exist finite constants Q_1 and Q_2 such that*

$$(24) \quad \max_{k=1, \dots, K} \sup_{t \in [0,1]} \sigma_{YX,k}^2(t) < Q_1 \quad \text{and} \quad \max_{r=1, \dots, p} \sigma_{YZ,r}^2 < Q_2,$$

the operator \mathcal{C}_{YW} is a compact operator.

Proof of Lemma 4.1 is provided in Appendix E.1.

LEMMA 4.2. *Define $\mathcal{C}_{WY} = \mathbb{E}[Y \otimes W] : \mathbb{R} \rightarrow \mathbb{H}$, which maps $d \in \mathbb{R}$ to a hybrid object in \mathbb{H} as follows:*

$$\mathcal{C}_{WY}d := \mathbb{E}[\langle Y_1, d \rangle W_1] = \mathbb{E}[Y_1 W_1 d] = d \Sigma_{YW}.$$

Then \mathcal{C}_{WY} is an adjoint operator of \mathcal{C}_{YW} . That is, $\mathcal{C}_{WY} = \mathcal{C}_{YW}^$.*

The proof of Lemma 4.2 is provided in Appendix E.2.

Based on these two operators, we define a positive operator as follows:

LEMMA 4.3 (Composite cross-covariance operator). *Define $\mathcal{U} := \mathcal{C}_{WY} \circ \mathcal{C}_{YW} : \mathbb{H} \rightarrow \mathbb{H}$ as an operator which performs the following mapping:*

$$\mathcal{U}h = \mathcal{C}_{WY}(\mathcal{C}_{YW}h) = \mathcal{C}_{YW}(\langle \Sigma_{YW}, h \rangle_{\mathbb{H}}) = \Sigma_{YW} \langle \Sigma_{YW}, h \rangle_{\mathbb{H}} = (\Sigma_{YW} \otimes \Sigma_{YW})h.$$

In other words, $\mathcal{U} = \Sigma_{YW} \otimes \Sigma_{YW}$. Then \mathcal{U} is a self-adjoint and positive-semidefinite operator. Under the conditions of Lemma 4.1, \mathcal{U} is a compact operator.

The proof of Lemma 4.3 is provided in Appendix E.3. By the Hilbert-Schmidt theorem (e.g., Theorem 4.2.4 in Hsing and Eubank, 2015), Lemma 4.3 guarantees the existence of a complete orthonormal system of eigenfunctions $\{\xi_{(u)}\}_{u \in \mathbb{N}}$ of \mathcal{U} in \mathbb{H} such that $\mathcal{U}\xi_{(u)} = \kappa_{(u)}\xi_{(u)}$, where $\{\kappa_{(u)}\}_{u \in \mathbb{N}}$ are the corresponding sequence of eigenvalues that goes to zero as $u \rightarrow \infty$, that is, $\kappa_{(1)} \geq \kappa_{(2)} \geq \dots \geq 0$.

The following theorem introduces the population-level PLS and Tucker's Criterion, adapted to our scalar-on-hybrid regression model setting, based on the aforementioned operators defined in the hybrid space.

THEOREM 4.4 (Tucker's criterion). *Under the conditions of Lemma 4.1, the constrained maximum*

$$\max_{\substack{\xi \in \mathbb{H} \\ \|\xi\|_{\mathbb{H}}=1}} \text{Cov}^2(\langle W_1, \xi \rangle_{\mathbb{H}}, Y_1)$$

is attained by the eigenfunction associated with the largest eigenvalue of the operator \mathcal{U} .

The proof of Theorem 4.4 is provided in Appendix E.4. In other words, l -th PLS component ξ_l can be obtained as the *first* eigenfunction of $\mathcal{U}^{[l]} = \Sigma_{YW}^{[l]} \otimes \Sigma_{YW}^{[l]}$ whose components are formulated using the l -th residualized response and predictors: $\Sigma_{YW}^{[l]} = E(Y^{[l]} W^{[l]})$.

THEOREM 4.5 (L^2 Convergence of Hybrid PLS). *Let Y_i be a scalar response and W_i be a hybrid predictor in the hybrid Hilbert space $\mathbb{H} = (\mathbb{L}^2[0, 1])^K \times \mathbb{R}^p$. Let $\hat{Y}_i = \langle \beta, W_i \rangle_{\mathbb{H}}$ denote the ordinary linear regression solution, where $\beta \in \mathbb{H}$ is the true regression coefficient. Let $\hat{Y}_{i,L} = \sum_{l=1}^L \hat{c}_i^{[l]} \rho_i^{[l]}$ be the Hybrid PLS approximation with L components, where $\rho_i^{[l]} = \langle W_i, \xi^{[l]} \rangle_{\mathbb{H}}$ is the l -th PLS score and $\xi^{[l]} \in \mathbb{H}$ is the l -th PLS direction. Then, the Hybrid PLS approximation converges to the ordinary linear regression solution in the mean-squared (L^2) sense:*

$$\lim_{L \rightarrow \infty} E \left[\|\hat{Y}_{i,L} - \hat{Y}_i\|_{L^2(\Omega)}^2 \right] = 0$$

PROOF. The proof proceeds in three main steps: first, establishing that the PLS directions form a complete orthonormal system (CONS) for the relevant predictor subspace; second, expressing the ordinary linear regression solution as an infinite series in this PLS basis; and third, showing that the mean-squared error of the PLS approximation vanishes as the number of components increases.

1. The PLS Directions Form a Complete Orthonormal System (CONS). The Hybrid PLS algorithm iteratively constructs a set of directions $\{\xi^{[l]}\}_{l=1}^{\infty}$ from the predictors. The l -th PLS direction, $\xi^{[l]}$, is obtained as the principal eigenfunction of the operator $\mathcal{U}^{[l]} = \Sigma_{YW}^{[l]} \otimes \Sigma_{YW}^{[l]}$, where the superscript $[l]$ denotes that the operator is constructed from the residualized response $Y^{[l]}$ and predictors $W^{[l]}$ from the previous step.

4.1.0.1. 1.1 Orthonormality through Residualization. The iterative nature of the PLS algorithm guarantees the orthonormality of the directions. The **residualization (or orthogonalization)** step is crucial. After the $(l-1)$ -th PLS direction $\xi^{[l-1]}$ is determined, the predictors and response are adjusted to remove the linear effect of this component. Specifically, for sample i :

$$W_i^{[l]} = W_i^{[l-1]} - \langle W_i^{[l-1]}, \xi^{[l-1]} \rangle_{\mathbb{H}} \xi^{[l-1]}$$

$$Y_i^{[l]} = Y_i^{[l-1]} - \hat{c}_i^{[l-1]} \rho_i^{[l-1]}$$

where $W_i^{[1]} = W_i$ and $Y_i^{[1]} = Y_i$. This process ensures that the residualized predictors $W_i^{[l]}$ are orthogonal to all previous directions $\xi^{[j]}$ for $j < l$. Consequently, the newly found direction $\xi^{[l]}$, being derived from $W_i^{[l]}$, will be orthogonal to $\{\xi^{[1]}, \dots, \xi^{[l-1]}\}$. By convention, PLS directions are normalized to have unit norm ($\|\xi^{[l]}\|_{\mathbb{H}} = 1$). Thus, the set $\{\xi^{[l]}\}_{l=1}^{\infty}$ is an orthonormal set.

4.1.0.2. 1.2 Completeness. The PLS algorithm, by repeatedly extracting the direction of maximum covariance from the residual space, exhaustively captures all linear information within the predictors that is correlated with the response. This iterative process terminates when the residualized predictors no longer exhibit covariance with the residualized response. Therefore, the span of the PLS directions, $\text{span}\{\xi^{[l]}\}_{l=1}^{\infty}$, precisely constitutes the subspace of \mathbb{H} that contains all the linear information in the predictors relevant to predicting the response. This implies that $\{\xi^{[l]}\}_{l=1}^{\infty}$ forms a CONS for this relevant predictor subspace.

2. The Ordinary Linear Regression Solution as a Series Expansion in the PLS Basis. Since $\{\xi^{[l]}\}_{l=1}^{\infty}$ is a CONS for the relevant predictor subspace, the true hybrid regression coefficient $\beta \in \mathbb{H}$ can be expressed as an infinite series in this basis:

$$\beta = \sum_{l=1}^{\infty} \langle \beta, \xi^{[l]} \rangle_{\mathbb{H}} \xi^{[l]}$$

Substituting this expansion into the definition of the ordinary linear regression solution \hat{Y}_i :

$$\hat{Y}_i = \langle \beta, W_i \rangle_{\mathbb{H}} = \left\langle \sum_{l=1}^{\infty} \langle \beta, \xi^{[l]} \rangle_{\mathbb{H}} \xi^{[l]}, W_i \right\rangle_{\mathbb{H}}$$

By the linearity of the inner product in a Hilbert space:

$$\hat{Y}_i = \sum_{l=1}^{\infty} \langle \beta, \xi^{[l]} \rangle_{\mathbb{H}} \langle \xi^{[l]}, W_i \rangle_{\mathbb{H}}$$

Let $c^{[l]} = \langle \beta, \xi^{[l]} \rangle_{\mathbb{H}}$ be the true coefficient for the l -th component and $\rho_i^{[l]} = \langle W_i, \xi^{[l]} \rangle_{\mathbb{H}}$ be the l -th PLS score for sample i . Then,

$$\hat{Y}_i = \sum_{l=1}^{\infty} c^{[l]} \rho_i^{[l]}$$

This shows that the exact regression solution is an infinite series involving the PLS scores. The PLS approximation $\hat{Y}_{i,L}$ is a partial sum of this series. Furthermore, it can be shown that the estimated coefficients $\hat{c}^{[l]}$ in the PLS approximation $\hat{Y}_{i,L}$ are consistent estimators of the true coefficients $c^{[l]}$.

