

OPTIMAL SUB-DATA SELECTION FOR NONPARAMETRIC FUNCTION ESTIMATION IN KERNEL LEARNING WITH LARGE-SCALE DATA

Anonymous authors

Paper under double-blind review

ABSTRACT

This paper considers estimating nonparametric functions in a reproducing kernel Hilbert space (RKHS) for kernel learning problems with large-scale data. Kernel learning with large-scale data is computationally intensive, particularly due to the high cost and complexity of tuning parameter selection. Existing sampling methods for scalable kernel learning, such as the leverage score-based sampling method and its variants, are designed to sketch the kernel matrix to minimize the expected global (in-sample or out-of-sample) prediction error. In complement to existing methods, this paper proposes an optimal informative sampling method to estimate nonparametric functions pointwisely when the subsample size is potentially small. Our method is tailored for scenarios where computational resources are limited, yet accurate pointwise prediction at each test location is desired. [It aims to select an informative sub-data, which is different from sketch methods that aim to select the columns of a kernel matrix.](#) Theoretical studies compare the efficiency of the proposed method to that based on the full data with optimally selected tuning parameters. [Furthermore, integrating our sub-data selection with sketching methods can improve computation time of full-data based sketching methods while maintaining the same statistical efficiency.](#) Numerical experiments demonstrate the statistical and computational efficiency of the proposed method.

1 INTRODUCTION

Large-scale data sets that have a large number of records with a great variety of resources are increasingly common in many applications such as genomics and genetics, neuro-imaging, and finance. While the increasing amount of data size brings tremendous potential for discoveries and makes it possible to fit complex models such as deep neural network models, it also brings tremendous challenges to many existing algorithms to process and analyze data quickly and efficiently. This paper is on large-scale data sets with large sample sizes. With increasing complexity and heterogeneity, nonparametric models are reliable and realistic because of their flexibility in the assumption and structure of the model. This paper focuses on nonparametric regularizations in an RKHS or the so-called kernel machine methods (Wahba, 1990; Wang, 2011; Gu, 2013; Liu et al., 2007).

The statistical properties of the kernel machine methods have been well-documented. However, the computation of the kernel machine method can be challenging for large-scale data sets. It is well known that the computational cost is at the order of $O(N^3)$ using a direct computation, where N is the sample size. To address the computational challenge, Zhang et al. (2015) developed a divide-and-conquer approach. The divide-and-conquer approach (Chen et al., 2021; Li et al., 2013) has been one of the most frequently used strategies. It first breaks down large-scale data into independent processable subsets, sending them to distributed machines for processing to obtain intermediate results, and these intermediate results are merged into final results. Besides the divide-and-conquer approach, there have been various approximation methods developed in the literature, including random Fourier features (Yang et al., 2012; Rahimi & Recht, 2007), the Nyström method (Williams & Seeger, 2001), FALKON (Rudi et al., 2017; Meanti et al., 2020) and EigenPro method (MA & Belkin, 2017; Abedsoltan et al., 2023). While random Fourier features approximate the functions in an RKHS using a smaller number of randomly sampled basis functions, Nyström method applies the idea of sketch to replace the empirical kernel matrix by a much smaller matrix with subsampled

054 columns. The FALKON and EigenPro method combine the preconditioning idea and the random
055 projection idea to speed up the computation.

056 The article aims to develop a sub-data selection-based method for nonparametric function estimation
057 in an RKHS to achieve the best statistical efficiency when the computational resource is limited
058 (Yao & Wang, 2021). The goal of our sub-data selection is to select the most informative subset of
059 observations, which is different from the existing sketch methods based on subsampling methods,
060 such as leverage score-based sampling (Alaoui & Mahoney, 2015; Rudi et al., 2015; 2018), which
061 have been developed to subsample columns of a kernel matrix. There exist abundant sub-data
062 selection approaches for data generated from parametric models. For example, Drineas et al. (2006);
063 Mahoney (2011) considered leverage sampling and Wang et al. (2019a) proposed an information-
064 based optimal subdata selection (IBOSS) to find subdata with the maximal information matrix under
065 the D- optimality for linear regression models. Ma & Sun (2015) developed local case-control
066 sampling and Cheng et al. (2020) generalized the idea of IBOSS for logistic regression models.
067 In addition, subdata selection approaches have been proposed for generalized linear models (Ai
068 et al., 2018) and quantile regression (Wang & Ma, 2020). A comprehensive review of these existing
069 subdata selection methods for parametric models may be found in Yao & Wang (2021); Chang (2024).
070 However, there is limited research on subdata selection methods targeting on nonparametric function
071 estimation. Recently, Chang (2024) developed a stratified subsampling approach for a supervised
072 learning in a nonparametric model setting. It is based on a partitioning estimate that is similar to
073 the regression tree or Nadaraya-Watson kernel estimator, which is different from the kernel machine
074 method discussed in this paper.

075 Inspired by the existing sub-data selection methods, we propose a new sub-data selection methodology
076 to estimate nonparametric functions in an RKHS. We first apply a clustering method such as k-means
077 or other clustering algorithms to select representative data points. The nonparametric function values
078 in each cluster are roughly approximated by the functional values at representative data points. To
079 decide the sampling weights for each cluster, we minimize the MSE of the nonparametric function
080 estimator that is constructed based on the selected representative data points. However, the sampling
081 weights depend on a tuning parameter which is unknown. To address this issue, we adopted a one-step
082 iteration procedure to choose a tuning parameter using representative data points via the BIC criterion
083 or cross-validation. The optimal weights depend on the cluster centers decided by the K-means
084 algorithm, the chosen kernel function, and the selected tuning parameter. After selecting multiple sub-
085 datasets using the optimal weights, we apply the kernel machine method to each selected sub-dataset
086 and aggregate the estimators to obtain the final estimate. [We further show that the sub-data selection
087 method can be integrated with sketching methods such as FALKON to improve the computation time
for pointwise estimation while maintaining the same statistical efficiency.](#)

088 From a theoretical perspective, the proposed method is designed to select a sub-sample (with a fixed
089 sample size) to minimize the expected prediction error at every test data point, hence it minimizes
090 the expected prediction error for every test data set. It is different from existing research for sketch
091 method in kernel learning, which has established its optimality to minimize the expected global
092 prediction error. For example, Musco & Musco (2017) considers in-sample prediction error and
093 Rudi et al. (2015; 2018); Alaoui & Mahoney (2015) consider the expected global generalization
094 (out-of-sample) prediction error. While these results are interesting and important, they do not directly
095 guarantee generalization performance for every test data point. Our contribution is to establish the
096 rate of convergence of the proposed estimator, and demonstrate its advantage in improving the
097 convergence rate of the RKHS estimator based on subsamples obtained from a Simple Random
098 Sampling (SRS). The results are interesting since they confirm that the proposed sub-data selection is
informative in making use of the full data to improve the prediction error.

099 From a numerical perspective, our numerical results show that our proposed method is computa-
100 tionally efficient when compared with a Simple Random Sample (SRS) approach and maintains
101 estimation accuracy when compared with the full data approach. The integrated mean squared errors
102 of our proposed method are comparable to that of using full data while computational time is as
103 good as the SRS based approach. More importantly, the proposed approach achieves a good balance
104 between statistical efficiency and computational efficiency when compared with the full data-based
105 and the SRS-based approaches. In addition, [we investigate a combination of proposed method with
106 sketching methods \(e.g., FALKON\) to improve computation time of existing sketching methods.](#)
107 The proposed method has shown better MSE in out-of-sample prediction than some existing sketch
algorithms, including Nyström (Williams & Seeger, 2001) and FALKON (Rudi et al., 2017; Meanti

et al., 2020). We also compare the performance of the proposed method with these sketch algorithms in the YearPredictionMSE data set.

This manuscript is organized in the following way. In Section 2, we provide the basic framework, an introduction of regularization in RKHS with full data sets, and the proposed method for regularization in RKHS with large-scale data, with the details of our algorithm. Theoretical justification is given in Section 3. The numerical and simulation studies are included in Section 4. Section 5 includes an application of the proposed method to a real data set and compares/combines it with other sketch algorithms. A brief discussion is provided in Section 6. [All the technical proofs and additional numerical results and details are included in the Appendix.](#)

2 REGULARIZATION IN REPRODUCING KERNEL HILBERT SPACES

Consider a continuous response y_i and a p -dimensional covariate vector X_i , modeled as

$$y_i = f_0(X_i) + \epsilon_i, \quad i = 1, \dots, N,$$

where ϵ_i are i.i.d. errors with mean zero and variance σ^2 . We assume f_0 belongs to a reproducing kernel Hilbert space (RKHS) \mathcal{H}_K induced by a symmetric, positive-definite kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$.

If μ is a finite measure on \mathcal{X} and $\int K(x, x)d\mu(x) < \infty$, [applying Mercer’s theorem \(Mercer, 1909\)](#), one can write $K(x, y) = \sum_{j=1}^{\infty} \lambda_j \psi_j(x) \psi_j(y)$, where $\{\psi_j\}$ form an orthonormal basis in $L^2(\mu)$, $\lambda_j > 0$, and $\sum_{j=1}^{\infty} \lambda_j < \infty$. The RKHS is $\mathcal{H}_K = \left\{ f(x) = \sum_{j=1}^{\infty} c_j \psi_j(x) : \sum_{j=1}^{\infty} c_j^2 / \lambda_j < \infty \right\}$, with norm $\|f\|_{\mathcal{H}_K}^2 = \sum_{j=1}^{\infty} c_j^2 / \lambda_j$. To estimate f_0 , we solve

$$\hat{f}_{N, \lambda^*} = \arg \min_{f \in \mathcal{H}_K} \frac{1}{N} \sum_{i=1}^N \{y_i - f(X_i)\}^2 + \lambda^* \|f\|_{\mathcal{H}_K}^2,$$

where $\lambda^* > 0$ is a tuning parameter. By [applying representer theorem \(Wahba, 1990; Kimeldorf & Wahba, 1971\)](#), the solution has the finite expansion $\hat{f}_{N, \lambda^*}(x) = \sum_{l=1}^N a_l K(x, X_l)$, with coefficients $\mathbf{a} = [a_1, \dots, a_N]^T$. Plugging this into the above objective for estimating f_0 , it yields an objective function for \mathbf{a} , which is $J(\mathbf{a}) = \frac{1}{N} \|\mathbf{y} - \mathbf{K}\mathbf{a}\|^2 + \lambda^* \mathbf{a}^T \mathbf{K}\mathbf{a}$, where $\mathbf{K} \in \mathbb{R}^{N \times N}$ has entries $K(X_i, X_j)$. The solution $\hat{\mathbf{a}}$ that minimizes $J(\mathbf{a})$ is $\hat{\mathbf{a}} = N^{-1} \lambda^{*-1} (\mathbf{I} + N^{-1} \lambda^{*-1} \mathbf{K})^{-1} \mathbf{y}$. Thus, for any x ,

$$\hat{f}_{N, \lambda^*}(x) = N^{-1} \lambda^{*-1} [K(x, X_1), \dots, K(x, X_N)] (\mathbf{I} + N^{-1} \lambda^{*-1} \mathbf{K})^{-1} \mathbf{y}.$$

This requires inverting a $N \times N$ matrix, which is computationally demanding for large N . Moreover, performance depends on selecting λ^* , adding to the computational cost.

2.1 PROPOSED METHOD FOR REGULARIZATION IN AN RKHS

Given data pairs $(X_1, y_1), \dots, (X_N, y_N)$ for a large N , the goal is to estimate a nonparametric function $f_0(x)$ for a given x . Because N is very large, a direct application of kernel machine method is time-consuming even not possible due to the inverse of an $N \times N$ matrix. Given a limited computational resource, we consider a sub-data selection approach by selecting a subsample $(X_{k_1}, y_{k_1}), \dots, (X_{k_n}, y_{k_n})$ with size n from the original sample with size N while maximizing the statistical efficiency of the resulting kernel machine estimator $\hat{f}_{n, \tilde{\lambda}^*}(x)$, [where \$\tilde{\lambda}^*\$ denotes the tuning parameter for the proposed method.](#)

Our proposed procedure makes use of the smoothness property of the nonparametric functions by approximating the functional values of $f_0(x)$ by a set of representative data points. To find the representative points, we firstly apply a clustering approach to cluster N data points into L representative clusters $\{\mathcal{C}_1, \dots, \mathcal{C}_L\}$, and then re-sampling with optimal weights $\{\omega_{x,1,C}, \dots, \omega_{x,L,C}\}$ from these L clusters to obtain these representative data points. [Euclidean k-means is a common choice for clustering, it is appropriate for kernels that preserve local Euclidean structure \(such as an RBF kernel\).](#) For RKHS spaces introduced by more general kernels, clustering algorithms such as kernel K-means or K-medians (Wang et al., 2019b) that preserve the geometry structure of the RKHS space are better choices.

Algorithm 1 Proposed Weighted Resampling RKHS Estimator

Require: Data $\{(X_i, y_i)\}_{i=1}^N$, subsample size n , number of clusters L , number of resamples B and kernel K

Ensure: $\hat{f}_{n, \tilde{\lambda}^*}(x)$

- 1: **Clustering.** Given X_1, \dots, X_N and subsample size n , apply a clustering method such as K-means (or random projected K-means or kernel K-means or K-medians) to partition the dataset into L clusters C_1, \dots, C_L . Let the cluster centers be $\{C_1, \dots, C_L\}$.
- 2: **Tuning $\tilde{\lambda}^*$.** Apply cross-validation or BIC method to find the tuning parameter $\tilde{\lambda}^*$ for the RKHS regression using the representative points $\{C_1, \dots, C_L\}$.
- 3: **Assigning weights.** Using $\tilde{\lambda}^*$, compute optimal cluster weights $\omega_{x,1,C}, \dots, \omega_{x,L,C}$ as defined in (2), and compute weights $\omega_i(x_0)$ for each data points defined in Section 3.2.
- 4: **for** $b = 1$ to B **do**
- 5: **Resampling with weights.** Using the weights $\omega_1(x), \dots, \omega_N(x)$ to sample n points from $\{X_1, \dots, X_N\}$ to obtain $\{X_1^*, \dots, X_n^*\}$ and corresponding outcomes $\mathbf{y}^* = [y_1^*, \dots, y_n^*]^\top$.
- 6: **Compute RKHS estimator.** Using the selected sample and $\tilde{\lambda}^*$ to compute $\hat{f}_{n, \tilde{\lambda}^*}^{(b)}(x)$ which is defined similarly to $\hat{f}_{N, \lambda^*}(x)$ except that it uses the b -th sampling data.
- 7: **end for**
- 8: **Aggregate:** The final estimate of $f_0(x)$ is $\hat{f}_{n, \tilde{\lambda}^*}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_{n, \tilde{\lambda}^*}^{(b)}(x)$.

The optimal weights $\{\omega_{x,1,C}, \dots, \omega_{x,L,C}\}$ are chosen to maximize the statistical efficiency of the proposed estimator. Let $\{C_1, \dots, C_L\}$ denote the clusters, where C_l contains N_l observations with $\sum_{l=1}^L N_l = N$. We denote the center of cluster C_l by C_l . Suppose we select a subdata set with size n , among them n_l data points are from the l -th cluster so that $n_l = n\omega_{x,l,C}$ and $\sum_{l=1}^L \omega_{x,l,C} = 1$. For all the data selected from the cluster C_l , we approximate them using their representative data centers C_l . With a slightly abuse of notations, we rearrange the data $\{y_1, \dots, y_N\}$ by clusters and denote y_{lj} the j -th data points selected from the l -th cluster. Then, the corresponding model for the l -th cluster mean $\bar{y}_l = \sum_{j=1}^{n_l} y_{lj}/n_l$ is

$$\bar{y}_l \approx f_0(C_l) + \bar{\epsilon}_l,$$

where $\bar{\epsilon}_l = \sum_{j=1}^{n_l} \epsilon_{lj}/n_l$ has mean zero and variance $\sigma^2/(n\omega_{x,l,C})$. We then consider the following RKHS nonparametric estimation of $f_0(x)$ given the above model:

$$\hat{f}_{n, \tilde{\lambda}^*}(x) = L^{-1} \tilde{\lambda}^{*-1} [K(x, C_1), \dots, K(x, C_L)] (\mathbf{I} + L^{-1} \tilde{\lambda}^{*-1} \mathbf{K}_c)^{-1} \bar{\mathbf{y}},$$

where \mathbf{K}_c is an $L \times L$ matrix with (i, j) -th element $K(C_i, C_j)$ and $\bar{\mathbf{y}} = [\bar{y}_1, \dots, \bar{y}_L]^\top$. Conditional on the centers $\{C_1, \dots, C_L\}$, the variance of the estimate $\hat{f}_{n, \tilde{\lambda}^*}(x)$ is therefore

$$\text{Var}\{\hat{f}_{n, \tilde{\lambda}^*}(x)\} = \sigma^2 L^{-3} \tilde{\lambda}^{*-2} \mathbf{K}_x(C)^\top (\mathbf{I} + L^{-1} \tilde{\lambda}^{*-1} \mathbf{K}_c)^{-1} \mathbf{D}_\omega (\mathbf{I} + L^{-1} \tilde{\lambda}^{*-1} \mathbf{K}_c)^{-1} \mathbf{K}_x(C),$$

where $\mathbf{K}_x(C) = [K(x, C_1), \dots, K(x, C_L)]^\top$, $\mathbf{D}_\omega = \text{diag}\{1/\omega_{x,1,C}, \dots, 1/\omega_{x,L,C}\}$. Because $\mathbb{E}(\bar{\mathbf{y}}) = [f_0(C_1), \dots, f_0(C_L)]^\top$, the conditional bias of the estimator $\hat{f}_{n, \tilde{\lambda}^*}(x)$ is independent of the choices of $\omega_{x,i,C}$'s. Therefore, to minimize the conditional MSE of $\hat{f}_{n, \tilde{\lambda}^*}(x)$, we could choose $\omega_{x,l,C}$'s that minimize the variance of $\hat{f}_{n, \tilde{\lambda}^*}(x)$:

$$\omega_{x,l,C} = \arg \min_{\omega_{x,l,C} \geq 0, \sum_{l=1}^L \omega_{x,l,C} = 1} \text{Var}\{\hat{f}_{n, \tilde{\lambda}^*}(x)\}. \quad (1)$$

It is not difficult to check that $\text{Var}\{\hat{f}_{n, \tilde{\lambda}^*}(x)\} = L^{-3} \sigma^2 \tilde{\lambda}^{*-2} \sum_{l=1}^L K_{xl}^2 / \omega_{x,l,C}$, where K_{xl} is the l -th element of the vector \mathbf{K}_x which is defined by

$$\mathbf{K}_x = (\mathbf{I} + L^{-1} \tilde{\lambda}^{*-1} \mathbf{K}_c)^{-1} \begin{bmatrix} K(x, C_1) \\ \vdots \\ K(x, C_L) \end{bmatrix} = \begin{bmatrix} K_{x1} \\ \vdots \\ K_{xL} \end{bmatrix}.$$

Therefore, the optimal weights $\omega_{x,l,C}$'s that minimize the variance is

$$\omega_{x,l,C} = |K_{xl}| / \sum_{l=1}^L |K_{xl}|. \quad (2)$$

Because the above optimal weights depend on the tuning parameter $\tilde{\lambda}^*$, an initial tuning parameter is needed to determine the optimal weights. We have compared an iterative algorithm with a one-step iterative algorithm, and we found that their performance were similar. To save computational time, we adapted the one-step iterative procedure in Algorithm 1.

