
Errors-in-variables Fréchet Regression with Low-rank Covariate Approximation

Anonymous Author(s)

Affiliation

Address

email

Abstract

Fréchet regression has emerged as a promising approach for regression analysis involving non-Euclidean response variables. However, its practical applicability has been hindered by its reliance on ideal scenarios with abundant and noiseless covariate data. In this paper, we present a novel estimation method that tackles these limitations by leveraging the low-rank structure inherent in the covariate matrix. Our proposed framework combines the concepts of global Fréchet regression and principal component regression, aiming to improve the efficiency and accuracy of the regression estimator. By incorporating the low-rank structure, our method enables more effective modeling and estimation, particularly in high-dimensional and errors-in-variables regression settings. We provide a theoretical analysis of the proposed estimator's large-sample properties, including a comprehensive rate analysis of bias, variance, and additional variations due to measurement errors. Furthermore, our numerical experiments provide empirical evidence that supports the theoretical findings, demonstrating the superior performance of our approach. Overall, this work introduces a promising framework for regression analysis of non-Euclidean variables, effectively addressing the challenges associated with limited and noisy covariate data, with potential applications in diverse fields.

1 Introduction

Regression analysis is a fundamental statistical methodology to model the relationship between response variables and explanatory variables (covariates). Linear regression, for example, models the (conditional) expected value of the response variable as a linear function of covariates. Regression models enable researchers and analysts to make predictions, gain insights into how input variables influence the outcomes of interest, and validate hypothetical associations between variables in inferential studies. As a result, regression is widely utilized across various scientific domains, including economics, psychology, biology, and engineering [31, 21, 29].

In recent decades, there has been a growing interest in developing statistical methods capable of handling random objects in non-Euclidean spaces. Examples of these include functional data analysis [42], statistical manifold learning [32], statistical network analysis [35], and object-oriented data analysis [40]. In such contexts, the response variable is defined in a metric space that may lack an algebraic structure, making it challenging to apply global, parametric approaches toward regression as in the classical Euclidean setting. To overcome this challenge, (global) Fréchet regression, which models the relationship by fitting the (conditional) barycenters of the responses as a function of covariates, has been introduced [41]. Notably, when the Euclidean metric is considered, Fréchet regression recovers classical Euclidean regression models. For more details on Fréchet regression and its recent developments, we refer readers to [30, 41, 22, 46, 27].

Nevertheless, most existing research on Fréchet regression has focused on ideal scenarios characterized by abundant covariate data that are accurately measured and free of noise. In practical applications, however, high-dimensional data often arise, which are also susceptible to measurement errors and other forms of contamination. These errors can stem from various sources, such as unreliable data collection methods (*e.g.*, low-resolution probes, subjective self-reports) or imperfect data storage and transmission. The high-dimensionality and the presence of measurement errors in covariates pose critical challenges for statistical inference, as regression analysis based on error-prone covariates may result in incorrect associations between variables, yielding misleading conclusions.

To address these limitations, it is crucial to extend the methodology and analysis of Fréchet regression to tackle high-dimensional errors-in-variables problems. In this work, we aim to leverage the low-rank structure in the covariates to enhance the estimation accuracy and computational efficiency of Fréchet regression. Specifically, we explore the extension of principal component regression to handle errors-in-variables regression problems with non-Euclidean response variables.

1.1 Contributions

This paper contributes to advancing the (global) Fréchet regression of non-Euclidean response variables, with a particular focus on high-dimensional, errors-in-variables regression.

Firstly, we propose a novel framework, called the regularized (global) Fréchet regression (Section 3) that combines the ideas from Fréchet regression [41] and the principal component regression [33]. This framework effectively utilizes the low-rank structure in the matrix of (Euclidean) covariates by extracting its principal components via low-rank matrix approximation. Our proposed method is straightforward to implement, not requiring any knowledge about the error-generating mechanism.

Furthermore, we provide a comprehensive theoretical analysis in three main theorems (Section 4) to establish the effectiveness of the proposed framework. Firstly, we prove the consistency of the proposed estimator for the true global Fréchet regression model (Theorem 1). Secondly, we investigate the convergence rate of the estimator’s bias and variance (Theorem 2). Lastly, we derive an upper bound for the distance between the estimates obtained using error-free covariates and those with errors-in-variables covariates (Theorem 3). Collectively, these results demonstrate that our approach effectively addresses model mis-specification and achieve more efficient model estimation by incorporating the low-rank structure of covariates, despite the presence of inherent bias due to unobserved measurement errors.

To validate our theoretical findings, we conduct numerical experiments on simulated datasets. Our results demonstrate that the proposed method provides more accurate estimates of the regression parameters, especially in high-dimensional settings. Our experiments emphasize the importance of incorporating the low-rank structure of covariates in Fréchet regression, and provide empirical evidence that aligns with our theoretical analysis.

1.2 Related work

Metric space-valued variables. Nonparametric regression models for Riemannian manifold-valued responses were proposed as a generalization of regression for multivariate outputs by Steinke *et al.* [49, 50]. These works provided a foundation for recent developments in regression analysis of non-Euclidean responses. Later, Hein [30] proposed a Nadaraya-Watson-type kernel estimation of regression model for general metric-space-valued outcomes. Since then, statistical properties of regression models for some special classes of metric space-valued outcomes, such as distribution functions [23, 53, 28] and matrix-valued responses [57, 20], have been investigated. Recently, many researchers have introduced further advances in Fréchet regression, including [41, 10, 38, 46]. In this study, we use the global Fréchet regression proposed by [41] as the basis for our proposed method.

Errors-in-variables regression. Much of earlier work on errors-in-variables problems in the statistical literature can be found in [13], which covers the simulation-extrapolation (SIMEX) [16, 11], the attenuation correction method [37], covariate-adjusted model [47, 19], and the deconvolution kernel method [25, 24, 18]. The regression calibration method [48], instrumental variable modeling [12, 44], and the two-phase study design [9, 4] were also proposed when additional data are available for correcting measurement errors. In the high-dimensional modeling literature, regularization

87 methods for recovering the true covariate structure can also be utilized [39, 7, 17]. However, most of
 88 these methods require prior knowledge about the measurement error distributions.

89 **Principal component regression.** The principal component regression (PCR) [33] is a statistical
 90 technique that regresses response variables on principal component scores of the covariate matrix.
 91 The conventional PCR selects a few principal components as the “new” regressors associated with
 92 the first leading eigenvalues to explain the highest proportion of variations observed in the original
 93 covariate matrix. In functional data analysis, PCR is known to have a shrinkage effect on the model
 94 estimate and produce robust prediction performance in functional regression [43, 34]. Recently,
 95 Agarwal *et al.* [2] investigated the robustness of PCR in the presence of measurement errors on
 96 covariates and the statistical guarantees for learning a good predictive model. Motivated by these
 97 findings, we will adopt the PCR framework to improve the estimation and prediction performance of
 98 the errors-in-variables Fréchet regression. in this study.

99 1.3 Organization

100 This paper is organized as follows. In Section 2, we introduce the notation used throughout the
 101 paper, and overview the global Fréchet regression framework. Section 3 presents the problem setup,
 102 objectives, and our proposed estimator, which we refer to as the regularized Fréchet regression
 103 (Definition 4). In Section 4, we discuss theoretical guarantees on the regularized Fréchet regression
 104 method in accurately estimating the global Fréchet regression function. Section 5 presents the
 105 results of numerical ‘proof-of-concept’ experiments that support the theoretical findings. Finally, we
 106 conclude this paper with discussions in Section 6. Due to space constraints, detailed proofs of the
 107 theorems as well as additional details and discussions of experiments are provided in the Appendix.

108 2 Preliminaries

109 2.1 Notation

110 Let \mathbb{N} denote the set of positive integers and \mathbb{R} denote the set of real numbers. Also, let $\mathbb{R}_+ :=$
 111 $\{x \in \mathbb{R} : x \geq 0\}$. For $n \in \mathbb{N}$, we let $[n] := \{1, \dots, n\}$. We mostly use plain letters to denote
 112 scalars, vectors, and random variables, but we also use boldface uppercase letters for matrices, and
 113 curly letters to denote sets when useful. Note that we may identify a vector with its column matrix
 114 representation. For a matrix \mathbf{X} , we let \mathbf{X}^{-1} denote its inverse (if exists) and \mathbf{X}^\dagger denote the Moore-
 115 Penrose pseudoinverse of \mathbf{X} . Also, we let $\text{rowsp}(\mathbf{X})$ and $\text{colsp}(\mathbf{X})$ denote the row and column
 116 spaces of \mathbf{X} , respectively. Furthermore, we let $\text{spec}(\mathbf{X})$ denote the set of non-zero singular values
 117 of \mathbf{X} , $\sigma_i(\mathbf{X})$ denote the i -th largest singular value of \mathbf{X} , and $\sigma^{(\lambda)}(\mathbf{X}) := \inf\{\sigma_i(\mathbf{X}) > \lambda : i \in \mathbb{N}\}$
 118 with the convention $\inf \emptyset = \infty$. We let $\mathbf{1}_n = (1, 1, \dots, 1) \in \mathbb{R}^d$ and let $\mathbb{1}$ denote the indicator
 119 function. We let $\|\cdot\|$ denote a norm, and set $\|\cdot\| = \|\cdot\|_2$ (the ℓ_2 -norm for vectors, and the spectral
 120 norm for matrices) by default, unless stated otherwise. For a finite set \mathcal{D} , we may identify \mathcal{D} with its
 121 empirical measure $\nu_{\mathcal{D}}^{\text{emp}} = \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \delta_x$, where δ_x denotes the Dirac measure supported on $\{x\}$.

122 Letting $f, g : \mathbb{R} \rightarrow \mathbb{R}$, we write $f(x) = O(g(x))$ as $x \rightarrow \infty$ if there exist $M > 0$ and $x_0 > 0$ such
 123 that $|f(x)| \leq M \cdot g(x)$ for all $x \geq x_0$. Likewise, we write $f(x) = \Omega(g(x))$ if $g(x) = O(f(x))$.
 124 Furthermore, we write $f(x) = o(g(x))$ as $x \rightarrow \infty$ if $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = 0$. For a sequence of random
 125 variables X_n , and a sequence a_n , we write $X_n = O_p(a_n)$ as $n \rightarrow \infty$ if for any $\varepsilon > 0$, there exists
 126 $M \in \mathbb{R}_+$ and $N \in \mathbb{N}$ such that $P(|\frac{X_n}{a_n}| > M) < \varepsilon$ for all $n \geq N$. Similarly, we write $X_n = o_p(a_n)$
 127 if $\lim_{n \rightarrow \infty} P(|\frac{X_n}{a_n}| > \varepsilon) = 0$ for all $\varepsilon > 0$.

128 2.2 Global Fréchet regression

129 Let (X, Y) be a random variable that has a joint distribution $F_{X,Y}$ supported on $\mathbb{R}^p \times \mathcal{M}$, where \mathbb{R}^p
 130 is the p -dimensional Euclidean space and $\mathcal{M} = (\mathcal{M}, d)$ is a metric space equipped with a distance
 131 function $d : \mathcal{M} \times \mathcal{M} \rightarrow \mathbb{R}$. We write the marginal distribution of X as F_X , and the conditional
 132 distribution of Y given X as $F_{Y|X}$.

133 **Definition 1** (Fréchet regression function). *Let (X, Y) be a random element that takes value in*
 134 *$\mathbb{R}^p \times \mathcal{M}$. The Fréchet regression function of Y on X is a function $\varphi^* : \mathbb{R}^p \rightarrow \mathcal{M}$ such that*

$$\varphi^*(x) = \arg \min_{y \in \mathcal{M}} \mathbb{E}[d^2(Y, y) | X = x], \quad \forall x \in \text{supp } F_X \subseteq \mathbb{R}^p. \quad (1)$$

135 We note that $\varphi^*(x)$ is the best predictor of Y given $X = x$, as it minimizes the marginal risk
 136 $\mathbb{E}[d^2(Y, \varphi^*(X))]$ under the squared-distance loss. In the literature, $\varphi^*(x)$ is also known as the
 137 conditional Fréchet mean [26] of Y given $X = x$. It is important to recognize that the existence
 138 and uniqueness of the Fréchet regression function are closely tied to the geometric characteristics of
 139 \mathcal{M} , and are not guaranteed in general [3, 8]. Nonetheless, extensive research has been conducted
 140 on the existence and uniqueness of Fréchet means in various metric spaces commonly encountered
 141 in practical applications. Examples include the unit circle in \mathbb{R}^2 [14], Riemannian manifolds [1, 5],
 142 Alexandrov spaces with non-positive curvature [52], metric spaces with upper bounded curvature
 143 [58], and Wasserstein space [59, 36].

144 While modeling and estimating the Fréchet regression function φ^* is often of interest, its global
 145 (parametric) modeling may not be straightforward, especially when \mathcal{M} lacks a useful algebraic
 146 structure, such as an inner product. For instance, in classical linear regression analysis with $\mathcal{M} = \mathbb{R}$,
 147 the distribution of $(Y|X = x)$ is normally distributed with a mean of $\varphi^*(x) = \alpha + \beta^\top x$ and variance
 148 σ_Y^2 , where α and β represent the regression coefficients. Similarly, when \mathcal{M} possesses a linear-
 149 algebraic structure, one can specify a class of regression functions that quantifies the association
 150 between the expected outcome and covariates in an additive and multiplicative manner. However,
 151 the lack of an algebraic structure in general metric spaces may prevent us from characterizing the
 152 barycenter $\varphi^*(x)$ in the same way classical regression analysis determines the expected value of
 153 outcomes with changing covariates.

154 To address this challenge, Petersen and Müller [41] recently proposed to exploit algebraic structures
 155 in the space of covariates, \mathbb{R}^p , instead of \mathcal{M} . Specifically, they consider a weighted Fréchet mean as

$$\varphi(x) = \arg \min_{y \in \mathcal{M}} \mathbb{E}[w(X, x) \cdot d^2(Y, y)], \quad (2)$$

156 where $w : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ is an arbitrary weight function such that $w(\xi, x)$ denotes the influence of
 157 ξ at x . In particular, we define the global Fréchet regression function with a specific choice of w ,
 158 following [41].

159 **Definition 2** (Global Fréchet regression function). *Let (X, Y) be a random variable in $\mathbb{R}^p \times \mathcal{M}$.*
 160 *Let $\mu = \mathbb{E}(X)$ and $\Sigma = \text{Var}(X)$. The global Fréchet regression function of Y on X is a function*
 161 *$\varphi_{\text{glo}} : \mathbb{R}^p \rightarrow \mathcal{M}$ such that*

$$\varphi_{\text{glo}}(x) = \arg \min_{y \in \mathcal{M}} \mathbb{E}[w_{\text{glo}}(X, x) \cdot d^2(Y, y)] \quad (3)$$

162 where $w_{\text{glo}}(X, x) = 1 + (X - \mu)^\top \Sigma^{-1}(x - \mu)$.

163 Note that when \mathcal{M} is an inner product space (e.g., $\mathcal{M} = \mathbb{R}$), the function φ_{glo} restores the standard
 164 least squares linear regression. For this reason, φ_{glo} is commonly referred to as the global Fréchet
 165 regression model for metric-space-valued outcomes in recent literature [41, 38, 54].

166 **What does it mean by “global” and where does it come from?** One might wonder why the
 167 term “global” is used to describe φ_{glo} as a Fréchet regression function. The use of the adjective
 168 “global” serves to emphasize its distinction from “local” nonparametric regression methods that
 169 interpolate data points. Notably, when \mathcal{M} is a Hilbert space, φ_{glo} reduces to the natural linear models.
 170 For instance, if $\mathcal{M} = \mathbb{R}$, then it follows that $\varphi_{\text{glo}}(x) = \mathbb{E}[w_{\text{glo}}(X, x) \cdot Y] = \alpha + \beta^\top (x - \mu)$,
 171 where $\alpha = \mathbb{E}[Y]$ and $\beta = \Sigma^{-1} \cdot \mathbb{E}[(X - \mu) \cdot Y]$. These linear models hold uniformly for the
 172 evaluation point x . Similarly, in the case of an L^2 space equipped with the squared-distance metric
 173 $d^2(y, y') = \|y - y'\|_2^2$ induced by the L^2 norm, φ_{glo} represents the linear regression model for
 174 functional responses. Thus, φ_{glo} establishes a globally defined model that spans the entire space.

3 Problem and methodology

3.1 Problem formulation

Let (X, Y) be a random variable in $\mathbb{R}^p \times \mathcal{M}$ and $F_{X,Y}$ be their joint distribution. Let $\mathcal{D}_n = \{(X_i, Y_i) : i \in [n]\}$ be an independent and identically distributed (IID) sample drawn from $F_{X,Y}$. Note that we may identify the set \mathcal{D}_n with its discrete measure (empirical distribution). We consider the problem of estimating the global Fréchet regression function φ_{glo} (see Definition 2) from data \mathcal{D}_n . In this setting, a natural estimator of φ_{glo} would be its sample-analogue estimator. With $\hat{\mu}_{\mathcal{D}_n} = \mathbb{E}_{(X,Y) \sim \mathcal{D}_n}(X) = \frac{1}{n} \sum_{i=1}^n X_i$ and $\hat{\Sigma}_{\mathcal{D}_n} = \text{Var}_{(X,Y) \sim \mathcal{D}_n}(X) = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_{\mathcal{D}_n}) \cdot (X_i - \hat{\mu}_{\mathcal{D}_n})^\top$, the sample-analogue estimator $\hat{\varphi}_{\mathcal{D}_n}$ is defined as

$$\hat{\varphi}_{\mathcal{D}_n}(x) = \arg \min_{y \in \mathcal{M}} \left\{ \frac{1}{n} \sum_{(X_i, Y_i) \in \mathcal{D}_n} \hat{w}_{\mathcal{D}_n}(X_i, x) \cdot d^2(Y_i, y) \right\} \quad (4)$$

where $\hat{w}_{\mathcal{D}_n}(X, x) = 1 + (X - \hat{\mu}_{\mathcal{D}_n})^\top \hat{\Sigma}_{\mathcal{D}_n}^{-1} (x - \hat{\mu}_{\mathcal{D}_n})$. The statistical properties of $\hat{\varphi}_{\mathcal{D}_n}$, including the asymptotic distribution, a ridge-type variable selection operation, and total variation regularization method have been investigated [41, 38, 54].

In practice, however, we may only be able to access $\tilde{\mathcal{D}}_n = \{(Z_i, Y_i) : i \in [n]\}$ instead of \mathcal{D}_n , where

$$Z_i = X_i + \varepsilon_i, \quad i = 1, \dots, n \quad (5)$$

denotes an error-prone observation of the covariates X by measurement error ε . This formulation corresponds to the classical errors-in-variables problem.

Objective. Given a dataset, either \mathcal{D}_n or $\tilde{\mathcal{D}}_n$, our aim is to produce an estimate $\hat{\varphi}$ of the global Fréchet regression function φ_{glo} so that the prediction error is minimized. Specifically, we evaluate the performance of $\hat{\varphi}$ by means of the distance in the response space, $d(\hat{\varphi}(x), \varphi_{\text{glo}}(x))$.

3.2 Fréchet regression with covariate principal components

Singular value thresholding. Among various low-rank matrix approximation methods, we consider the (hard) singular value thresholding (SVT). For any $\lambda \in \mathbb{R}_+$, we define the map $\text{SVT}^{(\lambda)} : \mathbb{R}^{n \times p} \rightarrow \mathbb{R}^{n \times p}$ that removes all singular values that are less than the threshold λ . To be precise, $\text{SVT}^{(\lambda)}$ can be expressed in terms of the singular value decomposition (SVD) as follows:

$$M = \sum_{i=1}^{\min\{n,p\}} s_i \cdot u_i v_i^\top \text{ is a SVD} \implies \text{SVT}^{(\lambda)}(M) = \sum_{i=1}^{\min\{n,p\}} s_i \cdot \mathbb{1}\{s_i > \lambda\} \cdot u_i v_i^\top. \quad (6)$$

Regularized Fréchet regression. We introduce a variant of the sample-analog estimator of the global Fréchet regression function based on principal components of the sample covariance. To facilitate the description of our proposed estimator, we introduce additional notation here.