3. Proof of L^2 Convergence. Now, we examine the mean-squared error between the PLS approximation and the ordinary linear regression solution:

$$E \left[\|\hat{Y}_{i,L} - \hat{Y}_i\|_{L^2(\Omega)}^2 \right] = E \left[\left\| \sum_{l=1}^L c^{[l]} \rho_i^{[l]} - \sum_{l=1}^{\infty} c^{[l]} \rho_i^{[l]} \right\|^2 \right]$$

This can be rewritten as the expectation of the squared norm of the "tail" of the series:

$$E \left[\|\hat{Y}_{i,L} - \hat{Y}_i\|_{L^2(\Omega)}^2 \right] = E \left[\left\| - \sum_{l=L+1}^{\infty} c^{[l]} \rho_i^{[l]} \right\|^2 \right]$$

Since the PLS scores $\{\rho_i^{[l]}\}_{l=1}^{\infty}$ form an orthogonal system in the space of random variables (due to the orthogonality of $\{\xi^{[l]}\}$ and properties of expectation), the expectation of the squared sum simplifies to the sum of the expectations of the squared individual terms:

$$E \left[\|\hat{Y}_{i,L} - \hat{Y}_i\|_{L^2(\Omega)}^2 \right] = \sum_{l=L+1}^{\infty} (c^{[l]})^2 E[(\rho_i^{[l]})^2]$$

The infinite series $\sum_{l=1}^{\infty} (c^{[l]})^2 E[(\rho_i^{[l]})^2]$ represents the total variance explained by the predictors and is a convergent series (related to the total variance of \hat{Y}_i). For any convergent series, the sum of its tail (i.e., from $L + 1$ to infinity) must approach zero as L tends to infinity.

$$\lim_{L \rightarrow \infty} \sum_{l=L+1}^{\infty} (c^{[l]})^2 E[(\rho_i^{[l]})^2] = 0$$

Therefore,

$$\lim_{L \rightarrow \infty} E \left[\|\hat{Y}_{i,L} - \hat{Y}_i\|_{L^2(\Omega)}^2 \right] = 0$$

This concludes the proof that the Hybrid PLS approximation converges to the ordinary linear regression solution in the mean-squared sense. \square

4.2. Geometric properties. A fundamental property of partial least squares is that between iterations, its derived directions are orthonormal and PLS scores are orthogonal. Our regularized estimates preserve this property, with respect to a modified inner product that incorporates the roughness penalty, defines as follows:

DEFINITION 4.6 (Roughness-sensitive inner product). Given two hybrid predictors $W_1 = (X_{11}, \dots, X_{1K}, \mathbf{Z}_1)$ and $W_2 = (X_{21}, \dots, X_{2K}, \mathbf{Z}_2)$, both elements of \mathbb{H} as defined in Definition 3.3, and a roughness penalty matrix $\Lambda = \text{blkdiag}(\lambda_1 I_M, \dots, \lambda_K I_M, 0_{p \times p})$, the roughness-sensitive inner product between W_1 and W_2 is defined as:

$$(25) \quad \langle W_1, W_2 \rangle_{\mathbb{H}, \Lambda} := \sum_{k=1}^K \int_0^1 X_{1k}(t) X_{2k}(t) dt + \sum_{k=1}^K \lambda_k \int_0^1 X_{1k}''(t) X_{2k}''(t) dt + \mathbf{Z}_1^\top \mathbf{Z}_2.$$

Based on this inner product, the following proposition states that the PLS component directions estimated from Proposition 3.5 are orthonormal.

PROPOSITION 4.7 (Orthonormality of estimated PLS component directions). *The PLS component directions $\hat{\xi}^{[1]}, \hat{\xi}^{[2]}, \dots, \hat{\xi}^{[L]}$, estimated via Proposition 3.5 with a roughness penalty matrix $\Lambda = \text{blkdiag}(\lambda_1 I_M, \dots, \lambda_K I_M, 0_{p \times p})$, are mutually orthonormal with respect to the inner product $\langle \cdot, \cdot \rangle_{\mathbb{H}, \Lambda}$. That is,*

$$\langle \hat{\xi}^{[l_1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}, \Lambda} = \mathbb{1}(l_1 = l_2), \quad l_1, l_2 = 1, \dots, L.$$

The proof of Proposition 4.7 is provided in Appendix E.5.

The next proposition states that the vectors of estimated PLS scores for different iteration numbers are mutually orthogonal.

PROPOSITION 4.8. *Recall from Lemma 3.7 that $\hat{\boldsymbol{\rho}}^{[l]}$ denote the n -dimensional vector whose elements consist of the l -th estimated PLS scores ($l = 1, \dots, L$) of n observations. The vectors $\hat{\boldsymbol{\rho}}^{[1]}, \hat{\boldsymbol{\rho}}^{[2]}, \dots, \hat{\boldsymbol{\rho}}^{[L]}$ are mutually orthogonal in the sense that*

$$\hat{\boldsymbol{\rho}}^{[l_1]^\top} \hat{\boldsymbol{\rho}}^{[l_2]} = 0 \quad \text{for } l_1, l_2 \in \{1, \dots, L\}, l_1 \neq l_2.$$

The proof of Proposition 4.8 is provided in Appendix E.6.

5. Simulation studies. (JM:) 08/16/2025: re-starting. the writing here is not completed.

To evaluate the superiority of our method under complex dependency structures — specifically, dependencies among functional predictors, among scalar predictors, and between scalar and functional predictors —

Matrix-normal setting. We begin by constructing predictors from a matrix-normal distribution, a framework that offers convenient and flexible control over the dependence structure across both rows and columns. Matrix-normal (MN) models, also known as Kronecker-separable covariance models, provide a principled approach to modeling multivariate data with structured covariance. Specifically, the matrix-normal distribution is defined as

$$\mathbf{X} \sim \mathcal{MN}_{m \times n}(\mathbf{M}; \mathbf{R}, \mathbf{C}),$$

and its log-density is given by

$$\log p(\mathbf{X} \mid \mathbf{M}, \mathbf{R}, \mathbf{C}) = -\frac{mn}{2} \log(2\pi) - \frac{m}{2} \log |\mathbf{C}| - \frac{n}{2} \log |\mathbf{R}| - \frac{1}{2} \text{Tr} \left[\mathbf{C}^{-1} (\mathbf{X} - \mathbf{M})^\top \mathbf{R}^{-1} (\mathbf{X} - \mathbf{M}) \right].$$

The key insight behind Kronecker separability is that if $\mathbf{Y} \sim \mathcal{MN}(\mathbf{M}, \mathbf{R}, \mathbf{C})$, then its vectorized form follows a multivariate normal distribution: $\text{vec}(\mathbf{Y}) \sim \mathcal{N}(\text{vec}(\mathbf{M}), \mathbf{C} \otimes \mathbf{R})$, where \otimes denotes the Kronecker product and vec is the vectorization operator.

Building on the functional graphical model simulation of [Zhu, Strawn and Dunson \(2016\)](#), we generate a mixed graphical model with five nodes, described as follows:

- Nodes F1 and F2: Two functional predictors modeled as Gaussian processes using a truncated Karhunen-Loève expansion, where the eigenbasis consists of Fourier basis functions with a fixed number of basis functions, $M = 9$.
- Nodes S1, S2, and S3: Three scalar predictors, each following an s -dimensional multivariate normal distribution. Unlike the functional predictors, these scalar predictors are modeled directly without basis expansion.

To capture dependencies among predictors, we introduce a graph structure that governs their conditional correlations. We consider two types of graph structures: a weakly connected graph and a strongly connected graph. In the Gaussian process framework, the precision matrix \mathbf{R}_0^{-1} encodes conditional independence relationships, while its inverse, \mathbf{R}_0 , represents marginal covariances. This structure extends to a blockwise correlation matrix $\mathbf{R} \in \mathbb{R}^{(2M+3s) \times (2M+3s)}$, where off-diagonal blocks represent correlations between FPC scores and scalar predictor values. Each block (i, j) of \mathbf{R} is given by $(R_0)_{ij} \mathbf{I}_{M_i, M_j}$, where \mathbf{I}_{M_i, M_j} is a rectangular identity matrix. Here, $M_i = 9$ if node i corresponds to a functional predictor and $M_j = s$ if node j corresponds to a scalar predictor.

For each functional predictor, we assign M reference eigenvalues (or FPC score) drawn independently from gamma distributions with decreasing means. These reference eigenvalues will later be multiplied by randomly multipliers drawn from correlated multivariate Normal distribution. These reference eigenvalues are fixed over all samples and all independent repetition of the experiment. We draw it one randomly just to ensure the differentiation between the two functional predictors. We independently draw for each functional predictor from . We then sample zero-mean multivariate normal data from the covariance matrix \mathbf{R} . The last $3s$ components are assigned as scalar predictors, while the first $2M$ components, scaled by their corresponding eigenvalues, serve as FPC scores. These scores are then expanded into functional data, evaluated over 100 equally spaced points on $[0, 1]$ using a Fourier basis.

Finally, we introduce structured variability by adding a common mean function, defined as a scaled and shifted sine function. To visualize the generated data, Figure 1 plots functional predictor realizations and heatmaps of sample correlation matrices for both graph structures, illustrating the distinct dependency patterns induced by the precision matrices.