For the tuning procedure in Algorithm 1, it requires $\hat{f}_{n,\tilde{\lambda}^*}(x)$ which is an RKHS estimator of $f_0(x)$ based on the centers selected. However, the estimator depends on the tuning parameter $\tilde{\lambda}^*$, embedded in the inverse operation, which needed to be computed for every possible candidate λ . This leads to a big computational burden. To mitigate computational burden for selecting the tuning parameter $\tilde{\lambda}^*$, we decompose the kernel matrix for the selected data with size n into $\mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^\top$, where \mathbf{Q} is an orthogonal matrix that is independent of $\tilde{\lambda}^*$, and $\mathbf{\Lambda} = \text{diag}\{\theta_1, \dots, \theta_n\}$ is an $n \times n$ diagonal matrix of eigenvalues. Using the notations defined in Algorithm 1, we can write $\hat{f}_{n,\tilde{\lambda}^*}(x) = [K(x, X_1^*), \dots, K(x, X_n^*)]^\top \mathbf{Q}\mathbf{\Lambda}_{\tilde{\lambda}^*}^{-1}\mathbf{Q}^\top \mathbf{y}^*$, where $\mathbf{y}^* = [y_1^*, \dots, y_n^*]^\top$ is an $n \times 1$ vector and $\mathbf{\Lambda}_{\tilde{\lambda}^*} = \text{diag}\{\tilde{\lambda}^* + \theta_1, \dots, \tilde{\lambda}^* + \theta_n\}$. Then, one only needs to compute the inverse of a diagonal matrix $\mathbf{\Lambda}_{\tilde{\lambda}^*}$ for different tuning parameters.

3 THEORETICAL JUSTIFICATION

3.1 ASYMPTOTIC RATE AND TUNING WITH FULL DATA

Theorem 1. Assume X has a probability density function $\pi(x)$ and $\mathbb{E}\{\psi_j^4(X)\} < \infty$. If λ_j decays at the order of j^{-k} for some $k > 1$ where λ_j is the j -th largest eigenvalue of the positive definite $K(\cdot, \cdot)$ and k is the decaying rate. For any $f \in \mathcal{H}_K$ with expansion $f(x) = \sum_{j=1}^{\infty} c_j \psi_j(x)$ and the coefficients satisfy $c_j^2 \lambda_j^{-1} \asymp j^{-a}$ ($a > 1$), the optimal rate of tuning parameter is: $\lambda^* \asymp N^{-\frac{k}{k+1+(a-1)}}$, and the corresponding asymptotic integrated mean squared error (IMSE) is $\|\hat{f}_{N,\lambda^*} - f_0\|^2 = O_p\{N^{-\frac{(a-1)+k}{(a-1)+k+1}}\}$, where $\|g\|$ is the L_2 norm of a function g .

The assumptions on eigenvalues and eigenfunctions of kernel are common in the literature (e.g., Yuan & Cai (2010); Zhang et al. (2015)). The details of the proof of Theorem 1 are given in the Appendix. Comparing with the results in the existing literature (e.g., Yuan & Cai (2010)), the results in Theorem 1 give a more specific rate of convergence for the functions $f \in \mathcal{H}_K$ with a specific form specified in Theorem 1. If we consider all the possible functions in \mathcal{H}_K , the rate of convergence is given by (letting $a \rightarrow 1$) $\|\hat{f}_{N,\lambda^*} - f_0\|^2 = O_p(N^{-\frac{k}{k+1}})$. For functions in a univariate Sobolev space of order m , the eigenvalue of the corresponding reproducing kernel is at the order of $\lambda_j \asymp j^{-2m}$. Yuan & Cai (2010), then the rate of the convergence is given by $\|\hat{f}_{N,\lambda^*} - f_0\|^2 = O_p(N^{-\frac{2m}{2m+1}})$. This rate is known to be optimal in the literature (e.g. Zhang et al. (2015)).

Theorem 1 has multi-purposes. First, the rate of convergence results of Theorem 1 is a basis to define an efficiency index, which is designed to measure the trade-off between computational and statistical efficiency. This would facilitate the comparison among different methods. The details are given in Section 4. Second, Theorem 1 provides a general guidance about the order of tuning parameter choices subject to a constant, which can be further selected by the BIC or cross-validation.

3.2 ASYMPTOTIC RATE AND TUNING WITH RESAMPLING UNDER PROPOSED METHOD

For a given x_0 , the proposed resampling weights for a sample $X_i \in \mathcal{C}_l$ where \mathcal{C}_l is the l -th cluster whose center is C_l with a cluster size N_l , the weight for X_i to be selected is: $\omega_i(x_0) = N_l^{-1} \omega_{x_0,l,C}$, where $\omega_{x_0,l,C}$ was defined in (2). For each x_0 and the j -th eigenfunctions, define $\omega_{x_0,j} = \mathbb{E}\left\{\sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i)\right\}$. Denote the indices of the sampled data points with size n for a given x_0 by $\mathcal{S}(x_0) = \{i_1^*, \dots, i_n^*\} \subset \{1, \dots, N\}$, and the sampled data points are $\{(X_i, y_i) : i \in \mathcal{S}(x_0)\}$.

The following Lemma studies the weight function $\omega_{x_0,j}$. The detailed proof of Lemma 1 can be found in the Appendix.

Lemma 1. Let $\omega_{x_0,l,C}$ be the cluster-level weights defined by (2), and assign each sample $X_i \in C_l$ the resampling weights $\omega_i(x_0) = N_l^{-1} \omega_{x_0,l,C}$. Define $\omega_{x_0,j} = \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i) \right\}$. If $\lambda_j \asymp j^{-k}$ and $\int K^2(x, y) d\pi(y) < \infty$, then $\omega_{x_0,j} \leq \lambda_j^{-2} \asymp j^{2k}$.

The proposed estimator of $f(x_0)$ is $\hat{f}_{n,\tilde{\lambda}^*,x_0}(x_0)$, which is given by:

$$\hat{f}_{n,\tilde{\lambda}^*,x_0} = \arg \min_{f \in \mathcal{H}_K} \left[\frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} \{y_i - f(X_i)\}^2 + \tilde{\lambda}^* \|f\|_{\mathcal{H}_K}^2 \right].$$

Define the functional $\hat{f}_{n,\tilde{\lambda}^*}$ of the proposed estimator by collecting the estimators at all $x_0 \in \mathcal{X}$ together $\hat{f}_{n,\tilde{\lambda}^*} = \left\{ \hat{f}_{n,\tilde{\lambda}^*,x_0}(x_0) : x_0 \in \mathcal{X} \right\}$, where \mathcal{X} is the space of the predictor/feature x .

Theorem 2. Assume X has a probability density function $\pi(x)$ and $\mathbb{E}\{\psi_j^4(X)\} < \infty$. If λ_j decays at the order of j^{-k} for some $k > 1$ where λ_j is the j -th largest eigenvalue of the positive definite $K(\cdot, \cdot)$ and k is the decaying rate. Assume $\omega_{x_0,j} \asymp j^\beta$ for $0 \leq \beta < k$. Then, for any function f_0 in RKHS, the optimal rate of the tuning parameter is: $\tilde{\lambda}^* \asymp n^{-\frac{k-\beta}{k+1-2\beta}}$ if $\beta \leq 1/2$ and $\tilde{\lambda}^* \asymp n^{-1/2}$ if $\beta > 1/2$, and the corresponding asymptotic IMSE of the estimator $\hat{f}_{n,\tilde{\lambda}^*}$ is

$$\|\hat{f}_{n,\tilde{\lambda}^*} - f_0\|^2 = \begin{cases} O_p\left(n^{-\frac{k}{k+1-2\beta}}\right) & \text{if } \beta \in [0, 1/2), \\ O_p\left(n^{-1} \log(n)\right) & \text{if } \beta = 1/2, \\ O_p\left(n^{-1}\right) & \text{if } \beta > 1/2. \end{cases}$$

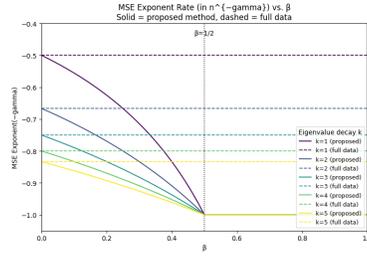


Figure 1: Asymptotic MSE for the proposed method for different values of β and k . Solid curves show the proposed method (using the optimal $\tilde{\lambda}^*$), and dashed curves show the rate the full-data estimator for comparison.

The assumption $\omega_{x_0,j} \asymp j^\beta$ in Theorem 2 is based on the result in Lemma 1. Detailed proof of Theorem 2 is given in the Appendix. Based on the results in Theorem 2, the rate of convergence of the proposed estimator is $\|\hat{f}_{n,\tilde{\lambda}^*} - f_0\|^2 = O_p\left(n^{-\frac{k}{k+1-2\beta}}\right)$ or $O_p(n^{-1})$ for any functions in the RKHS \mathcal{H}_K . Because $\beta \geq 0$, the rate of the proposed estimator is faster than the rate of the RKHS estimator sampled by the SRS with the same subsample size, which is given by $n^{-\frac{k}{k+1}}$. The proposed method leads to more substantial improvement when $\beta \geq 1/2$, where the optimum reaches the parametric rate n^{-1} with sample size n . The improvement in the estimator rate is due to informative sampling, where the sampling weights carry the information of the entire data set. An illustration of the rate improvement is provided in Figure 1.

4 SIMULATION STUDY

4.1 PROPOSED VS. SIMPLE RANDOM SAMPLE VS. FULL DATA

Starting with the simplest case, we evaluate the numerical performance of our proposed method in estimating the unknown function $f_0(x) = 5 \cos x$, and compare it with that based on the full data and data sampled using a simple random sampling (SRS) strategy. The full sample size is denoted by N and the subsample size is denoted by n . The random errors $\epsilon_1, \dots, \epsilon_n \sim N(0, 0.5)$ i.i.d.. The kernel function $K(\cdot, \cdot)$ was Gaussian. Results are based on 100 replications, and $B = 1$ was used for the proposed method. The simulation study was implemented in an R environment.

To evaluate each method, we compute the integrated mean squared error (IMSE) and record the computational time. To compare methods A and R , we use relative efficiency (RE):

$$\text{RE} = \frac{\text{Efficiency}_A}{\text{Efficiency}_R} = \frac{(\text{IMSE}_A)^{-\alpha} (\text{Time}_A)^{-\beta}}{(\text{IMSE}_R)^{-\alpha} (\text{Time}_R)^{-\beta}},$$

Table 1: IMSE, average timing (hr:min’sec”) and relative efficiency (RE) for $f_0(x) = 5 \cos(x)$.

N	n	Full Data		Proposed			SRS		
		IMSE	Time	IMSE	Time	RE	IMSE	Time	RE
1000	100	0.0034	0:0’2	0.0074	0:0’0	2.20	0.0297	0:0’0	0.68
1000	200	0.0034	0:0’2	0.0046	0:0’0	2.66	0.0152	0:0’0	0.89
5000	500	0.0008	0:3’35	0.0017	0:0’0	3.26	0.0067	0:0’0	0.85
5000	1000	0.0008	0:3’35	0.0011	0:0’3	2.72	0.0036	0:0’2	0.76
10000	1000	0.0005	0:22’16	0.0009	0:0’3	3.49	0.0036	0:0’2	0.76
10000	2000	0.0005	0:22’16	0.0006	0:0’26	2.98	0.0020	0:0’14	0.81

where R is the full-data estimator and we choose $\alpha = 5/4$, $\beta = 1/3$ so that the efficiency of the full data method is (asymptotically) a constant with respect to N because Theorem 1 suggests that the full data approach has $\text{IMSE} \asymp N^{-4/5}$ for functions in the Sobolev space with order $m = 2$ and computation time $\asymp N^3$ (using direct computation), and the RE of the method based on the SRS is approximately 1. $\text{RE} < 1$ indicates method A performed worse than the full data method; $\text{RE} > 1$ indicates method A performed better than the full data method.

4.2 SIMULATION STUDY: FALKON VS. PROPOSED VS. NYSTRÖM

We compare the proposed K-means-based resampling method with two fast kernel ridge regression baselines: FALKON methods including FALKON17 (Rudi et al., 2017) and FALKON20 (Meanti et al., 2020), and Nyström KRR (Williams & Seeger, 2001). We generate 20-dimensional predictors with five relevant variables and evaluate training sizes $N \in \{2000, 5000, 10000\}$ under matched computational budgets using subset sizes $n \in \{100, 500, 1000\}$. Across 100 replications, we report test MSE and total runtime (including tuning). Detailed experimental setup, tuning grids, and implementation notes are provided in Appendix E.1.

Table 2: Simulation results: mean squared error (MSE) and total time (seconds) averaged over 100 replications, n is the subset size (M for FALKON/Nyström and n for Proposed), $B = 3$ for Proposed. Time includes tuning and final training.

N	n	MSE				Time		
		FALKON17	FALKON20	Proposed	Nyström	FALKON17	Proposed	Nyström
2000	100	6.7196	6.7001	5.7745	6.5325	0.0615	0.1241	0.0166
2000	500	4.1357	4.1636	2.2412	2.5114	0.3186	0.7966	0.5946
2000	1000	3.6189	3.6094	1.7058	2.0939	0.8643	3.2982	1.2725
5000	100	6.4308	6.4577	5.7397	6.2693	0.1898	0.1805	0.0234
5000	500	2.9266	2.9216	1.8134	2.3298	0.6127	1.6280	0.6692
5000	1000	2.4558	2.4484	1.2729	1.9035	1.5446	7.1565	1.6389
10000	100	6.3071	6.3295	5.7606	6.2087	0.2486	0.3953	0.0365
10000	500	2.3723	2.3875	1.6411	2.0700	1.0347	3.0601	0.7498
10000	1000	2.0695	2.0586	1.0393	1.7595	2.3338	13.2765	1.9417

Table 2 shows that the proposed method attains the lowest MSE for all the sub-data considered in the simulation, reflecting the benefit of informative resampling. FALKON remains the fastest, especially at small n , while Nyström is competitive but trails the proposed estimator in accuracy under the similar computational budget.

5 REAL-DATA STUDY: YEARPREDICTIONMSD

5.1 METHOD COMPARISON: PROPOSED VS. FALKON VS. NYSTRÖM

We evaluate the proposed estimator on the YEARPREDICTIONMSD dataset (90 features, continuous response) and compare it with Nyström KRR and FALKON under matched computational budgets.

For each experiment we draw $N \in \{2000, 5000, 10000, 20000\}$ for training and use a fixed test set of size 1000, with all randomization seeded for reproducibility. We also apply a simple filter step by selecting the top- p features where $p \in \{30, 60, 90\}$. All methods use RBF kernels unless otherwise specified, with subset/landmark budget $n = M$ for Nyström and FALKON, and n as subsample size for the proposed method. Full implementation details, kernel settings, and runtime protocol are provided in Appendix E.2.

Table 3: YearPredictionMSD data using kernel learning with 1,000 testing data points. Proposed method uses random-projection- k -means centers, $B=3$, fixed $\lambda^*=1$, with three kernels: Prop.G = Gaussian RBF, Prop.M = Matérn-3/2, Prop.L = Laplace. Times for Proposed are averaged across these three kernels and shown as total with centers-time in parentheses. FALKON and Nyström use Gaussian RBF ($\sigma=6$). FAL17 and FAL20 refer to the implementations of Rudi et al. (2017) and Meanti et al. (2020), respectively.

N	n	p	MSE						Time (s)		
			Prop.G	Prop.M	Prop.L	FAL17	FAL20	Nyström	Proposed	FALKON	Nyström
2000	500	30	100.67	97.76	96.48	122.42	110.40	124.54	1.32 (0.16)	0.12	0.11
		60	96.26	93.87	94.14	105.75	97.93	109.25	1.22 (0.08)	0.06	0.05
		90	96.96	95.36	95.34	96.65	95.06	98.38	1.12 (0.07)	0.06	0.06
5000	1000	30	94.55	92.80	91.44	108.61	103.29	116.52	4.10 (0.17)	0.10	0.14
		60	95.37	92.56	92.98	105.11	106.85	105.20	3.95 (0.18)	0.09	0.16
		90	87.35	88.98	89.07	91.57	93.15	91.59	3.25 (0.16)	0.13	0.16
5000	2000	30	81.39	81.22	80.51	113.93	105.82	154.15	15.05 (0.20)	0.26	0.61
		60	78.45	78.79	78.42	95.62	103.94	99.60	14.19 (0.17)	0.26	0.63
		90	76.38	76.62	77.18	86.14	101.82	88.08	14.13 (0.16)	0.25	0.68
10000	1000	30	92.90	92.33	92.22	97.37	100.86	101.36	4.58 (0.41)	0.13	0.27
		60	91.16	89.99	91.14	97.00	98.03	93.86	4.50 (0.45)	0.13	0.18
		90	88.38	88.53	91.41	90.63	88.72	91.26	4.56 (0.44)	0.10	0.18
10000	2000	30	96.52	95.73	95.16	110.34	101.30	115.29	15.83 (0.48)	0.27	0.71
		60	88.77	89.86	90.75	94.10	99.89	95.96	15.48 (0.42)	0.31	0.72
		90	87.31	86.63	88.73	87.63	92.06	87.38	15.47 (0.50)	0.26	0.73
20000	1000	30	97.05	96.24	95.75	97.44	95.37	99.82	9.84 (1.45)	0.22	0.23
		60	96.93	96.82	95.11	96.95	95.19	98.38	10.38 (1.40)	0.15	0.25
		90	94.05	94.42	95.14	95.02	89.42	97.59	10.13 (1.38)	0.15	0.25
20000	2000	30	98.24	98.25	98.34	105.48	94.60	110.94	20.56 (1.42)	0.33	0.92
		60	98.50	99.32	97.92	101.90	94.21	103.28	20.38 (1.27)	0.33	1.00
		90	94.17	92.97	94.23	98.11	87.58	93.95	19.79 (1.34)	0.33	0.93

Table 3 summarizes the test MSEs and runtimes. The proposed method is consistently competitive with (and often outperforms) FALKON and Nyström in accuracy, with longer but manageable computation time. This is because our proposed method is pointwise and for each test data point, a sub-data is selected. The Laplace and Matérn variants provide additional robustness, with Gaussian RBF performing strongest in some settings. Although the proposed method is designed for sub-data selection for pointwise prediction, the reported run times show that the proposed approach scales well for $N=20k$ with $n=2000$, where total time cost remains at the order of 20 seconds. An investigation of the impact of clustering methods on the proposed method is given in the Appendix.

5.2 POINTWISE PREDICTION

Since our method is designed for pointwise prediction, we conduct an experiment to demonstrate how it can be used to improve the computational efficiency of FALKON while maintaining statistical accuracy. As discussed in Section 3.2, combining our sub-data selection with FALKON allows pointwise prediction to be performed on a much smaller but informative subset of the data, reducing training cost without sacrificing estimation quality. In this experiment, FALKON is trained once on the full training set of size N using M randomly chosen landmarks (Gaussian kernel with $\sigma = 6$, $\lambda^* = 10^{-6}$, 20 iterations), and then used to predict all 200 test points. In contrast, Prop+FALKON

begins with k -means clustering on all N points, computes sampling weights for each test point x_0 , draws a pointwise subdata of size $n \asymp 3N^{0.8}$, and fits FALKON with the same M landmarks on this subset to make the prediction at x_0 .

Table 4: Numerical comparison of pointwise prediction on YearPredictionMSD using full-data FALKON and the proposed+FALKON (Prop+F) method across different training sizes N , subdata sizes n , and landmark counts M . Timing columns give the average computation per test point. Clust is the k -means clustering time for Prop+F, Train is the training time for Prop+F, and Sel is the sampling/weighting time for Prop+F.