Definition 3 (Covariate mean/covariance). *For a probability distribution ν on $\mathbb{R}^p \times \mathcal{M}$, the covariate mean of ν , denoted by μ_ν , and the covariate covariance of ν , denoted by Σ_ν , are defined as*

$$\mu_\nu = \mathbb{E}_{(X,Y) \sim \nu}(X) \quad \text{and} \quad \Sigma_\nu = \text{Var}_{(X,Y) \sim \nu}(X). \quad (7)$$

Recall that a finite set $\mathcal{D} \subset \mathbb{R}^p \times \mathcal{M}$ may be identified with its empirical distribution; it follows that

$$\mu_{\mathcal{D}} = \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} x_i \quad \text{and} \quad \Sigma_{\mathcal{D}} = \frac{1}{|\mathcal{D}|} \sum_{(x_i, y_i) \in \mathcal{D}} (x_i - \mu_{\mathcal{D}}) \cdot (x_i - \mu_{\mathcal{D}})^\top. \quad (8)$$

Definition 4 (Regularized Fréchet regression). *Let ν be a probability distribution on $\mathbb{R}^p \times \mathcal{M}$ and $\lambda \in \mathbb{R}_+$. The λ -regularized Fréchet regression function for ν is a map $\varphi_\nu^{(\lambda)} : \mathbb{R}^p \rightarrow \mathcal{M}$ such that*

$$\varphi_\nu^{(\lambda)}(x) = \arg \min_{y \in \mathcal{M}} R_\nu^{(\lambda)}(y; x), \quad \text{where} \quad R_\nu^{(\lambda)}(y; x) = \mathbb{E}_{(X,Y) \sim \nu} \left[w_\nu^{(\lambda)}(X, x) \cdot d^2(Y, y) \right]$$

$$\text{and} \quad w_\nu^{(\lambda)}(x', x) = 1 + (x' - \mu_\nu)^\top \left[\text{SVT}^{(\lambda)}(\Sigma_\nu) \right]^\dagger (x - \mu_\nu). \quad (9)$$

When $\mathcal{D}_n = \{(X_i, Y_i) \in \mathbb{R}^p \times \mathcal{M} : i \in [n]\}$ is an IID sample from $F_{X,Y}$, the λ -regularized estimator $\varphi_{\mathcal{D}_n}^{(\lambda)}$ subsumes the sample-analogue estimator $\hat{\varphi}_{\mathcal{D}_n}$ in (4) as a special case where $\lambda = 0$.

Connection to principal component regression. Here we remark that when \mathcal{M} is a Euclidean space, the regularized Fréchet regression function $\varphi_{\nu}^{(\lambda)}$ effectively reduces to the principal component regression. Suppose that $\mathcal{M} = \mathbb{R}$ and $\mathcal{D}_n = \{(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R} : i \in [n]\}$ is a given dataset. Then $\varphi_{\mathcal{D}_n}^{(\lambda)}(x) = \bar{y} + \hat{\beta}_\lambda^\top (x - \mu_{\mathcal{D}_n})$ where $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and $\hat{\beta}_\lambda = [\text{SVT}^{(\lambda)}(\Sigma_{\mathcal{D}_n})]^\dagger \cdot [\frac{1}{n} \sum_{i=1}^n (x_i - \mu_{\mathcal{D}_n}) \cdot (y_i - \bar{y})]$. Observe that $\hat{\beta}_\lambda$ is exactly the regression coefficient of principal component regression applied to the centered dataset $\mathcal{D}_n^{\text{ctr}} = \{(x_i - \mu_{\mathcal{D}_n}, y_i - \bar{y}) : i \in [n]\}$ using k principal components where $k = \max_{a \in [p]} \{\sigma_a(\Sigma_{\mathcal{D}_n^{\text{ctr}}}) \geq \lambda\}$.

4 Main results

In this section, we investigate properties of $\varphi_{\nu}^{(\lambda)}$ for $\lambda \geq 0$, with a focus on two cases: $\nu = \mathcal{D}_n$ and $\nu = \tilde{\mathcal{D}}_n$, cf. Section 3.1. By denoting the true distribution that generates (X, Y) as ν^* , we can express φ_{glo} as $\varphi_{\nu^*}^{(0)}$. To analyze the discrepancy between the regularized global Fréchet regression estimators and $\varphi_{\text{glo}}(x)$, we examine the relationships depicted in the schematic in Figure 1. Our theoretical findings can be summarized as follows: even in the presence of covariate noises, $\varphi_{\tilde{\mathcal{D}}_n}^{(\lambda)}$ with a suitable $\lambda > 0$ can effectively eliminate the noise in Z to estimate X , thereby reducing the error in estimating φ_{glo} .

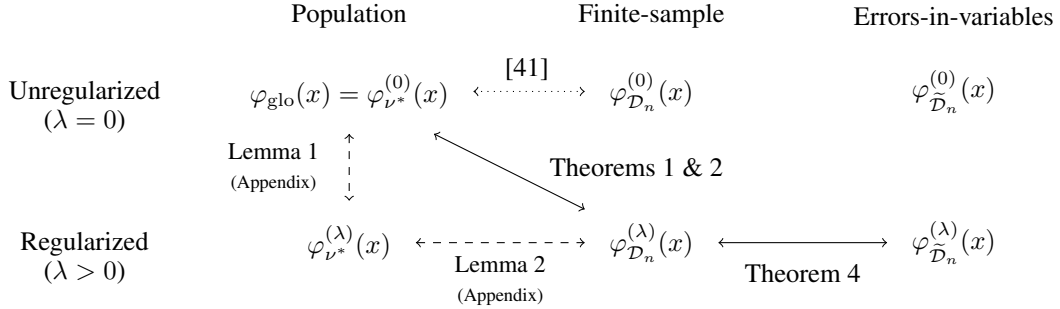


Figure 1: A schematic for the relationship between the regularized Fréchet regression estimators.

4.1 Model assumptions and examples

We impose the following assumptions for our analysis.

(C0) (Existence) For any probability distribution ν and any $\lambda \in \mathbb{R}_+$, the object $\varphi_{\nu}^{(\lambda)}(x)$ exists (almost surely) and is unique. In particular, $\inf_{y \in \mathcal{M}: d(y, \varphi_{\text{glo}}(x)) > \varepsilon} R(y; x) > R(\varphi_{\text{glo}}(x); x)$ for all $\varepsilon > 0$, where $R(y; x) := R_{\nu^*}^{(0)}(y; x)$.

(C1) (Growth) There exist $D_g > 0$, $C_g > 0$ and $\alpha > 1$, possibly depending on x , such that for any probability distribution ν and any $\lambda \in \mathbb{R}_+$,

$$\begin{cases} d(y, \varphi_{\nu}^{(\lambda)}(x)) < D_g & \implies R_{\nu}^{(\lambda)}(y; x) - R_{\nu}^{(\lambda)}(\varphi_{\nu}^{(\lambda)}(x); x) \geq C_g \cdot d(y, \varphi_{\nu}^{(\lambda)}(x))^\alpha, \\ d(y, \varphi_{\nu}^{(\lambda)}(x)) \geq D_g & \implies R_{\nu}^{(\lambda)}(y; x) - R_{\nu}^{(\lambda)}(\varphi_{\nu}^{(\lambda)}(x); x) \geq C_g \cdot D_g^\alpha. \end{cases} \quad (10)$$

(C2) (Bounded entropy) There exists $C_e > 0$, possibly depending on y , such that

$$\limsup_{\delta \rightarrow 0} \int_0^1 \sqrt{1 + \log \mathfrak{N}(B_d(y, \delta), \delta \varepsilon)} d\varepsilon \leq C_e, \quad (11)$$

where $B_d(y, \delta) := \{y' \in \mathcal{M} : d(y, y') \leq \delta\}$ and $\mathfrak{N}(S, \varepsilon)$ is the ε -covering number¹ of S .

¹A formal definition of covering number is provided in Appendix A; see Definition 6.

Assumption (C0) is common to establish the consistency of an M-estimator [55, Chapter 3.2]; in particular, it ensures the weak convergence of the empirical process $R_{\mathcal{D}_n}^{(\lambda)}$ to the population process $R_{\nu^*}^{(\lambda)}$ implying convergence of their minimizers. Furthermore, the conditions on the curvature (C1) and the covering number (C2) control the behavior of the objectives near the minimum in order to obtain rates of convergence; it is worth mentioning that (C2) corresponds to a (locally) bounded entropy for every $y \in \mathcal{M}$, while (P1) in [41] requires the same condition only with $y = \varphi_{\text{glo}}(x)$. These conditions arise from empirical process theory and are also commonly adopted [41, 45, 46].

Here we provide several examples of the space \mathcal{M} , in which the conditions (C0), (C1) and (C2) are satisfied. We verify the conditions in Appendix A; see Propositions 1, 2, 3, and 4.

Example 1. Let $\mathcal{M} = (\mathcal{H}, d_{\text{HS}})$ be a finite-dimensional Hilbert space \mathcal{H} equipped with the Hilbert-Schmidt metric $d_{\text{HS}}(y_1, y_2) = \langle y_1 - y_2, y_1 - y_2 \rangle^{1/2}$, e.g., $\mathcal{M} = (\mathbb{R}^r, d_2)$ where d_2 is the ℓ^2 -metric.

Example 2. Let \mathcal{M} be \mathcal{W} , the set of probability distributions G on \mathbb{R} such that $\int_{\mathbb{R}} x^2 dG(x) < \infty$, equipped with the Wasserstein metric d_W defined as

$$d_W(G_1, G_2)^2 = \int_0^1 (G_1^{-1}(t) - G_2^{-1}(t))^2 dt,$$

where G_1^{-1} and G_2^{-1} are the quantile functions of G_1 and G_2 , respectively. See [41, Section 6].

Example 3. Let $\mathcal{M} = \{M \in \mathbb{R}^{r \times r} : M = M^T, M \succeq 0 \text{ and } M_{ii} = 1, \forall i \in [r]\}$ be the set of correlation matrices of size r , equipped with the Frobenius metric, $d_F(M, M') = \|M - M'\|_F$.

Example 4. Let \mathcal{M} be a (bounded) Riemannian manifold of dimension r , and let d_g be the geodesic distance induced by the Riemannian metric.

4.2 Theorem statements

4.2.1 Noiseless covariate setting

First of all, we show the consistency of the λ -regularized Fréchet regression function.

Theorem 1 (Consistency). *Suppose that Assumption (C0) holds. If $\text{diam}(\mathcal{M}) < \infty$, then for any $\lambda \in \mathbb{R}$ such that $0 \leq \lambda < \min\{\sigma_i(\Sigma_{\nu^*}) : \sigma_i(\Sigma_{\nu^*}) > 0\}$, and any $x \in \mathbb{R}^p$,*

$$d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(0)}(x)) = o_P(1) \quad \text{as } n \rightarrow \infty. \quad (12)$$

If $\lambda < \sigma^{(0)}(\Sigma_{\nu^*}) = \min\{\sigma_i(\Sigma_{\nu^*}) : \sigma_i(\Sigma_{\nu^*}) > 0\}$, then the regularized estimator $\varphi_{\mathcal{D}_n}^{(\lambda)}(x)$ effectively reduces to the same as the sample-analog estimator $\widehat{\varphi}_{\mathcal{D}_n}(x)$ in (4) in the limit $n \rightarrow \infty$. Thus, $\varphi_{\mathcal{D}_n}^{(\lambda)}(x)$ inherits the consistency of $\widehat{\varphi}_{\mathcal{D}_n}$. We provide a detailed proof of Theorem 1 in Appendix B.

In addition to the consistency of $\varphi_{\mathcal{D}_n}^{(\lambda)}$ in the small λ limit, we show the rate of its convergence that holds for any fixed $\lambda \in \mathbb{R}_+$.

Definition 5. For a positive semidefinite matrix Σ , the Mahalanobis seminorm of x induced by Σ is

$$\|x\|_{\Sigma} := (x^{\top} \Sigma^{\dagger} x)^{1/2}. \quad (13)$$

Theorem 2 (Rate of convergence). *Suppose that Assumptions (C0)–(C2) hold. If $\text{diam}(\mathcal{M}) < \infty$, then for any $\lambda \in \mathbb{R}$ and $x \in \mathbb{R}^p$ such that $\|x - \mu_{\nu^*}\|_{\Sigma_{\nu^*}} \leq \frac{C_g \cdot D_g^{\alpha}}{\text{diam}(\mathcal{M})^2 \cdot \sqrt{\text{rank } \Sigma_{\nu^*}}}$,*

$$d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(0)}(x)) = O_P\left(b_{\lambda}(x)^{\frac{1}{\alpha-1}} + n^{-\frac{1}{2(\alpha-1)}}\right) \quad \text{as } n \rightarrow \infty, \quad (14)$$

where $b_{\lambda}(x) = \text{rank}(\Sigma_{\nu^*} - \Sigma_{\nu^*}^{(\lambda)})^{\frac{1}{2}} \cdot \|x - \mu_{\nu^*}\|_{\Sigma_{\nu^*} - \Sigma_{\nu^*}^{(\lambda)}}$.

We obtain Theorem 2 by showing a “bias” upper bound $d(\varphi_{\nu^*}^{(\lambda)}(x), \varphi_{\nu^*}^{(0)}(x)) = O(b_{\lambda}(x)^{\frac{1}{\alpha-1}})$ and a “variance” bound $d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(\lambda)}(x)) = O_P(n^{-\frac{1}{2(\alpha-1)}})$; see Lemmas 1 and 2 in Appendix C.

Here we remark that $b_{\lambda}(x)$ is a monotone non-decreasing function of λ , and if $\lambda < \sigma^{(0)}(\Sigma_{\nu^*})$ then $b_{\lambda}(x) = 0$. Also, the condition on $\|x - \mu_{\nu^*}\|_{\Sigma_{\nu^*}}$ is introduced for a technical reason, and can be removed when $D_g = \infty$.

Remark 1. Note that Condition (C1) holds with $D_g = \infty$ and $\alpha = 2$ for Examples 1, 2 and 3. Thus, we have $d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(0)}(x)) = O_P(b_{\lambda}(x) + n^{-\frac{1}{2}})$ as $n \rightarrow \infty$.

4.2.2 Error-prone covariate setting

Given a set $\mathcal{D}_n = \{(x_i, y_i) : i \in [n]\}$, let $\mathbf{X}_{\mathcal{D}_n} := [x_1 \ \cdots \ x_n]^\top \in \mathbb{R}^{n \times p}$. We let $\mathbf{X} = \mathbf{X}_{\mathcal{D}_n}$ and $\mathbf{Z} = \mathbf{X}_{\tilde{\mathcal{D}}_n}$ for shorthand, and further, we let $\mathbf{X}_{\text{ctr}} = (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \mathbf{X}$ and $\mathbf{Z}_{\text{ctr}} = (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \mathbf{Z}$ denote the ‘row-centered’ matrices.

Theorem 3 (De-noising covariates). *Suppose that Assumptions (C0) and (C1) hold. Then there exists a constant $C > 0$ such that for any $\lambda \in \mathbb{R}_+$, if*

$$x \in \mu_{\mathcal{D}_n} + \text{rowsp } \mathbf{X}_{\text{ctr}} \quad \text{and} \quad \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} \leq \frac{1}{2} \left(\frac{C_g \cdot D_g^\alpha}{2 \text{diam}(\mathcal{M})} \cdot \frac{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}) \wedge \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})}{\|\mathbf{Z} - \mathbf{X}\|} - 1 \right), \quad (15)$$

then

$$d\left(\varphi_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x), \varphi_{\mathcal{D}_n}^{(\lambda)}(x)\right) \leq C \cdot \left(\frac{\|\mathbf{Z} - \mathbf{X}\|}{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}) \wedge \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})} \cdot \frac{2 \cdot \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} + 1}{C_g} \right)^{\frac{1}{\alpha}}. \quad (16)$$

Again, we remark that the condition on $\|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}}$ in (15) can be removed when $D_g = \infty$. It is worth noting that the quantity $\frac{\|\mathbf{Z} - \mathbf{X}\|}{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}) \wedge \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})}$ serves as the reciprocal of the signal-to-noise ratio because $\|\mathbf{Z} - \mathbf{X}\|$ captures the magnitude of the ‘noise’ in the covariates, while $\min\{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})\}$ quantifies the strength of the ‘signal’ remaining in the λ -SVT of the design matrix. Additionally, we observe that the error bound (16) increases proportionally to the normalized deviation of x from the mean, $\mu_{\mathcal{D}_n}$, which is a reasonable outcome. For the complete version of Theorem 3 and its proof, please refer to Appendix D.

4.3 Proof sketches

We outline our proofs for the main theorems, whose details are presented in Appendices B, C and D.

Proof of Theorem 1. We show that $R_{\mathcal{D}_n}^{(\lambda)}(y; x)$ weakly converges to $R_{\nu^*}^{(0)}(y; x)$ in the $\ell^\infty(\mathcal{M})$ -sense. According to [55, Theorem 1.5.4], it suffices to show that (1) $R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(0)}(y; x) = o_p(1)$ for all $y \in \mathcal{M}$, and (2) $R_{\mathcal{D}_n}^{(\lambda)}$ is asymptotically equicontinuous in probability.

Proof of Theorem 2. We prove upper bounds for the bias and the variance separately.

To control the bias (Lemma 1 in Appendix C), we first show an upper bound for $R(\varphi^{(\lambda)}(x); x) - R(\varphi(x); x)$, and then convert it to an upper bound on the distance between the minimizers $d(\varphi^{(\lambda)}(x), \varphi(x))$ using the Growth condition (C1). We utilize the so-called ‘peeling technique’ in empirical process theory in this conversion.

To control the variance (Lemma 2 in Appendix C), we follow a similar strategy as in Lemma 1, but with additional technical considerations. We define the ‘fluctuation variable’ $Z_n^{(\lambda)}(y; x) := R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(\lambda)}(y; x)$ parametrized by $y \in \mathcal{M}$, and derive a probabilistic upper bound for $R_{\nu^*}^{(\lambda)}(\varphi_{\mathcal{D}_n}^{(\lambda)}(x); x) - R_{\nu^*}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x)$ by establishing a uniform upper bound for $Z_n^{(\lambda)}(y; x) - Z_n^{(\lambda)}(\varphi(x); x)$; here, the Entropy condition (C2) is used. Again, we use the Growth condition (C1) and the peeling technique to obtain a probabilistic upper bound for the distance $d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(\lambda)}(x))$.

Proof of Theorem 3. We express the difference in the objective functions $R_{\tilde{\mathcal{D}}_n}^{(\lambda)}(y; x) - R_{\mathcal{D}_n}^{(\lambda)}(y; x)$ using the difference in the weights $w_{\tilde{\mathcal{D}}_n}^{(\lambda)}(y; x) - w_{\mathcal{D}_n}^{(\lambda)}(y; x)$, which can be written in terms of \mathbf{X} and \mathbf{Z} . We use classical matrix perturbation theory to control $R_{\tilde{\mathcal{D}}_n}^{(\lambda)}(y; x) - R_{\mathcal{D}_n}^{(\lambda)}(y; x)$, and transform it to an upper bound on the distance $d(\varphi_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x), \varphi_{\mathcal{D}_n}^{(\lambda)}(x))$ using the Growth condition (C1).

5 Experiments

In this section, we present the results of our numerical simulations, which aim to validate and support the theoretical findings presented earlier. We focus on the problem of global Fréchet regression

analysis for one-dimensional distribution functions (Example 2) and conduct a comprehensive set of simulations under various conditions. Our simulations enable us to assess the performance of our proposed methodology and compare it to alternative approaches. Here we briefly summarize the results in Figure 2 and Table 1. See more details about the simulation settings, implementation details and evaluation metrics, as well as further discussions on the results, in Appendix E.