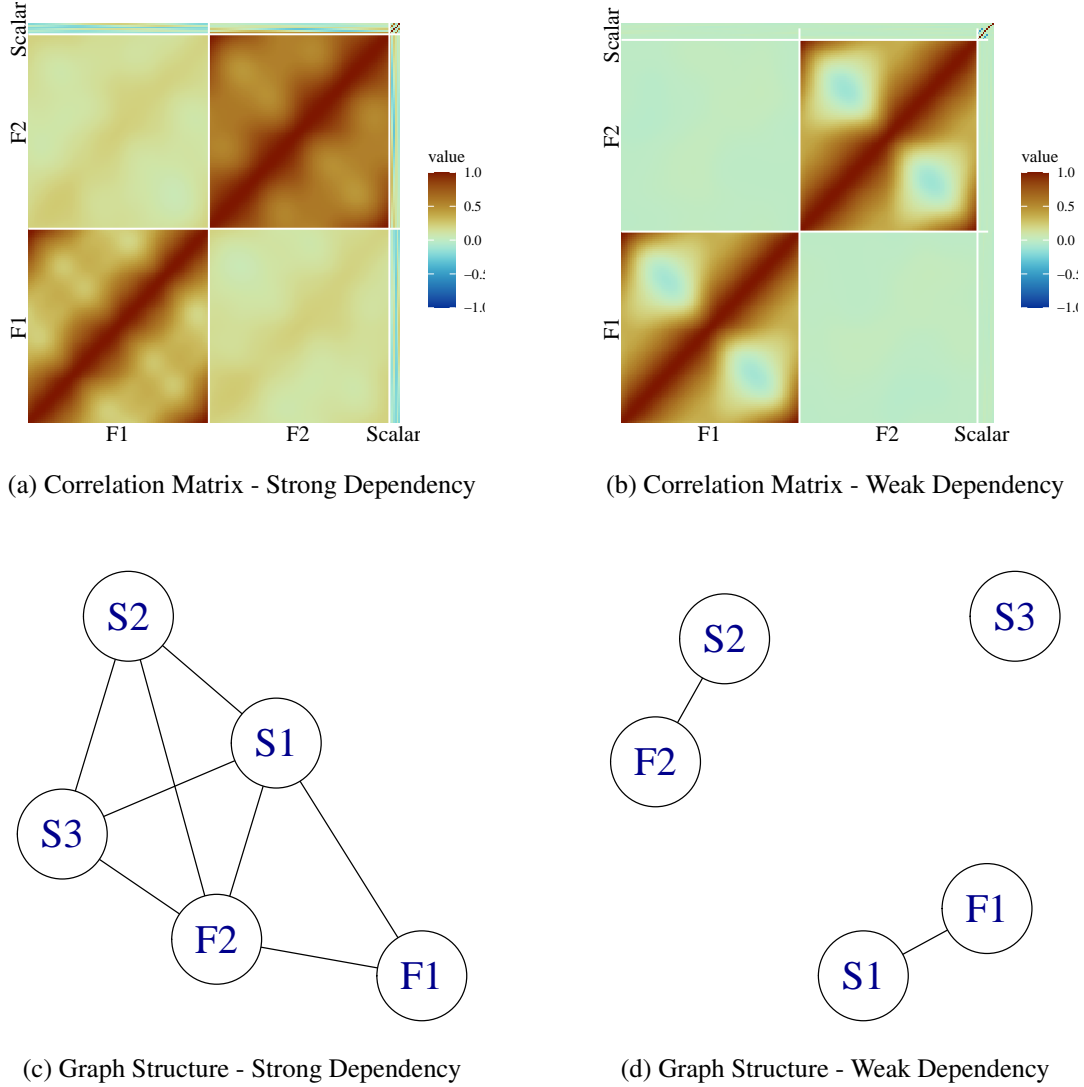


Fig 1: Comparison of Correlation Matrices and Graph Structures under Strong and Weak Dependencies.

5.1. *Simulation shown in february 2025 meeting.* To effectively manage the dependencies among functional predictors, scalar covariates, and between functional and scalar predictors, we utilize a Gaussian Markov random field (GMRF) within the framework of Gaussian undirected graphical models. In a GMRF, the off-diagonal elements of the precision matrix capture the conditional correlations between the corresponding components. Given that our setting involves both functional and scalar covariates, we adopt the simulation setup proposed by [Kolar, Liu and Xing \(2014\)](#), which focuses on mixed attribute Gaussian graphical models.

We construct a graph consisting of two functional predictor nodes and three vector predictor nodes. Below, we sequentially describe the graph structure and the generation process for each node.

Functional Predictors. We consider two functional predictors that share the same precision matrix. Denoted as $\Theta := (\theta_{ti}) \in \mathbb{R}^{p \times p}$, this precision matrix follows an AR(1) structure with

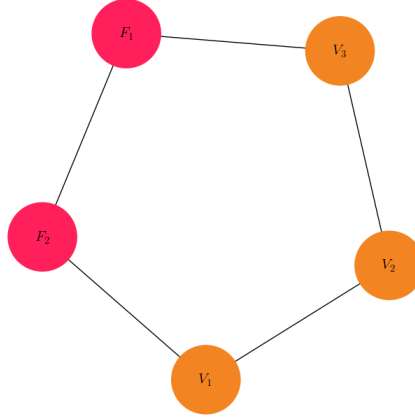


Fig 2

a white noise variance of 0.1^2 and an autoregressive coefficient of $\rho = 0.95$, ensuring a smooth functional trajectory.

More formally, each functional predictor is evaluated at p equally spaced points on $[0, 1]$, with values defined recursively as:

$$X_i^{(k)}(0) = 0, \quad X_i^{(k)}(t) = \rho X_i^{(k)}(t-1) + \varepsilon_{itk}, \quad \varepsilon_{itk} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2), \quad t = 1, \dots, p, \quad i = 1, \dots, n, \quad k = 1, 2.$$

In this setting, the precision matrix Θ exhibits a tridiagonal Toeplitz structure, where the diagonal entries are given by:

$$\theta_{tt} = \begin{cases} \frac{1}{\sigma^2(1 - 0.95^2)}, & t = 1, p, \\ \frac{1 + 0.95^2}{\sigma^2(1 - 0.95^2)}, & 2 \leq t \leq p - 1. \end{cases}$$

The off-diagonal entries are:

$$\theta_{t,t+1} = \theta_{t+1,t} = -\frac{\rho}{\sigma^2(1 - \rho^2)}, \quad 1 \leq t \leq p - 1.$$

Finally, the functional predictors are smoothed using a B-spline basis with 15 basis functions, as described in Section 3.3.1.

Vector Predictors. We consider s vector predictors, each following a d -dimensional zero-mean Gaussian distribution with a shared precision matrix. This precision matrix, denoted as $\Gamma := (\gamma_{ij}) \in \mathbb{R}^{d \times d}$, follows a Toeplitz structure with exponentially decaying entries, given by:

$$\gamma_{ij} := 0.5^{|i-j|}.$$

Graph Structure. We arrange five nodes in a chain structure, where each node follows a sequential order, and the last node connects back to the first, as illustrated in Figure 2. An edge in the graph indicates that the connected nodes remain correlated when the values of all other nodes are fixed. The resulting marginal dependence structure is significantly more complex than the chain structure itself. We designate nodes F_1 and F_2 as functional predictors and nodes V_1 , V_2 , and V_3 as vector predictors. To introduce conditional dependencies between the components, we set the off-diagonal blocks of the precision matrix to be 0.51 if the

Table 1: Gaussian Markov random field simulation results

Sample Size	Method	$p_{\text{scalar}} = 3$	$p_{\text{scalar}} = 6$	$p_{\text{scalar}} = 9$	$p_{\text{scalar}} = 12$	$p_{\text{scalar}} = 15$
100	fpcr	0.195 (0.028)	0.195 (0.030)	0.207 (0.034)	0.215 (0.036)	0.232 (0.041)
	hybridpls	0.116 (0.016)	0.130 (0.021)	0.162 (0.030)	0.185 (0.032)	0.215 (0.038)
	pfr	0.341 (0.104)	0.323 (0.100)	0.354 (0.102)	0.379 (0.105)	0.417 (0.110)
200	fpcr	0.175 (0.016)	0.171 (0.016)	0.175 (0.017)	0.178 (0.017)	0.182 (0.018)
	hybridpls	0.107 (0.010)	0.110 (0.010)	0.125 (0.014)	0.149 (0.017)	0.162 (0.017)
	pfr	0.201 (0.059)	0.186 (0.031)	0.193 (0.036)	0.204 (0.052)	0.215 (0.064)
300	fpcr	0.169 (0.013)	0.167 (0.012)	0.168 (0.013)	0.171 (0.012)	0.174 (0.013)
	hybridpls	0.104 (0.007)	0.106 (0.007)	0.113 (0.009)	0.137 (0.013)	0.150 (0.013)
	pfr	0.175 (0.018)	0.172 (0.013)	0.173 (0.015)	0.176 (0.013)	0.180 (0.018)
400	fpcr	0.167 (0.011)	0.164 (0.011)	0.167 (0.011)	0.167 (0.011)	0.170 (0.011)
	hybridpls	0.103 (0.006)	0.104 (0.006)	0.110 (0.007)	0.129 (0.011)	0.144 (0.012)
	pfr	0.172 (0.011)	0.169 (0.011)	0.172 (0.011)	0.173 (0.011)	0.175 (0.012)

Fig 3: Enter Caption

corresponding components are connected in the graph, and zero otherwise. Here, $\mathbf{1}$ denotes a matrix of appropriate dimensions where all elements are equal to 1. Overall, we generate a $2p + 3d$ -dimensional multivariate Gaussian distribution with mean zero and a precision matrix Ω , structured as follows:

$$\Omega = \begin{pmatrix} \Omega_F & 0.5\mathbf{1} & 0 & 0 & 0.5\mathbf{1} \\ 0.5\mathbf{1} & \Omega_F & 0.5\mathbf{1} & 0 & 0 \\ 0 & 0.5\mathbf{1} & \Omega_V & 0.5\mathbf{1} & 0 \\ 0 & 0 & 0.5\mathbf{1} & \Omega_V & 0.5\mathbf{1} \\ 0.5\mathbf{1} & 0 & 0 & 0.5\mathbf{1} & \Omega_V \end{pmatrix}$$

For each observation drawn from this multivariate Gaussian distribution, we process

The regression coefficients for the first functional predictor, β_{F_1} , are drawn from a multivariate normal distribution $N(0, 5\mathbf{I}_p)$, with a fixed random seed. The coefficients for the second functional predictor, β_{F_2} , are drawn independently from the same distribution and are also smoothed using a B-spline basis with 10 basis functions. For the vector covariates, the regression coefficients are sampled from $N(0, I_d)$. After calculating the inner product of the covariates and their corresponding coefficients, independent Gaussian noise from $N(0, 0.1^2)$ is added to the generated responses to simulate measurement noise.

The baseline methods are:

- Penalized functional regression (pfr) [Goldsmith et al. \(2011\)](#)
- Principal component regression (fpcr): run both of PCA for multiple functional predictors [Happ and Greven \(2018\)](#) and scalar PCA and run OLS on the PC scores.