Sample sizes			MSE		RMSE		Total Time (s)		Time dist. (Prop+F)		
N	M	n	FALKON	Prop+F	FALKON	Prop+F	FALKON	Prop+F	Clust	Train	Sel
10k	0.5k	4.76k	81.035	83.809	9.002	9.153	0.028	0.273	0.234	0.018	0.021
20k	1k	8.28k	88.075	87.974	9.383	9.382	0.107	0.851	0.713	0.067	0.071
50k	1k	17.23k	78.571	78.139	8.865	8.840	0.228	3.991	3.467	0.101	0.423
100k	5k	30.00k	80.573	85.890	8.977	9.273	3.522	11.765	8.781	1.844	1.140
100k	10k	30.00k	85.107	83.930	9.224	9.160	56.999	17.457	9.062	7.261	1.134
200k	5k	52.23k	82.603	77.079	9.090	8.778	47.972	26.735	21.687	2.162	2.886
200k	10k	52.23k	75.708	78.909	8.697	8.888	96.188	33.800	22.693	8.319	2.788

The pointwise results in Table 4 demonstrate that the proposed selection strategy can substantially reduce the effective training size while preserving predictive accuracy. Across all training sizes, the subdata-based estimator achieves MSE and RMSE comparable to (and occasionally better than) full FALKON, despite operating on significantly smaller subsets. The reduction in training points also leads to computation reduction, particularly for large $N = 200k$ without compromising performance. As for computational complexity, we compare the theoretical time and memory costs in Table 5.

Table 5: Time and memory complexity of kernel methods. N = full training size, M = number of landmarks, $n \asymp N^\gamma$ subdata size selected by the proposed method, where $\gamma = (k + 1 - 2\beta)/(k + 1)$ if $\beta \in [0, 1/2)$ and $\gamma = k/(k + 1)$ if $\beta > 1/2$, k is defined in Theorem 2. K = number of clusters, and t = number of iterations. Details in E.2 of the Appendix.

Method	Time Complexity	Memory Complexity
FALKON (Rudi et al. (2017); $M \asymp \sqrt{N}$)	$O(N^{3/2})$	$O(N)$
FALKON (Meanti et al. (2020))	$O(NM) + O(M^2)$	$O(NM + M^2)$
Proposed + FALKON	$O(N^\gamma M + N^\gamma Kt)$	$O(N^\gamma M + M^2 + K)$

6 DISCUSSION

A growing body of recent work has demonstrated theoretical equivalences and connections between deep neural networks (DNNs) and kernel learning methods (e.g., Jacot et al. (2018); Zhu et al. (2022); Zhang et al. (2024)). As highlighted in Belkin et al. (2018), developing a deeper understanding of more tractable kernel methods is an important step toward building a solid theoretical foundation for DNNs. This paper contributes to the understanding of kernel learning for large-scale data sets, both theoretically and empirically. We address the computational challenges of kernel learning by selecting an informative, prediction-oriented subset of the data, allowing kernel machines to operate on a much smaller effective sample size. Unlike sketching methods that approximate the kernel matrix, our approach selects informative samples by assigning sampling weights that minimize the prediction MSE.

A key finding is that the proposed method can improve the convergence rate by sampling from the full data in an informative manner. When integrated with FALKON, our method substantially reduces FALKON’s computational cost on the full data while maintaining its predictive accuracy. Numerical results show that the proposed method achieves lower IMSE in estimating the unknown nonparametric function than SRS-based methods and offers IMSE comparable to the full-data estimator. It also delivers competitive out-of-sample predictive performance relative to existing sketching algorithms, such as FALKON and Nyström methods.

REFERENCES

- 486
487
488 Amirhesam Abedsoltan, Mikhail Belkin, and Parthe Pandit. Toward large kernel models. In Andreas
489 Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan
490 Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume
491 202 of *Proceedings of Machine Learning Research*, pp. 61–78. PMLR, 23–29 Jul 2023. URL
492 <https://proceedings.mlr.press/v202/abedsoltan23a.html>.
- 493 Mingyao Ai, Jun Yu, Huiming Zhang, and Haiying Wang. Optimal subsampling algorithms for big
494 data regressions. *Statistica Sinica*, 2018. URL [https://api.semanticscholar.org/
495 CorpusID:198455923](https://api.semanticscholar.org/CorpusID:198455923).
- 496 Ahmed Alaoui and Michael W Mahoney. Fast randomized kernel ridge regression with
497 statistical guarantees. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett
498 (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Asso-
499 ciates, Inc., 2015. URL [https://proceedings.neurips.cc/paper_files/paper/
500 2015/file/f3f27a324736617f20abbf2ffd806f6d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2015/file/f3f27a324736617f20abbf2ffd806f6d-Paper.pdf).
- 501 Mikhail Belkin, Siyuan Ma, and Soumik Mandal. To understand deep learning we need to understand
502 kernel learning. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International
503 Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*,
504 pp. 541–549. PMLR, 10–15 Jul 2018. URL [https://proceedings.mlr.press/v80/
505 belkin18a.html](https://proceedings.mlr.press/v80/belkin18a.html).
- 506 Ming-Chung Chang. Supervised stratified subsampling for predictive analytics. *Journal of Computa-
507 tional and Graphical Statistics*, 33(3):1017–1036, 2024. doi: 10.1080/10618600.2024.2304075.
508 URL <https://doi.org/10.1080/10618600.2024.2304075>.
- 509 Xueying Chen, Jerry Q. Cheng, and Min-ge Xie. *Divide-and-Conquer Methods for Big Data Analysis*,
510 pp. 1–15. John Wiley Sons, Ltd, 2021. ISBN 9781118445112. doi: [https://doi.org/10.1002/
511 9781118445112.stat08298](https://doi.org/10.1002/9781118445112.stat08298). URL [https://onlinelibrary.wiley.com/doi/abs/10.
512 1002/9781118445112.stat08298](https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118445112.stat08298).
- 513 Qianshun Cheng, HaiYing Wang, and Min Yang. Information-based optimal subdata selection for
514 big data logistic regression. *Journal of Statistical Planning and Inference*, 209:112–122, 2020.
- 515 Petros Drineas, Michael W Mahoney, and Shan Muthukrishnan. Sampling algorithms for l_2
516 regression and applications. In *Proceedings of the seventeenth annual ACM-SIAM symposium on
517 Discrete algorithm*, pp. 1127–1136, 2006.
- 518 Chong Gu. *Smoothing spline ANOVA models*, volume 297. Springer, 2013.
- 519 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: convergence and
520 generalization in neural networks. In *Proceedings of the 32nd International Conference on Neural
521 Information Processing Systems, NIPS’ 18*, pp. 8580–8589, Red Hook, NY, USA, 2018. Curran
522 Associates Inc.
- 523 Leonard Kaufman and Peter J Rousseeuw. *Finding Groups in Data: An Introduction to Cluster
524 Analysis*. Wiley, 1990.
- 525 George Kimeldorf and Grace Wahba. Some results on tchebycheffian spline functions. *Journal of
526 mathematical analysis and applications*, 33(1):82–95, 1971.
- 527 Runze Li, Dennis K.J. Lin, and Bing Li. Statistical inference in massive data sets. *Applied Stochastic
528 Models in Business and Industry*, 29(5):399–409, 2013. doi: <https://doi.org/10.1002/asmb.1927>.
529 URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/asmb.1927>.
- 530 Dawei Liu, Xihong Lin, and Debashis Ghosh. Semiparametric regression of multidimensional
531 genetic pathway data: Least-squares kernel machines and linear mixed models. *Biometrics*, 63(4):
532 1079–1088, 2007.
- 533 Ping Ma and Xiaoxiao Sun. Leveraging for big data regression. *Wiley Interdisciplinary Reviews:
534 Computational Statistics*, 7(1):70–76, 2015.

- 540 SIYUAN MA and Mikhail Belkin. Diving into the shallows: a computational perspective on large-
541 scale shallow learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vish-
542 wanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30.
543 Curran Associates, Inc., 2017. URL [https://proceedings.neurips.cc/paper_](https://proceedings.neurips.cc/paper_files/paper/2017/file/bf424cb7b0dea050a42b9739eb261a3a-Paper.pdf)
544 [files/paper/2017/file/bf424cb7b0dea050a42b9739eb261a3a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/bf424cb7b0dea050a42b9739eb261a3a-Paper.pdf).
- 545 Michael W Mahoney. Randomized algorithms for matrices and data. *arXiv preprint arXiv:1104.5557*,
546 2011.
- 547
548 Giacomo Meanti, Luigi Carratino, Lorenzo Rosasco, and Alessandro Rudi. Kernel methods through
549 the roof: Handling billions of points efficiently. In *Advances in Neural Information Processing*
550 *Systems 32 (NeurIPS 2020)*, 2020.
- 551 James Mercer. Functions of positive and negative type and their connection with the theory of integral
552 equations. *Philos. Trans. R. Soc. Lond.*, 209:415–446, 1909.
- 553
554 Cameron Musco and Christopher Musco. Recursive sampling for the nystrom method. In
555 I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Gar-
556 nett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Asso-
557 ciates, Inc., 2017. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2017/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf)
558 [2017/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf).
- 559 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines.
560 In J. Platt, D. Koller, Y. Singer, and S. Roweis (eds.), *Advances in Neural In-*
561 *formation Processing Systems*, volume 20. Curran Associates, Inc., 2007. URL
562 [https://proceedings.neurips.cc/paper_files/paper/2007/file/](https://proceedings.neurips.cc/paper_files/paper/2007/file/013a006f03dbc5392effeb8f18fda7555-Paper.pdf)
563 [013a006f03dbc5392effeb8f18fda7555-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2007/file/013a006f03dbc5392effeb8f18fda7555-Paper.pdf).
- 564
565 Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Less is more: Nyström com-
566 putational regularization. In *Neural Information Processing Systems*, 2015. URL [https:](https://api.semanticscholar.org/CorpusID:384343)
567 [//api.semanticscholar.org/CorpusID:384343](https://api.semanticscholar.org/CorpusID:384343).
- 568 Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Falkon: An optimal large scale kernel
569 method. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- 570
571 Alessandro Rudi, Daniele Calandriello, Luigi Carratino, and Lorenzo Rosasco. On fast leverage
572 score sampling and optimal learning. In *Proceedings of the 32nd International Conference on*
573 *Neural Information Processing Systems, NIPS’18*, pp. 5677–5687, Red Hook, NY, USA, 2018.
574 Curran Associates Inc.
- 575 Tao Shi, Mikhail Belkin, and Bin Yu. Data spectroscopy: Eigenspaces of convolution operators and
576 clustering. *The Annals of Statistics*, 37(6B):3960 – 3984, 2009. doi: 10.1214/09-AOS700. URL
577 <https://doi.org/10.1214/09-AOS700>.
- 578
579 Grace Wahba. *Spline models for observational data*. SIAM, 1990.
- 580
581 Haiying Wang and Yanyuan Ma. Optimal subsampling for quantile regression in big data. *Biometrika*,
582 108(1):99–112, 07 2020. ISSN 0006-3444. doi: 10.1093/biomet/asaa043. URL [https://doi.](https://doi.org/10.1093/biomet/asaa043)
583 [org/10.1093/biomet/asaa043](https://doi.org/10.1093/biomet/asaa043).
- 584 HaiYing Wang, Min Yang, and John Stufken. Information-based optimal subdata selection for big
585 data linear regression. *Journal of the American Statistical Association*, 114(525):393–405, 2019a.
- 586
587 Shusen Wang, Alex Gittens, and Michael W. Mahoney. Scalable kernel k-means clustering with
588 nyström approximation: relative-error bounds. *J. Mach. Learn. Res.*, 20(1):431–479, January
589 2019b. ISSN 1532-4435.
- 590
591 Yuedong Wang. *Smoothing splines: methods and applications*. CRC press, 2011.
- 592
593 Christopher K. I. Williams and Matthias Seeger. Using the nyström method to speed up kernel
machines. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2001.

594 Tianbao Yang, Yu-feng Li, Mehrdad Mahdavi, Rong Jin, and Zhi-Hua Zhou. Nyström method vs ran-
595 dom fourier features: A theoretical and empirical comparison. In F. Pereira, C.J. Burges, L. Bottou,
596 and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 25. Cur-
597 ran Associates, Inc., 2012. URL [https://proceedings.neurips.cc/paper_files/
598 paper/2012/file/621bf66ddb7c962aa0d22ac97d69b793-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2012/file/621bf66ddb7c962aa0d22ac97d69b793-Paper.pdf).
599

600 Yaqiong Yao and HaiYing Wang. A review on optimal subsampling methods for massive datasets.
601 *Journal of Data Science*, 19(1):151–172, 2021.

602 Ming Yuan and T Tony Cai. A reproducing kernel hilbert space approach to functional linear
603 regression. *The Annals of Statistics*, 38(6):3412–3444, 2010.

604

605 Yuchen Zhang, John Duchi, and Martin Wainwright. Divide and conquer kernel ridge regression: a
606 distributed algorithm with minimax optimal rates. *J. Mach. Learn. Res.*, 16(1):3299–3340, January
607 2015. ISSN 1532-4435.

608 Zijian Zhang, Clark Hensley, and Zhiqian Chen. Improving node classification with neural tangent
609 kernel: A graph neural network approach. In *Proceedings of the International Conference on
610 Machine Learning, Pattern Recognition and Automation Engineering, MLPRAE '24*, pp. 93–97,
611 New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400709876. doi:
612 10.1145/3696687.3696703. URL <https://doi.org/10.1145/3696687.3696703>.

613 Libin Zhu, Chaoyue Liu, and Mikhail Belkin. Transition to linearity of general neural networks with
614 directed acyclic graph architecture. NIPS '22, Red Hook, NY, USA, 2022. Curran Associates Inc.
615 ISBN 9781713871088.
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A PROOF OF THEOREM 1

In this proof, we derive the asymptotic mean squared error of the nonparametric estimator for full data. First, we recall that, the objective function to estimate $f_0(x)$ is given by:

$$\hat{f}_{N,\lambda^*} = \arg \min_{f \in \mathcal{H}_K} \left\{ \frac{1}{N} \sum_{i=1}^N (y_i - f(X_i))^2 + \lambda^* \|f\|_{\mathcal{H}_K}^2 \right\},$$

where the norm of any function f in \mathcal{H}_K is:

$$\|f\|_{\mathcal{H}_K}^2 := \sum_{j=1}^{\infty} \frac{c_j^2}{\lambda_j} = \sum_{j=1}^{\infty} \lambda_j \langle f, \psi_j \rangle_{\mathcal{H}_K}^2 < \infty.$$

So the objective function is equivalent to:

$$\hat{f}_{N,\lambda^*} = \arg \min_{f \in \mathcal{H}_K} \left\{ \frac{1}{N} \sum_{i=1}^N (y_i - f(X_i))^2 + \lambda^* \sum_{j=1}^{\infty} \frac{c_j^2}{\lambda_j} \right\}.$$

To investigate the asymptotic mean squared error, we will decompose the MSE of the estimator into deterministic error and stochastic error. More specifically, the estimation error between \hat{f}_{N,λ^*} and the true function f_0 can be decomposed as

$$\hat{f}_{N,\lambda^*} - f_0 = (\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}) + (\bar{f}_{\infty,\lambda^*} - f_0)$$

where we refer $\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}$ as stochastic error and $\bar{f}_{\infty,\lambda^*} - f_0$ as deterministic error. Here $\bar{f}_{\infty,\lambda^*}$ is the solution of the following objective function:

$$\bar{f}_{\infty,\lambda^*} = \arg \min_{f \in \mathcal{H}_K} \{l_{\infty}(f) + \lambda^* \|f\|_{\mathcal{H}_K}^2\}$$

where the loss function $l_{\infty}(f)$ is the limit of $l_N(f) = \frac{1}{N} \sum_{i=1}^N (y_i - f(X_i))^2$. That is

$$\begin{aligned} l_{\infty}(f) &:= \mathbb{E}(l_N(f)) = \mathbb{E}\{y - f(X)\}^2 \\ &= \sigma^2 + \mathbb{E}\{f(X) - f_0(X)\}^2. \end{aligned}$$

Define the functional \hat{f}_{N,λ^*} of the population estimator by collecting the estimators at all $x_0 \in \mathcal{X}$ together $\hat{f}_{N,\lambda^*} = \{\hat{f}_{N,\lambda^*}(x_0) : x_0 \in \mathcal{X}\}$, where \mathcal{X} is the space of the predictor/feature x . Define the norm:

$$\begin{aligned} \|\hat{f}_{N,\lambda^*} - f\|^2 &= \int \langle \hat{f}_{N,\lambda^*} - f, K(x_0, \cdot) \rangle_{\mathcal{H}_K}^2 \pi(x_0) dx_0 \\ &= \int \langle \hat{f}_{N,\lambda^*} - f, \sum_{j=1}^{\infty} \lambda_j \psi_j(\cdot) \psi_j(x_0) \rangle_{\mathcal{H}_K}^2 \pi(x_0) dx_0 \\ &= \int \left(\sum_{j=1}^{\infty} \lambda_j \langle \hat{f}_{N,\lambda^*} - f, \psi_j(\cdot) \rangle_{\mathcal{H}_K} \psi_j(x_0) \right)^2 \pi(x_0) dx_0 \\ &= \sum_{j=1}^{\infty} \lambda_j^2 \langle \hat{f}_{N,\lambda^*} - f, \psi_j(\cdot) \rangle_{\mathcal{H}_K}^2 \\ &= \sum_{j=1}^{\infty} \lambda_j^2 \langle \sum_{l=1}^{\infty} (\hat{c}_l - c_l) \psi_l, \psi_j \rangle_{\mathcal{H}_K}^2 \\ &= \sum_{j=1}^{\infty} \lambda_j^2 (\hat{c}_j - c_j)^2 \langle \psi_j, \psi_j \rangle_{\mathcal{H}_K}^2 \\ &= \sum_{j=1}^{\infty} (\hat{c}_j - c_j)^2. \end{aligned}$$

We will prove the theorem using parts A.1-A.3 given below.

702 A.1 ASYMPTOTIC ORDER FOR DETERMINISTIC ERROR $\bar{f}_{\infty, \lambda^*} - f_0$.

703 Write $f_0(x) = \sum_{j=1}^{\infty} a_j \psi_j(x)$ and $f(x) = \sum_{j=1}^{\infty} c_j \psi_j(x)$, then we have

$$704 \quad l_{\infty}(f) = \sigma^2 + \sum_{j=1}^{\infty} (c_j - a_j)^2.$$

705 By the orthogonality of the eigenfunctions, we have $\int_R \psi_j(x)^2 \pi(x) dx = 1$ and
706 $\int_R \psi_j(x) \psi_l(x) \pi(x) dx = 0$, then we can derive

$$\begin{aligned} 707 \quad l_{\infty}(f) &= \sigma^2 + \mathbb{E}\{f(X) - f_0(X)\}^2 \\ 708 \quad &= \sigma^2 + \mathbb{E}\left\{\sum_{j=1}^{\infty} c_j \psi_j(X) - \sum_{l=1}^{\infty} a_l \psi_l(X)\right\}^2 \\ 709 \quad &= \sigma^2 + \mathbb{E}\left\{\sum_{j=1}^{\infty} \sum_{l=1}^{\infty} (c_j - a_j)(c_l - a_l) \psi_j(X) \psi_l(X)\right\} \\ 710 \quad &= \sigma^2 + \sum_{j=1}^{\infty} (c_j - a_j)^2 \mathbb{E}\{\psi_j(X)^2\} + \sum_{j \neq l}^{\infty} (c_j - a_j)(c_l - a_l) \mathbb{E}\{\psi_j(X) \psi_l(X)\} \\ 711 \quad &= \sigma^2 + \sum_{j=1}^{\infty} (c_j - a_j)^2 \int_R \psi_j(x)^2 \pi(x) dx + \sum_{j \neq l}^{\infty} (c_j - a_j)(c_l - a_l) \int_R \psi_j(x) \psi_l(x) \pi(x) dx \\ 712 \quad &= \sigma^2 + \sum_{j=1}^{\infty} (c_j - a_j)^2. \end{aligned}$$

713 Then the corresponding objective function can be expressed as:

$$714 \quad \bar{f}_{\infty, \lambda^*}(c_j) = \arg \min\{Q_{\infty}(c_j)\} = \arg \min\left\{\sigma^2 + \sum_{j=1}^{\infty} (c_j - a_j)^2 + \lambda^* \sum_{j=1}^{\infty} \frac{c_j^2}{\lambda_j}\right\}.$$

715 We then take derivative w.r.t. c_j and obtain

$$716 \quad Q'_{\infty}(c_j) = 2c_j - 2a_j + 2\lambda^* \lambda_j^{-1} c_j.$$

717 It follows that the minimizer of the above objective function can be written as:

$$718 \quad \bar{f}_{\infty, \lambda^*}(x) = \sum_{j=1}^{\infty} \bar{c}_j \psi_j(x) = \sum_{j=1}^{\infty} \frac{a_j}{1 + \lambda^* \lambda_j^{-1}} \psi_j(x),$$

719 where $\bar{c}_j = a_j / (1 + \lambda^* \lambda_j^{-1})$.

