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

Anonymous authorsPaper 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 pointwise 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 also serves as a complement to existing fast kernel learning algorithms, such as the Nyström method and FALKON, which rely on randomly selected sub-datasets. Theoretical studies compare the efficiency of the proposed method to that based on the full data with optimally selected tuning parameters. Numerical experiments demonstrate the statistical efficiency of the proposed method over some existing methods based on randomly sampled data.

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) 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 columns. The FALKON and

EigenPro method combine the preconditioning idea and the random projection idea to speed up the computation.

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

Inspired by the existing sub-data selection methods, we propose a new sub-data selection methodology to estimate nonparametric functions in an RKHS. We first apply a clustering method such as k-means or other clustering algorithms to select representative data points. The nonparametric function values in each cluster are roughly approximated by the functional values at representative data points. To decide the sampling weights for each cluster, we minimize the MSE of the nonparametric function estimator that is constructed based on the selected representative data points. However, the sampling weights depend on a tuning parameter which is unknown. To address this issue, we adopted a one-step iteration procedure to choose a tuning parameter using representative data points via the BIC criterion or cross-validation. The optimal weights depend on the cluster centers decided by the K-means algorithm, the chosen kernel function, and the selected tuning parameter. After selecting multiple sub-datasets using the optimal weights, we apply the kernel machine method to each selected sub-dataset and aggregate the estimators to obtain the final estimate.

From a theoretical perspective, the proposed method is designed to select a sub-sample (with a fixed sample size) to minimize the expected prediction error at every test data point, hence it minimizes the expected prediction error for every test data set. It is different from existing research for sketch method in kernel learning, which has established its optimality to minimize the expected global prediction error. For example, Musco & Musco (2017) consider in-sample prediction error and Rudi et al. (2015; 2018); Alaoui & Mahoney (2015) consider the expected global generalization (out-of-sample) prediction error. While these results are interesting and important, these results do not directly guarantee generalization performance for every test data point. Our contribution is to establish the rate of convergence of the proposed estimator, and demonstrate its advantage in improving the convergence rate of the RHKS estimator based on subsamples obtained from a Simple Random 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.

From a numerical perspective, our numerical results show that our proposed method is computationally efficient when compared with a Simple Random Sample (SRS) approach and maintains estimation accuracy when compared with the full data approach. The integrated mean squared errors of our proposed method are comparable to that of using full data while computational time is as good as the SRS based approach. More importantly, the proposed approach achieves a good balance between statistical efficiency and computational efficiency when compared with the full data-based and the SRS-based approaches. In addition, we show that the proposed method has better MSE in out-of-sample prediction than existing sketch algorithms, including Nyström Williams & Seeger (2001) and FALKON Rudi et al. (2017). We further compare the performance of the proposed method with these sketch algorithms in the YearPredictionMSE data set. The proposed method has demonstrated its superior performance in improving the prediction MSE.

This manuscript is organized in the following way. In Section 2, we provide the basic framework, an introduction of regularization in RHKS 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 it with other sketch algorithms. A brief discussion is provided in Section 6. All the technical proofs 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) \mathscr{H}_K induced by a symmetric, positive-definite kernel $K: \mathscr{X} \times \mathscr{X} \to \mathbb{R}$.

If μ is a finite measure on $\mathscr X$ and $\int K(x,x)d\mu(x)<\infty$, then $K(x,y)=\sum_{j=1}^\infty \lambda_j\psi_j(x)\psi_j(y)$, 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

$$\mathscr{H}_K = \left\{ f(x) = \sum_{j=1}^{\infty} c_j \psi_j(x) : \sum_{j=1}^{\infty} \frac{c_j^2}{\lambda_j} < \infty \right\},\,$$

with norm $||f||_{\mathcal{H}_K}^2 = \sum_{j=1}^{\infty} c_j^2 / \lambda_j$. To estimate f_0 , we solve

$$\hat{f}_{N,\lambda_T} = \arg\min_{f \in \mathscr{H}_K} \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2 + \lambda_T ||f||_{\mathscr{H}_K}^2,$$

where $\lambda_T > 0$ is a tuning parameter. By applying representer theorem, the solution has the finite expansion $\hat{f}_{N,\lambda_T}(x) = \sum_{l=1}^N a_l K(x,X_l)$, with coefficients $a = (a_1,\ldots,a_N)^{\top}$. Plugging this into the objective yields $J(a) = \frac{1}{N} \|\mathbf{y} - \mathbf{K}a\|^2 + \lambda_T a^{\top} \mathbf{K}a$, where $\mathbf{K} \in \mathbb{R}^{N \times N}$ has entries $K(X_i,X_j)$. The solution is $\hat{a} = N^{-1} \lambda_T^{-1} (\mathbf{I} + N^{-1} \lambda_T^{-1} \mathbf{K})^{-1} \mathbf{y}$. Thus, for any x,