We compare our method with these baseline methods using $p = 100$. We consider scenarios with $d = 1, 2, 3, 4, 5$ and $n = 100, 200, 300, 400$. For each scenario, we use 70% of the data for training and evaluate prediction performance on the remaining 30% test set, using the root prediction mean squared error as the evaluation metric. For our method and PCR, the maximum number of components are set as 20. The number of components is chosen by 5-fold cross validation. The results, summarized in Table ??, demonstrate that our method consistently outperforms the baseline methods across all scenarios.

6. Data Application. (JM:) 08/16/2025: re-starting. the writing here is not completed.

Renal study data. We applied our proposed hybrid functional PLS regression, along with other regression methods, to the Emory renal study data. The study collected data on 226

kidneys (left and right) from 113 subjects, including: (i) baseline renogram curves; (ii) post-furosemide renogram curves; (iii) ordinal ratings of kidney obstruction status (non-obstructed, equivocal, or obstructed) independently assessed by three nuclear medicine experts; (iv) eight kidney-level pharmacokinetic variables derived from radionuclide imaging; and (v) two subject-level variables (age and gender). The subjects had a mean age of 57.8 years (SD = 15.5; range = 18–83), with 54 males (48%) and 59 females (52%). The three experts unanimously classified 153 kidneys as non-obstructed, 5 as equivocal, and 40 as obstructed, while 28 kidneys had discrepant ratings.

The two renogram curves, (i) and (ii), were treated as functional predictors and smoothed using a B-spline basis of order 15. The remaining variables, excluding the diagnosis, were treated as scalar predictors. Given the nature of these variables, we assume they are correlated with the renogram curves but not entirely redundant, as they may contain additional useful information. Finally, the diagnoses provided by the three experts were averaged and transformed using a min-max logit transformation. We splitted the data into 70% of training data and 30% of testing data, and evaluated the prediction performance by root mean squared error on the test data, normalized by the range of the test data response.

6.1.

Acknowledgments. The authors would like to thank the anonymous referees, an Associate Editor and the Editor for their constructive comments that improved the quality of this paper.

Funding. The first author was supported by NSF Grant DMS-??-??????.

The second author was supported in part by NIH Grant ??????????.

SUPPLEMENTARY MATERIAL

Title of Supplement A

Short description of Supplement A.

Title of Supplement B

Short description of Supplement B.

REFERENCES

- AGUILERA, A. M., AGUILERA-MORILLO, M. C. and PREDÁ, C. (2016). Penalized Versions of Functional PLS Regression. *Chemometrics and Intelligent Laboratory Systems* **154** 80–92.
- AGUILERA, A. M., ESCABIAS, M., PREDÁ, C. and SAPORTA, G. (2010). Using Basis Expansions for Estimating Functional PLS Regression: Applications with Chemometric Data. *Chemometrics and Intelligent Laboratory Systems* **104** 289–305.
- BAZARAA, M. S., SHERALI, H. D. and SHETTY, C. M. (2006). *Nonlinear Programming: Theory and Algorithms*. Wiley-Interscience, Hoboken, N.J.
- BEYAZTAS, U. and LIN SHANG, H. (2022). A Robust Functional Partial Least Squares for Scalar-on-Multiple-Function Regression. *Journal of Chemometrics* **36** e3394.
- BEYAZTAS, U. and SHANG, H. L. (2020). On function-on-function regression: Partial least squares approach. *Environmental and Ecological Statistics* **27** 95–114.
- CAI, T. T. and HALL, P. (2006). Prediction in Functional Linear Regression. *The Annals of Statistics* **34** 2159–2179.
- CARDOT, H., FERRATY, F. and SARDA, P. (2003). Spline Estimators for the Functional Linear Model. *Statistica Sinica* **13** 571–591.
- CHANG, C., JANG, J. H., MANATUNGA, A., TAYLOR, A. T. and LONG, Q. (2020). A Bayesian Latent Class Model to Predict Kidney Obstruction in the Absence of Gold Standard. *Journal of the American Statistical Association* **115** 1645–1663.

- DELAIGLE, A. and HALL, P. (2012). Methodology and Theory for Partial Least Squares Applied to Functional Data. *The Annals of Statistics* **40** 322–352.
- FEBRERO-BANDE, M., GALEANO, P. and GONZÁLEZ-MANTEIGA, W. (2017). Functional Principal Component Regression and Functional Partial Least-Squares Regression: An Overview and a Comparative Study. *International Statistical Review* **85** 61–83.
- GENG, S., KOLAR, M. and KOYEJO, O. (2020). Joint Nonparametric Precision Matrix Estimation with Confounding. In *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference* 378–388. PMLR.
- GOCKENBACH, M. S. (2010). *Finite-Dimensional Linear Algebra*, 1st edition ed. CRC Press, Boca Raton, FL.
- GOLDSMITH, J., BOBB, J., CRAINICEANU, C. M., CAFFO, B. and REICH, D. (2011). Penalized Functional Regression. *Journal of Computational and Graphical Statistics* **20** 830–851. <https://doi.org/10.1198/jcgs.2010.10007>
- HALL, P. and HOROWITZ, J. L. (2007). Methodology and Convergence Rates for Functional Linear Regression. *The Annals of Statistics* **35** 70–91.
- HAPP, C. and GREVEN, S. (2018). Multivariate Functional Principal Component Analysis for Data Observed on Different (Dimensional) Domains. *Journal of the American Statistical Association* **113** 649–659.
- HSING, T. and EUBANK, R. (2015). *Theoretical Foundations of Functional Data Analysis, with an Introduction to Linear Operators*, 1st edition ed. Wiley.
- JANG, J. H. (2021). Principal Component Analysis of Hybrid Functional and Vector Data. *Statistics in Medicine* **40** 5152–5173. <https://doi.org/10.1002/sim.9117>
- KOLAR, M., LIU, H. and XING, E. P. (2014). Graph Estimation From Multi-Attribute Data. *Journal of Machine Learning Research* **15** 1713–1750.
- PRED, C. and SAPORTA, G. (2005). PLS Regression on a Stochastic Process. *Computational Statistics & Data Analysis* **48** 149–158.
- REISS, P. T. and OGDEN, R. T. (2007). Functional Principal Component Regression and Functional Partial Least Squares. *Journal of the American Statistical Association* **102** 984–996.
- RICE, J. A. and SILVERMAN, B. W. (1991). Estimating the mean and covariance structure nonparametrically when the data are curves. *Journal of the Royal Statistical Society: Series B* **53** 233–243.
- SARICAM, S., BEYAZTAS, U., ASIKGIL, B. and SHANG, H. L. (2022). On Partial Least-Squares Estimation in Scalar-on-Function Regression Models. *Journal of Chemometrics* **36** e3452.
- SILVERMAN, B. W. (1996). Smoothed functional principal components by choice of norm. *The Annals of Statistics* **24** 1–24.
- TUCKER, R. S. (1938). The reasons for price rigidity. *The American Economic Review* **28** 41–54.
- WANG, Y. (2018). Partial least squares methods for functional regression models, PhD thesis, University of North Carolina at Chapel Hill.
- ZHAO, Y., OGDEN, R. T. and REISS, P. T. (2012). Wavelet-Based LASSO in Functional Linear Regression. *Journal of Computational and Graphical Statistics* **21** 600–617.
- ZHU, H., STRAWN, N. and DUNSON, D. B. (2016). Bayesian Graphical Models for Multivariate Functional Data. *Journal of Machine Learning Research* **17** 1–27.

APPENDIX A: OVERVIEW OF APPENDIX

- Appendix B proves that our defined product space is a separable Hilbert space.
- Appendix C provides the proofs for the core problem formulations and computational schemes of the iterative steps for estimating the PLS direction and performing residualization.
- Appendix E provides proof of Theorem 4.4, demonstrating that the population-level version of the optimization problem solved at each step is well-defined.

APPENDIX B: PROOF OF LEMMA 3.2

PROOF. Both $\mathbb{L}^2[0, 1]$ and \mathbb{R}^p are Hilbert spaces, with their inner products, norms, and metrics corresponding to the terms defined in equation (2). A finite Cartesian product of Hilbert spaces, equipped with an ℓ_2 norm of the component metrics (which is our metric), is also a Hilbert space. Furthermore, both of $\mathbb{L}^2[0, 1]$ and \mathbb{R}^p are separable. For $\mathbb{L}^2[0, 1]$, the set of all polynomials with rational coefficients is a countable and dense subset. For \mathbb{R}^p , the set of all vectors with rational components forms a countable and dense subset. A finite Cartesian product of separable spaces is also separable. Therefore $\mathbb{H} = (\mathbb{L}^2[0, 1])^K \times \mathbb{R}^p$ is separable. \square

APPENDIX C: PROOF OF SECTION 3.3

This section provides the proofs for the core problem formulations and computational schemes of the iterative steps for estimating the PLS direction and performing residualization. Appendices C.1 and C.2 derive the eigenproblem formulations for unregularized and regularized PLS direction estimation, respectively. Appendices C.3 and C.4 provide proofs for the simple computation schemes for regularized PLS direction estimation and residualization, respectively.

C.1. Proof of Proposition 3.4.