720 To bound the deterministic error, assume $a_j^2 \lambda_j^{-1} = j^{-a}$ with $a > 1$ and $\lambda_j \asymp j^{-k}$ ($k > 1$), so
721 $a_j^2 \asymp j^{-(a+k)}$:

$$722 \quad \|\bar{f}_{\infty, \lambda^*} - f_0\|^2 := \sum_{j=1}^{\infty} (\bar{c}_j - a_j)^2 = \sum_{j=1}^{\infty} \left(\frac{\lambda^* \lambda_j^{-1}}{1 + \lambda^* \lambda_j^{-1}}\right)^2 a_j^2 = \sum_{j=1}^{\infty} \left(\frac{\lambda^*}{\lambda_j + \lambda^*}\right)^2 a_j^2.$$

723 Let J solve $\lambda_J \asymp \lambda^*$ so $J \asymp \lambda^{*-1/k}$. Split the sum at J :

$$724 \quad \sum_{j \leq J} \left(\frac{\lambda^*}{\lambda_j + \lambda^*}\right)^2 a_j^2 \lesssim \lambda^{*2} \sum_{j \leq J} \frac{a_j^2}{\lambda_j^2} \asymp \lambda^{*2} \sum_{j \leq J} j^{k-a} \asymp \lambda^{*2} J^{1+k-a} \asymp \lambda^{*\frac{a+k-1}{k}},$$

$$725 \quad \sum_{j > J} \left(\frac{\lambda^*}{\lambda_j + \lambda^*}\right)^2 a_j^2 \lesssim \sum_{j > J} a_j^2 \asymp \sum_{j > J} j^{-(a+k)} \asymp J^{-(a+k-1)} \asymp \lambda^{*\frac{a+k-1}{k}}.$$

726 Therefore,

$$727 \quad \|\bar{f}_{\infty, \lambda^*} - f_0\|^2 \asymp \lambda^{*\frac{a+k-1}{k}} \quad (a > 1, k > 1).$$

A.2 ASYMPTOTIC ORDER FOR STOCHASTIC ERROR $\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}$.

Recall that

$$l_N(f) = \frac{1}{N} \sum_{i=1}^N \{y_i - f(X_i)\}^2,$$

and the objective function for \hat{f}_{N,λ^*} can be written as:

$$\hat{f}_{N,\lambda^*} = \arg \min_f \{Q_N(f)\} = \arg \min_f \{l_N(f) + \lambda^* \|f\|_{\mathcal{H}_K}^2\}.$$

Note that the functional derivatives of the objective function with respect to f is stochastic, which leads to a stochastic denominator in the solution of the above objective function, and it is not straightforward to handle a stochastic denominator. To avoid such difficulty, we define an intermediate quantity

$$\tilde{f}_{N,\lambda^*} = \bar{f}_{\infty,\lambda^*} - \frac{1}{2} G_{\lambda^*}^{-1} D l_{N,\lambda^*}(\bar{f}_{\infty,\lambda^*})$$

where $G_{\lambda^*} = \frac{1}{2} D^2 l_{\infty,\lambda^*}(\bar{f}_{\infty,\lambda^*})$. Then we could write:

$$\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*} = (\hat{f}_{N,\lambda^*} - \tilde{f}_{N,\lambda^*}) + (\tilde{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}). \quad (3)$$

To find the orders of the above two terms, we need the functional derivatives given below. For functions $\eta, g \in \mathcal{H}_K$ and define the dot product $\eta \cdot g = \langle \eta, g \rangle_{\mathcal{H}_K}$

$$\begin{aligned} D l_N(f) \cdot \eta &= \frac{d}{dh} \frac{1}{N} \sum_{i=1}^N \{y_i - \langle f + h\eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K}\}^2 \Big|_{h=0} \\ &= \frac{d}{dh} \frac{1}{N} \sum_{i=1}^N \{y_i - \langle f, K(X_i, \cdot) \rangle_{\mathcal{H}_K} - h \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K}\}^2 \Big|_{h=0} \\ &= -\frac{2}{N} \sum_{i=1}^N \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \{y_i - \langle f, K(X_i, \cdot) \rangle_{\mathcal{H}_K}\} \\ &= -\frac{2}{N} \sum_{i=1}^N \eta(X_i) \{y_i - f(X_i)\}. \end{aligned}$$

$$\begin{aligned} D l_{\infty}(f) \cdot \eta &= -2 \int \eta(x) (f_0(x) - f(x)) \pi(x) dx \\ &= -2 \int \langle \eta, K(x, \cdot) \rangle_{\mathcal{H}_K} \langle f_0 - f, K(x, \cdot) \rangle_{\mathcal{H}_K} \pi(x) dx \end{aligned}$$

$$\begin{aligned} D^2 l_N(f) \cdot \eta \cdot g &= \frac{2}{N} \sum_{i=1}^N \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \langle g, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \\ &= \frac{2}{N} \sum_{i=1}^N \langle \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K} K(X_i, \cdot), g \rangle_{\mathcal{H}_K}. \end{aligned}$$

$$\begin{aligned} D^2 l_{\infty}(f) \cdot \eta \cdot g &= 2 \left\langle \langle \eta, K(x, \cdot) \rangle_{\mathcal{H}_K} K(x, \cdot) \pi(x) dx, g \right\rangle_{\mathcal{H}_K} \\ &= 2 \sum_{j=1}^{\infty} \lambda_j^2 \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K}. \end{aligned}$$

$$D \|f\|_{\mathcal{H}_K}^2 \cdot \eta = 2 \sum_{j=1}^{\infty} \lambda_j \langle f, \psi_j \rangle_{\mathcal{H}_K} \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \quad \text{and}$$

$$D^2 \|f\|_{\mathcal{H}_K}^2 \cdot \eta \cdot g = 2 \sum_{j=1}^{\infty} \lambda_j \langle g, \psi_j \rangle_{\mathcal{H}_K} \langle \eta, \psi_j \rangle_{\mathcal{H}_K}.$$

810 A.2.1 EVALUATE THE ORDER OF $\tilde{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}$.

811
812 Starting from the second term $\tilde{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}$ in equation (3), because the functional derivatives to the
813 penalty terms $\lambda^* \|f\|_{\mathcal{H}_K}^2$ are the same for $l_{N,\lambda^*}(f)$ and $l_{\infty,\lambda^*}(f)$, and \bar{f} is the minimizer of $l_{\infty,\lambda^*}(f)$
814 that satisfies that $Dl_{\infty,\lambda^*}(\bar{f}) = 0$. So, we have

$$815 \quad Dl_{N,\lambda^*}(\bar{f}) = Dl_{N,\lambda^*}(\bar{f}) - Dl_{\infty,\lambda^*}(\bar{f}) = Dl_N(\bar{f}) - Dl_{\infty}(\bar{f}).$$

816
817 For any function η , we have

$$818 \quad \begin{aligned} 819 \quad \mathbb{E}\{Dl_{N,\lambda^*}(\bar{f}) \cdot \eta\}^2 &= \mathbb{E}\{Dl_N(\bar{f}) \cdot \eta - Dl_{\infty}(\bar{f}) \cdot \eta\}^2 \\ 820 &= \mathbb{E}\left[-\frac{2}{N} \sum_{i=1}^N \{y_i - \bar{f}(X_i)\} \eta(X_i) + 2 \int \{\bar{f}(x) - f_0(x)\} \eta(x) \pi(x) dx\right]^2 \\ 821 &= \frac{4}{N^2} \mathbb{E}\left[\sum_{i=1}^N \left\{ (y_i - \bar{f}(X_i)) \eta(X_i) - \mathbb{E}[(\bar{f}(X) - f_0(X)) \eta(X)] \right\}\right]^2 \\ 822 &= \frac{4}{N} \text{Var}\left[\{y - \bar{f}(X)\} \eta(X)\right] \leq \frac{4}{N} \mathbb{E}[\{y - \bar{f}(X)\} \eta(X)]^2 \asymp \frac{4}{N}. \end{aligned}$$

823
824
825
826
827
828 By the definition of G_{λ^*} , we have

$$829 \quad \begin{aligned} 830 \quad G_{\lambda^*} \cdot \eta \cdot g &= \frac{1}{2} D_{\infty,\lambda^*}^2 \cdot \eta \cdot g = \frac{1}{2} Dl_{\infty}^2(f) \cdot \eta \cdot g + \frac{1}{2} D^2 \|f\|_{\mathcal{H}_K}^2 \cdot \eta \cdot g \\ 831 &= \sum_{j=1}^{\infty} \lambda_j^2 \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} + \sum_{j=1}^{\infty} \lambda^* \lambda_j \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} \\ 832 &= \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda_j^{-1} \lambda^*) \langle g, \psi_j \rangle_{\mathcal{H}_K} \langle \eta, \psi_j \rangle_{\mathcal{H}_K}. \end{aligned}$$

833
834
835
836
837
838 Let $\eta = \psi_m$ gives

$$839 \quad \langle G_{\lambda^*} g, \psi_m \rangle_{\mathcal{H}_K} = (\lambda_m + \lambda^*) \langle g, \psi_m \rangle_{\mathcal{H}_K},$$

840
841 which leads to

$$842 \quad \langle G_{\lambda^*}^{-1} g, \psi_m \rangle_{\mathcal{H}_K} = (\lambda_m + \lambda^*)^{-1} \langle g, \psi_m \rangle_{\mathcal{H}_K}.$$

843
844 Then, we can bound the second term in (3) by:

$$845 \quad \begin{aligned} 846 \quad \mathbb{E}\|\tilde{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|^2 &= \mathbb{E}\left\|\frac{1}{2} G_{\lambda^*}^{-1} Dl_{N,\lambda^*}(\bar{f}_{\infty,\lambda^*})\right\|^2 \\ 847 &= \frac{1}{4} \mathbb{E}\left\{\sum_{j=1}^{\infty} \lambda_j^2 \langle G_{\lambda^*}^{-1} Dl_{N,\lambda^*}(\bar{f}_{\infty,\lambda^*}), \psi_j \rangle_{\mathcal{H}_K}^2\right\} \\ 848 &= \frac{1}{4} \mathbb{E}\left\{\sum_{j=1}^{\infty} \lambda_j^2 \frac{1}{(\lambda_j + \lambda^*)^2} \langle Dl_{N,\lambda^*}(\bar{f}_{\infty,\lambda^*}), \psi_j \rangle_{\mathcal{H}_K}^2\right\} \\ 849 &= \frac{1}{4} \sum_{j=1}^{\infty} \frac{\lambda_j^2}{(\lambda_j + \lambda^*)^2} \mathbb{E}\left\{\langle Dl_{N,\lambda^*}(\bar{f}_{\infty,\lambda^*}), \psi_j \rangle_{\mathcal{H}_K}^2\right\} \\ 850 &\lesssim \frac{1}{N} \sum_{j=1}^{\infty} \frac{\lambda_j^2}{(\lambda_j + \lambda^*)^2} \quad (\text{using } \mathbb{E}\{\langle Dl_{N,\lambda^*}(\bar{f}), \psi_j \rangle_{\mathcal{H}_K}^2\} \asymp N^{-1}) \\ 851 &\asymp \frac{1}{N} \int_1^{\infty} \frac{1}{(1 + \lambda^* j^k)^2} dj \quad (\text{since } \lambda_j \asymp j^{-k}, k > 1) \\ 852 &\asymp \frac{1}{N} \lambda^{*-1/k}. \end{aligned}$$

864 A.2.2 EVALUATE THE ORDER OF $\hat{f}_{N,\lambda^*} - \tilde{f}_{N,\lambda^*}$.

865
866 Next we need to find the stochastic order of the second term $\hat{f}_{N,\lambda^*} - \tilde{f}_{N,\lambda^*}$ in the expression (3). For
867 brevity, write

$$868 \hat{f} := \hat{f}_{N,\lambda^*}, \quad \tilde{f} := \tilde{f}_{N,\lambda^*}, \quad \bar{f} := \bar{f}_{\infty,\lambda^*}.$$

869
870 Note that \hat{f} is the solution of the following first order equation

$$871 D l_{N,\lambda^*}(\hat{f}) = D l_{N,\lambda^*}(\bar{f}) + D^2 l_{N,\lambda^*}(\bar{f}) \cdot (\hat{f} - \bar{f}) = 0,$$

872
873 and by the definition of \tilde{f} , we have

$$874 D l_{N,\lambda^*}(\tilde{f}) + D^2 l_{\infty,\lambda^*}(\tilde{f}) \cdot (\tilde{f} - \bar{f}) = 0.$$

875
876 From the above two equations, we can find:

$$877 D^2 l_{N,\lambda^*}(\bar{f}) \cdot (\hat{f} - \bar{f}) = D^2 l_{\infty,\lambda^*}(\tilde{f}) \cdot (\hat{f} - \bar{f}).$$

880 Then we can write:

$$881 D^2 l_{\infty,\lambda^*}(\tilde{f}) \cdot (\hat{f} - \tilde{f}) = D^2 l_{\infty,\lambda^*}(\tilde{f}) \cdot (\hat{f} - \bar{f}) + D^2 l_{\infty,\lambda^*}(\tilde{f}) \cdot (\bar{f} - \tilde{f})$$

$$882 = D^2 l_{\infty,\lambda^*}(\tilde{f}) \cdot (\hat{f} - \bar{f}) - D^2 l_{N,\lambda^*}(\tilde{f}) \cdot (\hat{f} - \bar{f})$$

$$883 = D^2 l_{\infty}(\tilde{f}) \cdot (\hat{f} - \bar{f}) - D^2 l_N(\tilde{f}) \cdot (\hat{f} - \bar{f}).$$

884
885 Using the definition of G_{λ^*} , we have

$$886 \hat{f} - \tilde{f} = \frac{1}{2} G_{\lambda^*}^{-1} D^2 l_{\infty,\lambda^*}(\tilde{f}) \cdot (\hat{f} - \tilde{f}) = \frac{1}{2} G_{\lambda^*}^{-1} \left\{ D^2 l_{\infty}(\tilde{f}) \cdot (\hat{f} - \bar{f}) - D^2 l_N(\tilde{f}) \cdot (\hat{f} - \bar{f}) \right\}.$$

887
888 Using the similar steps in evaluating $\mathbb{E}\|\tilde{f} - \bar{f}\|^2$, we have

$$889 \|\hat{f} - \tilde{f}\|^2$$

$$890 = \left\| \frac{1}{2} G_{\lambda^*}^{-1} \left\{ D^2 l_{\infty}(\tilde{f})(\hat{f} - \bar{f}) - D^2 l_N(\tilde{f})(\hat{f} - \bar{f}) \right\} \right\|^2$$

$$891 = \frac{1}{4} \sum_{j=1}^{\infty} \lambda_j^2 \left\langle G_{\lambda^*}^{-1} \left\{ D^2 l_{\infty}(\tilde{f})(\hat{f} - \bar{f}) - D^2 l_N(\tilde{f})(\hat{f} - \bar{f}) \right\}, \psi_j \right\rangle_{\mathcal{H}_K}^2$$

$$892 = \frac{1}{4} \sum_{j=1}^{\infty} \lambda_j^2 \frac{1}{(\lambda_j + \lambda^*)^2} \left\{ \left\langle D^2 l_{\infty}(\tilde{f})(\hat{f} - \bar{f}), \psi_j \right\rangle_{\mathcal{H}_K} - \left\langle D^2 l_N(\tilde{f})(\hat{f} - \bar{f}), \psi_j \right\rangle_{\mathcal{H}_K} \right\}^2$$

$$893 = \frac{1}{4} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \left\{ \left\langle D^2 l_{\infty}(\tilde{f})(\hat{f} - \bar{f}), \psi_j \right\rangle_{\mathcal{H}_K} - \left\langle D^2 l_N(\tilde{f})(\hat{f} - \bar{f}), \psi_j \right\rangle_{\mathcal{H}_K} \right\}^2$$

$$894 = \frac{1}{4} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \left\{ 2 \sum_{l=1}^{\infty} \lambda_l^2 \left\langle \hat{f} - \bar{f}, \psi_l \right\rangle_{\mathcal{H}_K} \left\langle \psi_j, \psi_l \right\rangle_{\mathcal{H}_K} \right.$$

$$895 \left. - \frac{2}{N} \sum_{i=1}^N \sum_{l=1}^{\infty} \sum_{k=1}^{\infty} \lambda_l \lambda_k \left\langle \hat{f} - \bar{f}, \psi_l \right\rangle_{\mathcal{H}_K} \left\langle \psi_j, \psi_k \right\rangle_{\mathcal{H}_K} \psi_l(X_i) \psi_k(X_i) \right\}^2$$

$$896 = \frac{1}{4} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \left[\frac{2}{N} \sum_{i=1}^N \sum_{l=1}^{\infty} \lambda_l \left\langle \hat{f} - \bar{f}, \psi_l \right\rangle_{\mathcal{H}_K} \mathbb{E} \{ \psi_j(X_i) \psi_l(X_i) \} \right.$$

$$897 \left. - \frac{2}{N} \sum_{i=1}^N \sum_{l=1}^{\infty} \lambda_l \left\langle \hat{f} - \bar{f}, \psi_l \right\rangle_{\mathcal{H}_K} \psi_l(X_i) \psi_k(X_i) \right]^2$$

$$898 = \frac{1}{4} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \left[\frac{2}{N} \sum_{l=1}^{\infty} \lambda_l^{\frac{1}{2}} \left\langle \hat{f} - \bar{f}, \psi_l \right\rangle_{\mathcal{H}_K} \lambda_l^{\frac{1}{2}} \left\{ \sum_{i=1}^N \mathbb{E} \{ \psi_j(X_i) \psi_l(X_i) \} - \sum_{i=1}^N \psi_j(X_i) \psi_l(X_i) \} \right\} \right]^2$$

$$\begin{aligned}
&\leq N^{-2} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \left(\sum_{l=1}^{\infty} \lambda_l \langle \hat{f} - \bar{f}, \psi_l \rangle_{\mathcal{H}_K}^2 \right) \left[\sum_{l=1}^{\infty} \lambda_l \left\{ \sum_{i=1}^N \mathbb{E}(\psi_j(X_i) \psi_l(X_i)) - \sum_{i=1}^N \psi_j(X_i) \psi_l(X_i) \right\}^2 \right] \\
&= N^{-2} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \|\hat{f} - \bar{f}\|_{\mathcal{H}_K}^2 \underbrace{\left[\sum_{l=1}^{\infty} \lambda_l \left\{ \sum_{i=1}^N \mathbb{E}(\psi_j(X_i) \psi_l(X_i)) - \sum_{i=1}^N \psi_j(X_i) \psi_l(X_i) \right\}^2 \right]}_{\text{Term A}}.
\end{aligned}$$