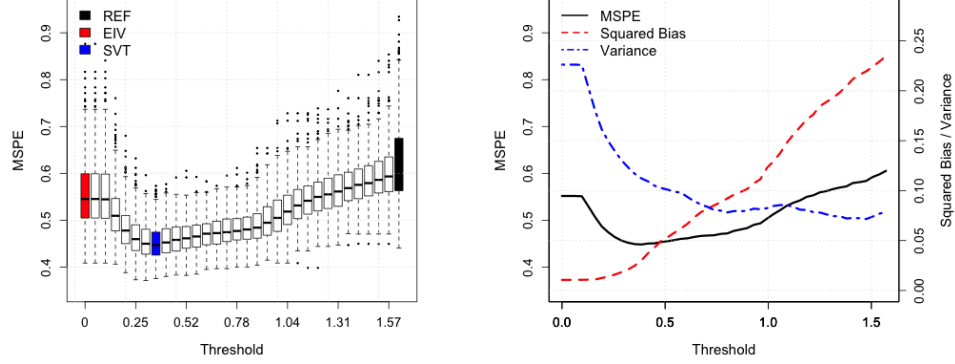


Figure 2: Comparison of the prediction performance of $\varphi_{\mathcal{D}_n}^{(0)}$ (REF), $\varphi_{\mathcal{D}_n}^{(0)}$ (EIV), and $\varphi_{\mathcal{D}_n}^{(\lambda)}$ (SVT) (left), and the trade-off between the bias and the variance (right) for $B = 500$, $p = 50$ and $n = 100$.

Table 1: Average performance of $\varphi_{\mathcal{D}_n}^{(0)}$ (REF), $\varphi_{\mathcal{D}_n}^{(0)}$ (EIV), and $\varphi_{\mathcal{D}_n}^{(\lambda)}$ (SVT) evaluated in four criteria via $B = 500$ Monte Carlo experiments under nine simulation scenarios (boldface=best).

p	Criterion	$n = 100$			$n = 200$			$n = 400$		
		REF	EIV	SVT	REF	EIV	SVT	REF	EIV	SVT
25	Bias	0.016	0.084	0.109	0.011	0.083	0.095	0.007	0.085	0.090
	$\sqrt{\text{Var}}$	0.332	0.290	0.252	0.214	0.187	0.172	0.145	0.126	0.120
	MSE	0.224	0.266	0.272	0.267	0.290	0.293	0.285	0.299	0.301
	MSPE	0.415	0.396	0.380	0.350	0.346	0.343	0.325	0.327	0.326
50	Bias	0.024	0.085	0.153	0.015	0.084	0.115	0.010	0.084	0.094
	$\sqrt{\text{Var}}$	0.567	0.495	0.350	0.327	0.287	0.246	0.211	0.186	0.176
	MSE	0.148	0.248	0.244	0.227	0.268	0.275	0.267	0.289	0.291
	MSPE	0.624	0.557	0.452	0.411	0.394	0.378	0.349	0.346	0.344
75	Bias	0.046	0.094	0.213	0.019	0.091	0.151	0.011	0.083	0.100
	$\sqrt{\text{Var}}$	1.000	0.884	0.410	0.436	0.384	0.292	0.270	0.237	0.215
	MSE	0.073	0.341	0.236	0.187	0.251	0.265	0.247	0.277	0.281
	MSPE	1.288	1.085	0.513	0.493	0.456	0.411	0.377	0.367	0.360

6 Discussion

This paper has addressed errors-in-variables regression of non-Euclidean response variables through the (global) Fréchet regression framework enhanced by low-rank approximation of covariates. Specifically, we introduce a novel regularized (global) Fréchet regression framework (Section 3), which combines the Fréchet regression with principal component regression. We also provide a comprehensive theoretical analysis in three main theorems (Section 4), and validate our theory through numerical experiments on simulated datasets.

We conclude this paper by proposing several promising directions for future research. First, it would be worthwhile to explore the large sample theory for selecting the optimal threshold parameter λ in the proposed SVT method, in order to characterize the theoretical phase transition of the bias-variance trade-off in the regularized (global) Fréchet regression. Second, we believe that our framework could be extended to errors-in-variables Fréchet regression for response variables in a broader class of metric spaces, e.g., by leveraging the quadruple inequality proposed by Schötz [45, 46]. Lastly, investigating the asymptotic distribution of the proposed SVT estimator would be highly appealing in the statistical literature, as it would enable us to make statistical inferences on the conditional Fréchet mean in non-Euclidean spaces [6, 8] with errors-in-variables covariates.

References

- [1] Bijan Afsari. Riemannian L^p center of mass: Existence, uniqueness, and convexity. Proceedings of the American Mathematical Society, 139(2):655–673, 2011.
- [2] Anish Agarwal, Devavrat Shah, Dennis Shen, and Dogyoon Song. On robustness of principal component regression. Journal of the American Statistical Association, 116(536):1731–1745, 2021.
- [3] Adil Ahidar-Coutrix, Thibaut Le Gouic, and Quentin Paris. Convergence rates for empirical barycenters in metric spaces: Curvature, convexity and extendable geodesics. Probability Theory and Related Fields, 177(1):323–368, 2020.
- [4] Gustavo Amorim, Ran Tao, Sarah Lotspeich, Pamela A Shaw, Thomas Lumley, and Bryan E Shepherd. Two-phase sampling designs for data validation in settings with covariate measurement error and continuous outcome. Journal of the Royal Statistical Society Series A: Statistics in Society, 184(4):1368–1389, 2021.
- [5] Marc Arnaudon and Laurent Miclo. Means in complete manifolds: Uniqueness and approximation. ESAIM: Probability and Statistics, 18:185–206, 2014.
- [6] Dennis Barden, Huiling Le, and Megan Owen. Central limit theorems for Fréchet means in the space of phylogenetic trees. Electronic Journal of Probability, 18(25):1–25, 2013.
- [7] Alexandre Belloni, Mathieu Rosenbaum, and Alexandre B Tsybakov. Linear and conic programming estimators in high dimensional errors-in-variables models. Journal of the Royal Statistical Society. Series B (Statistical Methodology), pages 939–956, 2017.
- [8] Rabi Bhattacharya and Lizhen Lin. Omnibus CLTs for Fréchet means and nonparametric inference on non-Euclidean spaces. Proceedings of the American Mathematical Society, 145(1):413–428, 2017.
- [9] Norman E Breslow and Thomas Lumley. Semiparametric models and two-phase samples: Applications to cox regression. IMS Collections, 9:65–77, 2013.
- [10] Louis Capitaine, Jérémie Bigot, Rodolphe Thiébaud, and Robin Genuer. Fréchet random forests for metric space valued regression with non Euclidean predictors. arXiv preprint arXiv:1906.01741, 2019.
- [11] Raymond J Carroll, Helmut Küchenhoff, F Lombard, and Leonard A Stefanski. Asymptotics for the SIMEX estimator in nonlinear measurement error models. Journal of the American Statistical Association, 91(433):242–250, 1996.
- [12] Raymond J Carroll, David Ruppert, Ciprian M Crainiceanu, Tor D Tosteson, and Margaret R Karagas. Nonlinear and nonparametric regression and instrumental variables. Journal of the American Statistical Association, 99(467):736–750, 2004.
- [13] Raymond J Carroll, David Ruppert, Leonard A Stefanski, and Ciprian M Crainiceanu. Measurement Error in Nonlinear Models: A Modern Perspective. Chapman and Hall/CRC, 2006.
- [14] Benjamin Charlier. Necessary and sufficient condition for the existence of a Fréchet mean on the circle. ESAIM: Probability and Statistics, 17:635–649, 2013.
- [15] Yan Mei Chen, Xiao Shan Chen, and Wen Li. On perturbation bounds for orthogonal projections. Numerical Algorithms, 73(2):433–444, 2016.
- [16] John R Cook and Leonard A Stefanski. Simulation-extrapolation estimation in parametric measurement error models. Journal of the American Statistical Association, 89(428):1314–1328, 1994.
- [17] Abhirup Datta and Hui Zou. CoCoLasso for high-dimensional error-in-variables regression. The Annals of Statistics, 45(6):2400–2426, 2017.

- [18] Aurore Delaigle, Jianqing Fan, and Raymond J Carroll. A design-adaptive local polynomial estimator for the errors-in-variables problem. Journal of the American Statistical Association, 104(485):348–359, 2009.
- [19] Aurore Delaigle, Peter Hall, and Wen-xin Zhou. Nonparametric covariate-adjusted regression. The Annals of Statistics, 44(5):2190–2220, 2016.
- [20] Shanshan Ding and R Dennis Cook. Matrix variate regressions and envelope models. Journal of the Royal Statistical Society Series B: Statistical Methodology, 80(2):387–408, 2018.
- [21] David Donoho. 50 years of data science. Journal of Computational and Graphical Statistics, 26(4):745–766, 2017.
- [22] Paromita Dubey and Hans-Georg Müller. Functional models for time-varying random objects. Journal of the Royal Statistical Society Series B: Statistical Methodology, 82(2):275–327, 2020.
- [23] Juan José Egozcue, Josep Daunis-I-Estadella, Vera Pawlowsky-Glahn, Karel Hron, and Peter Filzmoser. Simplicial regression. The normal model. Journal of Applied Probability and Statistics, 2012.
- [24] Jianqing Fan and Elias Masry. Multivariate regression estimation with errors-in-variables: Asymptotic normality for mixing processes. Journal of Multivariate Analysis, 43(2):237–271, 1992.
- [25] Jianqing Fan and Young K Truong. Nonparametric regression with errors in variables. The Annals of Statistics, pages 1900–1925, 1993.
- [26] Maurice Fréchet. Les éléments aléatoires de nature quelconque dans un espace distancié. In Annales de l’Institut Henri Poincaré, volume 10, pages 215–310, 1948.
- [27] Aritra Ghosal, Wendy Meiring, and Alexander Petersen. Fréchet single index models for object response regression. Electronic Journal of Statistics, 17(1):1074–1112, 2023.
- [28] Kyunghye Han, Hans-Georg Müller, and Byeong U Park. Additive functional regression for densities as responses. Journal of the American Statistical Association, 2019.
- [29] X He, C Madigan, J Wellner, and B Yu. Statistics at a crossroads: Who is for the challenge. In Workshop report: National Science Foundation. https://www.nsf.gov/mps/dms/documents/Statistics_at_a_Crossroads_Workshop_Report_2019.pdf, 2019.
- [30] Matthias Hein. Robust nonparametric regression with metric-space valued output. Advances in Neural Information Processing Systems, 22, 2009.
- [31] Anthony JG Hey, Stewart Tansley, Kristin Michele Tolle, et al. The fourth paradigm: Data-intensive scientific discovery, volume 1. Microsoft research Redmond, WA, 2009.
- [32] Alan J. Izenman. Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning. Springer, 2008.
- [33] Ian T Jolliffe. A note on the use of principal components in regression. Journal of the Royal Statistical Society Series C: Applied Statistics, 31(3):300–303, 1982.
- [34] Ioannis Kalogridis and Stefan Van Aelst. Robust functional regression based on principal components. Journal of Multivariate Analysis, 173:393–415, 2019.
- [35] Eric D. Kolaczyk. Statistical Analysis of Network Data. Springer, 2009.
- [36] Thibaut Le Gouic and Jean-Michel Loubes. Existence and consistency of Wasserstein barycenters. Probability Theory and Related Fields, 168:901–917, 2017.
- [37] Hua Liang, Wolfgang Härdle, and Raymond J Carroll. Estimation in a semiparametric partially linear errors-in-variables model. The Annals of Statistics, 27(5):1519–1535, 1999.

- 419 [38] Zhenhua Lin and Hans-Georg Müller. Total variation regularized Fréchet regression for metric-
420 space valued data. The Annals of Statistics, 49(6):3510–3533, 2021.
- 421 [39] Po-Ling Loh and Martin J Wainwright. High-dimensional regression with noisy and missing
422 data: Provable guarantees with non-convexity. Advances in Neural Information Processing
423 Systems, 24, 2011.
- 424 [40] James Stephen Marron and Ian L Dryden. Object Oriented Data Analysis. CRC Press, 2021.
- 425 [41] Alexander Petersen and Hans-Georg Müller. Fréchet regression for random objects with
426 Euclidean predictors. The Annals of Statistics, 47(2):691–719, 2019.
- 427 [42] J. O. Ramsay and B. W. Silverman. Functional Data Analysis. Springer, 2 edition, 2005.
- 428 [43] Philip T Reiss and R Todd Ogden. Functional principal component regression and functional
429 partial least squares. Journal of the American Statistical Association, 102(479):984–996, 2007.
- 430 [44] Susanne M Schennach. Instrumental variable estimation of nonlinear errors-in-variables models.
431 Econometrica, 75(1):201–239, 2007.
- 432 [45] Christof Schötz. Convergence rates for the generalized Fréchet mean via the quadruple inequality.
433 Electronic Journal of Statistics, 13:4280–4345, 2019.
- 434 [46] Christof Schötz. Nonparametric regression in nonstandard spaces. Electronic Journal of
435 Statistics, 16(2):4679–4741, 2022.
- 436 [47] Damla Şentürk and Hans-Georg Müller. Covariate-adjusted regression. Biometrika, 92(1):75–
437 89, 2005.
- 438 [48] Donna Spiegelman, Aidan McDermott, and Bernard Rosner. Regression calibration method
439 for correcting measurement-error bias in nutritional epidemiology. The American Journal of
440 Clinical Nutrition, 65(4):1179S–1186S, 1997.
- 441 [49] Florian Steinke and Matthias Hein. Non-parametric regression between manifolds. Advances
442 in Neural Information Processing Systems, 21, 2008.
- 443 [50] Florian Steinke, Matthias Hein, and Bernhard Schölkopf. Nonparametric regression between
444 general Riemannian manifolds. SIAM Journal on Imaging Sciences, 3(3):527–563, 2010.
- 445 [51] Gilbert W Stewart. On the perturbation of pseudo-inverses, projections and linear least squares
446 problems. SIAM Review, 19(4):634–662, 1977.
- 447 [52] Karl-Theodor Sturm. Probability measures on metric spaces of nonpositive curvature. Heat
448 Kernels and Analysis on Manifolds, Graphs, and Metric Spaces: Lecture Notes from a Quarter
449 Program on Heat Kernels, Random Walks, and Analysis on Manifolds and Graphs, 338:357,
450 2003.
- 451 [53] Renáta Talská, Alessandra Menafoglio, Jitka Machalová, Karel Hron, and E Fišerová. Com-
452 positional regression with functional response. Computational Statistics & Data Analysis,
453 123:66–85, 2018.
- 454 [54] Danielle C Tucker, Yichao Wu, and Hans-Georg Müller. Variable selection for global Fréchet
455 regression. Journal of the American Statistical Association, pages 1–15, 2021.
- 456 [55] Aad W Vaart and Jon A Wellner. Weak convergence. In Weak Convergence and Empirical
457 Processes, pages 16–28. Springer, 1996.
- 458 [56] Roman Vershynin. High-dimensional probability: An introduction with applications in data
459 science, volume 47. Cambridge University Press, 2018.
- 460 [57] Cinzia Viroli. On matrix-variate regression analysis. Journal of Multivariate Analysis, 111:296–
461 309, 2012.
- 462 [58] Takumi Yokota. Convex functions and barycenter on CAT(1)-spaces of small radii. Journal of
463 the Mathematical Society of Japan, 68(3):1297–1323, 2016.
- 464 [59] Y Zemel and V M Panaretos. Fréchet means and procrustes analysis in Wasserstein space.
465 Bernoulli, 25(2):932–976, 2019.

466 A Verification of the model assumptions

467 A.1 Additional background

468 **Definition 6** (ε -net and covering number). *Let (\mathcal{M}, d) be a metric space. Let $S \subseteq T$ be a subset and*
 469 *let $\varepsilon > 0$. A subset $\mathcal{N} \subseteq S$ is called an ε -net of S if every point in S is within distance ε of some*
 470 *point \mathcal{N} , i.e.,*

$$\forall x \in S, \exists x_0 \in \mathcal{N} \text{ such that } d(x, x_0) \leq \varepsilon.$$

471 *The ε -covering number of S , denoted by $\mathfrak{N}(S, \varepsilon)$, is the smallest possible cardinality of an ε -net of S ,*
 472 *i.e.,*

$$\mathfrak{N}(S, \varepsilon) := \min \left\{ k \in \mathbb{N} : \exists y_1, \dots, y_k \in \mathcal{M} \text{ such that } S \subseteq \bigcup_{i=1}^k B_d(y_i, \varepsilon) \right\}, \quad (17)$$

473 where $B_d(y, \varepsilon) = \{y' \in \mathcal{M} : d(y, y') \leq \varepsilon\}$ denotes the closed ε -ball centered at $y \in \mathcal{M}$.

474 Let $B_2^r(0, 1) := \{x \in \mathbb{R}^r : \|x\|_2 \leq 1\}$ denote the unit ℓ^2 -norm ball in \mathbb{R}^r . It is well known² that for
 475 any $\varepsilon > 0$,

$$\left(\frac{1}{\varepsilon}\right)^r \leq \mathfrak{N}(B_2^r(0, 1), \varepsilon) \leq \left(\frac{2}{\varepsilon} + 1\right)^r. \quad (18)$$

476 A.2 Example 1: Euclidean space

477 **Proposition 1.** *The space $(\mathcal{H}, d_{\text{HS}})$ defined in Example 1 satisfies Assumptions (C0), (C1), and (C2).*

478 *Proof of Proposition 1.* For any probability distribution ν and any $\lambda \in \mathbb{R}_+$, let $y_\nu^{(\lambda)} :=$
 479 $\mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \cdot Y]$. Then we observe that

$$\begin{aligned} R_\nu^{(\lambda)}(y; x) &= \mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \cdot d^2(Y, y)] \\ &= \mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \cdot \|Y - y\|^2] \\ &= \mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \cdot \|Y - y_\nu^{(\lambda)}\|^2] + \|y - y_\nu^{(\lambda)}\|_{\text{HS}}^2 \\ &\quad + 2 \underbrace{\left\langle \mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \cdot (Y - y_\nu^{(\lambda)})], y_\nu^{(\lambda)} - y \right\rangle}_{=0} \\ &= R_\nu^{(\lambda)}(y_\nu^{(\lambda)}; x) + \|y - y_\nu^{(\lambda)}\|_{\text{HS}}^2. \end{aligned}$$

480 As $R_\nu^{(\lambda)}(y; x)$ is a strictly convex and coercive function, there exists a unique minimizer, $\varphi_\nu^{(\lambda)}$. Thus,
 481 Condition (C0) is proved. Furthermore, Condition (C1) is also satisfied with $D_g = \infty$, $C_g = 1$, and
 482 $\alpha = 2$.

483 Lastly, for any $y \in \mathcal{H}$ and any $\varepsilon > 0$,

$$\mathfrak{N}(B_{d_{\text{HS}}}(y, \delta), \delta\varepsilon) = \mathfrak{N}(B_{d_{\text{HS}}}(y, 1), \varepsilon) \leq \left(\frac{2}{\varepsilon} + 1\right)^r \leq C \cdot \varepsilon^{-r}$$

484 where $r = \dim \mathcal{H}$ and $C > 1$ is a constant that depends on r only; see the covering number upper
 485 bound in (18). Thus, the integral (11) is bounded as follows:

$$\begin{aligned} \int_0^1 \sqrt{1 + \log \mathfrak{N}(B_d(\varphi(x), \delta), \delta\varepsilon)} \, d\varepsilon &\leq \int_0^1 \sqrt{1 + \log C - r \log \varepsilon} \, d\varepsilon \\ &\leq \sqrt{1 + \log C} + \sqrt{r} \int_0^1 \sqrt{-\log \varepsilon} \, d\varepsilon \\ &= \sqrt{1 + \log C} + \sqrt{r} \int_0^\infty e^{-t} \sqrt{t} \, dt \\ &= \sqrt{1 + \log C} + \frac{\sqrt{r\pi}}{2} \end{aligned}$$

²See [56, Corollary 4.2.13] for example.

486 using the change of variable $t = -\log \varepsilon$. Therefore, Assumption (C2) holds with $C_e = \sqrt{1 + \log C} +$
 487 $\frac{\sqrt{r\pi}}{2}$.

488 □

489 A.3 Example 2: set of probability distributions

490 **Proposition 2.** *The space (\mathcal{W}, d_W) defined in Example 2 satisfies Assumptions (C0), (C1), and (C2).*

491 *Proof of Proposition 2.* For a probability distribution function $y \in \mathcal{W}$ defined on \mathbb{R} , let $\mathcal{Q} =$
 492 $Q(\mathcal{W}) := \{Q(y) : y \in \mathcal{W}\}$ denote the collection of corresponding quantile functions, where
 493 $(Q(y))(u) = y^{-1}(u)$ for $u \in [0, 1]$.