PROOF. Let $\xi = (\xi_1, \dots, \xi_K, \zeta) \in \mathbb{H}$. Assume that each functional component $\xi_j \in \mathbb{L}^2([0, 1])$ lies in the span of the basis functions $b_1(t), \dots, b_M(t)$ and can be written as

$$\xi_j(t) := \sum_{m=1}^M d_{jm} b_m(t), \quad j = 1, \dots, K,$$

with coefficient vectors $\gamma_j := (d_{j1}, \dots, d_{jM})^\top \in \mathbb{R}^M$. Then, the empirical covariance is computed as follows:

$$\begin{aligned} \widehat{\text{Cov}}(\langle \widetilde{W}, \xi \rangle_{\mathbb{H}}, Y) &= \frac{1}{n} \sum_{i=1}^n y_i \langle \widetilde{W}_i, \xi \rangle_{\mathbb{H}} \\ &= \frac{1}{n} \sum_{i=1}^n y_i \left(\sum_{j=1}^K \langle \widetilde{X}_{ij}, \xi_j \rangle + \mathbf{z}_i^\top \zeta \right) \\ &= \frac{1}{n} \sum_{i=1}^n y_i \left(\sum_{k=1}^K \int_0^1 \widetilde{X}_{ik}(t) \xi_k(t) dt \right) + \frac{1}{n} \sum_{i=1}^n y_i (\mathbf{z}_i^\top \zeta) \\ &= \frac{1}{n} \sum_{i=1}^n y_i \left\{ \sum_{k=1}^K \int_0^1 \left(\sum_{m=1}^M \theta_{ikm} b_m(t) \right) \left(\sum_{m'=1}^M d_{km'} b_{m'}(t) \right) dt \right\} + \frac{1}{n} \mathbf{y}^\top Z \zeta \\ &= \frac{1}{n} \sum_{i=1}^n y_i \left\{ \sum_{k=1}^K \sum_{m=1}^M \sum_{m'=1}^M \theta_{ikm} d_{km'} \int_0^1 b_m(t) b_{m'}(t) dt \right\} + \frac{1}{n} \mathbf{y}^\top Z \zeta \\ &= \sum_{k=1}^K \left\{ \frac{1}{n} \sum_{i=1}^n y_i (\boldsymbol{\theta}_{ij}^\top B \gamma_j) \right\} + \frac{1}{n} \mathbf{y}^\top Z \zeta \\ &= \frac{1}{n} \left(\sum_{k=1}^K \mathbf{y}^\top \Theta_j B \gamma_j + \mathbf{y}^\top Z \zeta \right) \end{aligned}$$

Building on this computation, the squared empirical covariance is expressed as the quadratic form involving the matrix V defined in (15):

$$\begin{aligned} \widehat{\text{Cov}}^2(\langle \widetilde{W}, \xi \rangle_{\mathbb{H}}, Y) &= \frac{1}{n^2} \left(\sum_{k=1}^K \mathbf{y}^\top \Theta_j B \gamma_j + \mathbf{y}^\top Z \zeta \right)^2 \\ &= \frac{1}{n^2} (\mathbf{y}^\top \Theta \mathbb{B} \boldsymbol{\xi})^2 \\ &= \frac{1}{n^2} \boldsymbol{\xi}^\top (\mathbb{B} \Theta^\top \mathbf{y}) (\mathbb{B} \Theta^\top \mathbf{y})^\top \boldsymbol{\xi} \end{aligned}$$

$$(26) \quad = \boldsymbol{\xi}^\top V \boldsymbol{\xi}.$$

The squared norm of ξ in the hybrid Hilbert space \mathbb{H} is computed as follows:

$$\begin{aligned}
 \langle \xi, \xi \rangle_{\mathbb{H}} &= \sum_{j=1}^K \int_0^1 \xi_j(t) \xi_j(t) dt + \boldsymbol{\zeta}^\top \boldsymbol{\zeta} \\
 &= \sum_{j=1}^K \int_0^1 \left(\sum_{m=1}^M d_{jm} b_m(t) \right) \left(\sum_{m'=1}^M d_{jm'} b_{m'}(t) \right) dt + \boldsymbol{\zeta}^\top \boldsymbol{\zeta} \\
 &= \sum_{j=1}^K \sum_{m=1}^M \sum_{m'=1}^M d_{jm} d_{jm'} \int_0^1 b_m(t) b_{m'}(t) dt + \boldsymbol{\zeta}^\top \boldsymbol{\zeta} \\
 &= \sum_{j=1}^K \boldsymbol{\gamma}_j^\top B \boldsymbol{\gamma}_j + \boldsymbol{\zeta}^\top \boldsymbol{\zeta}, \\
 (27) \quad &= \boldsymbol{\xi}^\top \mathbb{B} \boldsymbol{\xi}
 \end{aligned}$$

Therefore, the covariance maximization problem (13) is equivalent to the generalized Raleigh quotient (14). This completes the proof of Proposition 3.4.

C.2. Proof of Proposition 3.5 . The penalization term $\sum_{j=1}^K \lambda_j \text{PEN}(\xi_K)$ can be written in matrix form as:

$$\begin{aligned}
 \sum_{j=1}^K \lambda_j \text{PEN}(\xi_j) &= \sum_{j=1}^K \lambda_j \int_0^1 \{\hat{\xi}_j''(t)\}^2 dt \\
 &= \sum_{j=1}^K \lambda_j \int_0^1 \left\{ \sum_{l=1}^M \theta_{ijl} b_l''(t) \right\}^2 dt \\
 &= \sum_{j=1}^K \lambda_j \sum_{l=1}^M \sum_{m=1}^M \theta_{ijl} \theta_{ijm} \int_0^1 b_l''(t) b_m''(t) dt \\
 &= \sum_{j=1}^K \lambda_j \boldsymbol{\gamma}_j^\top B'' \boldsymbol{\gamma}_j \\
 (28) \quad &= \boldsymbol{\xi}^\top \Lambda \mathbb{B}'' \boldsymbol{\xi},
 \end{aligned}$$

where $\boldsymbol{\xi}$ is the concatenated coefficient vector form of ξ as defined in (??), \mathbb{B}'' is defined in (10), and Λ is defined in (18).

Combining this with (27), we can compute

$$\|\xi\|_{\mathbb{H}} + \sum_{j=1}^K \lambda_j \text{PEN}(\xi_K) = \boldsymbol{\xi}^\top \mathbb{B} \boldsymbol{\xi} + \boldsymbol{\xi}^\top \Lambda \mathbb{B}'' \boldsymbol{\xi} = \boldsymbol{\xi}^\top (\mathbb{B} + \Lambda \mathbb{B}'') \boldsymbol{\xi}.$$

This computation along with (26) implies that the covariance maximization problem (17) is equivalent to the generalized Raleigh quotient (19). This completes the proof of Proposition 3.5. □

C.3. Proof of Proposition 3.6.

PROOF. The optimization problem (19) can be written as

$$\begin{aligned} & \max_{\gamma_1, \dots, \gamma_K \in \mathbb{R}^M, \zeta \in \mathbb{R}^p} \frac{1}{n^2} \left(\sum_{j=1}^K \mathbf{y}^\top \Theta_j B \gamma_j + \mathbf{y}^\top Z \zeta \right)^2, \\ & \text{subject to} \quad \sum_{j=1}^K \gamma_j^\top (B + \lambda_j B'') \gamma_j + \zeta^\top \zeta = 1. \end{aligned}$$

Let $\mathbf{u}_j := B \Theta_j^\top \mathbf{y} \in \mathbb{R}^M$ and $\mathbf{v} := Z^\top \mathbf{y} \in \mathbb{R}^p$. These are fixed problem data. Then the objective becomes

$$\frac{1}{n^2} \left(\sum_{j=1}^K \mathbf{u}_j^\top \gamma_j + \mathbf{v}^\top \zeta \right)^2.$$

The objective function is continuous. The constraint defines the boundary of an ellipsoid in a finite-dimensional Euclidean space, the feasible set is compact. By the Weierstrass Extreme Value Theorem (see, e.g., Theorem 2.3.1 in [Bazaraa, Sherali and Shetty 2006](#)), a global maximizer exists. The constraint function is continuously differentiable and its gradient vanishes only at the origin, which is not feasible. Thus the gradient is nonzero at all feasible points. Hence, the Linear Independence Constraint Qualification (LICQ) holds, and the Karush-Kuhn-Tucker (KKT) conditions are necessary for local optimality (see, e.g., Theorem 5.3.1 in [Bazaraa, Sherali and Shetty 2006](#)).

Define the Lagrangian:

$$\mathcal{L}(\{\gamma_j\}, \zeta, \mu) := \left(\sum_{j=1}^K \mathbf{u}_j^\top \gamma_j + \mathbf{v}^\top \zeta \right)^2 - \mu \left[\sum_{j=1}^K \gamma_j^\top (B + \lambda_j B'') \gamma_j + \zeta^\top \zeta - 1 \right].$$

Let $s := \frac{2}{n^2} \left(\sum_{j=1}^K \mathbf{u}_j^\top \gamma_j + \mathbf{v}^\top \zeta \right)$. The KKT conditions require that:

$$(29) \quad \nabla_{\gamma_j} \mathcal{L} = s \mathbf{u}_j - 2\mu (B + \lambda_j B'') \gamma_j = 0, \quad \text{for } j = 1, \dots, K,$$

$$(30) \quad \nabla_{\zeta} \mathcal{L} = s \mathbf{v} - 2\mu \zeta = 0.$$

From (29) and (30), we have

$$(B + \lambda_j B'') \gamma_j = \frac{s}{2\mu} \mathbf{u}_j, \quad \zeta = \frac{s}{2\mu} \mathbf{v}, \quad \text{for } j = 1, \dots, K.$$

This implies that for any local maximizer, there exists $c \neq 0$ such that

$$\gamma_j = c(B + \lambda_j B'')^{-1} \mathbf{u}_j, \quad \zeta = c \mathbf{v}, \quad \text{for } j = 1, \dots, K.$$

Since we assumed in Section 3.1 that the functions $\{b_m(t)\}$ and their second derivatives are linearly independent, both B and B'' are positive definite (see, for example, Theorem 273 of [Gockenbach, 2010](#)). As positive definiteness is preserved under conic combinations, $B + \lambda_j B''$ is also positive definite.