For Term A, we can expand it as:

$$\begin{aligned}
&\sum_{l=1}^{\infty} \lambda_l \left[\sum_{i=1}^N \left\{ \mathbb{E}(\psi_j(X_i) \psi_l(X_i)) - \psi_j(X_i) \psi_l(X_i) \right\} \right]^2 \\
&= \sum_{l=1}^{\infty} \lambda_l \left[\sum_{i=1}^N \left\{ \mathbb{E}(\psi_j(X_i) \psi_l(X_i)) - \psi_j(X_i) \psi_l(X_i) \right\}^2 \right] \\
&\quad + \sum_{l=1}^{\infty} \lambda_l \left[\sum_{i \neq k}^N \left\{ \mathbb{E}(\psi_j(X_i) \psi_l(X_i)) - \psi_j(X_i) \psi_l(X_i) \right\} \left\{ \mathbb{E}(\psi_j(X_k) \psi_l(X_k)) - \psi_j(X_k) \psi_l(X_k) \right\} \right] \\
&= \lambda_j \underbrace{\sum_{i=1}^N \left\{ \mathbb{E}(\psi_j^2(X_i)) - \psi_j^2(X_i) \right\}^2}_{\text{Term 1}} \\
&\quad + \underbrace{\sum_{l=1}^{\infty} \lambda_l \sum_{i \neq k}^N \left[\mathbb{E}\{\psi_j(X_i) \psi_l(X_i)\} - \psi_j(X_i) \psi_l(X_i) \right] \left[\mathbb{E}\{\psi_j(X_k) \psi_l(X_k)\} - \psi_j(X_k) \psi_l(X_k) \right]}_{\text{Term 2}}.
\end{aligned}$$

To find the orders of Term 1 and Term 2, we evaluate the expectation of Term 1 and the variance of Term 2, which are given below:

$$\begin{aligned}
&\mathbb{E} \left[\lambda_j \sum_{i=1}^N \left\{ \mathbb{E}(\psi_j^2(X_i)) - \psi_j^2(X_i) \right\}^2 \right] \\
&= \mathbb{E} \left[\lambda_j \sum_{i=1}^N \left\{ \psi_j^2(X_i) - 1 \right\}^2 + \sum_{l \neq j}^{\infty} \lambda_l \sum_{i=1}^N \left\{ \psi_j(X_i) \psi_l(X_i) \right\}^2 \right] \\
&= \lambda_j N \mathbb{E} \left\{ \psi_j^2(X) - 1 \right\}^2 + N \mathbb{E} \left\{ \sum_{l \neq j}^{\infty} \lambda_l \psi_j^2(X) \psi_l^2(X) \right\} \\
&\leq \lambda_j N \text{Var} \left\{ \psi_j^2(X) \right\} + N \left[\mathbb{E} \left\{ \sum_{l=1}^{\infty} \lambda_l \psi_l^2(X) \right\}^2 + \text{Var} \left\{ \sum_{l=1}^{\infty} \lambda_l \psi_l^2(X) \right\} \right]^{1/2} \left[\mathbb{E} \left\{ \psi_j^4(X) \right\} \right]^{1/2} \\
&= \lambda_j N \text{Var} \left\{ \psi_j^2(X) \right\} + N \left[\mathbb{E} \{ K(X, X) \}^2 + \text{Var} \{ K(X, X) \} \right]^{1/2} \left[\mathbb{E} \left\{ \psi_j^4(X) \right\} \right]^{1/2} \\
&\asymp N.
\end{aligned}$$

$$\begin{aligned}
&\text{Var} \left[\sum_{l=1}^{\infty} \lambda_l \sum_{i \neq k}^N \left\{ \mathbb{E}[\psi_j(X_i) \psi_l(X_i)] - \psi_j(X_i) \psi_l(X_i) \right\} \left\{ \mathbb{E}[\psi_j(X_k) \psi_l(X_k)] - \psi_j(X_k) \psi_l(X_k) \right\} \right] \\
&\leq \sum_{l=1}^{\infty} \lambda_l^2 (N-1)^2 \mathbb{E} \left[\left\{ \mathbb{E}[\psi_j(X_i) \psi_l(X_i)] - \psi_j(X_i) \psi_l(X_i) \right\}^2 \right] \mathbb{E} \left[\left\{ \mathbb{E}[\psi_j(X_k) \psi_l(X_k)] - \psi_j(X_k) \psi_l(X_k) \right\}^2 \right] \\
&\quad + \sum_{l \neq v}^{\infty} \lambda_l \lambda_v (N-1)^2 \text{Var} \left\{ \psi_j(X) \psi_l(X) \right\} \text{Var} \left\{ \psi_j(X) \psi_v(X) \right\}
\end{aligned}$$

$$\begin{aligned}
&= (N-1)^2 \sum_{l=1}^{\infty} \lambda_l^2 \text{Var}^2\{\psi_j(X)\psi_l(X)\} + (N-1)^2 \sum_{l \neq v}^{\infty} \lambda_l \lambda_v \text{Var}\{\psi_j(X)\psi_l(X)\} \text{Var}\{\psi_j(X)\psi_v(X)\} \\
&\leq (N-1)^2 \left[\sum_{l=1}^{\infty} \lambda_l^2 \text{Var}^2\{\psi_j(X)\psi_l(X)\} + \left[\sum_{l=1}^{\infty} \lambda_l^2 \text{Var}^2\{\psi_j(X)\psi_l(X)\} \right]^{1/2} \left[\sum_{v=1}^{\infty} \lambda_v^2 \text{Var}^2\{\psi_j(X)\psi_v(X)\} \right]^{1/2} \right] \\
&\asymp 2(N-1)^2 \int_1^{\infty} l^{-2k} dl \asymp N^2.
\end{aligned}$$

Now we found that Term 1 is at the order of N and Term 2 is at the order of $\sqrt{N^2} = N$, therefore we can write:

$$\begin{aligned}
\|\hat{f}_{N,\lambda^*} - \tilde{f}_{N,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 &= N^{-1} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 \\
&= N^{-1} \int_1^{\infty} \lambda_j^2 (1 + \lambda^* \lambda_j^{-1})^{-2} \|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 \\
&\asymp N^{-1} (\lambda^*)^{-1/k} \|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2.
\end{aligned}$$

Now, it follows that

$$\|\hat{f}_{N,\lambda^*} - \tilde{f}_{N,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 = O_p\left(N^{-1} (\lambda^*)^{-1/k} \|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2\right).$$

If $N^{-1} (\lambda^*)^{-1/k} \rightarrow 0$,

$$\|\hat{f}_{N,\lambda^*} - \tilde{f}_{N,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 = o_p\left(\|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}\right) = o_p(1) \|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}.$$

Observed that

$$\|\tilde{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 \geq \|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 - \|\hat{f}_{N,\lambda^*} - \tilde{f}_{N,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 = (1 - o_p(1)) \|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2,$$

then

$$\|\hat{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}^2 = O_p\left(\|\tilde{f}_{N,\lambda^*} - \bar{f}_{\infty,\lambda^*}\|_{\mathcal{H}_{\mathcal{X}}}\right) = O_p(N^{-1} (\lambda^*)^{-1/k}).$$

A.3 ASYMPTOTIC ORDER FOR MSE AND TUNING.

Combining the results in previous steps, we can express the order of MSE by:

$$\|\hat{f}_{N,\lambda^*} - f_0\|^2 = O_p(N^{-1} (\lambda^*)^{-1/k} + (\lambda^*)^{\frac{a+k-1}{k}}).$$

To find the optimal order of λ^* , set $M(d) = N^{-1} d^{-1/k} + d^{\frac{a+k-1}{k}}$. Then

$$M'(d) = -\frac{1}{k} N^{-1} d^{-1/k-1} + \frac{a+k-1}{k} d^{\frac{a+k-1}{k}-1} = 0,$$

which yields

$$d^{\frac{a+k}{k}} \asymp N^{-1} \implies \lambda^* \asymp N^{-\frac{k}{a+k}}.$$

Also check

$$N^{-1} (\lambda^*)^{-1/k} = N^{-1} (N^{-k/(a+k)})^{-1/k} = N^{-\frac{a+k-1}{a+k}} \rightarrow 0 \quad (a > 1, k > 1).$$

Therefore,

$$\|\hat{f}_{N,\lambda^*} - f_0\|^2 = O_p(N^{-1 + \frac{1}{a+k}}).$$

This completes the proof of Theorem 1.

B PROOF OF THEOREM 2

For any given x_0 , the proposed resampling weights and any sample $X_i \in \mathcal{C}_l$ where \mathcal{C}_l is the l -th cluster with its corresponding center C_l and its cluster size N_l , the weight/probability of X_i to be selected is:

$$\omega_i(x_0) = N_l^{-1} \omega_{x_0, l, C}.$$

So we have:

$$\sum_{i=1}^N \omega_i(x_0) = \sum_{l=1}^L \sum_{i \in \mathcal{C}_l} \frac{1}{N_l} \omega_{x_0, l, C} = \sum_{l=1}^L \omega_{x_0, l, C} = 1.$$

Denote the indicies of the resampling subset of size n for the given x_0 is $\mathcal{S}(x_0) = \{i_1^*, \dots, i_n^*\} \subset \{1, \dots, N\}$, with replacement using probabilities $\{\omega_i(x_0)\}$. Let $I_i(x_0)$ be the count of how many times index i appears in $\mathcal{S}(x_0)$ so $\mathbb{E}[I_i(x_0)] = n \omega_i(x_0)$. Define

$$l_{n, x_0}(f) = \frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} (y_i - f(X_i))^2 = \frac{1}{n} \sum_{i=1}^N I_i(x_0) (y_i - f(X_i))^2.$$

Recall that the proposed estimator of $f(x_0)$ is $\hat{f}_{n, \tilde{\lambda}^*, x_0}(x_0)$, which is given by:

$$\hat{f}_{n, \tilde{\lambda}^*, x_0} = \arg \min_{f \in \mathcal{H}_K} \left\{ \frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} (y_i - f(X_i))^2 + \tilde{\lambda}^* \|f\|_{\mathcal{H}_K}^2 \right\}.$$

Define the functional $\hat{f}_{n, \tilde{\lambda}^*}$ of the proposed estimator by collecting the estimators at all $x_0 \in \mathcal{X}$ together $\hat{f}_{n, \tilde{\lambda}^*} = \{\hat{f}_{n, \tilde{\lambda}^*, x_0}(x_0) : x_0 \in \mathcal{X}\}$, where \mathcal{X} is the space of the predictor/feature x . Define the norm:

$$\begin{aligned} \|\hat{f}_{n, \tilde{\lambda}^*} - f\|^2 &= \int \langle \hat{f}_{n, \tilde{\lambda}^*, x_0} - f, K(x_0, \cdot) \rangle_{\mathcal{H}_K}^2 \pi(x_0) dx_0 \\ &= \int \langle \hat{f}_{n, \tilde{\lambda}^*, x_0} - f, \sum_{j=1}^{\infty} \lambda_j \psi_j(\cdot) \psi_j(x_0) \rangle_{\mathcal{H}_K}^2 \pi(x_0) dx_0 \\ &= \int \left(\sum_{j=1}^{\infty} \lambda_j \langle \hat{f}_{n, \tilde{\lambda}^*, x_0} - f, \psi_j(\cdot) \rangle_{\mathcal{H}_K} \psi_j(x_0) \right)^2 \pi(x_0) dx_0. \end{aligned}$$

Here $l_{\infty}(f)$ is the limit of $l_n(f)$, with conditioning made explicit:

$$\begin{aligned} l_{\infty}(f) &:= \mathbb{E} \left[\frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} \{y_i - f(X_i)\}^2 \right] \\ &= \mathbb{E} \left[\mathbb{E} \left\{ \frac{1}{n} \sum_{i=1}^N I_i(x_0) (y_i - f(X_i))^2 \mid \{X_i\}_{i=1}^N \right\} \right] \\ &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^N \mathbb{E}(I_i(x_0) \mid \{X_i\}_{i=1}^N) (y_i - f(X_i))^2 \right] \\ &= \mathbb{E} \left[\sum_{i=1}^N \omega_i(x_0) (y_i - f(X_i))^2 \right] \quad (\mathbb{E}[I_i(x_0) \mid \{X_i\}] = n \omega_i(x_0)) \\ &= \mathbb{E} \left[\sum_{i=1}^N \omega_i(x_0) \{\sigma^2 + (f(X_i) - f_0(X_i))^2\} \right] \\ &= \sigma^2 + \mathbb{E} \left[\sum_{i=1}^N \omega_i(x_0) \left\{ \langle f - f_0, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \right\}^2 \right]. \end{aligned}$$

To evaluate the MSE of our proposed method, we evaluate the deterministic error and stochastic error separately by decomposing it in the following way:

$$\hat{f}_{n, \tilde{\lambda}^*, x_0} - f_0 = (\hat{f}_{n, \tilde{\lambda}^*, x_0} - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}) + (\bar{f}_{\infty, \tilde{\lambda}^*, x_0} - f_0),$$

where $\bar{f}_{\infty, \tilde{\lambda}^*, x_0}$ is the solution of the following objective function

$$\bar{f}_{\infty, \tilde{\lambda}^*, x_0} = \arg \min \{ l_{\infty}(f) + \tilde{\lambda}^* \|f\|_{\mathcal{H}_K}^2 \}.$$

First, we evaluate the functional derivatives for the empirical loss $l_{n, x_0}(f) = \frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} \{y_i - f(X_i)\}^2$ and for its population limit $l_{\infty, x_0}(f) = \sigma^2 + \mathbb{E} \left[\sum_{i=1}^N \omega_i(x_0) \{f(X_i) - f_0(X_i)\}^2 \right]$.

For any $\eta, g \in \mathcal{H}_K$, the derivatives of l_{∞, x_0} are

$$\begin{aligned} D l_{\infty, x_0}(f) \cdot \eta &= \frac{d}{dh} l_{\infty, x_0}(f + h\eta) \Big|_{h=0} \\ &= 2 \mathbb{E} \left[\sum_{i=1}^N \omega_i(x_0) \eta(X_i) \{f(X_i) - f_0(X_i)\} \right] \\ &= 2 \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \langle f - f_0, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \right\}. \end{aligned}$$

and

$$\begin{aligned} D^2 l_{\infty, x_0}(f) \cdot \eta \cdot g &= 2 \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \langle g, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \right\} \\ &= 2 \sum_{j=1}^{\infty} \lambda_j^2 \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} \underbrace{\mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j(X_i)^2 \right\}}_{=: \omega_{x_0, j}}. \end{aligned}$$

For the empirical loss l_{n, x_0} ,

$$\begin{aligned} D l_{n, x_0}(f) \cdot \eta &= \frac{d}{dh} \frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} (y_i - \langle f + h\eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K})^2 \Big|_{h=0} \\ &= -\frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K} (y_i - \langle f, K(X_i, \cdot) \rangle_{\mathcal{H}_K}), \end{aligned}$$

and

$$\begin{aligned} D^2 l_{n, x_0}(f) \cdot \eta \cdot g &= \frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} \langle \eta, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \langle g, K(X_i, \cdot) \rangle_{\mathcal{H}_K} \\ &= \frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} \sum_{j, k \geq 1} \lambda_j \lambda_k \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_k \rangle_{\mathcal{H}_K} \psi_j(X_i) \psi_k(X_i). \end{aligned}$$

For the RKHS penalty,

$$\begin{aligned} D \|f\|_{\mathcal{H}_K}^2 \cdot \eta &= 2 \langle f, \eta \rangle_{\mathcal{H}_K}, \\ D^2 \|f\|_{\mathcal{H}_K}^2 \cdot \eta \cdot g &= 2 \langle \eta, g \rangle_{\mathcal{H}_K} = 2 \sum_{j=1}^{\infty} \lambda_j \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K}. \end{aligned}$$

To evaluate the stochastic error, we use a similar method in the proof of Theorem 1 by defining an intermediate quantity

$$\tilde{f}_{n, \tilde{\lambda}^*, x_0} = \bar{f}_{\infty, \tilde{\lambda}^*, x_0} - \frac{1}{2} G_{\tilde{\lambda}^*, x_0}^{-1} D l_{n, x_0}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}),$$

where the local operator at x_0 is

$$G_{\tilde{\lambda}^*, x_0} = \frac{1}{2} D^2 l_{\infty, x_0}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}) + \tilde{\lambda}^* D^2 \|f\|_{\mathcal{H}_K}^2.$$

Then, we define the functional

$$\tilde{f}_{n, \tilde{\lambda}^*} = \{\tilde{f}_{n, \tilde{\lambda}^*, x_0} : x_0 \in \mathcal{X}\}.$$

Finally, we decompose the estimation error as

$$\hat{f}_{n, \tilde{\lambda}^*} - \bar{f}_{\infty, \tilde{\lambda}^*} = (\hat{f}_{n, \tilde{\lambda}^*} - \tilde{f}_{n, \tilde{\lambda}^*}) + (\tilde{f}_{n, \tilde{\lambda}^*} - \bar{f}_{\infty, \tilde{\lambda}^*}).$$

B.1 ASYMPTOTIC ORDER FOR STOCHASTIC ERROR.

To figure out the order for stochastic error, we finish the proof in 2 steps.

B.1.1 STEP 1: EVALUATE THE ORDER OF $\tilde{f}_{n, \tilde{\lambda}^*} - \bar{f}_{\infty, \tilde{\lambda}^*}$.