494 We note that $f \mapsto \mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \langle Q(Y), f \rangle]$ is a bounded linear functional defined on $L^2[0, 1]$
 495 because $\mathbb{E}_\nu|w_\nu^{(\lambda)}(X, x)|^2 \leq 2 + 2p \|(x - \mu_\nu)\|_{\Sigma_\nu}^2$ implies that $\mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \cdot \|Q(Y)\|_2]$ is bounded.
 496 It follows from the Riesz representation theorem that there exists $f_x^{(\lambda)} \in L^2[0, 1]$ such that

$$\mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \langle Q(Y), g \rangle_2] = \langle f_x^{(\lambda)}, g \rangle_2 \quad (19)$$

497 for all $g \in L^2[0, 1]$. Then, we have

$$R_\nu^{(\lambda)}(y; x) = \mathbb{E}_\nu[w_\nu^{(\lambda)}(X, x) \|Q(Y) - f_x^{(\lambda)}\|_2^2] + \|Q(y) - f_x^{(\lambda)}\|_2^2, \quad (20)$$

498 which yields that

$$\varphi_\nu^{(\lambda)}(x) = Q^{-1}\left(\arg \min_{Q \in \mathcal{Q}} \|Q - f_x^{(\lambda)}\|_2^2\right). \quad (21)$$

499 The condition (C0) follows from the convexity of $(\mathcal{Q}, \|\cdot\|_2)$. Moreover, the convexity also gives
 500 $\langle Q(y) - Q(\varphi_\nu^{(\lambda)}(x)), f_x^{(\lambda)}(x) - Q(\varphi_\nu^{(\lambda)}(x)) \rangle_2 \leq 0$ for all $y \in \mathcal{W}$, so that

$$\begin{aligned} R_\nu^{(\lambda)}(y; x) - R_\nu^{(\lambda)}(\varphi_\nu^{(\lambda)}(x); x) &= \|Q(y) - f_x^{(\lambda)}(x)\|_2^2 - \|Q(\varphi_\nu^{(\lambda)}(x)) - f_x^{(\lambda)}(x)\|_2^2 \\ &= \|Q(y) - Q(\varphi_\nu^{(\lambda)}(x))\|_2^2 - 2\langle Q(y) - Q(\varphi_\nu^{(\lambda)}(x)), f_x^{(\lambda)}(x) - Q(\varphi_\nu^{(\lambda)}(x)) \rangle_2 \\ &\geq \|Q(y) - Q(\varphi_\nu^{(\lambda)}(x))\|_2^2 \\ &= d_W^2(y, \varphi_\nu^{(\lambda)}(x)). \end{aligned} \quad (22)$$

501 Therefore, the condition (C1) holds for any arbitrary constant $D_g > 0$ with $C_g = 1$ and $\alpha = 2$.

502 Finally, we refer to [41, Proposition 1] to ensure that for any $\delta > 0$ and any $\varepsilon > 0$,

$$\sup_{y \in \mathcal{W}} \log \mathfrak{N}(B_{d_W}(y, \delta), D_e \varepsilon) \leq \sup_{Q \in \mathcal{Q}} \log \mathfrak{N}(B_{d_2}(Q, \delta), \delta \varepsilon) \leq C \cdot \varepsilon^{-1} \quad (23)$$

503 holds with an absolute constant $C > 0$. Technically, this fact comes from the covering number bound
 504 for a class of uniformly bounded and monotone functions in L^2 . This confirms that the entropy
 505 condition (C2) holds. □

506 A.4 Example 3: set of correlation matrices

507 **Proposition 3.** *The space (\mathcal{M}, d_F) defined in Example 3 satisfies Assumptions (C0), (C1), and (C2).*

508 *Proof of Proposition 3.* First of all, we note that \mathcal{M} is a convex subset of $\mathcal{S}^r := \{X \in \mathbb{R}^{r \times r} : X =$
 509 $X^\top\}$, which is the set of $r \times r$ symmetric matrices. It is because $\mathcal{M} = \mathcal{S}_+^r \cap H$ where \mathcal{S}_+^r denotes
 510 the cone of $r \times r$ positive semidefinite matrices and $H := \{X \in \mathcal{S}^r : X_{ii} = 1, \forall i \in [r]\}$ denotes an
 511 affine subspace of \mathcal{S}^r , both of which are convex.

512 Next, we observe that \mathcal{S}^r equipped with the Frobenius metric d_F is isometrically isomorphic to
 513 $\mathbb{R}^{r(r+1)/2}$ equipped with the ℓ^2 -metric. Hence, (\mathcal{M}, d_F) satisfies Assumptions (C0), (C1), and (C2),
 514 inheriting these properties from the ambient space \mathcal{S}^r , which is established by Proposition 1. We
 515 note that the inheritance of (C0), (C1) relies on the convexity of \mathcal{M} , while (C2) is inherited simply
 516 based on the inclusion $\mathcal{M} \subset \mathcal{S}^r$. □

517 A.5 Example 4: bounded Riemannian manifold

518 **Proposition 4.** *The space (\mathcal{M}, d_g) defined in Example 4 satisfies Assumption (C2) provided that the*
 519 *Riemannian metric is equivalent to the ambient Euclidean metric.*

520 *Furthermore, let $T_y\mathcal{M}$ be the tangent space of \mathcal{M} at y , and $\text{Exp}_y : T_y\mathcal{M} \rightarrow \mathcal{M}$ be the manifold*
 521 *exponential map at y . Let $g_\nu^{(\lambda)}(u; y, x) := R_\nu^{(\lambda)}(\mathbb{E}_y(u), x)$ for $u \in T_y\mathcal{M}$. If (C0) holds and the*
 522 *Hessian of $g_\nu^{(\lambda)}(u; \varphi_\nu^{(\lambda)}(x), x)$ is positive definite, then (C1) for some $D_g > 0$.*

523 *Proof of Proposition 3.* Since \mathcal{M} has finite dimension and is bounded, the bounded entropy condition
 524 (C2) follows from the metric equivalence.

525 Suppose that (C0) holds, and let $\delta > 0$ be the injectivity radius at $\varphi_\nu^{(\lambda)}(x)$. Consider $y \in \mathcal{M}$ such
 526 that $d(y, \varphi_\nu^{(\lambda)}(x)) < \delta$, and let $u_y = \text{Log}_{\varphi_\nu^{(\lambda)}(x)}(y)$. Then we have

$$R_\nu^{(\lambda)}(y; x) - R_\nu^{(\lambda)}(\varphi_\nu^{(\lambda)}(x); x) = g_\nu^{(\lambda)}(u_y; \varphi_\nu^{(\lambda)}(x), x) - g_\nu^{(\lambda)}(0; \varphi_\nu^{(\lambda)}(x), x) = u_y^\top \nabla^2 g_\nu^{(\lambda)}(\bar{u}_y) u_y$$

527 for some \bar{u}_y between 0 and u_y . Since $u_y^\top u_y = d(y, \varphi_\nu^{(\lambda)}(x))^2$ and $g_\nu^{(\lambda)}$ is continuous, the positive
 528 definiteness of $\nabla^2 g_\nu^{(\lambda)}(\bar{u}_y)$ implies (C1) with $\alpha = 1$. \square

529 B Proof of Theorem 1

530 *Proof of Theorem 1.* Recall from Definition 4, cf. (9), that for any probability distribution ν on \mathbb{R}^p ,
 531 any $\lambda \in \mathbb{R}_+$, and any $x \in \mathbb{R}^p$, the λ -regularized Fréchet regression function evaluated at x is given
 532 as the minimizer of a function $R_\nu^{(\lambda)}$ as

$$\varphi_\nu^{(\lambda)}(x) = \arg \min_{y \in \mathcal{M}} R_\nu^{(\lambda)}(y; x)$$

533 where

$$R_\nu^{(\lambda)}(y; x) = \mathbb{E}_{(X,Y) \sim \nu} \left[w_\nu^{(\lambda)}(X, x) \cdot d^2(Y, y) \right] \quad \text{and}$$

$$w_\nu^{(\lambda)}(x', x) = 1 + (x' - \mu_\nu)^\top \left[\text{SVT}^{(\lambda)}(\Sigma_\nu) \right]^\dagger (x - \mu_\nu).$$

534 In this proof, we follow a similar strategy to that in the proof of [41, Theorem 1]. Specifically,
 535 it suffices to show $\sup_{y \in \mathcal{M}} |R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(0)}(y; x)|$ converges to zero in probability, due to [55,
 536 Corollary 3.2.3]. To this end, we show $R_{\mathcal{D}_n}^{(\lambda)}(y; x)$ weakly converges to $R_{\nu^*}^{(0)}(y; x)$ in the $\ell^\infty(\mathcal{M})$ -
 537 sense, and then apply [55, Theorem 1.3.6]. Again, according to [55, Theorem 1.5.4], this weak
 538 convergence can be proved by showing that

539 (S1) $R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(0)}(y; x) = o_p(1)$ for all $y \in \mathcal{M}$, and

540 (S2) $R_{\mathcal{D}_n}^{(\lambda)}$ is asymptotically equicontinuous in probability, i.e., for any $\varepsilon, \eta > 0$, there exists
 541 $\delta > 0$ such that

$$\limsup_n P \left(\sup_{y_1, y_2 \in \mathcal{M}: d(y_1, y_2) < \delta} \left| R_{\mathcal{D}_n}^{(\lambda)}(y_1; x) - R_{\mathcal{D}_n}^{(\lambda)}(y_2; x) \right| > \varepsilon \right) < \eta.$$

542 In what follows, we prove these two statements, (S1) and (S2), thereby completing the proof of
 543 Theorem 1.

544 **Step 1: proof of (S1).** First of all, we observe that

$$R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(0)}(y; x) = \underbrace{\left(R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\mathcal{D}_n}^{(0)}(y; x) \right)}_{=: T_1} + \underbrace{\left(R_{\mathcal{D}_n}^{(0)}(y; x) - R_{\nu^*}^{(0)}(y; x) \right)}_{=: T_2}. \quad (24)$$

545 We separately analyze the two terms T_1 and T_2 below to show $T_1 = o_p(1)$ and $T_2 = o_p(1)$ as
 546 $n \rightarrow \infty$.

547

(i) $\underline{T_1 = o_p(1)}$.

548

Let $\mathcal{D}_n = \{(X_i, Y_i) : i \in [n]\}$, and we re-write

$$T_1 = \frac{1}{n} \sum_{i=1}^n \left(w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) - w_{\mathcal{D}_n}^{(0)}(X_i, x) \right) \cdot d^2(Y_i, y).$$

549

Letting $\hat{\mu}_n = \mu_{\mathcal{D}_n}$, $\hat{\Sigma}_n = \Sigma_{\mathcal{D}_n}$, and $\hat{\Sigma}_n^{(\lambda)} = \text{SVT}^{(\lambda)}(\hat{\Sigma}_n)$ for shorthand, we observe that

$$w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) - w_{\mathcal{D}_n}^{(0)}(X_i, x) = (X_i - \hat{\mu}_n)^\top \left[\hat{\Sigma}_n^{(\lambda), \dagger} - \hat{\Sigma}_n^\dagger \right] (x - \hat{\mu}_n).$$

550

Let $\mathbf{X} = [X_1 \ \cdots \ X_n]^\top \in \mathbb{R}^{n \times p}$, and note that $\hat{\Sigma}_n = \frac{1}{n} (\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top)^\top (\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top)$. Then it follows that

551

$$\frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_n)^\top \left[\hat{\Sigma}_n^{(\lambda), \dagger} - \hat{\Sigma}_n^\dagger \right] = \frac{1}{n} \mathbf{1}_n^\top (\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top) \left[\hat{\Sigma}_n^{(\lambda), \dagger} - \hat{\Sigma}_n^\dagger \right]$$

552

Consider a singular value decomposition of $\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top$, namely,

$$\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top = \sum_{i=1}^{\min\{n, p\}} s_i \cdot u_i v_i^\top,$$

553

and observe that $\hat{\Sigma}_n = \sum_{i=1}^{\min\{n, p\}} s_i^2 \cdot v_i v_i^\top$ is an eigenvalue decomposition of $\hat{\Sigma}_n$. Letting

554

 $\mathcal{V}_n^{(\lambda)} := \text{span} \left\{ v_i : i \in [p], 0 < s_i \leq \sqrt{\lambda} \right\}$ be a subspace of \mathbb{R}^p spanned by the eigenvectors

555

of $\hat{\Sigma}_n$ corresponding to the nonzero eigenvalues no greater than λ , we have

$$\begin{aligned} \hat{\Sigma}_n^{(\lambda), \dagger} - \hat{\Sigma}_n^\dagger &= \sum_{i=1}^p \frac{1}{s_i^2} \cdot \mathbb{1}\{s_i > \sqrt{\lambda}\} \cdot v_i v_i^\top - \sum_{i=1}^p \frac{1}{s_i^2} \cdot \mathbb{1}\{s_i > 0\} \cdot v_i v_i^\top \\ &= \sum_{i=1}^p \frac{1}{s_i^2} \cdot \mathbb{1}\{0 < s_i \leq \sqrt{\lambda}\} \cdot v_i v_i^\top \\ &= \hat{\Sigma}_n^\dagger \cdot \Pi_{\mathcal{V}_n^{(\lambda)}} \\ &= n \cdot (\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top)^\dagger (\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top)^{\dagger, \top} \cdot \Pi_{\mathcal{V}_n^{(\lambda)}} \end{aligned} \quad (25)$$

556

where $\Pi_{\mathcal{V}_n^{(\lambda)}}$ denotes the projection matrix onto the subspace $\mathcal{V}_n^{(\lambda)}$. Note that $\Pi_{\mathcal{V}_n^{(\lambda)}} = 0$ if

557

and only if $\min \{i \in [p] : 0 < s_i \leq \sqrt{\lambda}\} = \emptyset$.

558

Therefore, we have

$$\begin{aligned} T_1 &= \frac{1}{n} \sum_{i=1}^n \left(w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) - w_{\mathcal{D}_n}^{(0)}(X_i, x) \right) \cdot d^2(Y_i, y) \\ &\leq \frac{\text{diam}(\mathcal{M})^2}{n} \mathbf{1}_n^\top (\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top) \left[\hat{\Sigma}_n^{(\lambda), \dagger} - \hat{\Sigma}_n^\dagger \right] (x - \hat{\mu}_n) \\ &= \text{diam}(\mathcal{M})^2 \cdot \mathbf{1}_n^\top (\mathbf{X} - \mathbf{1}_n \hat{\mu}_n^\top)^{\dagger, \top} \cdot \Pi_{\mathcal{V}_n^{(\lambda)}} \cdot (x - \hat{\mu}_n) \quad \because (25) \\ &= o_p(1). \end{aligned}$$

559

The last line follows from the fact that $\sup_{i \in [p]} (\sigma_i(\hat{\Sigma}_n) - \sigma_i(\Sigma_{\nu^*})) \rightarrow 0$ in probability, and thus, $\Pi_{\mathcal{V}_n^{(\lambda)}} \rightarrow 0$ in probability.

560

561

(ii) $\underline{T_2 = o_p(1)}$.

562 Letting $\tilde{R}_n(y; x) = \frac{1}{n} \sum_{i=1}^n w_{\nu^*}^{(0)}(X_i, x) \cdot d^2(Y_i, y)$, we decompose T_2 as follows:

$$\begin{aligned}
T_2 &= R_{\mathcal{D}_n}^{(0)}(y; x) - \tilde{R}_n(y; x) + \tilde{R}_n(y; x) - R_{\nu^*}^{(0)}(y; x) \\
&= \underbrace{\frac{1}{n} \sum_{i=1}^n \left\{ w_{\mathcal{D}_n}^{(0)}(X_i, x) - w_{\nu^*}^{(0)}(X_i, x) \right\} \cdot d^2(Y_i, y)}_{=: T_{2A}} \\
&\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \left\{ w_{\nu^*}^{(0)}(X_i, x) \cdot d^2(Y_i, y) - \mathbb{E} \left[w_{\nu^*}^{(0)}(X_i, x) \cdot d^2(Y_i, y) \right] \right\}}_{=: T_{2B}}
\end{aligned}$$

563 Note that T_{2B} converges to 0 in probability by the weak law of large numbers.

564 Now it remains to show $T_{2A} = o_p(1)$. To this end, we note that

$$\begin{aligned}
w_{\mathcal{D}_n}^{(0)}(X_i, x) - w_{\nu^*}^{(0)}(X_i, x) &= V_n(x) + X_i^\top W_n(x) \\
\text{where } \begin{cases} V_n(x) = -\hat{\mu}_n^\top \hat{\Sigma}_n^\dagger (x - \hat{\mu}_n) + \mu^\top \Sigma^\dagger (x - \mu), \\ W_n(x) = \hat{\Sigma}_n^\dagger (x - \hat{\mu}_n) - \Sigma^\dagger (x - \mu). \end{cases} & \quad (26)
\end{aligned}$$

565 Since $\hat{\mu}_n$ and $\hat{\Sigma}_n$ respectively converge to μ and Σ in probability, it is possible to verify
566 that $|V_n(x)|, \|W_n(x)\|$ converge to 0 in probability. As a result, T_2 also converges to 0 in
567 probability.

568 All in all, we have $R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(0)}(y; x) = o_p(1)$, and thus, proved (S1).

569 **Step 2: proof of (S2).** For any $y_1, y_2 \in \mathcal{M}$,

$$\begin{aligned}
\left| R_{\mathcal{D}_n}^{(\lambda)}(y_1; x) - R_{\mathcal{D}_n}^{(\lambda)}(y_2; x) \right| &= \left| \frac{1}{n} \sum_{i=1}^n w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) \cdot \{d^2(Y_i, y_1) - d^2(Y_i, y_2)\} \right| \\
&\leq \frac{1}{n} \sum_{i=1}^n \left| w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) \right| \cdot |d(Y_i, y_1) + d(Y_i, y_2)| \cdot |d(Y_i, y_1) - d(Y_i, y_2)| \\
&\leq 2 \operatorname{diam}(\mathcal{M}) \cdot d(y_1, y_2) \cdot \left(\frac{1}{n} \sum_{i=1}^n \left| w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) \right| \right) \\
&= O_p(d(y_1, y_2))
\end{aligned}$$

570 where the O_p term is independent of $y_1, y_2 \in \mathcal{M}$. Therefore,

$$\sup_{y_1, y_2 \in \mathcal{M}: d(y_1, y_2) < \delta} \left| R_{\mathcal{D}_n}^{(\lambda)}(y_1; x) - R_{\mathcal{D}_n}^{(\lambda)}(y_2; x) \right| = O_p(\delta),$$

571 which proves (S2).

572 □

573 C Proof of Theorem 2

574 In this section, we prove the two claims in Theorem 2. Specifically, in Section C.1, we present and
575 prove a lemma that controls the bias in the population estimator (Lemma 1), and in Section C.2, we
576 present and prove a lemma that controls the variance of the empirical estimator (Lemma 2).

577 C.1 Bias in the population estimator

578 We recall the definition of Mahalanobis seminorm from Definition 5: $\|x\|_\Sigma := (x^\top \Sigma^\dagger x)^{1/2}$.

579 **Lemma 1.** Suppose that Assumptions (C0) and (C1) hold. If

$$\|x - \mu_{\nu^*}\|_{\Sigma_{\nu^*}} \leq \frac{C_g \cdot D_g^\alpha}{\text{diam}(\mathcal{M})^2 \cdot \sqrt{\text{rank} \Sigma_{\nu^*}}}, \quad (27)$$

580 then for any $\lambda \in \mathbb{R}_+$,

$$d(\varphi^{(\lambda)}(x), \varphi(x)) \leq 2^{K_0} \cdot b_\lambda(x)^{\frac{1}{\alpha-1}} = O(b_\lambda(x)^{\frac{1}{\alpha-1}}) \quad (28)$$

581 where

$$K_0 = \left\lfloor \frac{1}{(\alpha-1) \log 2} \cdot \log \left(\frac{4 \text{diam}(\mathcal{M})}{C_g \cdot (1 - 2^{-(\alpha-1)})} \right) \right\rfloor + 1 \quad \text{and}$$

$$b_\lambda(x) = \sqrt{\text{rank}(\Sigma_{\nu^*} - \Sigma_{\nu^*}^{(\lambda)})} \cdot \|x - \mu_{\nu^*}\|_{\Sigma_{\nu^*} - \Sigma_{\nu^*}^{(\lambda)}}.$$

582 *Proof of Lemma 1.* For the sake of brevity, we write $\varphi^{(\lambda)}(x) = \varphi_{\nu^*}^{(\lambda)}(x)$ and $\varphi(x) = \varphi_{\nu^*}^{(0)}(x)$
 583 throughout this proof, dropping the subscript ν^* . Likewise, we simply write $\mu = \mu_{\nu^*}$ and $\Sigma = \Sigma_{\nu^*}$.