A local maximizer must also be primal feasible. Substituting the conditions above in to the primal constraint:

$$\sum_{j=1}^K (c(B + \lambda_j B'')^{-1} \mathbf{u}_j)^\top (B + \lambda_j B'') (c(B + \lambda_j B'')^{-1} \mathbf{u}_j) + c^2 \mathbf{v}^\top \mathbf{v} = 1$$

$$\Rightarrow c^2 \left(\sum_{j=1}^K \mathbf{u}_j^\top (B + \lambda_j B'')^{-1} \mathbf{u}_j + \mathbf{v}^\top \mathbf{v} \right) = 1.$$

Solving for c^2 gives:

$$c^2 = \left(\sum_{j=1}^K \mathbf{u}_j^\top (B + \lambda_j B'')^{-1} \mathbf{u}_j + \mathbf{v}^\top \mathbf{v} \right)^{-1}.$$

Hence, the unique maximizer (up to sign) is:

$$\gamma_j^* = \frac{1}{\sqrt{q}} (B + \lambda_j B'')^{-1} \mathbf{u}_j, \quad \zeta^* = \frac{1}{\sqrt{q}} \mathbf{v},$$

where

$$q := \sum_{j=1}^K \mathbf{u}_j^\top (B + \lambda_j B'')^{-1} \mathbf{u}_j + \mathbf{v}^\top \mathbf{v}.$$

This completes the proof of Proposition C.3. □

C.4. Proof of Lemma 3.7 . For notational simplicity, we omit the iteration superscript $[l]$. Recall from (20) that $\hat{\rho}_i = \langle \widetilde{W}_i, \hat{\xi} \rangle_{\mathbb{H}}$. Using the basis expansion coefficient notation from (4), (6), (7), and (9), the full vector of PLS scores $\hat{\rho} = (\hat{\rho}_1, \dots, \hat{\rho}_n)^\top$ is computed through the following matrix multiplication:

$$(31) \quad \hat{\rho}^\top = \sum_{k=1}^K (\hat{\gamma}_k)^\top B \Theta_k^\top + \zeta^\top \mathbf{Z}^\top.$$

The least square criterion can be decomposed as follows:

$$\begin{aligned} \text{SSE}(\delta) &= \sum_{i=1}^n \langle \widetilde{W}_i^{[l]} - \hat{\rho}_i^{[l]} \delta, \widetilde{W}_i^{[l]} - \hat{\rho}_i^{[l]} \delta \rangle_{\mathbb{H}} \\ &= \underbrace{\sum_{i=1}^n \langle \widetilde{W}_i^{[l]}, \widetilde{W}_i^{[l]} \rangle_{\mathbb{H}}}_A - 2 \underbrace{\sum_{i=1}^n \hat{\rho}_i^{[l]} \langle \delta, \widetilde{W}_i^{[l]} \rangle_{\mathbb{H}}}_{B_1(\delta)} + \underbrace{(\hat{\rho}^\top \hat{\rho}) \langle \delta, \delta \rangle_{\mathbb{H}}}_{B_2(\delta)} \end{aligned}$$

Here, part A does not contain δ , so it does not contribute to the minimization problem.

Let $\delta := (\pi_1^\top, \dots, \pi_K^\top, \chi^\top)^\top$ be the $MK + p$ -dimensional concatenated vector of basis coefficients and the scalar part of δ . Abusing notation, we treat SSE as a function of δ . Next, We will now demonstrate that the combined term $B_1(\delta) + B_2(\delta)$ is a quadratic function of δ , by expanding its component functions:

$$\begin{aligned} B_1(\delta) &= -2 \sum_{i=1}^n \hat{\rho}_i \langle \delta, \widetilde{W}_i \rangle_{\mathbb{H}} = -2 \hat{\rho}^\top \left(\sum_{j=1}^K \Theta_j B \pi_j + \mathbf{Z} \chi \right), \\ B_2(\delta) &= (\hat{\rho}^\top \hat{\rho}) \langle \delta, \delta \rangle_{\mathbb{H}} = (\hat{\rho}^\top \hat{\rho}) \left(\sum_{j=1}^K \pi_j^\top B \pi_j + \chi^\top \chi \right), \end{aligned}$$

Now we compute the gradients:

$$\begin{aligned}\nabla_{\pi_j} B_1(\delta) &= -2B\Theta_j^\top \hat{\rho}, \quad j = 1, \dots, K, \quad \nabla_{\chi} B_1(\delta) = -2\mathbf{Z}^\top \hat{\rho}, \\ \nabla_{\pi_j} B_2(\delta) &= 2(\hat{\rho}^\top \hat{\rho}) B \pi_j, \quad j = 1, \dots, K, \quad \nabla_{\chi} B_2(\delta) = 2(\hat{\rho}^\top \hat{\rho}) \chi,\end{aligned}$$

The Hessian of $\text{SSE}(\delta)$ is then given by $2(\hat{\rho}^\top \hat{\rho})B$. Since B is positive definite, $\text{SSE}(\delta)$ is convex. Consequently, the gradient vanishes at its unique minimizer:

$$\begin{aligned}\nabla_{\pi_j} \text{SSE}(\hat{\delta}) &= -2B\Theta_j^\top \hat{\rho} + 2(\hat{\rho}^\top \hat{\rho}) B \hat{\pi}_j = 0, \quad j = 1, \dots, K, \\ \nabla_{\chi} \text{SSE}(\hat{\delta}) &= -2\mathbf{Z}^\top \hat{\rho} + 2(\hat{\rho}^\top \hat{\rho}) \hat{\chi} = 0,\end{aligned}$$

providing the following closed-form solution:

$$\hat{\pi}_j = \{(\hat{\rho}^\top \hat{\rho})B\}^{-1} B \Theta_j^\top \hat{\rho} = \frac{1}{\hat{\rho}^\top \hat{\rho}} \Theta_j^\top \hat{\rho} \quad j = 1, \dots, K, \quad \hat{\chi} = \frac{1}{(\hat{\rho}^\top \hat{\rho})} \mathbf{Z}^\top \hat{\rho}.$$

Expanding with respect to these coefficients, we obtain

$$\delta^{[l]} = \frac{1}{\|\hat{\rho}^{[l]}\|_2^2} \sum_{i=1}^n \hat{\rho}_i^{[l]} \widetilde{W}_i^{[l]}.$$

On the other hand since $\mathbf{y}^{[l]}$ and $\hat{\rho}^{[l]}$ are zero-mean, it is well known that the least square estimate of linear regression coefficient is computed as $\hat{\nu}^{[l]} = \frac{\mathbf{y}^{[l]\top} \hat{\rho}^{[l]}}{\|\hat{\rho}^{[l]}\|_2^2}$. This completes the proof of Lemma 3.7.

APPENDIX D: PROOF OF LEMMA 3.8

PROOF. We prove the lemma using mathematical induction.

Base case ($l = 1$). For $l = 1$, the lemma states $\hat{\rho}_i^{[1]} = \langle W_i, \hat{\iota}^{[1]} \rangle_{\mathbb{H}}$. Since $\hat{\iota}^{[1]} := \hat{\xi}^{[1]}$ and $W_i^{[1]} = W_i$, the base case holds.

Inductive steps. Assume $\hat{\rho}_i^{[u]} = \langle W_i, \hat{\iota}^{[u]} \rangle_{\mathbb{H}}$ for $u = 1, \dots, l$. We want to show that $\hat{\rho}_i^{[l+1]} = \langle W_i, \hat{\iota}^{[l+1]} \rangle_{\mathbb{H}}$. Recall from (20) and Lemma 3.7 that we have

$$(32) \quad \hat{\rho}_i^{[l+1]} = \langle \widetilde{W}_i^{[l+1]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}}, \quad \widetilde{W}_i^{[l+1]} = \widetilde{W}_i^{[l]} - \hat{\rho}_i^{[l]} \hat{\delta}^{[l]}$$

Combining these two, we have

$$(33) \quad \hat{\rho}_i^{[l+1]} = \langle \widetilde{W}_i^{[l]} - \hat{\rho}_i^{[l]} \hat{\delta}^{[l]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} = \langle \widetilde{W}_i^{[l]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} - \hat{\rho}_i^{[l]} \langle \hat{\delta}^{[l]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}}.$$

Next, we can express $\widetilde{W}_i^{[l]}$ in terms of the original predictor W_i by recursively applying Lemma 3.7:

$$\widetilde{W}_i^{[l]} = \widetilde{W}_i^{[l-1]} - \hat{\rho}_i^{[l-1]} \hat{\delta}^{[l-1]} = \dots = W_i - \sum_{u=1}^{l-1} \hat{\rho}_i^{[u]} \hat{\delta}^{[u]}$$

Substituting this into our equation (33), we have:

$$\begin{aligned}\hat{\rho}_i^{[l+1]} &= \left\langle W_i - \sum_{u=1}^{l-1} \hat{\rho}_i^{[u]} \hat{\delta}^{[u]}, \hat{\xi}^{[l+1]} \right\rangle_{\mathbb{H}} - \hat{\rho}_i^{[l]} \langle \hat{\delta}^{[l]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} \\ &= \langle W_i, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} - \sum_{u=1}^{l-1} \hat{\rho}_i^{[u]} \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} - \hat{\rho}_i^{[l]} \langle \hat{\delta}^{[l]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}}\end{aligned}$$

$$\begin{aligned}
&= \langle W_i, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} - \sum_{u=1}^l \hat{\rho}_i^{[u]} \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} \\
&\stackrel{(i)}{=} \langle W_i, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} - \sum_{u=1}^l \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} \langle W_i, \hat{t}^{[u]} \rangle_{\mathbb{H}} \\
&= \langle W_i, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} - \left\langle W_i, \sum_{u=1}^l \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} \hat{t}^{[u]} \right\rangle_{\mathbb{H}} \\
&= \left\langle W_i, \hat{\xi}^{[l+1]} - \sum_{u=1}^l \langle \hat{\delta}^{[u]}, \hat{\xi}^{[l+1]} \rangle_{\mathbb{H}} \hat{t}^{[u]} \right\rangle_{\mathbb{H}} \\
&= \hat{\rho}_i^{[l+1]}.
\end{aligned}$$

where step (i) uses the inductive hypothesis $\langle W_i, \hat{t}^{[u]} \rangle_{\mathbb{H}} = \hat{\rho}_i^{[u]}$ for all $u \leq l$. By the principle of mathematical induction, the lemma holds for all $l \geq 1$. This completes the proof of Lemma 3.8. \square

APPENDIX E: PROOF OF SECTION 4.1

This section provides proof of Theorem 4.4, demonstrating that the population-level version of the optimization problem solved at each step is well-defined, and proofs of the intermediate lemmas that support the main theorem.

E.1. Proof of Lemma 4.1 . *(JM:) An alternative proof to the compactness of \mathcal{U} ; rather than directly showing \mathcal{U} is compact, we show \mathcal{C}_{YW} is compact*

PROOF. Any bounded linear functional from a Hilbert space to \mathbb{R} is a compact operator. This is because the image of any bounded set under such a functional is a bounded set in \mathbb{R} , and by Bolzano-Weierstrass Theorem, every bounded sequence of real numbers has a convergent subsequence. Therefore, to show that \mathcal{C}_{YW} is a compact operator, it suffices to show that \mathcal{C}_{YW} is a bounded linear functional.