Assume that $\omega_{x_0, j} = \mathbb{E}\{\sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i)\}$. To this end, we first obtain the eigenvalues of the operator $G_{\tilde{\lambda}^*}$.

$$\begin{aligned} G_{\tilde{\lambda}^*, x_0} \cdot \eta \cdot g &= \frac{1}{2} D^2 l_{\infty, x_0}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}) \cdot \eta \cdot g + \tilde{\lambda}^* D^2 \|f\|_{\mathcal{H}_K}^2 \cdot \eta \cdot g \\ &= \sum_{j=1}^{\infty} \lambda_j^2 \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} \mathbb{E}\left\{\sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i)\right\} \\ &\quad + \sum_{j=1}^{\infty} \tilde{\lambda}^* \lambda_j \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} \\ &= \sum_{j=1}^{\infty} \omega_{x_0, j} \lambda_j^2 \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} + \sum_{j=1}^{\infty} \tilde{\lambda}^* \lambda_j \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} \\ &= \sum_{j=1}^{\infty} (\omega_{x_0, j} \lambda_j^2 + \tilde{\lambda}^* \lambda_j) \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K} \\ &= \sum_{j=1}^{\infty} \lambda_j^2 (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1}) \langle \eta, \psi_j \rangle_{\mathcal{H}_K} \langle g, \psi_j \rangle_{\mathcal{H}_K}. \end{aligned}$$

Let $\eta = \psi_m$ gives

$$\langle G_{\tilde{\lambda}^*, x_0} g, \psi_m \rangle_{\mathcal{H}_K} = \lambda_m (\omega_{x_0, m} + \tilde{\lambda}^* \lambda_m^{-1}) \langle g, \psi_m \rangle_{\mathcal{H}_K}$$

which leads to

$$\langle G_{\tilde{\lambda}^*, x_0}^{-1} g, \psi_m \rangle_{\mathcal{H}_K} = \lambda_m^{-1} (\omega_{x_0, m} + \tilde{\lambda}^* \lambda_m^{-1})^{-1} \langle g, \psi_m \rangle_{\mathcal{H}_K}.$$

Using the expression of the norm and the expression of the operator $G_{\tilde{\lambda}^*, x_0}$ and assuming $\lambda_j \asymp j^{-k}$ with $k > 1$, and $\omega_{x_0, j} \asymp j^\beta$ uniformly in x_0 for some $\beta < k$ (so that $k - \beta > 0$), we have the following:

$$\begin{aligned} \|\tilde{f}_{n, \tilde{\lambda}^*} - \bar{f}_{\infty, \tilde{\lambda}^*}\|^2 &= \frac{1}{4} \int \left(\sum_{j=1}^{\infty} \lambda_j \langle G_{\tilde{\lambda}^*, x_0}^{-1} D l_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right)^2 \pi(x_0) dx_0 \\ &= \frac{1}{4} \int \left[\left\{ \sum_{j=1}^{\infty} \lambda_j \lambda_j^{-1} (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \langle D l_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right\}^2 \right] \pi(x_0) dx_0 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{4} \int \left[\left\{ \sum_{j=1}^{\infty} (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right\}^2 \right] \pi(x_0) dx_0 \\
&= \frac{1}{4} \int \sum_{j=1}^{\infty} (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1})^{-2} \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K}^2 \psi_j^2(x_0) \pi(x_0) dx_0 \\
&\quad + \frac{1}{4} \int \sum_{j \neq l}^{\infty} (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} (\omega_{x_0, l} + \tilde{\lambda}^* \lambda_l^{-1})^{-1} \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K} \\
&\quad \cdot \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_l \rangle_{\mathcal{H}_K} \psi_j(x_0) \psi_l(x_0) \pi(x_0) dx_0 \\
&= \frac{1}{4} \int \sum_{j=1}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2} \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K}^2 \psi_j^2(x_0) \pi(x_0) dx_0 \\
&\quad + \frac{1}{4} \int \sum_{j \neq l}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)} \frac{1}{(l^\beta + \tilde{\lambda}^* l^k)} \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K} \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_l \rangle_{\mathcal{H}_K} \\
&\quad \times \psi_j(x_0) \psi_l(x_0) \pi(x_0) dx_0 \\
&\asymp \frac{1}{4} \sum_{j=1}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2} \mathbb{E} \{ \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K}^2 \} \int \psi_j^2(x_0) \pi(x_0) dx_0 \\
&\quad + \frac{1}{4} \sum_{j \neq l}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)} \frac{1}{(l^\beta + \tilde{\lambda}^* l^k)} \mathbb{E} \{ \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K} \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_l \rangle_{\mathcal{H}_K} \} \\
&\quad \times \int \psi_j(x_0) \psi_l(x_0) \pi(x_0) dx_0 \\
&= \frac{1}{4} \sum_{j=1}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2} \mathbb{E} \{ \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K}^2 \}
\end{aligned}$$

To find the asymptotic order of $\|\tilde{f}_n - \bar{f}_n\|^2$, we need to derive the coefficients $\sum_{j=1}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2}$ and

$$\mathbb{E} \{ \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K}^2 \}.$$

Derivation for $\mathbb{E} \{ \langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \rangle_{\mathcal{H}_K}^2 \}$.

$$\begin{aligned}
&\mathbb{E} \left\{ \left\langle Dl_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \right\rangle_{\mathcal{H}_K}^2 \right\} \\
&= \mathbb{E} \left\{ \left\langle Dl_n(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}) - Dl_{\infty}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \right\rangle_{\mathcal{H}_K} \right\}^2 \\
&= \mathbb{E} \left[-\frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} \{y_i - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)\} \psi_j(X_i) + 2 \sum_{i=1}^N \omega_i(x_0) \mathbb{E} \{ (f_0(X_i) - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)) \psi_j(X_i) \mid X_i \} \right]^2 \\
&= \mathbb{E} \left[-\frac{2}{n} \sum_{i=1}^N I_i(x_0) \{y_i - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)\} \psi_j(X_i) + \frac{2}{n} \mathbb{E} \left\{ \sum_{i=1}^N I_i(x_0) (f_0(X_i) - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)) \psi_j(X_i) \mid X_i, y_i \right\} \right]^2 \\
&= \frac{4}{n^2} \mathbb{E} \left\{ \sum_{i=1}^N (I_i(x_0) - n \omega_i(x_0)) (y_i - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)) \psi_j(X_i) \right\}^2 \quad (\text{since } \mathbb{E}[I_i(x_0) \mid X_i, y_i] = n \omega_i(x_0)) \\
&= \frac{4}{n^2} \mathbb{E} \left[\mathbb{E} \left\{ \left(\sum_{i=1}^N (I_i(x_0) - n \omega_i(x_0)) (y_i - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)) \psi_j(X_i) \right)^2 \mid X_i, y_i \right\} \right] \\
&= \frac{4}{n^2} \mathbb{E} \left[\text{Var} \left\{ \sum_{i=1}^N I_i(x_0) a_i \mid X_i, y_i \right\} + \left\{ \mathbb{E} \left\{ \sum_{i=1}^N I_i(x_0) a_i \mid X_i, y_i \right\} - n \sum_{i=1}^N \omega_i(x_0) a_i \right\}^2 \right]
\end{aligned}$$

$$= \frac{4}{n^2} \mathbb{E} \left[\text{Var} \left\{ \sum_{i=1}^N I_i(x_0) a_i \mid X_i, y_i \right\} \right]$$

where we denoted $a_i := (y_i - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)) \psi_j(X_i)$. Using the multinomial covariance, $\text{Var}(I_i \mid X_i) = n\omega_i(1 - \omega_i)$ and $\text{Cov}(I_i, I_k \mid x) = -n\omega_i\omega_k$ for $i \neq k$, we get

$$\begin{aligned} \text{Var} \left\{ \sum_{i=1}^N I_i(x_0) a_i \mid X_i, y_i \right\} &= \sum_{i=1}^N a_i^2 \text{Var}(I_i \mid X_i) + 2 \sum_{1 \leq i < k \leq N} a_i a_k \text{Cov}(I_i, I_k \mid x) \\ &= n \sum_{i=1}^N \omega_i a_i^2 - n \sum_{i=1}^N \sum_{k=1}^N \omega_i \omega_k a_i a_k \\ &= n \left\{ \sum_{i=1}^N \omega_i a_i^2 - \left(\sum_{i=1}^N \omega_i a_i \right)^2 \right\} \leq n \sum_{i=1}^N \omega_i a_i^2. \end{aligned}$$

Therefore,

$$\begin{aligned} &\mathbb{E} \left\{ \left\langle D l_{n, \tilde{\lambda}^*}(\bar{f}_{\infty, \tilde{\lambda}^*, x_0}), \psi_j \right\rangle_{\mathcal{H}_K}^2 \right\} \\ &\leq \frac{4}{n} \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) a_i^2 \right\} \\ &= \frac{4}{n} \sum_{i=1}^N \omega_i(x_0) \mathbb{E} \left[\{y_i - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)\}^2 \psi_j(X_i)^2 \right] \\ &= \frac{4}{n} \sum_{i=1}^N \omega_i(x_0) \mathbb{E} \left[\{(f_0(X_i) - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)) + \epsilon_i\}^2 \psi_j(X_i)^2 \right] \\ &= \frac{4}{n} \sum_{i=1}^N \omega_i(x_0) \mathbb{E} \left[\{(f_0(X_i) - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i))^2 + \sigma^2\} \psi_j(X_i)^2 \right] \\ &\leq \frac{4}{n} \left[\sigma^2 \sum_{i=1}^N \omega_i(x_0) \mathbb{E} \{\psi_j(X_i)^2\} + \sum_{i=1}^N \omega_i(x_0) \mathbb{E} \left[\{f_0(X_i) - \bar{f}_{\infty, \tilde{\lambda}^*, x_0}(X_i)\}^2 \psi_j(X_i)^2 \right] \right] \asymp \frac{1}{n}. \end{aligned}$$

Derivation for the coefficients $\sum_{j=1}^{\infty} (j^\beta + \tilde{\lambda}^* j^k)^{-2}$. Assumed $\lambda_j \asymp j^{-k}$ with $k > 1$, and $\omega_{x_0, j} \asymp j^\beta$ uniformly in x_0 for some $\beta < k$ (so that $k - \beta > 0$). Consider

$$S_1(x_0) := \sum_{j=1}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2}.$$

Let J solve $j^\beta = \tilde{\lambda}^* j^k$, i.e. $J \asymp (\tilde{\lambda}^*)^{-1/(k-\beta)}$.

We split

$$S_1(x_0) = \sum_{j \leq J} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2} + \sum_{j > J} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2}.$$

For $j \leq J$, $j^\beta \gg \tilde{\lambda}^* j^k$ and $(j^\beta + \tilde{\lambda}^* j^k)^2 \asymp j^{2\beta}$, giving

$$\sum_{j \leq J} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2} \asymp \sum_{j \leq J} j^{-2\beta} \asymp \begin{cases} J^{1-2\beta}, & \beta < \frac{1}{2}, \\ \log J, & \beta = \frac{1}{2}, \\ 1, & \frac{1}{2} < \beta < k. \end{cases} \quad (4)$$

For $j > J$, $\tilde{\lambda}^* j^k \gg j^\beta$ and $(j^\beta + \tilde{\lambda}^* j^k)^2 \asymp (\tilde{\lambda}^*)^2 j^{2k}$, giving

$$\sum_{j > J} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2} \asymp \frac{1}{(\tilde{\lambda}^*)^2} \sum_{j > J} j^{-2k} \asymp \frac{1}{(\tilde{\lambda}^*)^2} J^{-(2k-1)} \asymp (\tilde{\lambda}^*)^{\frac{2\beta-1}{k-\beta}}. \quad (5)$$

Using $J \asymp (\tilde{\lambda}^*)^{-1/(k-\beta)}$ to rewrite equation 4 in terms of $\tilde{\lambda}^*$,

$$\sum_{j \leq J} j^{-2\beta} \asymp \begin{cases} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}, & \beta < \frac{1}{2}, \\ \log(1/\tilde{\lambda}^*), & \beta = \frac{1}{2}, \\ 1, & \frac{1}{2} < \beta < k. \end{cases}$$

Comparing with equation 5, we obtain

$$S_1(x_0) \asymp \begin{cases} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}, & \beta < \frac{1}{2}, \\ \log(1/\tilde{\lambda}^*), & \beta = \frac{1}{2}, \\ 1, & \frac{1}{2} < \beta < k. \end{cases}$$

Since $\mathbb{E}[\langle Dl_{n,\tilde{\lambda}^*}(\bar{f}_{\infty,\tilde{\lambda}^*,x_0}), \psi_j \rangle_{\mathcal{H}_K}^2] \asymp n^{-1}$, we obtain

$$\|\tilde{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 = \begin{cases} O_p(n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}), & \beta < \frac{1}{2}, \\ O_p(n^{-1} \log(1/\tilde{\lambda}^*)), & \beta = \frac{1}{2}, \\ O_p(n^{-1}), & \frac{1}{2} < \beta < k. \end{cases} \quad (\text{A})$$

B.1.2 STEP 2: EVALUATE THE ORDER OF $\hat{f}_{n,\tilde{\lambda}^*} - \tilde{f}_{n,\tilde{\lambda}^*}$.

For brevity, write

$$\hat{f}_{x_0} := \hat{f}_{n,\tilde{\lambda}^*,x_0}, \quad \bar{f}_{x_0} := \bar{f}_{\infty,\tilde{\lambda}^*,x_0}, \quad \tilde{f}_{x_0} := \tilde{f}_{n,\tilde{\lambda}^*,x_0}.$$

By definition of \hat{f}_{x_0} , we know $Dl_{n,\tilde{\lambda}^*}(\hat{f}_{x_0}) = 0$. We now check the following equation:

$$Dl_{n,\tilde{\lambda}^*}(\hat{f}_{x_0}) = Dl_{n,\tilde{\lambda}^*}(\bar{f}_{x_0}) + D^2l_{n,\tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) = 0. \quad (6)$$

Using the functional derivatives derived above, we have

$$Dl_{n,\tilde{\lambda}^*} = -\frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} K(X_i, \cdot) (y_i - \langle \bar{f}_{x_0}, K(X_i, \cdot) \rangle_{\mathcal{H}_K}) + 2\tilde{\lambda}^* \langle \bar{f}_{x_0}, \cdot \rangle_{\mathcal{H}_K} \quad (7)$$

$$D^2l_{n,\tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) = \frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} \langle \hat{f}_{x_0} - \bar{f}_{x_0}, K(X_i, \cdot) \rangle_{\mathcal{H}_K} K(X_i, \cdot) + 2\tilde{\lambda}^* \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \cdot \rangle_{\mathcal{H}_K} \quad (8)$$

Adding (7) and (8) together, we obtain:

$$\begin{aligned} (7) + (8) &= -\frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} K(X_i, \cdot) (y_i - \langle \hat{f}_{x_0}, K(X_i, \cdot) \rangle_{\mathcal{H}_K}) + 2\tilde{\lambda}^* \langle \hat{f}_{x_0}, \cdot \rangle_{\mathcal{H}_K} \\ &= Dl_{n,\tilde{\lambda}^*}(\hat{f}_{x_0}) = 0. \end{aligned}$$

Using the definition of \tilde{f}_{x_0} , we have

$$Dl_{n,\tilde{\lambda}^*}(\bar{f}_{x_0}) + D^2l_{\infty,\tilde{\lambda}^*}(\bar{f}_{x_0})(\tilde{f}_{x_0} - \bar{f}_{x_0}) = 0.$$

Combining with the equation $Dl_{n,\tilde{\lambda}^*}(\hat{f}_{x_0}) = Dl_{n,\tilde{\lambda}^*}(\bar{f}_{x_0}) + D^2l_{n,\tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) = 0$, we get

$$D^2l_{\infty,\tilde{\lambda}^*}(\bar{f}_{x_0})(\tilde{f}_{x_0} - \bar{f}_{x_0}) = D^2l_{n,\tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}).$$

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

Then, we can derive:

$$\begin{aligned} D^2 l_{\infty, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \tilde{f}_{x_0}) &= D^2 l_{\infty, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) + D^2 l_{\infty, \tilde{\lambda}^*}(\bar{f}_{x_0})(\bar{f}_{x_0} - \tilde{f}_{x_0}) \\ &= D^2 l_{\infty, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) - D^2 l_{n, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) \\ &= D^2 l_{\infty}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) - D^2 l_n(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}), \end{aligned}$$

where in the last line we used that the penalty part is the same in $D^2 l_{\infty, \tilde{\lambda}^*}$ and $D^2 l_{n, \tilde{\lambda}^*}$.

Then $\hat{f}_{x_0} - \tilde{f}_{x_0}$ can be expressed as

$$\hat{f}_{x_0} - \tilde{f}_{x_0} = \frac{1}{2} G_{\tilde{\lambda}^*}^{-1} \left\{ D^2 l_{\infty}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) - D^2 l_n(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) \right\}.$$

So we write:

$$\begin{aligned} &\|\hat{f}_{n, \tilde{\lambda}^*} - \tilde{f}_{n, \tilde{\lambda}^*}\|^2 \\ &:= \int \langle \hat{f}_{n, \tilde{\lambda}^*, x_0} - \tilde{f}_{n, \tilde{\lambda}^*, x_0}, K(x_0, \cdot) \rangle_{\mathcal{H}_K}^2 \pi(x_0) dx_0 \\ &= \int \left(\sum_{j=1}^{\infty} \lambda_j \langle \hat{f}_{n, \tilde{\lambda}^*, x_0} - \tilde{f}_{n, \tilde{\lambda}^*, x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right)^2 \pi(x_0) dx_0 \\ &= \frac{1}{4} \int \left\{ \sum_{j=1}^{\infty} \lambda_j \left\langle G_{\tilde{\lambda}^*}^{-1} \left(D^2 l_{\infty, \tilde{\lambda}^*}(\bar{f}_{x_0}) - D^2 l_{n, \tilde{\lambda}^*}(\bar{f}_{x_0}) \right) [\hat{f}_{x_0} - \bar{f}_{x_0}], \psi_j \right\rangle_{\mathcal{H}_K} \psi_j(x_0) \right\}^2 \pi(x_0) dx_0 \\ &= \frac{1}{4} \int \left[\sum_{j=1}^{\infty} (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \left\{ \langle D^2 l_{\infty, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}), \psi_j \rangle_{\mathcal{H}_K} \right\} \psi_j(x_0) \right]^2 \pi(x_0) dx_0 \\ &\quad - \langle D^2 l_{n, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}), \psi_j \rangle_{\mathcal{H}_K} \\ &=: \frac{1}{4} \int \left\{ \sum_{j=1}^{\infty} (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} A_j(x_0) \psi_j(x_0) \right\}^2 \pi(x_0) dx_0, \tag{2.1} \end{aligned}$$

where

$$\begin{aligned} A_j(x_0) &= \langle D^2 l_{\infty, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}), \psi_j \rangle_{\mathcal{H}_K} - \langle D^2 l_{n, \tilde{\lambda}^*}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}), \psi_j \rangle_{\mathcal{H}_K} \\ &= 2 \sum_{\ell=1}^{\infty} \lambda_{\ell}^2 \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_{\ell} \rangle_{\mathcal{H}_K} \langle \psi_j, \psi_{\ell} \rangle_{\mathcal{H}_K} \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_{\ell}^2(X_i) \right\} \\ &\quad - \frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} \sum_{k=1}^{\infty} \lambda_k^2 \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_k \rangle_{\mathcal{H}_K} \langle \psi_j, \psi_k \rangle_{\mathcal{H}_K} \psi_k^2(X_i) \\ &= 2 \lambda_j^2 \frac{1}{\lambda_j} \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i) \right\} \\ &\quad - \frac{2}{n} \sum_{i \in \mathcal{S}(x_0)} \lambda_j^2 \frac{1}{\lambda_j} \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j^2(X_i) \\ &= 2 \lambda_j \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \left[\mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i) \right\} - \frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} \psi_j^2(X_i) \right]. \end{aligned}$$

plug it back into 2.1

$$= \frac{1}{4} \int \left[\sum_{j=1}^{\infty} (\omega_{x_0, j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \left\{ \left(2 \lambda_j \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \right) \left(\omega_{x_0, j} - \frac{1}{n} \sum_{i \in \mathcal{S}(x_0)} \psi_j^2(X_i) \right) \psi_j(x_0) \right\} \right]^2 \pi(x_0) dx_0$$

$$\begin{aligned}
&= \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} \lambda_j (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \left\{ n \omega_{x_0,j} - \sum_{i \in \mathcal{S}(x_0)} \psi_j^2(X_i) \right\} \right]^2 \pi(x_0) dx_0 \\
&= \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} \lambda_j (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right. \\
&\quad \left. \times \left\{ n \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i) \right\} - \sum_{i=1}^N I_i(x_0) \psi_j^2(X_i) \right\} \right]^2 \pi(x_0) dx_0. \tag{2.2}
\end{aligned}$$

Set

$$\Delta_j(x_0) := n \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i) \right\} - \sum_{i=1}^N I_i(x_0) \psi_j^2(X_i).$$

Known that $\mathbb{E}\{I_i(x_0) \mid X_1, \dots, X_N\} = n \omega_i(x_0)$ and $\sum_{i=1}^N \omega_i(x_0) = 1$, $\mathbb{E}\{\Delta_j(x_0)\} = 0$. For the second moment, condition on $\{X_i\}_{i=1}^N$, then

$$\begin{aligned}
\mathbb{E}\{\Delta_j(x_0)^2\} &= \mathbb{E} \left[\text{Var} \left\{ \sum_{i=1}^N I_i(x_0) \psi_j^2(X_i) \mid X_1, \dots, X_N \right\} \right] \\
&= \mathbb{E} \left[n \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^4(X_i) - \left(\sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i) \right)^2 \right\} \right] \\
&\leq \mathbb{E} \left\{ n \sum_{i=1}^N \omega_i(x_0) \psi_j^4(X_i) \right\} \\
&= n \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^4(X_i) \right\} = n \mathbb{E} \{ \psi_j^4(X) \}.
\end{aligned}$$

Therefore,

$$\Delta_j(x_0) \asymp O_p \left(\sqrt{n \mathbb{E} \{ \psi_j^4(X) \}} \right).$$