584 **Step 1: A naïve upper bound.** Observe that for any $\lambda \in \mathbb{R}_+$, $x \in \mathbb{R}^p$, and $y \in \mathcal{M}$,

$$\begin{aligned} & |R(y; x) - R^{(\lambda)}(y; x)| \\ &= \left| \mathbb{E}_{\nu^*} \left[(X - \mu)^\top \cdot (\Sigma^\dagger - \Sigma^{(\lambda), \dagger}) \cdot (x - \mu) \cdot d^2(Y, y) \right] \right| \\ &\leq \text{diam}(\mathcal{M})^2 \cdot \mathbb{E}_{\nu^*} [\|X - \mu\|_{\Sigma - \Sigma^{(\lambda)}}] \cdot \|x - \mu\|_{\Sigma - \Sigma^{(\lambda)}} \quad \because \text{Cauchy-Schwarz inequality} \\ &\leq \text{diam}(\mathcal{M})^2 \cdot \left(\mathbb{E}_{\nu^*} \|X - \mu\|_{\Sigma - \Sigma^{(\lambda)}}^2 \right)^{1/2} \cdot \|x - \mu\|_{\Sigma - \Sigma^{(\lambda)}} \quad \because \text{Jensen's inequality} \\ &= \text{diam}(\mathcal{M})^2 \cdot \sqrt{\text{rank}(\Sigma - \Sigma^{(\lambda)})} \cdot \|x - \mu\|_{\Sigma - \Sigma^{(\lambda)}}, \end{aligned} \quad (29)$$

585 where the last inequality follows from $\mathbb{E}_{\nu^*} \|X - \mu\|_{\Sigma - \Sigma^{(\lambda)}}^2 = \text{rank}(\Sigma - \Sigma^{(\lambda)})$.

586 We observe that the upper bound in (29) is monotone non-decreasing with respect to $\lambda \in \mathbb{R}_+$, and it
 587 converges to 0 as $\lambda \rightarrow 0$. To see this, for any $\lambda \in \mathbb{R}_+$, we let

$$\mathcal{V}^{(\lambda)} := \text{span} \{v_i : i \in [p], 0 < \lambda_i \leq \lambda\}$$

588 where $\Sigma = \sum_{i=1}^p \lambda_i \cdot v_i v_i^\top$ is an eigendecomposition of Σ . Letting $\Pi_{\mathcal{V}^{(\lambda)}}$ denote the projection
 589 matrix onto the subspace $\mathcal{V}^{(\lambda)}$, we note that $\Sigma - \Sigma^{(\lambda)} = \Pi_{\mathcal{V}^{(\lambda)}} \Sigma \Pi_{\mathcal{V}^{(\lambda)}}$, and that $(\Sigma - \Sigma^{(\lambda)})^\dagger =$
 590 $\Pi_{\mathcal{V}^{(\lambda)}} \Sigma^\dagger \Pi_{\mathcal{V}^{(\lambda)}}$. Thus, $\text{rank}(\Sigma - \Sigma^{(\lambda)}) = \dim \mathcal{V}^{(\lambda)}$, and furthermore, we notice that $\mathcal{V}^{(\lambda)} = \{0\}$ if
 591 and only if $\lambda < \lambda_{\min} := \min\{\lambda_i : \lambda_i > 0\}$. Therefore,

$$\lambda < \lambda_{\min} \implies R^{(\lambda)}(y; x) - R(y; x) = 0 \implies \varphi^{(\lambda)}(x) = \varphi(x), \quad \forall x. \quad (30)$$

592 The observation (30), together with Assumption (C0), implies that $d(\varphi^{(\lambda)}(x), \varphi(x)) = o(1)$ as
 593 $\lambda \rightarrow 0$.

594 **Step 2: Controlling risk difference.** Next, we move on to determine the order of $d(\varphi^{(\lambda)}(x), \varphi(x))$
 595 — as a function of $b_\lambda(x)$ — for a fixed $\lambda \in \mathbb{R}$. We may assume $\lambda > \lambda_{\min}$ for the proof because the
 596 lemma is trivial otherwise, cf. (30). Assuming $\lambda > \lambda_{\min}$, we may decompose the difference in the
 597 population objective at $\varphi^{(\lambda)}(x)$ and $\varphi(x)$ as follows:

$$\begin{aligned} R(\varphi^{(\lambda)}(x); x) - R(\varphi(x); x) &= \underbrace{\left\{ R(\varphi^{(\lambda)}(x); x) - R^{(\lambda)}(\varphi^{(\lambda)}(x); x) + R^{(\lambda)}(\varphi(x); x) - R(\varphi(x); x) \right\}}_{=:\mathfrak{R}_1} \\ &\quad - \underbrace{\left\{ R^{(\lambda)}(\varphi(x); x) - R^{(\lambda)}(\varphi^{(\lambda)}(x); x) \right\}}_{=:\mathfrak{R}_2}. \end{aligned}$$

598 We observe that both \mathfrak{R}_1 and \mathfrak{R}_2 are non-negative, due to the optimality of $\varphi(x)$ and $\varphi^{(\lambda)}(x)$. Then,
 599 we obtain an upper bound for \mathfrak{R}_1 using a similar argument as in (29). Specifically,

$$\begin{aligned} R(\varphi^{(\lambda)}(x); x) - R(\varphi(x); x) &\leq \mathfrak{R}_1 \\ &= \mathbb{E}_{\nu^*} \left[\left\{ w_{\nu^*}^{(0)}(X, x) - w_{\nu^*}^{(\lambda)}(X, x) \right\} \cdot \left\{ d^2(Y, \varphi^{(\lambda)}(x)) - d^2(Y, \varphi(x)) \right\} \right] \\ &\leq 2 \operatorname{diam}(\mathcal{M}) \cdot \mathbf{b}_\lambda(x) \cdot d(\varphi^{(\lambda)}(x), \varphi(x)). \end{aligned} \quad (31)$$

600 **Step 3: Converting risk difference to bias.** Lastly, we convert the upper bound (31) to an upper
 601 bound on the distance $d(\varphi^{(\lambda)}(x), \varphi(x))$ using Assumption (C1). To this end, we begin by confirming
 602 that

$$\begin{aligned} R(\varphi^{(\lambda)}(x); x) - R(\varphi(x); x) &= \mathbb{E}_{\nu^*} \left[(X - \mu)^\top \cdot \Sigma^\dagger \cdot (x - \mu) \cdot \left\{ d^2(Y, \varphi^{(\lambda)}(x)) - d^2(Y, \varphi(x)) \right\} \right] \\ &\leq \operatorname{diam}(\mathcal{M})^2 \cdot \left(\mathbb{E}_{\nu^*} \|X - \mu\|_\Sigma^2 \right)^{1/2} \cdot \|x - \mu\|_\Sigma \\ &= \operatorname{diam}(\mathcal{M})^2 \cdot \sqrt{\operatorname{rank} \Sigma} \cdot \|x - \mu\|_\Sigma \\ &\leq C_g \cdot D_g^\alpha. \end{aligned}$$

603 Thereafter, we choose an arbitrary $K \in \mathbb{N}$ and $r \in \mathbb{R}_+$ whose values will be determined later in this
 604 proof. Then we obtain the following inequality using the so-called peeling technique:

$$\begin{aligned} &\mathbb{1} \left\{ d(\varphi^{(\lambda)}(x), \varphi(x)) > 2^K \cdot \mathbf{b}_\lambda(x)^r \right\} \\ &= \sum_{k=K}^{\infty} \mathbb{1} \left\{ 2^k \cdot \mathbf{b}_\lambda(x)^r < d(\varphi^{(\lambda)}(x), \varphi(x)) \leq 2^{k+1} \cdot \mathbf{b}_\lambda(x)^r \right\} \\ &\leq \sum_{k=K}^{\infty} \mathbb{1} \left\{ 2^k \cdot \mathbf{b}_\lambda(x)^r < d(\varphi^{(\lambda)}(x), \varphi(x)) \leq 2^{k+1} \cdot \mathbf{b}_\lambda(x)^r \right\} \\ &\leq \sum_{k=K}^{\infty} \frac{R(\varphi^{(\lambda)}(x); x) - R(\varphi(x); x)}{C_g \cdot (2^k \cdot \mathbf{b}_\lambda(x)^r)^\alpha} \cdot \mathbb{1} \left\{ d(\varphi^{(\lambda)}(x), \varphi(x)) \leq 2^{k+1} \cdot \mathbf{b}_\lambda(x)^r \right\}. \quad \because (C1) \end{aligned} \quad (32)$$

605 Moreover, we decompose the numerator in the fraction appearing in the upper bound (32) as follows:

606 Combining (31) with (32), we have

$$\begin{aligned} &\mathbb{1} \left\{ d(\varphi^{(\lambda)}(x), \varphi(x)) > 2^K \cdot \mathbf{b}_\lambda(x)^r \right\} \\ &\leq \sum_{k=K}^{\infty} \frac{2 \operatorname{diam}(\mathcal{M}) \cdot \mathbf{b}_\lambda(x) \cdot d(\varphi^{(\lambda)}(x), \varphi(x))}{C_g \cdot (2^k \cdot \mathbf{b}_\lambda(x)^r)^\alpha} \cdot \mathbb{1} \left\{ d(\varphi^{(\lambda)}(x), \varphi(x)) \leq 2^{k+1} \cdot \mathbf{b}_\lambda(x)^r \right\} \\ &\leq \frac{4 \operatorname{diam}(\mathcal{M})}{C_g} \cdot \mathbf{b}_\lambda(x)^{1-r(\alpha-1)} \sum_{k=K}^{\infty} \frac{1}{2^{k(\alpha-1)}}. \end{aligned} \quad (33)$$

607 Note that $C := \frac{4 \operatorname{diam}(\mathcal{M})}{C_g} > 0$ is a constant independent of λ . Let $r = 1/(\alpha - 1)$, and observe that
 608 the upper bound in (33) becomes smaller than 1 for a sufficiently large K . Specifically,

$$K \geq \left\lceil \frac{1}{(\alpha - 1) \log 2} \cdot \log \left(\frac{4 \operatorname{diam}(\mathcal{M})}{C_g \cdot (1 - 2^{-(\alpha-1)})} \right) \right\rceil + 1 \quad \implies \quad \frac{4 \operatorname{diam}(\mathcal{M})}{C_g} \cdot \sum_{k=K}^{\infty} \frac{1}{2^{k(\alpha-1)}} < 1.$$

609 As a result, the inequality “ $d(\varphi^{(\lambda)}(x), \varphi(x)) > 2^{K_0} \cdot \mathbf{b}_\lambda(x)^r$ ” in the indicator function must be false,
 610 and we conclude that

$$d(\varphi^{(\lambda)}(x), \varphi(x)) \leq 2^{K_0} \cdot \mathbf{b}_\lambda(x)^{\frac{1}{\alpha-1}}.$$

611

□

612 C.2 Variance of the empirical estimator

613 **Lemma 2.** Suppose that Assumptions (C0), (C1) and (C2) hold. For any $\lambda \in \mathbb{R}_+$ such that
 614 $\lambda \notin \text{spec}(\Sigma_{\nu^*})$, it holds that

$$d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(\lambda)}(x)) = O_P\left(n^{-\frac{1}{2(\alpha-1)}}\right).$$

615 *Proof of Lemma 2.* Recall from the definition of λ -regularized Fréchet regression (Definition 4) and
 616 (9) that

$$R_{\mathcal{D}_n}^{(\lambda)}(y; x) = \frac{1}{n} \sum_{i=1}^n w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) \cdot d^2(Y_i, y) \quad \text{and} \quad R_{\nu^*}^{(\lambda)}(y; x) = \mathbb{E}_{(X, Y) \sim \nu^*} \left[w_{\nu^*}^{(\lambda)}(X, x) \cdot d^2(Y, y) \right].$$

617 Additionally, we define an auxiliary function $\tilde{R}_n(y; x)$ as the “empirical risk with population weight”
 618 such that

$$\tilde{R}_n(y; x) := \frac{1}{n} \sum_{i=1}^n w_{\nu^*}^{(\lambda)}(X_i, x) \cdot d^2(Y_i, y).$$

619 We present the rest of this proof in three steps, outlined as follows. In Step 1, we show the consistency
 620 of $\varphi_{\mathcal{D}_n}^{(\lambda)}(x)$, i.e., $d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(\lambda)}(x)) = o_P(1)$ as $n \rightarrow \infty$. In Step 2, we define the discrepancy vari-
 621 able $Z_n^{(\lambda)}(y; x) := R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(\lambda)}(y; x)$ between the finite-sample and the population objectives,
 622 cf. (36), and prove a uniform upper bound for $Z_n^{(\lambda)}(y; x)$ that holds in a neighborhood of $\varphi_{\nu^*}^{(\lambda)}(y; x)$.
 623 Lastly, in Step 3, we utilize the peeling technique from empirical process theory to obtain the desired
 624 rate of convergence.

625 **Step 1: Consistency.** We first claim that $d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(\lambda)}(x)) = o_P(1)$ by an argument similar to
 626 that used in the proof of Theorem 1. Specifically, it suffices to show that

627 (S1') $R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(\lambda)}(y; x) = o_P(1)$, and

628 (S2') $R_{\mathcal{D}_n}^{(\lambda)}(\cdot; x) : \mathcal{M} \rightarrow \mathbb{R}$ is asymptotically equicontinuous in probability.

629 Note that we already showed the asymptotic equicontinuity in the proof of Theorem 1; see (S2).
 630 Thus, it remains to show the pointwise convergence in probability. To show (S1'), we decompose
 631 $R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(\lambda)}(y; x)$ as follows.

$$\begin{aligned} R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(\lambda)}(y; x) &= \{R_{\mathcal{D}_n}^{(\lambda)}(y; x) - \tilde{R}_n(y; x)\} + \{\tilde{R}_n(y; x) - R_{\nu^*}^{(\lambda)}(y; x)\} \\ &= \underbrace{\frac{1}{n} \sum_{i=1}^n \{w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) - w_{\nu^*}^{(\lambda)}(X_i, x)\} \cdot d^2(Y_i, y)}_{:= A_n^{(\lambda)}(y; x)} \\ &\quad + \underbrace{\frac{1}{n} \sum_{i=1}^n \left(w_{\nu^*}^{(\lambda)}(X_i, x) \cdot d^2(Y_i, y) - \mathbb{E}_{\nu^*} [w_{\nu^*}^{(\lambda)}(X_i, x) \cdot d^2(Y_i, y)] \right)}_{:= B_n^{(\lambda)}(y; x)}. \end{aligned}$$

632 Next, we show that $A_n^{(\lambda)}(y; x)$ and $B_n^{(\lambda)}(y; x)$ respectively converge to 0 in probability.

633 • Letting $\hat{\mu}_n = \mu_{\mathcal{D}_n}$, $\hat{\Sigma}_n = \Sigma_{\mathcal{D}_n}$, and $\hat{\Sigma}_n^{(\lambda)} = \text{SVT}^{(\lambda)}(\hat{\Sigma}_n)$ for shorthand, we can write

$$w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) - w_{\nu^*}^{(\lambda)}(X_i, x) = V_n^{(\lambda)}(x) + X_i^\top W_n^{(\lambda)}(x),$$

634

similarly to (26), where

$$\begin{aligned}
V_n^{(\lambda)}(x) &= -\widehat{\mu}_n^\top [\widehat{\Sigma}_n^{(\lambda)}]^\dagger (x - \widehat{\mu}_n) + \mu^\top [\Sigma^{(\lambda)}]^\dagger (x - \mu), \\
W_n^{(\lambda)}(x) &= [\widehat{\Sigma}_n^{(\lambda)}]^\dagger (x - \widehat{\mu}_n) - [\Sigma^{(\lambda)}]^\dagger (x - \mu).
\end{aligned} \tag{34}$$

635

Since $\|\widehat{\mu}_n - \mu\|_2 = O_P(n^{-1/2})$ and $\|\widehat{\Sigma}_n^{(\lambda)} - \Sigma^{(\lambda)}\| = O_P(n^{-1/2})$ (if $\lambda \notin \text{spec } \Sigma$) independent of $\lambda > 0$, we also have $|V_n^{(\lambda)}(x)| = O_P(n^{-1/2})$ and $\|W_n^{(\lambda)}(x)\|_2 = O_P(n^{-1/2})$. This implies that $A_n^{(\lambda)}(y; x) = o_P(1)$.

637

638

- Moreover, we note that if $\|x - \mu\|_\Sigma < \infty$, then the random variable $w_{\nu^*}^{(\lambda)}(X, x)$ has finite second moment

639

$$\begin{aligned}
\mathbb{E}_{\nu^*} [w_{\nu^*}^{(\lambda)}(X, x)^2] &\leq 2 \left(1 + \mathbb{E}_{\nu^*} \left[\left| (X - \mu)^\top [\Sigma^{(\lambda)}]^\dagger (x - \mu) \right|^2 \right] \right) \\
&\leq 2 \left(1 + \mathbb{E}_{\nu^*} \left[\|X - \mu\|_{\Sigma^{(\lambda)}}^2 \cdot \|x - \mu\|_{\Sigma^{(\lambda)}}^2 \right] \right) \\
&\leq 2 \{ 1 + p \|x - \mu\|_\Sigma^2 \},
\end{aligned} \tag{35}$$

640

regardless of the value of $\lambda > 0$. When $\text{diam}(\mathcal{M}) < \infty$, the product $w_{\nu^*}^{(\lambda)}(X, x) \cdot d^2(Y, y)$ also has finite second moment. Since $B_n^{(\lambda)}(y; x)$ is the sample mean of IID random variables with mean zero and finite variance, it follows that

641

642

$$B_n^{(\lambda)}(y; x) = O_P \left(\sqrt{\frac{\text{Var}[w_{\nu^*}^{(\lambda)}(X_1, x) \cdot d^2(Y_1, y)]}{n}} \right) = O_P(n^{-1/2}).$$

643

Step 2: Uniform control of the fluctuation in objective discrepancy. For any $\lambda \in \mathbb{R}_+$ and any

644

$(x, y) \in \mathbb{R}^p \times \mathcal{M}$, we let $Z_n^{(\lambda)}(y; x)$ denote the random variable defined as

$$Z_n^{(\lambda)}(y; x) := R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(\lambda)}(y; x) \tag{36}$$

645

We observed that

$$\begin{aligned}
&Z_n^{(\lambda)}(y; x) - Z_n^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) \\
&= \left\{ R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\nu^*}^{(\lambda)}(y; x) \right\} - \left\{ R_{\mathcal{D}_n}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) - R_{\nu^*}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) \right\} \\
&= \left[\left\{ R_{\mathcal{D}_n}^{(\lambda)}(y; x) - \tilde{R}_n(y; x) \right\} - \left\{ R_{\mathcal{D}_n}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) - \tilde{R}_n(\varphi_{\nu^*}^{(\lambda)}(x); x) \right\} \right] \\
&\quad + \left[\left\{ \tilde{R}_n(y; x) - R_{\nu^*}^{(\lambda)}(y; x) \right\} - \left\{ \tilde{R}_n(\varphi_{\nu^*}^{(\lambda)}(x); x) - R_{\nu^*}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) \right\} \right] \\
&= \frac{1}{n} \sum_{i=1}^n \left\{ w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) - w_{\nu^*}^{(\lambda)}(X_i, x) \right\} \cdot \ell_i^{(\lambda)}(y; x) \\
&\quad \underbrace{\hspace{10em}}_{=: \mathfrak{A}_n^{(\lambda)}(y; x)} \\
&\quad + \frac{1}{n} \sum_{i=1}^n \left(w_{\nu^*}^{(\lambda)}(X_i, x) \cdot \ell_i^{(\lambda)}(y; x) - \mathbb{E}_{\nu^*} \left[w_{\nu^*}^{(\lambda)}(X_i, x) \cdot \ell_i^{(\lambda)}(y; x) \right] \right) \\
&\quad \underbrace{\hspace{10em}}_{=: \mathfrak{B}_n^{(\lambda)}(y; x)}
\end{aligned} \tag{37}$$

646

where $\ell_i^{(\lambda)}(y; x) := d^2(Y_i, y) - d^2(Y_i, \varphi_{\nu^*}^{(\lambda)}(x))$.