Linearity. The operator $\mathcal{C}_{YW} : \mathbb{H} \rightarrow \mathbb{R}$ is defined as $\mathcal{C}_{YW}h = \langle \Sigma_{YW}, h \rangle_{\mathbb{H}}$. For any $h_1, h_2 \in \mathbb{H}$ and scalar $c \in \mathbb{R}$, the linearity of the inner product implies:

$$\begin{aligned}
\mathcal{C}_{YW}(h_1 + h_2) &= \langle \Sigma_{YW}, h_1 + h_2 \rangle_{\mathbb{H}} = \langle \Sigma_{YW}, h_1 \rangle_{\mathbb{H}} + \langle \Sigma_{YW}, h_2 \rangle_{\mathbb{H}} = \mathcal{C}_{YW}h_1 + \mathcal{C}_{YW}h_2, \\
\mathcal{C}_{YW}(ch) &= \langle \Sigma_{YW}, ch \rangle_{\mathbb{H}} = c \langle \Sigma_{YW}, h \rangle_{\mathbb{H}} = c \mathcal{C}_{YW}h.
\end{aligned}$$

Thus, \mathcal{C}_{YW} is a linear functional.

Boundedness. To show that \mathcal{C}_{YW} is bounded, we need to find a finite constant M such that $|\mathcal{C}_{YW}h| \leq M \|h\|_{\mathbb{H}}$ for all $h \in \mathbb{H}$. By the Cauchy-Schwarz inequality, we have:

$$|\mathcal{C}_{YW}h| = |\langle \Sigma_{YW}, h \rangle_{\mathbb{H}}| \leq \|\Sigma_{YW}\|_{\mathbb{H}} \|h\|_{\mathbb{H}}.$$

Now, we must verify that $\|\Sigma_{YW}\|_{\mathbb{H}}$ is finite. Recall that

$$\Sigma_{YW} = (\sigma_{YX,1}(t), \dots, \sigma_{YX,K}(t), \sigma_{YZ,1}, \dots, \sigma_{YZ,p}).$$

Let $\mu([0, 1])$ denote a Lebesgue measure of $[0, 1]$, and let $T = \max_{k=1, \dots, K} \mu([0, 1])$. The norm of Σ_{YW} in \mathbb{H} is given by:

$$\|\Sigma_{YW}\|_{\mathbb{H}}^2 = \sum_{k=1}^K \int_0^1 \sigma_{YX,k}(t)^2 dt + \sum_{r=1}^p \sigma_{YZ,r}^2 < KTQ_1 + pQ_2 < \infty,$$

where the last inequality uses the condition (24). Let $M = \sqrt{KTQ_1 + pQ_2}$. Thus, we have shown that $|\mathcal{C}_{YW}h| \leq M\|h\|_{\mathbb{H}}$ for a finite constant M . This completes the proof of Lemma 4.1. \square

E.2. Proof of Lemma 4.2.

PROOF. For $h \in \mathbb{H}$ and $d \in \mathbb{R}$, we have:

$$\langle \mathcal{C}_{YW}h, d \rangle = \mathbb{E}[\langle W_1, h \rangle_{\mathbb{H}} Y_1] d = \mathbb{E}\langle Y_1 W_1 d, h \rangle_{\mathbb{H}} = \langle h, \mathbb{E}[Y_1 W_1 d] \rangle_{\mathbb{H}} = \langle h, \mathcal{C}_{WY}d \rangle_{\mathbb{H}}$$

Thus, we have $\langle \mathcal{C}_{YW}h, d \rangle = \langle h, \mathcal{C}_{WY}d \rangle_{\mathbb{H}}$, which implies $\mathcal{C}_{WY} = \mathcal{C}_{YW}^*$. This completes the proof of Lemma 4.2. \square

E.3. Proof of Lemma 4.3.

PROOF. We show that \mathcal{U} is self-adjoint, positive-semidefinite, and compact, in turn.

Self-adjoint. For any $h_1, h_2 \in \mathbb{H}$, we have

$$\begin{aligned} \langle \mathcal{U}h_1, h_2 \rangle_{\mathbb{H}} &= \langle \langle h_1, \Sigma_{YW} \rangle_{\mathbb{H}} \Sigma_{YW}, h_2 \rangle_{\mathbb{H}} = \langle h_1, \Sigma_{YW} \rangle_{\mathbb{H}} \langle \Sigma_{YW}, h_2 \rangle_{\mathbb{H}} \\ &= \langle h_1, \langle \Sigma_{YW}, h_2 \rangle_{\mathbb{H}} \Sigma_{YW} \rangle_{\mathbb{H}} \\ &= \langle h_1, \mathcal{U}h_2 \rangle_{\mathbb{H}}. \end{aligned}$$

Positive-semidefinite. For every $h \in \mathbb{H}$, we have

$$\langle \mathcal{U}h, h \rangle_{\mathbb{H}} = \langle \langle \Sigma_{YW}, h \rangle_{\mathbb{H}} \Sigma_{YW}, h \rangle_{\mathbb{H}} = \langle \Sigma_{YW}, h \rangle_{\mathbb{H}}^2 \geq 0.$$

Compact. By Lemma 4.1, \mathcal{C}_{YW} is a compact operator. By Theorem 4.1.3 of Hsing and Eubank (2015), the composition of two operators is compact if either operator is compact. Therefore $\mathcal{U} := \mathcal{C}_{WY} \circ \mathcal{C}_{YW}$ is a compact operator. This completes the proof of Lemma 4.3. \square

E.4. Proof of Theorem 4.4.

PROOF. by Theorem 4.3.1 of Hsing and Eubank (2015), since \mathcal{C}_{YW} is compact by Lemma 4.1, it has the singular value decomposition of \mathcal{C}_{YW} , given by:

$$\mathcal{C}_{YW} = \sum_{j=1}^{\infty} \iota_j (f_{1j} \otimes f_{2j}).$$

Let $\|\cdot\|_{\text{op}}$ denote an operator norm. By 4.3.4 in Hsing and Eubank (2015), we have

$$\|\mathcal{C}_{YW}\|_{\text{op}} = \sup_{\substack{h \in \mathbb{H} \\ \|h\|_{\mathbb{H}}=1}} |\mathcal{C}_{YW}h|^2 = \sup_{\substack{h \in \mathbb{H} \\ \|h\|_{\mathbb{H}}=1}} |E(\langle W, h \rangle_{\mathbb{H}} Y)|^2 = \sup_{\substack{h \in \mathbb{H} \\ \|h\|_{\mathbb{H}}=1}} \text{Cov}^2(\langle W, h \rangle_{\mathbb{H}}, Y) = \kappa_1^2,$$

with maximum attained at $h = f_{11}$, which is an eigenfunction of $\mathcal{C}_{YW}^* \circ \mathcal{C}_{YW} = \mathcal{C}_{WY} \circ \mathcal{C}_{YW} = \mathcal{U}$ corresponding to the largest eigenvalue κ_1^2 . This completes the proof of Theorem 4.4. \square

E.5. Proof of Proposition 4.7.

PROOF. The unit norm condition is trivially met by the constraint $1 = \hat{\xi}_l^\top (\mathbb{B} + \Lambda \mathbb{B}'') \hat{\xi}_l$ enforced in (19), because we have:

$$\begin{aligned}
 \hat{\xi}_l^\top (\mathbb{B} + \Lambda \mathbb{B}'') \hat{\xi}_l &= \hat{\xi}_l^\top \mathbb{B} \hat{\xi}_l + \hat{\xi}_l^\top \Lambda \mathbb{B}'' \hat{\xi}_l \\
 &\stackrel{(i)}{=} \sum_{j=1}^K \int_0^1 \hat{\xi}_l(t) \hat{\xi}_l(t) dt + \hat{\zeta}^\top \hat{\zeta} + \hat{\xi}^\top \Lambda \mathbb{B}'' \hat{\xi} \\
 &\stackrel{(ii)}{=} \sum_{j=1}^K \int_0^1 \hat{\xi}_l(t) \hat{\xi}_l(t) dt + \hat{\zeta}^\top \hat{\zeta} + \sum_{j=1}^K \lambda_l \int_0^1 \{\hat{\xi}_l''(t)\}^2 dt \\
 (34) \quad &= \langle \hat{\xi}_l, \hat{\xi}_l \rangle_{\mathbb{H}, \Lambda},
 \end{aligned}$$

where step (i) uses (27) and step (ii) uses (28).

Now to switch our gears toward the the orthogonality. For $l_1 > l_2$, we have

$$\begin{aligned}
 \langle \hat{\xi}_{l_1}, \hat{\xi}_{l_2} \rangle_{\mathbb{H}, \Lambda} &\stackrel{(i)}{=} \hat{\xi}_{l_1}^\top (\mathbb{B} + \Lambda \mathbb{B}'') \hat{\xi}_{l_2} \\
 &\stackrel{(ii)}{=} \frac{1}{\kappa_{l_1}} (V^{[l_1]} \hat{\xi}_{l_1})^\top \hat{\xi}_{l_2} \\
 &= \frac{1}{\kappa_{l_1}} \hat{\xi}_{l_1}^\top V^{[l_1]} \hat{\xi}_{l_2} \\
 &\stackrel{(iii)}{=} \frac{1}{n^2 \kappa_{l_1}} \hat{\xi}_{l_1}^\top (\mathbb{B} \Theta^{[l_1]^\top} \mathbf{y}) (\mathbb{B} \Theta^{[l_1]^\top} \mathbf{y})^\top \hat{\xi}_{l_2} \\
 &= \frac{1}{n^2 \kappa_{l_1}} \hat{\xi}_{l_1}^\top (\mathbb{B} \Theta^{[l_1]^\top} \mathbf{y}) \mathbf{y}^\top \Theta^{[l_1]} \mathbb{B} \hat{\xi}_{l_2} \\
 &= c \mathbf{y}^\top (\Theta^{[l_1]} \mathbb{B} \hat{\xi}_{l_2}),
 \end{aligned}$$

where step (i) uses (34), step (ii) uses the fact that $\hat{\xi}_{l_1}$ is a generalized eigenvector (κ_{l_1} denotes the corresponding generalized eigenvalue), as presented in Proposition 3.5, and step (iii) uses the definition of V matrix given in (15). Here, c represents a scalar that condenses all multiplicative terms preceding \mathbf{y}^\top .