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

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

2.1 Proposed Method for Regularization in an RKHS

Given data pairs $(X_1,y_1),\cdots,(X_N,y_N)$ for a large N, the goal is to estimate a nonparametric function f(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 $(X_{k_1},y_{k_1}),\cdots,(X_{k_n},y_{k_n})$ with size n from the original sample with size n while maximizing the statistical efficiency of the resulting kernel machine method estimator $\hat{f}_{s\lambda}^*(x)$.

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 $\{\mathcal{C}_1,\cdots,\mathcal{C}_L\}$, and then re-sampling with optimal weights $\{w_{x,1,C},\cdots,w_{x,L,C}\}$ from these L clusters to obtain these representative data points.

The optimal weights $\{w_{x,1,C},\cdots,w_{x,L,C}\}$ are chosen by maximizing the statistical efficiency of the proposed estimator. Assume $\{C_1,\cdots,C_L\}$ are centers of the clusters $\{\mathcal{C}_1,\cdots,\mathcal{C}_L\}$ and the i-th cluster has N_i data points so that $\sum_{i=1}^n N_i = N$, the centers of those clusters are denoted by $\{C_1,\cdots,C_L\}$. Suppose we select a subdata set with size n, among them n_i data points are from the i-th cluster so that $n_i = nw_{x,i,C}$ and $\sum_{i=1}^L w_{x,i,C} = 1$. For all the data selected from the cluster \mathcal{C}_i ,

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: $f_{n\lambda'_{T}}^{*}(x)$

- 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) to partition the dataset into L clusters C_1, \dots, C_L . Let the cluster centers be $\{C_1, \dots, C_L\}$.
- 2: **Tuning** λ^* . Apply cross-validation or BIC method to find the tuning parameter λ^* for the RKHS regression using the representative points $\{C_1, \dots, C_L\}$.
- 3: Assigning weights. Using λ^* , compute optimal cluster weights w_1, \dots, w_L as defined in (1) for selecting clusters C_1, \dots, C_L . For the *i*-th cluster C_i with size N_i , assign the weight for the *j*-th data point in C_i as $w_{ij} = w_{x,i,C}/n_i$ for $j = 1, \dots, n_i$. Denote the resulting pointwise weights for $\{X_1, \dots, X_N\}$ by $w_1(x), \dots, w_N(x)$.
- 4: **for** b = 1 to B **do**
- 5: **Resampling with weights.** Using the weights $w_1(x), \dots, w_N(x)$, 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: **Fit RKHS estimator.** Using the selected sample $\{X_1^*, \dots, X_n^*\}$, compute the estimator $\hat{f}_{n\lambda'_n}^{*(b)}(x)$ for the *b*-th sampling with λ^* selected by a BIC approach or cross-validation, where

$$\hat{f}_{n\lambda_T'}^{*(b)}(x) = n^{-1}\lambda_T'^{-1} \left[K(x, X_1^*), \dots, K(x, X_n^*) \right] \left(\mathbf{I} + n^{-1}\lambda_T'^{-1} \mathbf{K}^* \right)^{-1} \mathbf{y}^*,$$

and \mathbf{K}^* is the $n \times n$ matrix with entries $\{K(X_i^*, X_j^*)\}_{i,j=1}^n$.

- 7: end for
- 8: Aggregate: The final estimate of f(x) is $\hat{f}_{n\lambda'_T}^*(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_{n\lambda'_T}^{*(b)}(x)$.

we approximate them as replications of the representative data centers C_i . Then, the corresponding mean model for $\bar{y}_i = \sum_{j=1}^{n_i} y_{ij}/n_i$ is

$$\bar{y}_i \approx f_0(C_i) + \bar{\epsilon}_i,$$

where $\bar{\epsilon}_i = \sum_{j=1}^{n_i} y_{ij}/n_i$ has mean zero and variance $\sigma^2/(nw_{x,i,C})$. We then consider the following RKHS nonparametric estimation of $f_0(x)$ given the above model:

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

where \mathbf{K}_c is an $L \times L$ matrix with (i, j)-th element $\{K(C_i, C_j)\}_{i,j=1}^L$ and $\bar{\mathbf{y}} = (\bar{y}_1, \dots, \bar{y}_L)$. Conditional on the centers $\{C_1, \dots, C_L\}$, the variance of the estimate $\hat{f}(x)$ is therefore

$$\operatorname{var}\{\hat{f}_{s\lambda}(x)\} = \sigma^2 L^{-3} \lambda^{-2} \left\{ K(x, C_1), \dots, K(x, C_L) \right\} \left(\mathbf{I} + L^{-1} \lambda^{-1} \mathbf{K}_c \right)^{-1} \mathbf{D}_w \left(\mathbf{I} + L^{-1} \lambda^{-1} \mathbf{K}_c \right)^{-1} \times \left\{ K(x, C_1), \dots, K(x, C_L) \right\}',$$

where $D_w = \text{diag}(1/w_{x,1,C}, \dots, 1/w_{x,L,C})'$. Because $E(\bar{\mathbf{y}}) = \{f_0(C_1), \dots, f_0(C_L)\}'$, the bias of the estimator $\hat{f}_{s\lambda}(x)$ is independent of the choices of $w_{x,i,C}$'s. Therefore, to minimize the conditional MSE of $\hat{f}_{s\lambda}(x)$, we could choose $w_{x,i,C}$'s that minimize the variance of $\hat{f}_{s\lambda}(x)$:

$$\hat{w}_{x,i,C} = \arg\min_{w_{x,i,C} \ge 0, \sum w_{x,i,C} = 1} \text{var}\{\hat{f}_{s\lambda}(x)\}.$$
 (1)

It is not difficult to check that $\operatorname{var}\{\hat{f}_{s\lambda}(x)\}=L^{-3}\sigma^2\lambda^{-2}\sum_{i=1}^LK_{xi}^2/w_{x,i,C}$. Therefore, the optimal weights $w_{x,i,C}$'s that minimize the variance is

$$w_{x,i,C} = |K_{xi}| / \sum_{i=1}^{L} |K_{xi}|, \tag{2}$$

where K_{xi} is the *i*-th element of the vector \boldsymbol{K}_x which is defined by

$$\boldsymbol{K}_{x} = \left(\boldsymbol{I} + L^{-1}\lambda^{-1}\boldsymbol{K}_{c}\right)^{-1} \begin{bmatrix} K(x,C_{1}) \\ \vdots \\ K(x,C_{L}) \end{bmatrix} = \begin{bmatrix} K_{x1} \\ \vdots \\ K_{xL} \end{bmatrix}.$$

217

218

219

220

221

222

223

224

225

226

227

228

229 230 231

232 233

234 235

236

237

238 239

240

241 242

243 244

245

246

247

248

249 250

251 252

253

254 255

256

257

258

259

260 261

262

263

264

265

266 267

268

269

Because the above optimal weights depend on the tuning parameter λ , 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. More specifically, our one-step iterative procedure is shown in Algorithm 1.

For the tuning procedure in Algorithm 1, it requires $\hat{f}_{s\lambda^*}(x)$ which is an RKHS estimator of $f_0(x)$ based on the centers selected. However, the estimator depends on the tuning parameter λ^* , which means the inverse operation in the estimator needed to computed for every possible candidates λ . This leads to a big computational burden. To mitigate computational burden for selecting the tuning parameter λ^* , we apply eigen-decomposition to the kernel matrix $K^* = Q\Lambda Q^T$, where Q is an orthogonal matrix that is independent of λ , and $\Lambda = \operatorname{diag}(\theta_i)_{n * n}$ is a diagonal matrix. Then, we can write

$$\hat{f}_{n\lambda_T^*}(x) = \boldsymbol{K}_x^T Q \Lambda_{\lambda^*}^{-1} Q^T \boldsymbol{y}^*,$$

where $K_x = (K(x, X_i^*))_{i=1}^n$ is an $n \times 1$ vector and $\Lambda_{\lambda^*} = \operatorname{diag}(\lambda^* + \theta_i)_{i=1}^n$.

THEORECTICAL JUSTIFICATION

ASYMPTOTIC RATE AND TUNING WITH FULL DATA

Theorem 1. Assume X has a probability density function $\pi(x)$ and $E[\psi_i^4(X)] < \infty$. If λ_i 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 function in RKHS $\|f\|_{\mathcal{H}_K}^2 = \sum_{j=1}^\infty a_j^2 \lambda_j$, $a_j^2 \lambda_j = j^{-a}(a>1)$, then the optimal rate of tuning parameter is: $\lambda_T \asymp N^{-\frac{k}{k+1+(a-1)}}$, and the corresponding asymptotic

integrated mean squared error is

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

where ||g|| is the L_2 norm of a function g.