647

Next, we analyze the asymptotic behavior of the two terms, $\mathfrak{A}_n^{(\lambda)}(y; x)$ and $\mathfrak{B}_n^{(\lambda)}(y; x)$. Specifically, we establish upper bounds on their magnitudes that hold uniformly over a δ -neighborhood of

648

$\varphi^{(\lambda)}(x) = \varphi_{\nu^*}^{(\lambda)}(x)$, which will be used later in Step 3 of this proof.

649

650

- Firstly, we observe that for any $\delta > 0$,

$$\begin{aligned}
& \sup_{y \in B_d(\varphi_{\nu^*}^{(\lambda)}(x); \delta)} |\mathfrak{A}_n^{(\lambda)}(y; x)| \\
& \leq \frac{1}{n} \sum_{i=1}^n |w_{\mathcal{D}_n}^{(\lambda)}(X_i, x) - w_{\nu^*}^{(\lambda)}(X_i, x)| \cdot \sup_{y \in B_d(\varphi_{\nu^*}^{(\lambda)}(x); \delta)} |d^2(Y_i, y) - d^2(Y_i, \varphi_{\nu^*}^{(\lambda)}(x))| \\
& \leq 2 \operatorname{diam}(\mathcal{M}) \cdot \left\{ \frac{1}{n} \sum_{i=1}^n \left\{ |V_n^{(\lambda)}(x)| + \|X_i\|_2 \|W_n^{(\lambda)}(x)\|_2 \right\} \right\} \\
& \quad \times \sup_{y \in B_d(\varphi_{\nu^*}^{(\lambda)}(x); \delta)} d(y, \varphi_{\nu^*}^{(\lambda)}(x)) \\
& = O_P\left(\delta \cdot n^{-1/2}\right), \tag{38}
\end{aligned}$$

651

where we used the property of $V_n^{(\lambda)}(x)$ and $W_n^{(\lambda)}(x)$ discussed in the paragraph following

652

(34). Since the stochastic magnitudes of $V_n^{(\lambda)}(x)$ and $W_n^{(\lambda)}(x)$ are independent of δ , (38)

653

implies that there exists $C_1^{(\lambda)} = C_1^{(\lambda)}(x) > 0$ such that for any $\delta > 0$,

$$\liminf_{n \rightarrow \infty} P\left(\sup_{y \in \mathcal{M}} \left\{ |\mathfrak{A}_n^{(\lambda)}(y; x)| : d(y, \varphi_{\nu^*}^{(\lambda)}(x)) < \delta \right\} \leq C_1^{(\lambda)} \cdot \delta \cdot n^{-1/2}\right) = 1. \tag{39}$$

654

Furthermore, for any $\gamma, \delta \in \mathbb{R}_+$ such that $0 \leq \gamma < \delta$, let $\mathfrak{E}_n^{(\lambda)}(\gamma, \delta; x)$ be defined as an event such that

655

$$\mathfrak{E}_n(\gamma, \delta; x) = \left(\sup_{y \in \mathcal{M}} \left\{ |\mathfrak{A}_n^{(\lambda)}(y; x)| : d(y, \varphi_{\nu^*}^{(\lambda)}(x)) \in [\gamma, \delta] \right\} \leq C_1^{(\lambda)} \cdot \delta \cdot n^{-1/2} \right). \tag{40}$$

656

For any $\gamma \in [0, \delta]$, we have $\mathfrak{E}_n(0, \delta; x) \subseteq \mathfrak{E}_n(\gamma, \delta; x)$, and thus,

657

$$\liminf_{n \rightarrow \infty} P(\mathfrak{E}_n(\gamma, \delta; x)) = 1.$$

658

- Next, we note that

$$|w_{\nu^*}^{(\lambda)}(X_i, x) \cdot \ell_i^{(\lambda)}(y; x)| \leq 2 \operatorname{diam}(\mathcal{M}) \cdot d(y, \varphi_{\nu^*}^{(\lambda)}(x)) \cdot |w_{\nu^*}^{(\lambda)}(X_i, x)|.$$

659

Observe that $d(y, \varphi_{\nu^*}^{(\lambda)}(x)) \leq \operatorname{diam}(\mathcal{M}) < \infty$ and recall that $\mathbb{E}_{\nu^*} \left[w_{\nu^*}^{(\lambda)}(X, x)^2 \right] \leq 2\{1 + p \|x - \mu\|_\Sigma^2\}$ as shown in Step 1 of this proof, cf. (35). It follows from the uniform entropy condition (C2), Theorem 2.7.11, and Theorem 2.14.2 in [55] that there exists $D_e = D_e(x) > 0$ such that for all $\delta \in [0, D_e]$,

662

$$\begin{aligned}
& \mathbb{E} \left[\sup_{y \in \mathcal{M}} \left\{ |\mathfrak{B}_n^{(\lambda)}(y; x)| : d(y, \varphi_{\nu^*}^{(\lambda)}(x)) < \delta \right\} \right] \\
& \leq 2 \operatorname{diam}(\mathcal{M}) \cdot \delta \cdot n^{-1/2} \sqrt{1 + p \|x - \mu\|_\Sigma^2} \int_0^1 \sqrt{1 + \log \mathfrak{N}(B_d(\varphi^{(\lambda)}(x); \delta), \delta \epsilon)} d\epsilon \\
& \leq C_2^{(\lambda)} \cdot \delta \cdot n^{-1/2} \tag{41}
\end{aligned}$$

663

where $C_2^{(\lambda)} = 2(C_e + 1) \cdot \operatorname{diam}(\mathcal{M}) \cdot \sqrt{1 + p \|x - \mu\|_\Sigma^2}$ is independent of $\delta > 0$ and $n \geq 1$.

664

Step 3: Concluding the proof. Lastly, we combine the results from Steps 1-2 to show that, for any $\eta > 0$, there exist $K = K(\eta) > 0$ and $N = N(\eta) \geq 1$ such that $P\left(d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(\lambda)}(x)) > 2^K n^{-\beta}\right) < \eta$ for any $n \geq N$, where $\beta > 0$ is an absolute constant that will be determined later in this proof. We prove this claim using the peeling technique, in a similar manner as we did in the proof of Lemma 1. To avoid cluttered notation, we let $\Delta(x) = d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\nu^*}^{(\lambda)}(x))$ in the rest of this proof.

670

671 For any fixed $K \in \mathbb{N}$ and a sufficiently large $n = n(K) \geq 1$ satisfying $2^K n^{-\beta} < D_* := D_g \wedge D_e$,
 672 we observe that

$$P\left(\Delta(x) > 2^K n^{-\beta}\right) = P\left(\Delta(x) \geq D_*\right) + P\left(2^K n^{-\beta} \leq \Delta(x) < D_*\right) \quad (42)$$

673 where we used $P(A) \leq P(B^c) + P(A \cap B)$ to get the inequality. As we know that $P\left(\Delta(x) \geq\right.$
 674 $D_*) = o(1)$ by Step 1 of this proof, we focus on showing an upper bound for the other term,
 675 $P(2^K n^{-\beta} \leq \Delta(x) < D_*)$.

676 Step 3-A: Decomposition of $P(2^K n^{-\beta} \leq \Delta(x) < D_*)$. For each $n, k \in \mathbb{N}$, we define

$$\begin{aligned} \mathfrak{F}_{n,k} &= \bigcap_{k'=K}^k \mathfrak{E}_n^{(\lambda)}(2^{k'} n^{-\beta}, 2^{k'+1} n^{-\beta} \wedge D_*; x), \\ \mathfrak{G}_{n,k} &= \left(\bigcap_{k'=K}^{k-1} \mathfrak{E}_n^{(\lambda)}(2^{k'} n^{-\beta}, 2^{k'+1} n^{-\beta} \wedge D_*; x) \right) \cap \mathfrak{E}_n^{(\lambda)}(2^k n^{-\beta}, 2^{k+1} n^{-\beta} \wedge D_*; x)^c, \end{aligned} \quad (43)$$

677 where we set $\mathfrak{F}_{n,K-1}$ to be the entire event space so that $\mathfrak{G}_{n,K} = (\mathfrak{F}_{n,K})^c$. It is worth mentioning
 678 that $\mathfrak{G}_{n,k}$ and $\mathfrak{G}_{n,k'}$ are mutually exclusive for any $k \neq k' \geq K$, and we will use this property when
 679 concluding the proof in Step 3-C below.

680 Now, we observe that

$$\begin{aligned} P\left(2^K n^{-\beta} \leq \Delta(x) < D_*\right) &\leq P\left(\mathfrak{E}_n^{(\lambda)}(2^K n^{-\beta}, 2^{K+1} n^{-\beta} \wedge D_*; x)^c\right) \\ &\quad + P\left(\left(2^K n^{-\beta} \leq \Delta(x) < D_*\right) \cap \mathfrak{E}_n^{(\lambda)}(2^K n^{-\beta}, 2^{K+1} n^{-\beta} \wedge D_*; x)\right) \\ &= P(\mathfrak{G}_{n,K}) + P\left(\left(2^K n^{-\beta} \leq \Delta(x) < D_*\right) \cap \mathfrak{F}_{n,K}\right) \\ &= P(\mathfrak{G}_{n,K}) + P\left(\left(2^K n^{-\beta} \leq \Delta(x) < 2^{K+1} n^{-\beta} \wedge D_*\right) \cap \mathfrak{F}_{n,K}\right) \\ &\quad + P\left(\left(2^{K+1} n^{-\beta} \leq \Delta(x) < D_*\right) \cap \mathfrak{F}_{n,K}\right) \end{aligned}$$

681 and that for every $k \geq K$,

$$P\left(\left(2^{k+1} n^{-\beta} \leq \Delta(x) < D_*\right) \cap \mathfrak{F}_{n,k}\right) \leq P\left(\left(2^{k+1} n^{-\beta} \leq \Delta(x) < D_*\right) \cap \mathfrak{F}_{n,k+1}\right) + P(\mathfrak{G}_{n,k+1}).$$

682 As a result, we have

$$P\left(2^K n^{-\beta} \leq \Delta(x) < D_*\right) = \sum_{k=K}^{\infty} P(\mathfrak{G}_{n,k}) + \underbrace{\sum_{k=K}^{\infty} P\left(\left(2^k n^{-\beta} \leq \Delta(x) < 2^{k+1} n^{-\beta} \wedge D_*\right) \cap \mathfrak{F}_{n,k}\right)}_{=: \mathfrak{C}_{n,k}}. \quad (44)$$

683 **Step 3-B: Controlling $\mathfrak{C}_{n,k}$.** Next, we show an upper bound for $\mathfrak{C}_{n,k}$. Suppose that $2^k n^{-\beta} \leq \Delta(x) <$
684 $2^{k+1} n^{-\beta} \wedge D_*$ and the event $\mathfrak{F}_{n,k}$ occurs. Then it follows from Assumption (C1) that

$$\begin{aligned}
& C_g \cdot \Delta(x)^\alpha \\
& \leq R_{\nu^*}^{(\lambda)}(\varphi_{\mathcal{D}_n}^{(\lambda)}(x); x) - R_{\nu^*}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) \\
& \leq \left\{ R_{\nu^*}^{(\lambda)}(\varphi_{\mathcal{D}_n}^{(\lambda)}(x); x) - R_{\nu^*}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) \right\} + \underbrace{\left\{ R_{\mathcal{D}_n}^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) - R_{\mathcal{D}_n}^{(\lambda)}(\varphi_{\mathcal{D}_n}^{(\lambda)}(x); x) \right\}}_{\geq 0} \\
& = Z_n^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) - Z_n^{(\lambda)}(\varphi_{\mathcal{D}_n}^{(\lambda)}(x); x) \quad \text{cf. (36)} \\
& \leq \left| \mathfrak{A}_n^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) \right| + \left| \mathfrak{B}_n^{(\lambda)}(\varphi_{\nu^*}^{(\lambda)}(x); x) \right| \quad \therefore (37) \\
& \leq \sup_{y \in \mathcal{M}} \left\{ \left| \mathfrak{A}_n^{(\lambda)}(y; x) \right| + \left| \mathfrak{B}_n^{(\lambda)}(y; x) \right| : 2^k n^{-\beta} \leq d(y, \varphi_{\nu^*}^{(\lambda)}(x)) < 2^{k+1} n^{-\beta} \wedge D_* \right\} \\
& \leq C_1^{(\lambda)} \cdot (2^{k+1} n^{-\beta} \wedge D_*) \cdot n^{-1/2} + \sup_{y \in \mathcal{M}} \left\{ \left| \mathfrak{B}_n^{(\lambda)}(y; x) \right| : d(y, \varphi_{\nu^*}^{(\lambda)}(x)) < 2^{k+1} n^{-\beta} \wedge D_* \right\}. \quad \therefore (40) \\
& \quad \quad \quad (45)
\end{aligned}$$

685 Therefore, we obtain that for each $k \geq K$,

$$\begin{aligned}
\mathfrak{C}_{n,k} &= P\left(\left(2^k n^{-\beta} \leq \Delta(x) < 2^{k+1} n^{-\beta} \wedge D_*\right) \cap \mathfrak{F}_{n,k}\right) \\
&\leq P\left(\left(\Delta(x)^\alpha \geq (2^k n^{-\beta})^\alpha\right) \cap \mathfrak{F}_{n,k}\right) \\
&\leq \frac{C_1^{(\lambda)} \cdot (2^{k+1} n^{-\beta} \wedge D_*) \cdot n^{-1/2} + \mathbb{E}\left[\sup_{y \in \mathcal{M}} \left\{ \left| \mathfrak{B}_n^{(\lambda)}(y; x) \right| : d(y, \varphi_{\nu^*}^{(\lambda)}(x)) < 2^{k+1} n^{-\beta} \wedge D_* \right\}\right]}{C_g \cdot (2^k n^{-\beta})^\alpha} \\
&\quad \quad \quad \therefore (45) \text{ \& Markov's inequality} \\
&\leq \frac{(C_1^{(\lambda)} + C_2^{(\lambda)}) \cdot (2^{k+1} n^{-\beta} \wedge D_*) \cdot n^{-1/2}}{C_g \cdot (2^k n^{-\beta})^\alpha} \quad \therefore (41) \quad (46)
\end{aligned}$$

686 **Step 3-C: Concluding Step 3.** Combining (42), (44), and (46), we have

$$\begin{aligned}
P\left(\Delta(x) > 2^K n^{-\beta}\right) &\leq \frac{2(C_1^{(\lambda)} + C_2^{(\lambda)})}{C_g} n^{-\frac{1}{2} + \beta(\alpha-1)} \sum_{k=K}^{\infty} 2^{-k(\alpha-1)} \\
&\quad + \underbrace{P\left(\Delta(x) \geq D_*\right)}_{=o(1) \therefore \text{Step 1 of this proof}} + \sum_{k=K}^{\infty} P\left(\mathfrak{G}_{n,k}\right).
\end{aligned}$$

687 Moreover, $\mathfrak{G}_{n,k}$ are mutually exclusive, and thus,

$$\sum_{k=K}^{\infty} P\left(\mathfrak{G}_{n,k}\right) = P\left(\bigcup_{k=K}^{\infty} \mathfrak{G}_{n,k}\right) = P\left(\left(\bigcup_{k=K}^{\infty} \mathfrak{C}_n^{(\lambda)}(2^k n^{-\beta}, 2^{k+1} n^{-\beta} \wedge D_*; x)\right)^c\right) \rightarrow 0 \quad \therefore (40)$$

688 Finally, we obtain the desired result by letting $\beta = \frac{1}{2(\alpha-1)}$.

689 □

690 D Proof of Theorem 3

691 In this section, we prove Theorem 4 that establishes an upper bound on $d(\varphi_{\mathcal{D}_n}^{(\lambda)}(x), \varphi_{\mathcal{D}_n}^{(\lambda)}(x))$. This
692 section is organized as follows. Firstly, in Section D.1, we present several useful results from matrix
693 perturbation theory as lemmas. Next, in Section D.2, we provide a key lemma (Lemma 6) that
694 establishes the stability of the weight function when there is covariate noise. Lastly, in Section D.3,
695 we state and prove Theorem 4, from which Theorem 3 can be easily derived.

696 D.1 Useful lemmas

697 **Definition 7.** Let $n, p \in \mathbb{N}$ and let $M \in \mathbb{R}^{n \times p}$. The row projection matrix for M , denoted by
698 $\Pi_M^{\text{row}} \in \mathbb{R}^{p \times p}$, is a matrix such that

$$\Pi_M^{\text{row}} := M^\dagger \cdot M. \quad (47)$$

699 and the column projection matrix for M , denoted by $\Pi_M^{\text{col}} \in \mathbb{R}^{n \times n}$, is a matrix such that

$$\Pi_M^{\text{col}} := M \cdot M^\dagger. \quad (48)$$

700 We recall from (6) that for any $\lambda \in \mathbb{R}_+$, the singular value thresholding (SVT) operator $\text{SVT}^{(\lambda)}$ is
701 defined such that

$$M = \sum_{i=1}^{\min\{n,p\}} s_i \cdot u_i v_i^\top \text{ is a SVD} \quad \mapsto \quad \text{SVT}^{(\lambda)}(M) = \sum_{i=1}^{\min\{n,p\}} s_i \cdot \mathbb{1}\{s_i > \lambda\} \cdot u_i v_i^\top.$$

702 In the rest of this section, we let $M^{(\lambda)} := \text{SVT}^{(\lambda)}(M)$ for shorthand.

703 **Lemma 3** (Properties of the row/column projection matrices). Let $n, p \in \mathbb{N}$, and $M \in \mathbb{R}^{n \times p}$. For
704 any $\lambda \in \mathbb{R}_+$, the following statements are true.

- 705 1. $\Pi_{M^{(\lambda)}}^{\text{row}}$ defines a projection in \mathbb{R}^p and $\text{rank } \Pi_{M^{(\lambda)}}^{\text{row}} = \text{rank } M^{(\lambda)}$.
- 706 2. $\Pi_{M^{(\lambda)}}^{\text{col}}$ defines a projection in \mathbb{R}^n and $\text{rank } \Pi_{M^{(\lambda)}}^{\text{col}} = \text{rank } M^{(\lambda)}$.
- 707 3. $M \Pi_{M^{(\lambda)}}^{\text{row}} M^\dagger = \Pi_{M^{(\lambda)}}^{\text{col}}$ and $M^\dagger \Pi_{M^{(\lambda)}}^{\text{col}} M = \Pi_{M^{(\lambda)}}^{\text{row}}$.

708 *Proof.* Let $r = \text{rank } M$ and consider a compact singular value decomposition (SVD) of M :

$$M = \sum_{i=1}^r s_i \cdot u_i v_i^\top$$

709 where s_1, \dots, s_r are non-zero singular values of M . Noticing that

$$M^{(\lambda)} = \text{SVT}^{(\lambda)}(M) = \sum_{i=1}^r \mathbb{1}\{s_i > \lambda\} \cdot u_i v_i^\top$$

710 and that $M^\dagger = \sum_{i=1}^r s_i^{-1} \cdot v_i u_i^\top$, the three conclusions of the lemma follow straightforwardly from
711 the orthonormality of singular vectors.

- 712 • $\Pi_{M^{(\lambda)}}^{\text{row}} = \sum_{i=1}^r v_i v_i^\top \cdot \mathbb{1}\{s_i > \lambda\}$ is the projection onto the row space of $M^{(\lambda)}$.
- 713 • $\Pi_{M^{(\lambda)}}^{\text{col}} = \sum_{i=1}^r u_i u_i^\top \cdot \mathbb{1}\{s_i > \lambda\}$ is the projection onto the column space of $M^{(\lambda)}$.
- 714 • Due to the orthonormality of singular vectors,

$$\begin{aligned} M \Pi_{M^{(\lambda)}}^{\text{row}} M^\dagger &= \left(\sum_{i=1}^r s_i \cdot u_i v_i^\top \right) \left(\sum_{i=1}^r v_i v_i^\top \cdot \mathbb{1}\{s_i > \lambda\} \right) \left(\sum_{i=1}^r s_i^{-1} \cdot v_i u_i^\top \right) \\ &= \sum_{i=1}^r u_i u_i^\top \cdot \mathbb{1}\{s_i > \lambda\} \\ &= \Pi_{M^{(\lambda)}}^{\text{col}}, \end{aligned}$$

715 and likewise, $M^\dagger \Pi_{M^{(\lambda)}}^{\text{col}} M = \Pi_{M^{(\lambda)}}^{\text{row}}$.