Orthogonality can be demonstrated by showing that $\Theta^{[l_1]} \mathbb{B} \hat{\xi}_{l_2} = \mathbf{0} \in \mathbb{R}^n$. From the construction in (8), (7) and (10), the i th entry of $\Theta^{[l_1]} \mathbb{B} \hat{\xi}_{l_2}$ is

$$\begin{aligned}
 (\Theta^{[l_1]} \mathbb{B} \hat{\xi}_{l_2})_i &= (\theta_{i11}^{[l_1]}, \dots, \theta_{i1M}^{[l_1]}, \dots, \theta_{iK1}^{[l_1]}, \dots, \theta_{iKM}^{[l_1]}, Z_{i1}^{[l_1]}, \dots, Z_{ip}^{[l_1]}) \mathbb{B} \hat{\xi}_{l_2} \\
 &= \sum_{k=1}^K \theta_{ik}^{[l_1]} B \hat{\gamma}_k^{[l_2]} + \mathbf{z}_i^{[l_1]^\top} \hat{\zeta}^{[l_2]} \\
 (35) \quad &= \langle \widetilde{W}_i^{[l_1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}},
 \end{aligned}$$

where the last equality uses the computation of the i th entry in (31). To expand the last quantity, we use a recursive relationship derived as follows. Recall from Lemma 3.7 that we have

$$(36) \quad \widetilde{W}_i^{[l_1]} := \widetilde{W}_i^{[l_1-1]} - \frac{\hat{\rho}_{i l_1-1}}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \sum_{i=1}^n \hat{\rho}_{i l_1-1} \widetilde{W}_i^{[l_1-1]}.$$

Let us denote, with some abuse of notation:

$$(37) \quad \widetilde{\mathbf{W}}^{[l_1]} := (\widetilde{W}_1^{[l_1-1]}, \dots, \widetilde{W}_n^{[l_1-1]})^\top \in \widetilde{\mathbb{H}}^{\otimes n},$$

and

$$(38) \quad \langle \widetilde{\mathbf{W}}^{[l_1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} := \left(\langle \widetilde{W}_1^{[l_1-1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}}, \dots, \langle \widetilde{W}_n^{[l_1-1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} \right)^\top \in \mathbb{R}^n.$$

Recalling (35), to achieve our goal of showing $\Theta^{[l_1]} \mathbb{B} \hat{\xi}_{l_2} = \mathbf{0} \in \mathbb{R}^n$, it suffices to show $\langle \widetilde{\mathbf{W}}^{[l_1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} = \mathbf{0}$. Using these notations and (36), we have, with some abuse of notation,

$$(39) \quad \widetilde{\mathbf{W}}^{[l_1]} = \widetilde{\mathbf{W}}^{[l_1-1]} - \left(\frac{1}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \hat{\rho}^{[l_1-1]} \hat{\rho}^{[l_1-1]\top} \right) \widetilde{\mathbf{W}}^{[l_1-1]},$$

and thus

$$\langle \widetilde{\mathbf{W}}^{[l_1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} = \langle \widetilde{\mathbf{W}}^{[l_1-1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} - \left(\frac{1}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \hat{\rho}^{[l_1-1]} \hat{\rho}^{[l_1-1]\top} \right) \langle \widetilde{\mathbf{W}}^{[l_1-1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}}.$$

Clearly, the right-hand side represents a linear operator from \mathbb{R}^n to \mathbb{R}^n , which we denote as $P^{[l_1-1]}$. Thus, we have $\langle \widetilde{\mathbf{W}}^{[l_1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} = P^{[l_1-1]} \langle \widetilde{\mathbf{W}}^{[l_1-1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}}$. Repeatedly using this relationship, we have

$$\begin{aligned} \langle \widetilde{\mathbf{W}}^{[l_1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} &= P^{[l_1-1]} \langle \widetilde{\mathbf{W}}^{[l_1-1]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} \\ &= P^{[l_1-1]} P^{[l_1-2]} \langle \widetilde{\mathbf{W}}^{[l_1-2]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} \\ &= \underbrace{P^{[l_1-1]} P^{[l_1-2]} \dots P^{[l_2+1]}}_{:=P} P^{[l_2]} \langle \widetilde{\mathbf{W}}^{[l_2]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} \\ &= P \left\{ \langle \widetilde{\mathbf{W}}^{[l_2]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} - \left(\frac{1}{\|\hat{\rho}^{[l_2]}\|_2^2} \hat{\rho}^{[l_2]} \hat{\rho}^{[l_2]\top} \right) \langle \widetilde{\mathbf{W}}^{[l_2]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} \right\} \\ &= P \left\{ \hat{\rho}^{[l_2]} - \left(\frac{1}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \hat{\rho}^{[l_2]} \hat{\rho}^{[l_2]\top} \right) \hat{\rho}^{[l_2]} \right\} \\ &= P \left\{ \hat{\rho}^{[l_2]} - \left(\frac{\hat{\rho}^{[l_2]\top} \hat{\rho}^{[l_2]}}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \right) \hat{\rho}^{[l_2]} \right\} \\ &= P \left\{ \hat{\rho}^{[l_2]} - \hat{\rho}^{[l_2]} \right\} = P \mathbf{0} = \mathbf{0}. \end{aligned}$$

This completes the proof of Proposition 4.7. □

E.6. Proof of Proposition 4.8.

PROOF. For $l_1 < l_2$, we have

$$\begin{aligned} \hat{\rho}^{[l_1]\top} \hat{\rho}^{[l_2]} &= \sum_{i=1}^n \hat{\rho}_i^{[l_1]} \hat{\rho}_i^{[l_2]} \\ &= \sum_{i=1}^n \hat{\rho}_i^{[l_1]} \langle \widetilde{W}_i^{[l_2]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}} \\ &= \left\langle \sum_{i=1}^n \hat{\rho}_i^{[l_1]} \widetilde{W}_i^{[l_2]}, \hat{\xi}^{[l_2]} \right\rangle_{\mathbb{H}} \end{aligned}$$

Therefore, to show that $\hat{\rho}^{[l_1]\top} \hat{\rho}^{[l_2]} = 0$, it suffices to show that $\sum_{i=1}^n \hat{\rho}_i^{[l_1]} \widetilde{W}_i^{[l_2]} = 0 \in \mathbb{H}$, where this zero element represents an ordered pair of K zero functions and a zero matrix. Using Lemma 3.7, notations from (37) and (38), and equation (39), we have:

$$\begin{aligned} \widetilde{W}^{[l_1]} &= \widetilde{W}^{[l_1-1]} - \left(\frac{1}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \hat{\rho}^{[l_1-1]} \hat{\rho}^{[l_1-1]\top} \right) \widetilde{W}^{[l_1-1]} \\ &= \widetilde{W}^{[l_1-1]} - \left(\frac{1}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \langle \widetilde{W}^{[l_1-1]}, \hat{\xi}^{[l_1-1]} \rangle_{\mathbb{H}} \hat{\rho}^{[l_1-1]\top} \right) \widetilde{W}^{[l_1-1]}. \end{aligned}$$

Using this relationship and with some abuse of notation, we have:

$$\begin{aligned} \sum_{i=1}^n \hat{\rho}_i^{[l_1]} \widetilde{W}_i^{[l_2]} &= \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1]} \\ &= \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1-1]} - \hat{\rho}^{[l_2]\top} \left(\frac{1}{\|\hat{\rho}^{[l_1-1]}\|_2^2} \langle \widetilde{W}^{[l_1-1]}, \hat{\xi}^{[l_1-1]} \rangle_{\mathbb{H}} \hat{\rho}^{[l_1-1]\top} \right) \widetilde{W}^{[l_1-1]} \\ &= \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1-1]} - \frac{\langle \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1-1]}, \hat{\xi}^{[l_1-1]} \rangle_{\mathbb{H}}}{\|\hat{\rho}^{[l_1-1]}\|_2^2} (\hat{\rho}^{[l_1-1]\top} \widetilde{W}^{[l_1-1]}) \\ &= h^{[l_1-1]} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1-1]}) \text{ (say)}. \end{aligned}$$

Here, the function $h^{[l_1-1]} : \mathbb{H} \mapsto \mathbb{H}$ maps the zero element of \mathbb{H} to itself, where the zero element represents an ordered pair of K zero functions and a zero matrix. Using this relationship, we have

$$\begin{aligned} \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1]} &= h^{[l_1-1]} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1-1]}) \\ &= h^{[l_1-1]} \circ h^{[l_1-2]} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_1-2]}) \\ &= h^{[l_1-1]} \circ h^{[l_1-2]} \circ \dots \circ h^{[l_2]} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]}). \end{aligned}$$

Therefore, to show that $\hat{\rho}^{[l_1]\top} \hat{\rho}^{[l_2]} = 0$, it suffices to show $h^{[l_2]} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]}) = 0$:

$$\begin{aligned} h^{[l_2]} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]}) &= \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]} - \frac{\langle \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}}}{\|\hat{\rho}^{[l_2]}\|_2^2} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]}) \\ &= \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]} - \frac{\hat{\rho}^{[l_2]\top} \langle \widetilde{W}^{[l_2]}, \hat{\xi}^{[l_2]} \rangle_{\mathbb{H}}}{\|\hat{\rho}^{[l_2]}\|_2^2} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]}) \\ &= \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]} - \frac{\hat{\rho}^{[l_2]\top} \hat{\rho}^{[l_2]}}{\|\hat{\rho}^{[l_2]}\|_2^2} (\hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]}) \\ &= \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]} - \hat{\rho}^{[l_2]\top} \widetilde{W}^{[l_2]} \\ &= 0. \end{aligned}$$

This completes the proof of Proposition 4.8. □