Plug it back to (2.2):

$$\begin{aligned}
&\asymp \frac{1}{n^2} \int \left\{ \sum_{j=1}^{\infty} \lambda_j (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right\}^2 \pi(x_0) dx_0 \\
&= \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-1} \left\{ \lambda_j \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right\} \right]^2 \pi(x_0) dx_0 \\
&\leq \frac{1}{n^2} \int \left\{ \sum_{j=1}^{\infty} (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-2} \right\} \left\{ \sum_{j=1}^{\infty} \lambda_j^2 \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K}^2 \psi_j(x_0)^2 \right\} \pi(x_0) dx_0 \\
&\leq \frac{1}{n^2} \int \underbrace{\left\{ \sum_{j=1}^{\infty} (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-2} \right\}}_{S_1(x_0)} \left\{ \sum_{j=1}^{\infty} \lambda_j \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right\}^2 \pi(x_0) dx_0.
\end{aligned}$$

Assume $\lambda_j \asymp j^{-k}$ with $k > 1$ and, uniformly in x_0 , $\omega_{x_0,j} \asymp j^\beta$ for some $\beta < k$. Recall

$$S_1(x_0) := \sum_{j=1}^{\infty} (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-2} \asymp \sum_{j=1}^{\infty} \frac{1}{(j^\beta + \tilde{\lambda}^* j^k)^2} \asymp \begin{cases} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}, & \beta < \frac{1}{2}, \\ \log(1/\tilde{\lambda}^*), & \beta = \frac{1}{2}, \\ 1, & \frac{1}{2} < \beta < k. \end{cases}$$

Consequently,

$$\begin{aligned} \|\hat{f}_{n,\tilde{\lambda}^*} - \tilde{f}_{n,\tilde{\lambda}^*}\|^2 &\leq \frac{1}{n^2} \int \underbrace{\sum_{j=1}^{\infty} (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1})^{-2}}_{= S_1(x_0)} \left\{ \sum_{j=1}^{\infty} \lambda_j \langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right\}^2 \pi(x_0) dx_0 \\ &= \frac{S_1(x_0)}{n^2} \sum_{j=1}^{\infty} \lambda_j^2 \langle \hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}, \psi_j \rangle_{\mathcal{H}_K}^2 \\ &= \frac{S_1(x_0)}{n^2} \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2, \end{aligned}$$

with $S_1(x_0)$ as above. Therefore,

$$\|\hat{f}_{n,\tilde{\lambda}^*} - \tilde{f}_{n,\tilde{\lambda}^*}\|^2 = \begin{cases} O_p\left(n^{-2} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}} \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2\right), & \beta < \frac{1}{2}, \\ O_p\left(n^{-2} \log(1/\tilde{\lambda}^*) \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2\right), & \beta = \frac{1}{2}, \\ O_p\left(n^{-2} \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2\right), & \frac{1}{2} < \beta < k. \end{cases} \quad (\text{B})$$

Recall

$$\|\tilde{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 = \begin{cases} O_p\left(n^{-1} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}\right), & \beta < \frac{1}{2}, \\ O_p\left(n^{-1} \log(1/\tilde{\lambda}^*)\right), & \beta = \frac{1}{2}, \\ O_p(n^{-1}), & \frac{1}{2} < \beta < k. \end{cases} \quad (\text{A})$$

From (B), write

$$\|\hat{f}_{n,\tilde{\lambda}^*} - \tilde{f}_{n,\tilde{\lambda}^*}\|^2 = \eta_n(\tilde{\lambda}^*, \beta) \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2,$$

where

$$\eta_n(\tilde{\lambda}^*, \beta) = \begin{cases} O_p\left(n^{-2} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}\right), & \beta < \frac{1}{2}, \\ O_p\left(n^{-2} \log(1/\tilde{\lambda}^*)\right), & \beta = \frac{1}{2}, \\ O_p(n^{-2}), & \frac{1}{2} < \beta < k. \end{cases}$$

By the triangle inequality,

$$\|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 \leq \|\tilde{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 + \|\hat{f}_{n,\tilde{\lambda}^*} - \tilde{f}_{n,\tilde{\lambda}^*}\|^2. \quad (\text{C})$$

If $\eta_n(\tilde{\lambda}^*, \beta) \xrightarrow{P} 0$, then

$$\begin{cases} n^{-2} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}} \rightarrow 0, & \beta < \frac{1}{2}, \\ n^{-2} \log(1/\tilde{\lambda}^*) \rightarrow 0, & \beta = \frac{1}{2}, \\ n^{-2} \rightarrow 0, & \frac{1}{2} < \beta < k, \end{cases} \quad (\text{D})$$

1512 then from (C) we get

$$1513 (1-\eta_n) \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 \leq \|\tilde{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 \quad \text{and} \quad \|\tilde{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 \geq (1-\eta_n) \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2,$$

1514 hence, under (D),

$$1515 \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 = \{1 + o_p(1)\} \|\tilde{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2.$$

1516 Combining with (A) yields

$$1517 \|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 = \begin{cases} O_p\left(n^{-1} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}\right), & \beta < \frac{1}{2}, \\ O_p(n^{-1} \log(1/\tilde{\lambda}^*)), & \beta = \frac{1}{2}, \\ O_p(n^{-1}), & \frac{1}{2} < \beta < k, \end{cases} \quad (\text{E})$$

1518 when (D) holds.

1519 **Condition (D) and corresponding constraints.**

1520 For $\eta_n(\tilde{\lambda}^*, \beta) \rightarrow 0$ in the case $\omega_{x_0,j} \asymp j^\beta$, the corresponding rate constraints on $\tilde{\lambda}^*$ are:

- 1521 • Case $\beta < \frac{1}{2}$: We require

$$1522 n^{-2} (\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}} \rightarrow 0,$$

1523 equivalently,

$$1524 (\tilde{\lambda}^*)^{\frac{1-2\beta}{k-\beta}} \gg n^{-2} \iff \tilde{\lambda}^* \gg n^{-\frac{2(k-\beta)}{1-2\beta}}.$$

- 1525 • Case $\beta = \frac{1}{2}$: We require

$$1526 n^{-2} \log(1/\tilde{\lambda}^*) \rightarrow 0,$$

1527 equivalently,

$$1528 \log(1/\tilde{\lambda}^*) = o(n^2).$$

- 1529 • Case $\frac{1}{2} < \beta < k$: Here $n^{-2} \rightarrow 0$ automatically; no further condition on $\tilde{\lambda}^*$ is needed.

1530 B.2 ASYMPTOTIC ORDER FOR DETERMINISTIC ERROR.

1531 Recall

$$1532 l_{\infty,x_0}(f) := \mathbb{E}\{l_{n,x_0}(f)\} = \sigma^2 + \mathbb{E}\left\{\sum_{i=1}^N \omega_i(x_0) \langle f - f_0, K(X_i, \cdot) \rangle_{\mathcal{H}_K}^2\right\}$$

$$1533 = \sigma^2 + \sum_{i=1}^N \mathbb{E}\left[\omega_i(x_0) \{f(X_i) - f_0(X_i)\}^2\right]$$

$$1534 = \sigma^2 + \sum_{i=1}^N \mathbb{E}\left[\omega_i(x_0) \left\{\sum_{j=1}^{\infty} (c_j - a_j) \psi_j(X_i)\right\}^2\right]$$

$$1535 = \sigma^2 + \sum_{j,\ell \geq 1} (c_j - a_j)(c_\ell - a_\ell) \sum_{i=1}^N \mathbb{E}\{\omega_i(x_0) \psi_j(X_i) \psi_\ell(X_i)\}$$

$$1536 = \sigma^2 + \sum_{j=1}^{\infty} \omega_{x_0,j} (c_j - a_j)^2,$$

1537 where $\omega_{x_0,j} := \mathbb{E}\{\sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i)\}$.

1538 The objective function about c_j at x_0 becomes

$$1539 Q_{x_0}(c) = \sum_{j \geq 1} \left\{ \omega_{x_0,j} (c_j - a_j)^2 + \tilde{\lambda}^* \frac{c_j^2}{\lambda_j} \right\}.$$

1566

Taking derivatives,

1567

1568

$$2\omega_{x_0,j}(c_j - a_j) + 2\tilde{\lambda}^* \frac{c_j}{\lambda_j} = 0 \implies (\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1}) c_j = \omega_{x_0,j} a_j,$$

1569

1570

hence

1571

1572

$$\bar{c}_j(x_0) = \frac{\omega_{x_0,j}}{\omega_{x_0,j} + \tilde{\lambda}^* \lambda_j^{-1}} a_j \quad \text{and} \quad \bar{c}_j(x_0) - a_j = -\frac{\tilde{\lambda}^*}{\omega_{x_0,j} \lambda_j + \tilde{\lambda}^*} a_j.$$

1573

1574

1575

Assume $\lambda_j \asymp j^{-k}$ with $k > 1$, $a_j^2 \lambda_j^{-1} = j^{-a}$ with $a > 1$ (so $a_j^2 \asymp j^{-(a+k)}$), and fix $\beta < k$ such that, uniformly in x_0 ,

1576

$$\omega_{x_0,j} \asymp j^\beta \quad (\text{equivalently, } \lambda_j \omega_{x_0,j} \asymp j^{\beta-k}).$$

1577

1578

Now we evaluate $\|\bar{f}_{n,\tilde{\lambda}^*} - f_0\|^2$:

1579

1580

1581

1582

1583

1584

1585

1586

1587

1588

1589

1590

1591

1592

1593

1594

1595

1596

$$\begin{aligned} \|\bar{f}_{n,\tilde{\lambda}^*} - f_0\|^2 &:= \int \langle \bar{f}_{n,\tilde{\lambda}^*,x_0} - f_0, K(x_0, \cdot) \rangle_{\mathcal{H}_K}^2 \pi(x_0) dx_0 \\ &= \int \left\{ \sum_{j=1}^{\infty} \lambda_j \langle \bar{f}_{n,\tilde{\lambda}^*,x_0} - f_0, \psi_j \rangle_{\mathcal{H}_K} \psi_j(x_0) \right\}^2 \pi(x_0) dx_0 \\ &= \int \left[\sum_{j=1}^{\infty} \{\bar{c}_j(x_0) - a_j\} \psi_j(x_0) \right]^2 \pi(x_0) dx_0 \\ &= \int \left\{ \sum_{j=1}^{\infty} \frac{\tilde{\lambda}^*}{\omega_{x_0,j} \lambda_j + \tilde{\lambda}^*} a_j \psi_j(x_0) \right\}^2 \pi(x_0) dx_0 \\ &\asymp \int \left\{ \sum_{j=1}^{\infty} \frac{\tilde{\lambda}^*}{j^{\beta-k} + \tilde{\lambda}^*} a_j \psi_j(x_0) \right\}^2 \pi(x_0) dx_0 \\ &= \sum_{j=1}^{\infty} \left(\frac{\tilde{\lambda}^*}{j^{\beta-k} + \tilde{\lambda}^*} a_j \right)^2 \asymp \sum_{j=1}^{\infty} \frac{(\tilde{\lambda}^*)^2 j^{-(a+k)}}{(j^{\beta-k} + \tilde{\lambda}^*)^2}. \end{aligned}$$

1597

1598

Let $J \asymp (\tilde{\lambda}^*)^{-1/(k-\beta)}$ so that $J^{\beta-k} \asymp \tilde{\lambda}^*$. Split at J :

1599

1600

1601

1602

$$\sum_{j \leq J} \frac{(\tilde{\lambda}^*)^2 j^{-(a+k)}}{(j^{\beta-k})^2} \asymp (\tilde{\lambda}^*)^2 \sum_{j \leq J} j^{-(a+2\beta-k)}, \quad \sum_{j > J} \frac{(\tilde{\lambda}^*)^2 j^{-(a+k)}}{(\tilde{\lambda}^*)^2} = \sum_{j > J} j^{-(a+k)}.$$

1603

1604

1605

1606

Using the integral test directly on the exponent $a + 2\beta - k$ (with $a + k > 1$), we have

1607

1608

1609

1610

1611

1612

1613

With $J \asymp (\tilde{\lambda}^*)^{-1/(k-\beta)}$, this gives

$$\sum_{j \leq J} j^{-(a+2\beta-k)} \asymp \begin{cases} (\tilde{\lambda}^*)^{\frac{a+k-1}{k-\beta}}, & a + 2\beta - k < 1, \\ (\tilde{\lambda}^*)^2 \log(1/\tilde{\lambda}^*), & a + 2\beta - k = 1, \\ (\tilde{\lambda}^*)^2, & a + 2\beta - k > 1, \end{cases} \quad \sum_{j > J} j^{-(a+k)} \asymp (\tilde{\lambda}^*)^{\frac{a+k-1}{k-\beta}}.$$

1614

1615

1616

1617

1618

1619

Comparing the two parts yields, uniformly in x_0 and for any $\beta < k$,

$$\|\bar{f}_{n,\tilde{\lambda}^*} - f_0\|^2 \asymp \begin{cases} (\tilde{\lambda}^*)^{\frac{a+k-1}{k-\beta}}, & a + 2\beta - k < 1, \\ (\tilde{\lambda}^*)^2 \log(1/\tilde{\lambda}^*), & a + 2\beta - k = 1, \\ (\tilde{\lambda}^*)^2, & a + 2\beta - k > 1. \end{cases} \quad (\text{F})$$

B.3 ASYMPTOTIC ORDER FOR MSE AND TUNING.

Assume $a > 1$, $k > 1$, and $\beta < k$.

Deterministic error (F) (split by $a + 2\beta - k$).

$$\|\bar{f}_{n,\tilde{\lambda}^*} - f_0\|^2 \asymp \begin{cases} (\tilde{\lambda}^*)^{\frac{a+k-1}{k-\beta}}, & a + 2\beta - k < 1, \\ (\tilde{\lambda}^*)^2 \log(1/\tilde{\lambda}^*), & a + 2\beta - k = 1, \\ (\tilde{\lambda}^*)^2, & a + 2\beta - k > 1. \end{cases}$$

Stochastic error (E) (split by $\beta = \frac{1}{2}$). When condition (D) holds and $\beta < k$,

$$\|\hat{f}_{n,\tilde{\lambda}^*} - \bar{f}_{n,\tilde{\lambda}^*}\|^2 = \begin{cases} O_p\left(n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}\right), & \beta < \frac{1}{2}, \\ O_p(n^{-1} \log(1/\tilde{\lambda}^*)), & \beta = \frac{1}{2}, \\ O_p(n^{-1}), & \frac{1}{2} < \beta < k. \end{cases}$$

Condition (D). We require $\eta_n(\tilde{\lambda}^*, \beta) \rightarrow 0$:

$$\begin{cases} n^{-2}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}} \rightarrow 0, & \beta < \frac{1}{2}, \\ n^{-2} \log(1/\tilde{\lambda}^*) \rightarrow 0, & \beta = \frac{1}{2}, \\ n^{-2} \rightarrow 0, & \frac{1}{2} < \beta < k. \end{cases}$$

For $\beta < \frac{1}{2}$, this is equivalent to $(\tilde{\lambda}^*) \gg n^{-2(k-\beta)/(1-2\beta)}$; for $\beta = \frac{1}{2}$, any polynomial decay of $\tilde{\lambda}^*$ ensures $n^{-2} \log(1/\tilde{\lambda}^*) \rightarrow 0$; for $\beta > \frac{1}{2}$, (D) reduces to $n^{-2} \rightarrow 0$ and imposes no additional restriction.

We minimize $MSE(\tilde{\lambda}^*) := \text{Det}(\tilde{\lambda}^*) + \text{Stoch}(\tilde{\lambda}^*)$ by differentiation.

Regimes (by deterministic split $a + 2\beta - k$ and stochastic split $\beta = \frac{1}{2}$, with $\beta < k$). Let

$$p := \frac{a+k-1}{k-\beta}.$$

(i) $a + 2\beta - k < 1$ (Det $\asymp (\tilde{\lambda}^*)^p$)

- $\beta < \frac{1}{2}$ (Stoch $\asymp n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}$).

$$\text{MSE}(\tilde{\lambda}^*) = (\tilde{\lambda}^*)^p + n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}, \quad p(\tilde{\lambda}^*)^{p+\frac{1-2\beta}{k-\beta}} = n^{-1}.$$

Hence

$$\tilde{\lambda}^* \asymp n^{-\frac{k-\beta}{a+k-2\beta}}, \quad \text{MSE} \asymp n^{-\frac{a+k-1}{a+k-2\beta}}.$$

- $\beta = \frac{1}{2}$ (Stoch $\asymp n^{-1} \log(1/\tilde{\lambda}^*)$).

$$\text{MSE}(\tilde{\lambda}^*) = (\tilde{\lambda}^*)^p + n^{-1} \log(1/\tilde{\lambda}^*), \quad p(\tilde{\lambda}^*)^p \asymp n^{-1}.$$

Thus

$$\tilde{\lambda}^* \asymp n^{-\frac{k-\beta}{a+k-1}}, \quad \text{MSE} \asymp n^{-1} \log n.$$

- $\frac{1}{2} < \beta < k$ (Stoch $\asymp n^{-1}$). Choose $\tilde{\lambda}^* \rightarrow 0$ with $(\tilde{\lambda}^*)^p = o(n^{-1})$; e.g.

$$\tilde{\lambda}^* \ll n^{-\frac{k-\beta}{a+k-1}} \Rightarrow \text{MSE} \asymp n^{-1}.$$

(ii) $a + 2\beta - k = 1$ (Det $\asymp (\tilde{\lambda}^*)^2 \log(1/\tilde{\lambda}^*)$)

1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727

- $\beta < \frac{1}{2}$ (Stoch $\asymp n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}$).

$$\text{MSE}(\tilde{\lambda}^*) = (\tilde{\lambda}^*)^2 \log(1/\tilde{\lambda}^*) + n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}.$$

It yields

$$(\tilde{\lambda}^*)^{2+\frac{1-2\beta}{k-\beta}} (2 \log(1/\tilde{\lambda}^*) - 1) \asymp n^{-1}.$$

Hence

$$\tilde{\lambda}^* \asymp n^{-\frac{k-\beta}{2k+1-4\beta}} (\log n)^{\frac{k-\beta}{2k+1-4\beta}}, \quad \text{MSE} \asymp n^{-\frac{2(k-\beta)}{2k+1-4\beta}} (\log n)^{\frac{2(k-\beta)}{2k+1-4\beta}}.$$

- $\beta = \frac{1}{2}$ (Stoch $\asymp n^{-1} \log(1/\tilde{\lambda}^*)$). Solving $2\tilde{\lambda}^* - \frac{1}{n\tilde{\lambda}^*} = 0$ gives

$$\tilde{\lambda}^* \asymp n^{-1/2}, \quad \text{MSE} \asymp n^{-1} \log n.$$

- $\frac{1}{2} < \beta < k$ (Stoch $\asymp n^{-1}$). Choose $\tilde{\lambda}^* \rightarrow 0$ with $(\tilde{\lambda}^*)^2 \log(1/\tilde{\lambda}^*) = o(n^{-1})$; e.g.

$$\tilde{\lambda}^* \ll n^{-1/2} \Rightarrow \text{MSE} \asymp n^{-1}.$$

(iii) $a + 2\beta - k > 1$ (Det $\asymp (\tilde{\lambda}^*)^2$)

- $\beta < \frac{1}{2}$ (Stoch $\asymp n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}$). Balancing $(\tilde{\lambda}^*)^2$ and $n^{-1}(\tilde{\lambda}^*)^{-\frac{1-2\beta}{k-\beta}}$ gives

$$\tilde{\lambda}^* \asymp n^{-\frac{k-\beta}{2k+1-4\beta}}, \quad \text{MSE} \asymp n^{-\frac{2(k-\beta)}{2k+1-4\beta}}.$$

- $\beta = \frac{1}{2}$ (Stoch $\asymp n^{-1} \log(1/\tilde{\lambda}^*)$).