716 □

717 In addition, we collect two classical results from matrix perturbation theory and state them as lemmas.

718 **Lemma 4** ([51, Theorem 3.2]). Let $\mathbf{X}, \mathbf{Z} \in \mathbb{R}^{n \times p}$. Then the following equation is true:

$$\mathbf{Z}^\dagger - \mathbf{X}^\dagger = -\mathbf{Z}^\dagger \Pi_{\mathbf{Z}}^{\text{col}} (\mathbf{Z} - \mathbf{X}) \Pi_{\mathbf{X}}^{\text{row}} \mathbf{X}^\dagger + \mathbf{Z}^\dagger \Pi_{\mathbf{Z}}^{\text{col}} \Pi_{\mathbf{X}}^{\text{col}^\perp} - \Pi_{\mathbf{Z}}^{\text{row}^\perp} \Pi_{\mathbf{X}}^{\text{row}} \mathbf{X}^\dagger \quad (49)$$

719 where $\Pi_{\mathbf{X}}^{\text{col}^\perp} = \mathbf{I}_n - \Pi_{\mathbf{X}}^{\text{col}}$ and $\Pi_{\mathbf{Z}}^{\text{row}^\perp} = \mathbf{I}_p - \Pi_{\mathbf{Z}}^{\text{row}}$.

720 **Lemma 5** ([15, Theorems 2.4 & 2.5]). Let $\mathbf{X}, \mathbf{Z} \in \mathbb{R}^{n \times p}$. Then

$$\|\Pi_{\mathbf{Z}}^{\text{col}} - \Pi_{\mathbf{X}}^{\text{col}}\| \leq \max \left\{ \left\| (\mathbf{Z} - \mathbf{X}) \mathbf{X}^\dagger \right\|, \left\| (\mathbf{Z} - \mathbf{X}) \mathbf{Z}^\dagger \right\| \right\}. \quad (50)$$

721 Moreover, if $\text{rank } \mathbf{X} = \text{rank } \mathbf{Z}$, then

$$\|\Pi_{\mathbf{Z}}^{\text{col}} - \Pi_{\mathbf{X}}^{\text{col}}\| \leq \min \left\{ \left\| (\mathbf{Z} - \mathbf{X}) \mathbf{X}^\dagger \right\|, \left\| (\mathbf{Z} - \mathbf{X}) \mathbf{Z}^\dagger \right\| \right\}. \quad (51)$$

722 D.2 Stability of the weights under (small) perturbation in covariates

723 Let $\mathcal{D}_n = \{(x_i, y_i) \in \mathbb{R}^p \times \mathcal{M} : i \in [n]\}$ and $\tilde{\mathcal{D}}_n = \{(z_i, y_i) \in \mathbb{R}^p \times \mathcal{M} : i \in [n]\}$ be two sets in
 724 $\mathbb{R}^p \times \mathcal{M}$. We may identify these sets with their empirical distributions. Recall the definition of $w_\nu^{(\lambda)}$
 725 from (9): for any probability measure ν on $\mathbb{R}^p \times \mathcal{M}$, any $\lambda \in \mathbb{R}_+$, and any $x, x' \in \mathbb{R}^p$,

$$w_\nu^{(\lambda)}(x', x) = 1 + (x' - \mu_\nu)^\top \left[\text{SVT}^{(\lambda)}(\Sigma_\nu) \right]^\dagger (x - \mu_\nu)$$

726 where $\mu_\nu = \mathbb{E}_{(X,Y) \sim \nu}(X)$ and $\Sigma_\nu = \text{Var}_{(X,Y) \sim \nu}(X)$, cf. (7). We define the weight vectors induced
 727 by \mathcal{D}_n and $\tilde{\mathcal{D}}_n$ as follows: for any $\lambda \in \mathbb{R}_+$ and any $x \in \mathbb{R}^p$,

$$\begin{aligned} \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) &:= \begin{bmatrix} w_{\mathcal{D}_n}^{(\lambda)}(x_1, x) & \cdots & w_{\mathcal{D}_n}^{(\lambda)}(x_n, x) \end{bmatrix} \in \mathbb{R}^n, \\ \vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) &:= \begin{bmatrix} w_{\tilde{\mathcal{D}}_n}^{(\lambda)}(z_1, x) & \cdots & w_{\tilde{\mathcal{D}}_n}^{(\lambda)}(z_n, x) \end{bmatrix} \in \mathbb{R}^n. \end{aligned} \quad (52)$$

728 **Lemma 6** (Stability of weights). Let $\mathcal{D}_n = \{(x_i, y_i) \in \mathbb{R}^p \times \mathcal{M} : i \in [n]\}$ and $\tilde{\mathcal{D}}_n =$
 729 $\{(z_i, y_i) \in \mathbb{R}^p \times \mathcal{M} : i \in [n]\}$. Let $\mathbf{X} = [x_1 \cdots x_n]^\top \in \mathbb{R}^{n \times p}$ and $\mathbf{Z} = [z_1 \cdots z_n]^\top \in$
 730 $\mathbb{R}^{n \times p}$. For any $\lambda \in \mathbb{R}_+$, if $x \in \mathbb{R}^p$ satisfies $x - \mu_{\mathcal{D}_n} \in \text{rowsp}(\mathbf{X}_{\text{ctr}})$, then

$$\|\vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) - \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x)\| \leq \frac{\sqrt{n} \cdot \|\mathbf{Z} - \mathbf{X}\|}{\min \{ \sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}}) \}} \cdot \left(2 \cdot \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} + 1 \right) \quad (53)$$

731 where $\mathbf{X}_{\text{ctr}} = (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \mathbf{X}$ and $\sigma^{(\lambda)}(\mathbf{X}) := \inf \{ \sigma_i(\mathbf{X}) > \lambda : i \in \mathbb{N} \}$ (likewise for \mathbf{Z}).

732 *Proof of Lemma 6.* This proof consists of three steps. In Step 1, we express the weight discrepancy
 733 $\vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) - \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x)$ as a sum of matrix products using projections. In Step 2, we establish upper
 734 bounds on the norm of the expression obtained in Step 1. In Step 3, we collect intermediate results
 735 together and conclude the proof.

736 **Step 1: Decomposition of the weight discrepancy.** First of all, we rewrite $\vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) - \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x)$ in
 737 a compact matrix representation that is presented in (61) at the end of this step. To this end, we begin
 738 by observing that

$$\mu_{\mathcal{D}_n} = \frac{1}{n} \mathbf{X}^\top \mathbf{1}_n, \quad \text{and} \quad \Sigma_{\mathcal{D}_n} = \frac{1}{n} (\mathbf{X} - \mathbf{1}_n \mu_{\mathcal{D}_n}^\top)^\top (\mathbf{X} - \mathbf{1}_n \mu_{\mathcal{D}_n}^\top) = \frac{1}{n} \mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}}. \quad (54)$$

739 For given $\lambda \in \mathbb{R}_+$, we let $\mathbf{X}_{\text{ctr}}^{(\lambda)} := \text{SVT}^{(\lambda)}(\mathbf{X}_{\text{ctr}})$, and observe that

$$\Sigma_{\mathcal{D}_n}^{(\lambda)} = \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \left(\frac{1}{n} \mathbf{X}_{\text{ctr}}^\top \mathbf{X}_{\text{ctr}} \right) \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} = \frac{1}{n} \cdot \mathbf{X}_{\text{ctr}}^{(\lambda)\top} \cdot \mathbf{X}_{\text{ctr}}^{(\lambda)}. \quad (55)$$

740 Then it follows that

$$\left[\Sigma_{\mathcal{D}_n}^{(\lambda)} \right]^\dagger = n \cdot \left[\mathbf{X}_{\text{ctr}}^{(\lambda)\top} \cdot \mathbf{X}_{\text{ctr}}^{(\lambda)} \right]^\dagger = n \cdot \left[\mathbf{X}_{\text{ctr}}^{(\lambda)} \right]^\dagger \cdot \left[\mathbf{X}_{\text{ctr}}^{(\lambda)\top} \right]^\dagger = n \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \mathbf{X}_{\text{ctr}}^\dagger \cdot \left(\mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}}.$$

Therefore, we have

$$\begin{aligned}
\vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) &= \mathbf{1}_n + (\mathbf{X} - \mathbf{1}_n \mu_{\mathcal{D}_n}^\top) \cdot \left[\Sigma_{\mathcal{D}_n}^{(\lambda)} \right]^\dagger \cdot (x - \mu_{\mathcal{D}_n}) \\
&= \mathbf{1}_n + n \cdot \mathbf{X}_{\text{ctr}} \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \mathbf{X}_{\text{ctr}}^\dagger \cdot \left(\mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot (x - \mu_{\mathcal{D}_n}) \\
&= \mathbf{1}_n + n \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot (x - \mu_{\mathcal{D}_n}), \tag{56}
\end{aligned}$$

where the equality in the last line follows from Lemma 3: $\mathbf{X}_{\text{ctr}} \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \mathbf{X}_{\text{ctr}}^\dagger = \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}}$.

Likewise, we repeat the above for $\tilde{\mathcal{D}}_n$ and \mathbf{Z} to write

$$\mu_{\tilde{\mathcal{D}}_n} = \frac{1}{n} \mathbf{Z}^\top \mathbf{1}_n \quad \text{and} \quad \Sigma_{\tilde{\mathcal{D}}_n} = \frac{1}{n} \mathbf{Z}_{\text{ctr}}^\top \mathbf{Z}_{\text{ctr}}.$$

Then, we obtain an expression for $\vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x)$ in a similar form to (56), namely,

$$\vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) = \mathbf{1}_n + n \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{Z}_{\text{ctr}}^\top \right)^\dagger \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot (x - \mu_{\tilde{\mathcal{D}}_n}). \tag{57}$$

Thereafter, we define $c_x, \tilde{c}_x \in \mathbb{R}^{n \times 1}$ so that

$$\begin{aligned}
c_x &= \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} = \left(\frac{1}{\sqrt{n}} \mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot (x - \mu_{\mathcal{D}_n}) \quad \text{and} \\
\tilde{c}_x &= \|x - \mu_{\tilde{\mathcal{D}}_n}\|_{\Sigma_{\tilde{\mathcal{D}}_n}} = \left(\frac{1}{\sqrt{n}} \mathbf{Z}_{\text{ctr}}^\top \right)^\dagger \cdot (x - \mu_{\tilde{\mathcal{D}}_n}). \tag{58}
\end{aligned}$$

Then we observe that for any $x \in \mathbb{R}^p$,

$$n \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot (x - \mu_{\mathcal{D}_n}) = n \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \frac{1}{\sqrt{n}} \mathbf{X}_{\text{ctr}}^\top \cdot \left(\frac{1}{\sqrt{n}} \mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot (x - \mu_{\mathcal{D}_n}) = \sqrt{n} \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \mathbf{X}_{\text{ctr}}^\top \cdot c_x. \tag{59}$$

Likewise,

$$n \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot (x - \mu_{\tilde{\mathcal{D}}_n}) = n \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \frac{1}{\sqrt{n}} \mathbf{Z}_{\text{ctr}}^\top \cdot \left(\frac{1}{\sqrt{n}} \mathbf{Z}_{\text{ctr}}^\top \right)^\dagger \cdot (x - \mu_{\tilde{\mathcal{D}}_n}) = \sqrt{n} \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \mathbf{Z}_{\text{ctr}}^\top \cdot \tilde{c}_x. \tag{60}$$

Consequently, for any $x \in \mathbb{R}^p$, we obtain from (56) and (57) with aid of (59) and (60) that

$$\begin{aligned}
\vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) - \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) &= \sqrt{n} \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{Z}_{\text{ctr}}^\top \right)^\dagger \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \mathbf{Z}_{\text{ctr}}^\top \cdot \tilde{c}_x - \sqrt{n} \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{row}} \cdot \mathbf{X}_{\text{ctr}}^\top \cdot c_x \\
&= \sqrt{n} \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \tilde{c}_x - \sqrt{n} \cdot \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot c_x \quad \quad \quad \because \text{Lemma 3} \\
&= \sqrt{n} \cdot \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot (\tilde{c}_x - c_x) + \sqrt{n} \cdot \left(\Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} - \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \right) \cdot c_x. \tag{61}
\end{aligned}$$

By triangle inequality, we obtain the following upper bound:

$$\left\| \vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) - \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) \right\| \leq \sqrt{n} \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot (\tilde{c}_x - c_x) \right\| + \sqrt{n} \cdot \left\| \left(\Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} - \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \right) \cdot c_x \right\|. \tag{62}$$

Step 2: Upper bounding the norm. Next, we establish separate upper bounds for the two terms in (62).

(1) The first term in (62). First of all, we observe from the definition of c_x and \tilde{c}_x , cf. (58), that

$$\begin{aligned}
\tilde{c}_x - c_x &= \left(\frac{1}{\sqrt{n}} \mathbf{Z}_{\text{ctr}}^\top \right)^\dagger \cdot (x - \mu_{\tilde{\mathcal{D}}_n}) - \left(\frac{1}{\sqrt{n}} \mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot (x - \mu_{\mathcal{D}_n}) \\
&= \sqrt{n} \cdot \left(\mathbf{Z}_{\text{ctr}}^\top - \mathbf{X}_{\text{ctr}}^\top \right)^\dagger \cdot (x - \mu_{\mathcal{D}_n}) + \sqrt{n} \cdot \left[\mathbf{Z}_{\text{ctr}}^\top \right]^\dagger \cdot (\mu_{\tilde{\mathcal{D}}_n} - \mu_{\mathcal{D}_n}).
\end{aligned}$$

753 Then we can upper bound the first term in (62) as follows:

$$\left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot (\tilde{c}_x - c_x) \right\| = \sqrt{n} \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{Z}_{\text{ctr}}^{\top \dagger} - \mathbf{X}_{\text{ctr}}^{\top \dagger} \right) \cdot (x - \mu_{\mathcal{D}_n}) \right\| + \sqrt{n} \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left[\mathbf{Z}_{\text{ctr}}^{\top} \right]^{\dagger} \cdot \left(\mu_{\tilde{\mathcal{D}}_n} - \mu_{\mathcal{D}_n} \right) \right\|. \quad (63)$$

754 Next, we consider the orthogonal decomposition of $x - \mu_{\mathcal{D}_n}$:

$$x - \mu_{\mathcal{D}_n} = \Pi_{\mathbf{X}_{\text{ctr}}}^{\text{row}} (x - \mu_{\mathcal{D}_n}) + \Pi_{\mathbf{X}_{\text{ctr}}}^{\text{row} \perp} (x - \mu_{\mathcal{D}_n}) = \frac{1}{\sqrt{n}} \mathbf{X}_{\text{ctr}}^{\top} \cdot c_x + \Pi_{\mathbf{X}_{\text{ctr}}}^{\text{row} \perp} (x - \mu_{\mathcal{D}_n}). \quad (64)$$

755 If $x - \mu_{\mathcal{D}_n} \in \text{rowsp}(\mathbf{X}_{\text{ctr}})$, then we obtain the following upper bound for the first term in (63):

$$\begin{aligned} & \sqrt{n} \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{Z}_{\text{ctr}}^{\top \dagger} - \mathbf{X}_{\text{ctr}}^{\top \dagger} \right) \cdot (x - \mu_{\mathcal{D}_n}) \right\| \\ & \leq \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{Z}_{\text{ctr}}^{\top \dagger} - \mathbf{X}_{\text{ctr}}^{\top \dagger} \right) \cdot \mathbf{X}_{\text{ctr}}^{\top} \cdot c_x \right\| \\ & \quad + \underbrace{\sqrt{n} \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left(\mathbf{Z}_{\text{ctr}}^{\top \dagger} - \mathbf{X}_{\text{ctr}}^{\top \dagger} \right) \cdot \Pi_{\mathbf{X}_{\text{ctr}}}^{\text{row} \perp} \cdot (x - \mu_{\mathcal{D}_n}) \right\|}_{=0} \quad \because (64) \\ & \leq \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left\{ -\mathbf{Z}_{\text{ctr}}^{\top \dagger} \cdot \Pi_{\mathbf{Z}_{\text{ctr}}}^{\text{row}} \cdot \left(\mathbf{Z}_{\text{ctr}}^{\top} - \mathbf{X}_{\text{ctr}}^{\top} \right) \cdot \Pi_{\mathbf{X}_{\text{ctr}}}^{\text{col}} \cdot \mathbf{X}_{\text{ctr}}^{\top \dagger} \right\} \cdot \mathbf{X}_{\text{ctr}}^{\top} \cdot c_x \right\| \quad \because \text{Lemma 4} \\ & \leq \left\| \left[\mathbf{Z}_{\text{ctr}}^{(\lambda) \top} \right]^{\dagger} \right\| \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}}^{\text{row}} \cdot (\mathbf{Z} - \mathbf{X})^{\top} \cdot \Pi_{\mathbf{1}_n^{\perp}}^{\text{row}} \cdot \Pi_{\mathbf{X}_{\text{ctr}}}^{\text{col}} \right\| \cdot \|c_x\| \\ & \leq \frac{\|\mathbf{Z} - \mathbf{X}\|}{\sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})} \cdot \|c_x\|. \end{aligned}$$

756 Similarly, the second term in (63) can be bounded by

$$\begin{aligned} \sqrt{n} \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot \left[\mathbf{Z}_{\text{ctr}}^{\top} \right]^{\dagger} \cdot \left(\mu_{\tilde{\mathcal{D}}_n} - \mu_{\mathcal{D}_n} \right) \right\| & \leq \frac{1}{\sqrt{n}} \cdot \left\| \left[\mathbf{Z}_{\text{ctr}}^{(\lambda) \top} \right]^{\dagger} \right\| \cdot \left\| \mathbf{1}_n^{\top} \cdot (\mathbf{Z} - \mathbf{X}) \right\| \\ & \leq \frac{1}{\sqrt{n}} \cdot \frac{\|\mathbf{Z} - \mathbf{X}\|}{\sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})} \cdot \|\mathbf{1}_n\|. \end{aligned}$$

757 All in all, we obtain

$$\sqrt{n} \cdot \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \cdot (\tilde{c}_x - c_x) \right\| \leq \frac{\|\mathbf{Z} - \mathbf{X}\|}{\sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})} \cdot \left(\sqrt{n} \cdot \|c_x\| + \|\mathbf{1}_n\| \right) \quad (65)$$

758 **(2) The second term in (62).** Letting $\mathbf{E}^{(\lambda)} := \mathbf{Z}_{\text{ctr}}^{(\lambda)} - \mathbf{X}_{\text{ctr}}^{(\lambda)}$, we observe that

$$\begin{aligned} \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} - \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \right\| & \leq \max \left\{ \left\| \mathbf{E}^{(\lambda)} \cdot \mathbf{X}_{\text{ctr}}^{(\lambda) \dagger} \right\|, \left\| \mathbf{E}^{(\lambda)} \cdot \mathbf{Z}_{\text{ctr}}^{(\lambda) \dagger} \right\| \right\} \quad \because \text{Lemma 5} \\ & \leq \left\| \mathbf{E}^{(\lambda)} \right\| \cdot \max \left\{ \left\| \mathbf{X}_{\text{ctr}}^{(\lambda) \dagger} \right\|, \left\| \mathbf{Z}_{\text{ctr}}^{(\lambda) \dagger} \right\| \right\} \\ & \leq \frac{\|\mathbf{Z} - \mathbf{X}\|}{\min \{ \sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}}) \}}. \end{aligned}$$

759 All in all, we obtain the following upper bound:

$$\sqrt{n} \left\| \left(\Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} - \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \right) \cdot c_x \right\| \leq \left\| \Pi_{\mathbf{Z}_{\text{ctr}}^{(\lambda)}}^{\text{col}} - \Pi_{\mathbf{X}_{\text{ctr}}^{(\lambda)}}^{\text{col}} \right\| \cdot \|c_x\| \leq \frac{\|\mathbf{Z} - \mathbf{X}\|}{\min \{ \sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}}) \}} \cdot \sqrt{n} \cdot \|c_x\|. \quad (66)$$

760 **Step 3: Concluding the proof.** We conclude this proof by inserting the upper bounds (65) and (66)
761 from Step 2 into the upper bound (62) in Step 1. Specifically, we obtain

$$\begin{aligned} \left\| \tilde{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) - \tilde{w}_{\mathcal{D}_n}^{(\lambda)}(x) \right\| & \leq \frac{\|\mathbf{Z} - \mathbf{X}\|}{\sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})} \cdot \left(\sqrt{n} \cdot \|c_x\| + \|\mathbf{1}_n\| \right) + \frac{\|\mathbf{Z} - \mathbf{X}\|}{\min \{ \sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}}) \}} \cdot \sqrt{n} \cdot \|c_x\| \\ & \leq \frac{\|\mathbf{Z} - \mathbf{X}\|}{\min \{ \sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}}) \}} \cdot (2\sqrt{n} \cdot \|c_x\| + \|\mathbf{1}_n\|). \end{aligned}$$

762 Lastly, we note that $\|c_x\| = \sqrt{(x - \mu_{\mathcal{D}_n})^{\top} \Sigma_{\mathcal{D}_n}^{\dagger} (x - \mu_{\mathcal{D}_n})} = \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}}$ and $\|\mathbf{1}_n\| = \sqrt{n}$. \square

763 **D.3 Completing the proof of Theorem 3**

764 Recall that given a set $\mathcal{D}_n = \{(x_i, y_i) : i \in [n]\}$, we let $\mathbf{X}_{\mathcal{D}_n} := [x_1 \ \cdots \ x_n]^\top \in \mathbb{R}^{n \times p}$. In
765 addition, we let

$$\forall y \in \mathcal{M}, \quad \vec{d}_{\mathcal{D}_n}^2(y) := [d^2(y_1, y) \ \cdots \ d^2(y_n, y)] \in \mathbb{R}^n. \quad (67)$$

766 Recall that we let $\mathbf{X} = \mathbf{X}_{\mathcal{D}_n}$ and $\mathbf{Z} = \mathbf{X}_{\tilde{\mathcal{D}}_n}$ for shorthand, and further, we let $\mathbf{X}_{\text{ctr}} = (\mathbf{I}_n -$
767 $\frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \mathbf{X}$ and $\mathbf{Z}_{\text{ctr}} = (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top) \mathbf{Z}$ denote the ‘row-centered’ matrices. Here we present and
768 prove the complete version of Theorem 3.