The 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 $f \in \mathscr{H}_K$ 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 $a \to 1$)

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

For functions in a univariate Sobolev space of order m, the eigenvalue of the corresponding reproducing kernel is at the order of $\lambda_i \approx j^{-2m}$ Yuan & Cai (2010), then the rate of the convergence is given

$$\|\hat{f}_{N,\lambda_T} - f_0\|^2 = O_p(N^{-\frac{2m}{2m+1}}).$$

The above rate is known to be optimal in the literature (e.g. Zhang et al. (2015)).

Theorem 1 serves for multi-purposes. First, the rate of convergence results of Theorem 1 can be used as a basis for defining an efficiency index to compare different methods. The efficiency index is design to measure the trade-off between computational and statistical efficiency for the full data method, so that the efficiency index is roughly constant for all the sample sizes N. This would facilitate the comparison among different methods. Based on the results in Theorem 1, we may define the Efficiency Index as IMSE $^{-(2m+1)/(2m)}N^{-1/3}$. Details are given in Section 4. Second, Theorem 1 provides a general guidance about the order of tuning parameter choices subject to a constant. Given the optimal order of the tuning parameter was given in Theorem 1, we might specified a range for the tuning parameter λ_T . However, it is subject to a constant. To pin-down the constant, we adapted a commonly used method based on the BIC criterion Liu et al. (2007). Specifically, for a full data set, we might choose λ_T by minimizing the following BIC criterion:

$$\mathrm{BIC}(\lambda_T) = N \log \left[(\boldsymbol{y} - \hat{\mathbf{f}}_{N\lambda_T})^T (\boldsymbol{y} - \hat{\mathbf{f}}_{N\lambda_T}) \right] + \log(N) \mathrm{tr} \left[N^{-1} \lambda_T^{-1} \boldsymbol{K} (\boldsymbol{I} + N^{-1} \lambda_T^{-1} \boldsymbol{K})^{-1} \right],$$

where $\mathbf{f}_{N\lambda_T} = (f_{N\lambda_T}(x_1), \cdots, f_{N\lambda_T}(x_N))'$. Similar to the full data method, we can also apply BIC criterion to select λ^* for the proposed method. Third, the rate of convergence in Theorem 1 enables us to compare between the estimator based on the full data and the proposed estimator in Theorem 2.

3.2 ASYMPTOTIC RATE AND TUNING WITH RESAMPLING UNDER PROPOSED METHOD

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: $\omega_i(x_0) = n_l^{-1}\omega_{x_0,l,C}$, where $\omega_{x_0,l,C}$ was defined in (2). For each x_0 , define $\omega_{x_0,j} := \mathbb{E}\Big[\sum_{i=1}^N \omega_i(x_0)\,\psi_j^2(x_i)\Big]$. Denote the resampling subset of size n for the given x_0 is $\mathscr{S}(x_0) = \{x_1^*, \dots, x_n^*\}$. Then the proposed estimator of $f(x_0)$ is $\hat{f}_{n\lambda_T',x_0}(x_0)$, which is given by:

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

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

Theorem 2. Assume X has a probability density function $\pi(x)$ and $E[\psi_j^4(X)] < \infty$. 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 function in RKHS $\|f\|_{\mathscr{H}_K}^2 = \sum_{j=1}^\infty a_j^2 \lambda_j, \ a_j^2 \lambda_j = j^{-a}(a>1)$. Assume $\omega_{x_0,j} \ \asymp \ j^{2\beta}$ for $2\beta \le k \le 4\beta$. Then, the optimal rate of tuning parameter is: $\lambda_T' \ \asymp n^{-\frac{k-2\beta}{a+2k-4\beta}}$, and the corresponding asymptotic IMSE of the estimator $\hat{f}_{n,\lambda_T'}$ is

$$\|\hat{f}_{n,\lambda'_T} - f_0\|^2 = O_p\left(n^{-\frac{a-1+k}{(a-1)+(k+1)+(k-4\beta)}}\right).$$

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 $\left\|\hat{f}_{n,\lambda_T'} - f_0\right\|^2 = O_p\left(n^{-\frac{k}{(k+1)+(k-4\beta)}}\right)$ for any functions in the RKHS \mathscr{H}_K (by letting $a \to 1$). Because $k \le 4\beta$, 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 improvement of the estimator rate is due to the informative sampling, where the sampling weights carry the information of the full data set.