$$\tilde{\lambda}^* \asymp n^{-1/2}, \quad \text{MSE} \asymp n^{-1} \log n.$$

- $\frac{1}{2} < \beta < k$ (Stoch $\asymp n^{-1}$). Choose $\tilde{\lambda}^* \rightarrow 0$ with $(\tilde{\lambda}^*)^2 = o(n^{-1})$; e.g.

$$\tilde{\lambda}^* \ll n^{-1/2} \Rightarrow \text{MSE} \asymp n^{-1}.$$

Result. The optimal tuning $\tilde{\lambda}^*$ and the resulting MSE are:

- (i) $k - 2\beta > a - 1$ (equivalently $a + 2\beta - k < 1$):

$$\tilde{\lambda}^* \asymp \begin{cases} n^{-\frac{k-\beta}{a+k-2\beta}}, & \beta < \frac{1}{2}, \\ n^{-\frac{k-\beta}{a+k-1}}, & \beta = \frac{1}{2}, \\ \ll n^{-\frac{k-\beta}{a+k-1}}, & \frac{1}{2} < \beta < k, \end{cases}$$

$$\text{MSE} \asymp \begin{cases} n^{-\frac{a-1+k}{(a-1)+(k+1)-2\beta}}, & \beta < \frac{1}{2}, \\ n^{-1} \log n, & \beta = \frac{1}{2}, \\ n^{-1}, & \frac{1}{2} < \beta < k. \end{cases}$$

- (ii) $k - 2\beta = a - 1$ (equivalently $a + 2\beta - k = 1$):

$$\tilde{\lambda}^* \asymp \begin{cases} n^{-\frac{k-\beta}{2k+1-4\beta}} (\log n)^{\frac{k-\beta}{2k+1-4\beta}}, & \beta < \frac{1}{2}, \\ n^{-1/2}, & \beta = \frac{1}{2}, \\ \ll n^{-1/2}, & \frac{1}{2} < \beta < k, \end{cases}$$

$$\text{MSE} \asymp \begin{cases} n^{-\frac{2(k-\beta)}{2k+1-4\beta}} (\log n)^{\frac{2(k-\beta)}{2k+1-4\beta}}, & \beta < \frac{1}{2}, \\ n^{-1} \log n, & \beta = \frac{1}{2}, \\ n^{-1}, & \frac{1}{2} < \beta < k. \end{cases}$$

(iii) $k - 2\beta < a - 1$ (equivalently $a + 2\beta - k > 1$):

$$\tilde{\lambda}^* \asymp \begin{cases} n^{-\frac{k-\beta}{2k+1-4\beta}}, & \beta < \frac{1}{2}, \\ n^{-1/2}, & \beta = \frac{1}{2}, \\ \ll n^{-1/2}, & \frac{1}{2} < \beta < k, \end{cases} \quad \text{MSE} \asymp \begin{cases} n^{-\frac{2(k-\beta)}{2k+1-4\beta}}, & \beta < \frac{1}{2}, \\ n^{-1} \log n, & \beta = \frac{1}{2}, \\ n^{-1}, & \frac{1}{2} < \beta < k. \end{cases}$$

To consider all the function in the RKHS, we can set $a \rightarrow 1$ to find the rate of all the functions in the RHKS. Assume $k > 1$ and $\beta < k$. Then the asymptotic orders of the optimal tuning $\tilde{\lambda}^*$ and the resulting mean squared error are:

$$\begin{cases} \beta < \frac{1}{2} : \tilde{\lambda}^* \asymp n^{-\frac{k-\beta}{k+1-2\beta}}, & \text{MSE} \asymp n^{-\frac{k}{k+1-2\beta}}, \\ \beta = \frac{1}{2} : \tilde{\lambda}^* \asymp n^{-\frac{k-\frac{1}{2}}{k}}, & \text{MSE} \asymp n^{-1} \log n, \\ \beta > \frac{1}{2} : \tilde{\lambda}^* \lesssim n^{-1/2}, & \text{MSE} \asymp n^{-1}. \end{cases}$$

C PROOF OF LEMMA 1

Proof. Recall the definition of \mathbf{K}_c , an $L \times L$ kernel matrix and the kernel vector evaluated at the cluster centers are

$$\mathbf{K}_c = \begin{pmatrix} K(C_1, C_1) & \cdots & K(C_1, C_L) \\ \vdots & \ddots & \vdots \\ K(C_L, C_1) & \cdots & K(C_L, C_L) \end{pmatrix}, \quad \mathbf{K}(x_0, C) = \begin{bmatrix} K(x_0, C_1) \\ \vdots \\ K(x_0, C_L) \end{bmatrix}.$$

Let ψ_1, \dots, ψ_L denote the first L eigenfunctions of the kernel function $K(\cdot, \cdot)$ evaluated on the centers and define

$$\boldsymbol{\eta}_\ell := \begin{bmatrix} \psi_\ell(C_1) \\ \vdots \\ \psi_\ell(C_L) \end{bmatrix}, \quad \ell = 1, \dots, L.$$

There exists an orthogonal matrix \mathbf{Q} $\mathbf{Q} = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_L)$, where $\mathbf{Q}^T \mathbf{Q} = \mathbf{Q} \mathbf{Q}^T = \mathbf{I}_L$, such that

$$\mathbf{Q}^T \mathbf{K}_c \mathbf{Q} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_L \end{pmatrix} =: \Lambda.$$

The vector \mathbf{K}_{x_0} satisfies

$$\mathbf{K}_{x_0} = (\mathbf{I} + L^{-1} \tilde{\lambda}^{*-1} \mathbf{K}_c)^{-1} \mathbf{K}(x_0, C) = \begin{bmatrix} K_{x_0 1} \\ \vdots \\ K_{x_0 L} \end{bmatrix}.$$

Using the eigen decomposition of \mathbf{K}_c ,

$$(\mathbf{I} + L^{-1} \tilde{\lambda}^{*-1} \mathbf{K}_c)^{-1} = \mathbf{Q} (\mathbf{I} + L^{-1} \tilde{\lambda}^{*-1} \Lambda)^{-1} \mathbf{Q}^T,$$

1782 the transformed kernel vector at x_0 is

$$1783 \mathbf{K}_{x_0} = \mathbf{Q}(\mathbf{I} + L^{-1}\tilde{\lambda}^{*-1}\Lambda)^{-1}\mathbf{Q}^\top \mathbf{K}(x_0, C).$$

1786 Since

$$1787 \mathbf{Q}^\top \mathbf{K}(x_0, C) = \begin{bmatrix} \boldsymbol{\eta}_1^\top \mathbf{K}(x_0, C) \\ \vdots \\ \boldsymbol{\eta}_L^\top \mathbf{K}(x_0, C) \end{bmatrix},$$

1789 the l -th coordinate of $\mathbf{K}_{x_0} = (K_{x_01}, \dots, K_{x_0L})^\top$ is

$$1791 K_{x_0l} = \sum_{r=1}^L \frac{\boldsymbol{\eta}_r^\top \mathbf{K}(x_0, C)}{1 + L^{-1}\tilde{\lambda}^{*-1}\lambda_r} \psi_r(C_l), \quad l = 1, \dots, L. \quad (9)$$

1796 The cluster-level resampling weight is defined as

$$1797 \omega_{x_0, l, C} = \frac{|K_{x_0l}|}{\sum_{m=1}^L |K_{x_0m}|},$$

1800 and each sample $x_i \in C_l$ receives

$$1802 \omega_i(x_0) = N_l^{-1} \omega_{x_0, l, C}.$$

1804 Then we can write

$$1805 \omega_{x_0, l, C} = \frac{\sum_{r=1}^L \frac{L^{-1}\boldsymbol{\eta}_r^\top \mathbf{K}(x_0, C)}{1 + L^{-1}\tilde{\lambda}^{*-1}\lambda_r} \psi_r(C_l)}{\sum_{m=1}^L \sum_{r=1}^L \frac{L^{-1}\boldsymbol{\eta}_r^\top \mathbf{K}(x_0, C)}{1 + L^{-1}\tilde{\lambda}^{*-1}\lambda_r} \psi_r(C_m)}, \quad l = 1, \dots, L.$$

1811 Note that

$$1812 L^{-1}\boldsymbol{\eta}_r^\top \mathbf{K}(x_0, C) = L^{-1} \sum_{q=1}^L \psi_r(C_q) K(x_0, C_q) \approx \int \psi_r(t) K(x_0, t) dt = \lambda_r \psi_r(x_0),$$

1815 and substituting this approximation into the numerator of $\omega_{x_0, l, C}$ gives

$$1817 \sum_{r=1}^L \frac{L^{-1}\boldsymbol{\eta}_r^\top \mathbf{K}(x_0, C)}{1 + L^{-1}\tilde{\lambda}^{*-1}\lambda_r} \psi_r(C_l) \approx \frac{1}{L} \sum_{r=1}^L \frac{\lambda_r}{1 + L^{-1}\tilde{\lambda}^{*-1}\lambda_r} \psi_r(x_0) \psi_r(C_l).$$

1821 Similarly, the denominator becomes

$$1822 \sum_{m=1}^L \sum_{r=1}^L \frac{L^{-1}\boldsymbol{\eta}_r^\top \mathbf{K}(x_0, C)}{1 + L^{-1}\tilde{\lambda}^{*-1}\lambda_r} \psi_r(C_m) \approx \frac{1}{L} \sum_{m=1}^L \sum_{r=1}^L \frac{\lambda_r}{1 + L^{-1}\tilde{\lambda}^{*-1}\lambda_r} \psi_r(x_0) \psi_r(C_m).$$

1825 So

$$1826 \omega_{x_0, l, C} = \frac{\sum_{r=1}^L \frac{\lambda_r}{L\tilde{\lambda}^* + \lambda_r} \psi_r(x_0) \psi_r(C_l)}{\sum_{m=1}^L \sum_{r=1}^L \frac{\lambda_r}{L\tilde{\lambda}^* + \lambda_r} \psi_r(x_0) \psi_r(C_m)}, \quad l = 1, \dots, L.$$

1832 Define a new kernel \tilde{K} by the following

$$1833 \tilde{K}(x, y) = \sum_{r=1}^{\infty} \frac{\lambda_r}{L\tilde{\lambda}^* + \lambda_r} \psi_r(x) \psi_r(y),$$

1835

1836 so that

$$1837 \omega_{x_0, l, C} \asymp \frac{1838 \tilde{K}(x_0, C_l)}{1839 \sum_{m=1}^L \tilde{K}(x_0, C_m)}.$$

1841 Each sample $x_i \in C_l$ receives weight

$$1842 \omega_i(x_0) = N_l^{-1} \omega_{x_0, l, C},$$

1845 so

$$1846 \omega_{x_0, j} = \mathbb{E} \left\{ \sum_{i=1}^N \omega_i(x_0) \psi_j^2(X_i) \right\} = \mathbb{E} \left\{ \sum_{l=1}^L \sum_{i \in C_l} N_l^{-1} \omega_{x_0, l, C} \psi_j^2(X_i) \right\}.$$

1850 Because $\omega_{x_0, l, C}$ is constant within C_l , we have

$$1851 \omega_{x_0, j} = \mathbb{E} \left[\sum_{l=1}^L \omega_{x_0, l, C} \mathbb{E} \left\{ N_l^{-1} \sum_{i \in C_l} \psi_j^2(X_i) | C \right\} \right].$$

1856 The within-cluster average may be approximated by the value at the center:

$$1857 \mathbb{E} \left\{ N_l^{-1} \sum_{i \in C_l} \psi_j^2(X_i) \right\} \approx \psi_j^2(C_l),$$

1861 hence

$$1862 \omega_{x_0, j} \asymp \mathbb{E} \left\{ \sum_{l=1}^L \omega_{x_0, l, C} \psi_j^2(C_l) \right\}. \quad (10)$$

1866 As $L \rightarrow \infty$ and the centers $\{C_l\}$ form a dense design from the marginal π , the sums in (equation 10) converge to integrals:

$$1867 \omega_{x_0, j} \asymp \frac{\int \tilde{K}(x_0, c) \psi_j^2(c) d\pi_c}{\int \tilde{K}(x_0, c) d\pi_c}.$$

1874 Applying the Cauchy–Schwarz inequality (e.g. Theorem 1 in Shi et al. (2009)), for kernels $K(x, y)$ such that $\int K^2(x, y) d\pi(y) < \infty$, the j -th eigenfunction of $K(x, y)$ is bounded by $|\psi_j(x)| \leq C/\lambda_j$ for some constant C . If $\lambda_j \asymp j^{-k}$, then it implies that $\omega_{x_0, j} \leq \lambda_j^{-2} \asymp j^{2k}$. \square

1880 D ADDITIONAL NUMERICAL RESULT

1883 D.1 CLUSTERING METHOD ANALYSIS

1885 To investigate the effect of clustering methods on the proposed estimator, we conduct additional
1886 experiments by testing different clustering methods while holding all other components fixed. Using
1887 the same dataset and retaining all 90 acoustic features under a Gaussian RBF kernel, we evaluate three
1888 clustering methods strategies for each (N, n) configuration: K-means (Lloyd’s algorithm with k-
1889 means++ initialization), random projection K-means, and K-medoids, a PAM-style medoid method
that is more robust to outliers and non-spherical clusters (Kaufman & Rousseeuw, 1990).

Table 6: Comparison the effect of clustering methods including K-means, random projection K-means (RP-kmeans) and K-mediods, on the proposed method using all the 90 features of the YearPredictionMSD dataset. Each entry reports (RMSE / MSE). Best MSE within each (N, n) configuration is highlighted in bold.

N	n	K-means		RP-kmeans		K-medoids	
		RMSE	MSE	RMSE	MSE	RMSE	MSE
2000	500	9.789	95.83	9.856	97.15	9.734	94.75
5000	1000	9.734	94.75	9.739	94.75	9.761	95.29
5000	2000	9.812	96.20	9.811	96.26	9.728	94.51
10000	1000	9.246	85.45	9.356	87.51	9.325	86.96
10000	2000	9.149	83.70	9.090	82.62	9.073	82.31
20000	1000	9.374	87.89	9.409	88.89	9.289	86.28
20000	2000	8.872	80.48	8.976	80.55	9.024	81.43
50000	1000	9.360	87.61	9.499	90.24	9.339	87.20
50000	2000	9.196	84.56	9.117	83.12	9.126	83.28

E NUMERICAL EXPERIMENT SETUP

E.1 ADDITIONAL DETAILS FOR SECTION 4.2

Data generation. Let $X = (X_1, \dots, X_{20})^\top$ with $X_j \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$. The response is

$$Y = f_0(X_1, \dots, X_5) + \epsilon, \quad f_0(x) = 10 \sin(\pi x_1 x_2) + 20(x_3 - 0.5)^2 + 10x_4 + 5x_5,$$

with $\epsilon \sim \mathcal{N}(0, 0.1^2)$. Only the first five coordinates influence the signal. For each replicate we draw a training set of size $N \in \{2000, 5000, 10000\}$ and evaluate on a fixed test set of size 100. We run 100 replications for each (N, n) , with $n \in \{100, 500, 1000\}$.

Hyperparameter tuning. All methods use the RBF kernel and an 8% hold-out for tuning.

FALKON (2017). Nyström landmarks are chosen uniformly with $M = n$. Bandwidth is selected from $\{\gamma_0/2, \gamma_0, 2\gamma_0\}$, where γ_0 is the median heuristic from the hold-out set. The ridge parameter is chosen from $\{0.5, 1.0, \dots, 4.0\}/N$. We run 300 PCG iterations.

FALKON (2020). The PyTorch/KeOps implementation of Meanti et al. (2020) uses the same landmark scheme, bandwidth grid, ridge grid, and 8% hold-out. Only MSE is reported due to cross-platform runtime differences.

Proposed method. We resample n points using the proposed weights, and average predictions over $B = 3$ subsets. Each replicate applies MATLAB k-means (k-means++ initialization) to compute n cluster centers. The center bandwidth is the median squared inter-center distance; the subset bandwidth is set to twice this value. We tune λ^* over $\{10^{-4}, 10^{-3}, \dots, 10^1\}$.

Nyström KRR. We use $M = n$ random landmarks and tune the ridge parameter over $\{10^{-4}, 10^{-3}, \dots, 10^1\}$ using the same hold-out split.

E.2 ADDITIONAL DETAILS FOR THE REAL-DATA STUDY

Data and preprocessing. We use the standard split of the YEARPREDICTIONMSD dataset: the first 463,715 rows form the training pool and the remaining observations form the test pool. For each experiment we sample $N \in \{2000, 5000, 10000, 20000\}$ from the training pool and use a fixed test set of size 1000. All randomization is seeded (rng(42)).

Feature selection. We rank the 90 features by absolute Pearson correlation with the response and keep the top $p \in \{30, 60, 90\}$, where $p=90$ means using all features.

Proposed method. We compute 500 clustering centers using RP- k -means (random projection to 32 dimensions followed by k -means). At estimation time we draw $B = 3$ independent subsets of size n , with bandwidths set by the median-distance rule. We use $\lambda_{\text{tuning}} = 1$ and report results for three kernels: Gaussian RBF (Prop(G)), Matérn-3/2 (Prop(M)), and Laplace (Prop(L)). Runtime includes center computation and the sum of the three kernel runs.

FALKON. We draw M uniform Nyström landmarks and run the PCG solver for 20 iterations with regularization $\lambda = 10^{-6}$. The kernel is always Gaussian RBF with $\sigma = 6$ following Rudi et al. (2017).

Nyström KRR. We sample M uniform landmarks, compute a low-rank decomposition of K_{ZZ} , and solve ridge regression using the same (M, γ, λ) as in FALKON, again using only the Gaussian RBF kernel.

Evaluation. For all methods we report test MSE and total wall-clock time, including tuning, subsampling, clustering, feature preprocessing, and model training.

E.3 TIME AND MEMORY COMPLEXITIES

We briefly sketch how the orders in Table 5 are obtained.

FALKON (Rudi et al., 2017). Standard Nyström KRR with M landmarks has time cost $O(NM) + O(M^3)$ and memory cost $O(NM + M^2)$. The FALKON analysis in Rudi et al. (2017) chooses the theoretically optimal $M \asymp \sqrt{N}$, which yields

$$O(NM) + O(M^3) = O(N^{3/2}) + O(N^{3/2}) = O(N^{3/2}),$$

with memory $O(NM + M^2) = O(N)$.

FALKON (Meanti et al., 2020). The large-scale implementation of Meanti et al. (2020) treats M as a tuning parameter and uses an iterative solver with per-iteration cost $O(NM) + O(M^2)$ and memory $O(NM + M^2)$. Since M is not tied to \sqrt{N} in theory, the generic $O(NM) + O(M^2)$ form is reported in Table 5.

Proposed + FALKON. Let n be the subdata size selected by the proposed method. The clustering-based selection uses K clusters and t iterations of k -means. At each iteration the cost is $O(nK)$, so the total time is $O(nKt)$ and memory is $O(n + K)$ for storing the data assignments and centers. Applying a Nyström-type solver (such as FALKON) to the subdata has computational cost $O(nM) + O(M^2)$ in time and $O(nM + M^2)$ in memory. Theorem 2 shows that the proposed estimator achieves the optimal full-data rate in Theorem 1 when

$$n \asymp N^{\frac{k+1-2\beta}{k+1}} \quad (\beta \in [0, 1/2)), \quad n \asymp N^{\frac{k}{k+1}} \quad (\beta > 1/2).$$

Writing $n \asymp N^\gamma$ for these two regimes, the combined selection + solver cost becomes

$$\text{Time: } O(N^\gamma M + N^\gamma Kt), \quad \text{Memory: } O(N^\gamma M + M^2 + K),$$

which matches the expressions reported for the Proposed + FALKON method in Table 5.