769 **Theorem 4** (De-noising covariates). *Suppose that Assumptions (C0) and (C1) hold. For any $\lambda \in \mathbb{R}_+$,*
770 *if $x \in \mu_{\mathcal{D}_n} + \text{rowsp } \mathbf{X}_{\text{ctr}}$ and*

$$\|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} \leq \frac{1}{2} \left(\frac{C_g \cdot D_g^\alpha}{2 \text{diam}(\mathcal{M})} \cdot \frac{\min \{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})\}}{\|\mathbf{Z} - \mathbf{X}\|} - 1 \right), \quad (68)$$

771 *then*

$$\begin{aligned} & d\left(\varphi_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x), \varphi_{\mathcal{D}_n}^{(\lambda)}(x)\right) \\ & \leq \left(\frac{\|\mathbf{Z} - \mathbf{X}\|}{\min \{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})\}} \cdot \frac{2 \cdot \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} + 1}{C_g} \cdot \frac{\|\vec{d}_{\mathcal{D}_n}^2(\tilde{\varphi}_n)\| + \|\vec{d}_{\mathcal{D}_n}^2(\varphi_n)\|}{\sqrt{n}} \right)^{\frac{1}{\alpha}}. \end{aligned} \quad (69)$$

772 *Proof of Theorem 4.* First of all, we recall from (52) that

$$\vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) = [w_{\mathcal{D}_n}^{(\lambda)}(x_1, x) \ \cdots \ w_{\mathcal{D}_n}^{(\lambda)}(x_n, x)] \quad \text{and} \quad \vec{w}_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x) = [w_{\tilde{\mathcal{D}}_n}^{(\lambda)}(z_1, x) \ \cdots \ w_{\tilde{\mathcal{D}}_n}^{(\lambda)}(z_n, x)].$$

773 In addition, recall that we let for any $y \in \mathcal{M}$,

$$\vec{d}_{\mathcal{D}_n}^2(y) = [d^2(y_1, y) \ \cdots \ d^2(y_n, y)] \in \mathbb{R}^n.$$

774 Thereafter, we observe that for any $y \in \mathcal{M}$ and any $x \in (\mu_{\mathcal{D}_n} + \text{rowsp } \mathbf{X}_{\text{ctr}})$,

$$\begin{aligned} \left| R_{\mathcal{D}_n}^{(\lambda)}(y; x) - R_{\mathcal{D}_n}^{(\lambda)}(y; x) \right| &= \frac{1}{n} \left| \sum_{i=1}^n \left(w_{\mathcal{D}_n}^{(\lambda)}(z_i, x) - w_{\mathcal{D}_n}^{(\lambda)}(x_i, x) \right) \cdot d^2(y_i, y) \right| \\ &= \frac{1}{n} \left| \left\langle \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) - \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x), \vec{d}_{\mathcal{D}_n}^2(y) \right\rangle \right| \\ &\stackrel{(a)}{\leq} \frac{1}{n} \left\| \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) - \vec{w}_{\mathcal{D}_n}^{(\lambda)}(x) \right\| \cdot \left\| \vec{d}_{\mathcal{D}_n}^2(y) \right\| \\ &\stackrel{(b)}{\leq} \frac{\|\mathbf{Z} - \mathbf{X}\|}{\min \{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})\}} \cdot (2 \cdot \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} + 1) \cdot \frac{\|\vec{d}_{\mathcal{D}_n}^2(y)\|}{\sqrt{n}} \end{aligned} \quad (70)$$

775 where (a) is due to Cauchy-Schwarz inequality, and (b) follows from Lemma 6.

776 Using shorthand notation $R_n = R_{\mathcal{D}_n}^{(\lambda)}$, $\tilde{R}_n = R_{\tilde{\mathcal{D}}_n}^{(\lambda)}$, $\varphi_n = \varphi_{\mathcal{D}_n}^{(\lambda)}(x)$, and $\tilde{\varphi}_n = \varphi_{\tilde{\mathcal{D}}_n}^{(\lambda)}(x)$, we observe
777 that

$$\begin{aligned} & R_n(\tilde{\varphi}_n) - R_n(\varphi_n) \\ &= R_n(\tilde{\varphi}_n) - \tilde{R}_n(\tilde{\varphi}_n) + \tilde{R}_n(\tilde{\varphi}_n) - R_n(\varphi_n) \\ &\stackrel{(a)}{\leq} R_n(\tilde{\varphi}_n) - \tilde{R}_n(\tilde{\varphi}_n) + \tilde{R}_n(\varphi_n) - R_n(\varphi_n) \\ &\stackrel{(b)}{\leq} \frac{\|\mathbf{Z} - \mathbf{X}\|}{\min \{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})\}} \cdot (2 \cdot \|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} + 1) \cdot \frac{\|\vec{d}_{\mathcal{D}_n}^2(\tilde{\varphi}_n)\| + \|\vec{d}_{\mathcal{D}_n}^2(\varphi_n)\|}{\sqrt{n}} \end{aligned} \quad (71)$$

778 where (a) follows from the optimality of $\tilde{\varphi}_n$, i.e., $\tilde{R}_n(\varphi_n) \geq \tilde{R}_n(\tilde{\varphi}_n)$, and (b) is due to (70).

779 Finally, we note that if

$$\|x - \mu_{\mathcal{D}_n}\|_{\Sigma_{\mathcal{D}_n}} \leq \frac{1}{2} \left(\frac{C_g \cdot D_g^\alpha}{2 \text{diam}(\mathcal{M})} \cdot \frac{\min\{\sigma^{(\lambda)}(\mathbf{X}_{\text{ctr}}), \sigma^{(\lambda)}(\mathbf{Z}_{\text{ctr}})\}}{\|\mathbf{Z} - \mathbf{X}\|} - 1 \right),$$

780 then the upper bound in (71) certifies that $R_n(\tilde{\varphi}_n) - R_n(\varphi_n) < C_g \cdot D_g^\alpha$. Thus, we can use
781 Assumption (C1) to convert the risk bound (71) to derive a distance bound between the minimizers:

$$d(\tilde{\varphi}_n, \varphi_n) \leq \left(\frac{R_n(\tilde{\varphi}_n) - R_n(\varphi_n)}{C_g} \right)^{\frac{1}{\alpha}},$$

782 which completes the proof. \square

783 E Details on the experiments

784 **Experimental setup.** We consider combinations of $p \in \{25, 50, 75\}$ and $n \in \{100, 200, 400\}$.
785 The datasets $\mathcal{D}_n = \{(X_i, Y_i) : i \in [n]\}$ and $\tilde{\mathcal{D}}_n = \{(Z_i, Y_i) : i \in [n]\}$ are generated as follows.

786 (True covariate X) Let $X_i \sim \mathcal{N}_p(\mathbf{0}_p, \Sigma)$ be IID multivariate Gaussian with mean $\mathbf{0}_p$ and covariance
787 Σ such that $\text{spec}(\Sigma) = \{\lambda_j > 0 : j \in [p]\}$ is an exponentially decreasing sequence such that
788 $\text{tr}(\Sigma) = \sum_{j=1}^p \lambda_j = p$. In particular, for each p , we consider an exponentially decreasing sequence
789 $1 = a_1 > \dots > a_p = 10^{-3}$, and then set $\lambda_j = p \cdot a_j / (\sum_{j'=1}^p a_{j'})$ for each $j \in [p]$. Note that
790 $\sum_{j=1}^{\lfloor p/3 \rfloor} \lambda_j / \sum_{j'=1}^p \lambda_{j'} \approx 0.9$ for all $p \in \{25, 50, 75\}$, and thus, Σ is effectively low-rank.

791 (Noisy covariate Z) Let $Z = X + \varepsilon$, where $\varepsilon \sim \mathcal{N}_p(\mathbf{0}_p, \sigma_\varepsilon^2 \cdot \text{diag}(\mathbf{1}_p))$. Note that in this setting, we
792 have the signal-to-noise ratio $\mathbb{E}(\|X\|_2^2) / \mathbb{E}(\|\varepsilon\|_2^2) = 1 / \sigma_\varepsilon^2$. We set $\sigma_\varepsilon^2 = 0.05^2$.

793 (Response Y) Given $X = x$, let Y be the distribution function of $\mathcal{N}(\mu_{\alpha, \beta}(x) + \eta, \tau^2)$, where

- 794 • $\mu_{\alpha, \beta}(x) = \alpha + \beta^\top x$ with $\alpha = 1$ and $\beta = p^{-1/2} \cdot \mathbf{1}_p$,
- 795 • $\eta \sim \mathcal{N}(0, \sigma_\eta^2)$,
- 796 • $\tau^2 \sim \mathcal{IG}(s_1, s_2)$, an inverse gamma distribution with shape s_1 and scale s_2 .

797 We note that $\mathbb{E}(\tau^2) = \frac{s_2}{s_1 - 1}$ and $\text{Var}(\tau^2) = \frac{s_2^2}{(s_1 - 1)^2(s_1 - 2)}$. In particular, when $\tau^2 = 0$, this setting
798 corresponds to the classical linear regression model for scalar responses. We set $\sigma_\eta^2 = 0.5^2$, and
799 $(s_1, s_2) = (18, 17)$. In this setting, we have

- 800 • $\mathbb{E}(\mu_{\alpha, \beta}(X)) = 1$ and $\text{Var}(\mu_{\alpha, \beta}(X)) = \beta^\top \Sigma \beta \approx 1$ for all $p \in \{25, 50, 75\}$,
- 801 • $\mathbb{E}(\tau^2) = 1$ and $\text{Var}(\tau^2) = 0.25^2$.

802 **Evaluation metrics.** For the assessment of simulation results, we perform $B = 500$ Monte
803 Carlo experiments, i.e., we draw $\mathcal{D}_n^{(b)}$ and $\tilde{\mathcal{D}}_n^{(b)}$ independent copies of \mathcal{D}_n and $\tilde{\mathcal{D}}_n$, respectively, for
804 $b = 1, \dots, B$.

805 (Model estimation) Being motivated by the standard regression analysis, we evaluate the accuracy
806 and efficiency of the Fréchet regression function estimator with

$$\text{Bias}^2(\varphi_\nu^{(\lambda)}(x)) = d_W(\overline{\varphi}_\nu^{(\lambda)}(x), \varphi_{\nu^*}^{(0)}(x))^2 \quad \text{and} \quad \text{Var}(\varphi_\nu^{(\lambda)}(x)) = \frac{1}{B} \sum_{b=1}^B d_W(\varphi_{\nu^{(b)}}^{(\lambda)}(x), \overline{\varphi}_\nu^{(\lambda)}(x))^2,$$

807 where $\nu \in \{\mathcal{D}_n, \tilde{\mathcal{D}}_n\}$. We note that the above representation is a generalization of the standard bias
808 and variance of the regression function estimator in Euclidean spaces. For the global assessment of
809 the estimation performance, we use the average criterion

$$\overline{\text{Bias}}^2(\varphi_\nu^{(\lambda)}) = \frac{1}{M} \sum_{m=1}^M \text{Bias}^2(\varphi_\nu^{(\lambda)}(x_m)) \quad \text{and} \quad \overline{\text{Var}}(\varphi_\nu^{(\lambda)}) = \frac{1}{M} \sum_{m=1}^M \text{Var}(\varphi_\nu^{(\lambda)}(x_m)),$$

where $\mathcal{G}_M = \{x_m : m = 1, \dots, M\}$ a set of fixed evaluation points. In our simulation, we generated x_1, \dots, x_M from $\mathcal{N}_p(\mathbf{0}_p, \Sigma)$ with $M = 500$ and the same evaluation set was used throughout the Monte Carlo experiments. In Table 1, we have reported $|\text{Bias}|(\varphi_\nu^{(\lambda)}) = [\overline{\text{Bias}}^2(\varphi_\nu^{(\lambda)})]^{1/2}$ and $\sqrt{\text{Var}}(\varphi_\nu^{(\lambda)}) = [\overline{\text{Var}}(\varphi_\nu^{(\lambda)})]^{1/2}$ to have them on the same scale of the metric distance.

(In-sample regression fit) In addition to the above bias and variance, we assess the model error by validating the global Fréchet regression fits of the estimated model with the mean squared error

$$\text{MSE}(\varphi_\nu^{(\lambda)}) = \frac{1}{n} \sum_{i=1}^n d_W(Y_i, \varphi_\nu^{(\lambda)}(X_i))^2.$$

The MSE is the average of squared metric-distance residuals from the observed responses, which is often unitized to measure the model adequacy in the classical regression analysis. Similarly, the overall performance $\overline{\text{MSE}}(\varphi_\nu^{(\lambda)}) = B^{-1} \sum_{b=1}^B \text{MSE}(\varphi_{\nu^{(b)}}^{(\lambda)})$ is reported in Table 1.

(Out-of-sample prediction) For $N = 1000$, generate $(X_1^{\text{new}}, Y_1^{\text{new}}, Z_1^{\text{new}}), \dots, (X_N^{\text{new}}, Y_N^{\text{new}}, Z_N^{\text{new}})$ from (X, Y, Z) , and set $\mathcal{D}_N^{\text{new}} = \{(X_i^{\text{new}}, Y_i^{\text{new}}) : i = 1, \dots, N\}$ and $\tilde{\mathcal{D}}_N^{\text{new}} = \{(Z_i^{\text{new}}, Y_i^{\text{new}}) : i = 1, \dots, N\}$ which independent of \mathcal{D}_n and $\tilde{\mathcal{D}}_n$, respectively. We measure the our-of-sample prediction performance with the mean squared prediction error

$$\text{MSPE}(\varphi_\nu^{(\lambda)}) = \frac{1}{N} \sum_{i=1}^N d_W(Y_i^{\text{new}}, \varphi_\nu^{(\lambda)}(X_i^{\text{new}}))^2,$$

where $\nu \in \{\mathcal{D}_n, \tilde{\mathcal{D}}_n\}$. We evaluate the average performance with $\overline{\text{MSPE}}(\varphi_\nu^{(\lambda)}) = B^{-1} \sum_{b=1}^B \text{MSPE}(\varphi_{\nu^{(b)}}^{(\lambda)})$.

(The choice of threshold) For simplicity, we chose a universal threshold value as

$$\hat{\lambda}_n = \arg \min_{\lambda \in \Lambda} \overline{\text{MSPE}}(\varphi_{\tilde{\mathcal{D}}_n}^{(\lambda)}),$$

where Λ is a fine grid on $(0, \sqrt{\lambda_1 \cdot p/n})$. Then the same threshold $\hat{\lambda}_n$ was used to evaluate $\text{Bias}^2(\varphi_{\tilde{\mathcal{D}}^{(b)}}^{(\lambda)}(x))$, $\text{Var}(\varphi_{\tilde{\mathcal{D}}^{(b)}}^{(\lambda)}(x))$, and $\text{MSE}(\varphi_{\tilde{\mathcal{D}}^{(b)}}^{(\lambda)}(x))$ for all $b = 1, \dots, B$. Therefore, we claim that the performance of the SVT estimator reported in Table 1 has further room for improvement if one substitute $\hat{\lambda}_n^{(b)} = \arg \min_{\lambda \in \Lambda} \text{MSPE}(\varphi_{\nu^{(b)}}^{(\lambda)})$ for each Monte Carlo experiment. Although suboptimal results are reported, we note that the proposed SVT outperforms both the oracle estimator and the naive EIV estimator in our simulation study. In practice, one may employ cross-validation for better performance. For the MSPE in Table 1, we reported $\min_{\lambda \in \Lambda} \overline{\text{MSPE}}(\varphi_{\tilde{\mathcal{D}}_n}^{(\lambda)})$.

Discussion on the simulation results. The results of our numerical study demonstrate that the proposed SVT method consistently improves the estimation and prediction performance, particularly in the errors-in-variables setting. Figure 2 illustrates how the proposed SVT estimator outperforms the naive errors-in-variables (EIV) estimator that corresponds to the SVT with zero thresholding. The naive EIV has an intrinsic model bias, known as the attenuation effect [13], because it regresses responses on error-prone covariates. We note that the naive EIV misspecifies the association structure between responses and the true covariates, and eventually, it leads to statistical inference on the misspecified model. Although the naive EIV analysis has the same efficiency as the global Fréchet regression analysis attains [41], this is not the interest of the original study designed by the error-free sample.

As shown in Theorem 2, the proposed SVT estimator is biased from thresholding singular values in the covariate matrix. However, unlike the naive EIV approach, the SVT estimator benefits from a shrinkage estimation effect such that the error-prone covariates are projected on a low-rank space and the global Fréchet regression model has a reduced dimension of effective covariates. Therefore, the SVT approach gains estimation efficiency by having a smaller estimation variance in the finite sample.

Motivated by these observations, we conducted a finite-sample study to evaluate the estimation and prediction performance of the SVT estimator. Table 1 summarizes our numerical experiments.

851 As discussed earlier, the EIV consistently showed intrinsic bias, which cannot be improved by
852 increasing the sample size. The SVT method has a greater bias than the naive EIV, but the variance is
853 always smaller. This bias-variance trade-off, as a consequence, significantly improved the prediction
854 performance of the SVT compared to the naive EIV.

855 Remarkably, it turned out that the SVT estimator achieved a smaller mean squared prediction error
856 (MSPE) than even the oracle estimator (REF) obtained from the error-free sample. The REF estimator
857 showed the smallest mean squared error (MSE) because it had a considerably small bias. However,
858 in our simulation study, the REF overfitted the training sample and showed poor performance
859 in prediction. It is also worth mentioning that the naive EIV estimator showed better prediction
860 performance than the REF estimator. This phenomenon happened because the true covariate matrix
861 is nearly singular in our simulation setup, and the REF suffered from the multicollinearity between
862 covariate components. In addition, measurement errors introduced non-ignorable minimum singular
863 values in the errors-in-variables covariate matrix with the magnitude of the variance of measurement
864 errors. As a result, the naive EIV could unintentionally avoid the multicollinearity issue and have a
865 shrinkage effect like a ridge regression.

866 These findings provide empirical evidence of the effectiveness and superiority of our approach,
867 reinforcing its practical relevance and potential impact in non-Euclidean regression analysis.