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. We consider N = 1000, 5000, 10000 and n = N/10, N/5 with kmeans centers of size L = n to compare estimation performance as the subsample size varies.

The random errors $\epsilon_1, \ldots, \epsilon_n \sim N(0, 0.5)$ i.i.d.. The kernel function $K(\cdot, \cdot)$ was Gaussian, with its dispersion parameter set to the sample standard deviation. Tuning parameters for the nonparametric estimation were selected using the BIC criterion. Results are based on 100 replications, and B=1 was used for the proposed method.

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):

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

where R is the full-data estimator and we choose $\alpha=5/4,\,\beta=1/3$. RE <1 indicates method A performed worse than the full data method; RE >1 indicates method A performed better than the full data method.

Table 1 summarizes IMSE, computation time, and RE for three methods: full-data, proposed, and SRS. Times are averaged per replication; EI is baseline-scaled to full data. The simulation study was implemented in an R environment.

Full Data

324 325

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

Proposed

SRS

326
327
328
329
330
331
000

332333334335

335 336

337 338

339 340

341342343

344 345

346 347 348

349 350 351

356 357 358

359 360 361

362

363

364 365 366

367

372373374

375

376

377

N**IMSE** Time **IMSE** Time RE **IMSE** Time RE n1000 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 0.0002.66 0.0152 0:0'0 0.89 0.0008 0:3'35 0.0017 0.0067 0:0'0 5000 500 0.003.26 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 0:0'20.76 10000 2000 0.0005 0:22'16 0.0006 0:0'26 2.98 0.0020 0:0'14 0.81

4.2 SIMULATION STUDY: FALKON vs. Proposed vs. Nyström

To assess the performance of our proposed K-means-based resampling method in higher dimensional case, we conducted a simulation study comparing it with two existing methods: (i) FALKON Rudi et al. (2017), a fast kernel ridge regression (KRR) solver combining Nyström approximation with a preconditioned conjugate gradient (PCG) scheme; (ii) a Nyström KRR baseline Williams & Seeger (2001).

Data generation. Let $X = (X_1, \dots, X_{20})^T$ be a 20-dimensional predictor with $X_j \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1]$. The response is generated by

$$Y = f(X_1, ..., X_5) + \varepsilon,$$
 $f(x) = 10\sin(\pi x_1 x_2) + 20(x_3 - 0.5)^2 + 10x_4 + 5x_5,$

with independent noise $\varepsilon \sim \mathcal{N}(0,0.1^2)$. Only the first five variables are relevant; the remaining 15 are nuisances. For each replicate we draw a training set of size $N \in \{2000,5000,10000\}$ and evaluate on a single fixed test set of size 100 (shared across all reps and settings). We run 100 replications per (N,n). All methods use the RBF kernel. To enforce an equal computational budget, we choose subsample size or the number of sketch columns $n \in \{100,500,1000\}$.

We report test MSE in prediction and the average computational time (in seconds), including both tuning and final training. Note that all three methods in this simulation study are implemented in MATLAB to make sure the comparison among computational cost is fair, since the FALKON code provided by Rudi et al. (2017) is in MATLAB.

Hyperparameter tuning. To keep the comparison fair, each method uses an 8% hold-out split from the training data for tuning to achieve the best performance.

- FALKON. We choose the Nyström landmarks by uniform subsampling with M=n. The bandwidth is tuned over $\{\gamma_0 \times 0.5, \ \gamma_0, \ \gamma_0 \times 2\}$, where γ_0 is the median-heuristic computed from a training subsample. The ridge parameter is tuned over $\lambda \in \{0.5, 1.0, \dots, 4.0\}/N$ to match FALKON's 1/N loss scaling. We run 300 PCG iterations.
- **Proposed.** We resample n data points with proposed weights, and final predictions average over B=3. For each replicate we run MATLAB's builtin k-means on the full training inputs with k-means++ initialization to obtain n centers. We set the center bandwidth ρ_{centers} to the median squared distance among the centers and, on each resampled subset, use a slightly broader kernel $\rho_{\text{sub}}=2\times (\text{median squared distance})$. We tune the parameter λ_T over $\{10^{-4},10^{-3},\ldots,10^1\}$ on the 8% hold-out.
- Nyström KRR. We use M=n random landmarks and tune the unscaled ridge over $\{10^{-4},10^{-3},\ldots,10^1\}$ on the same 8% hold-out.

Results. Table 2 reports the test MSE and total time averaged over 100 replications for each (N,n). 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 Typically trails the Proposed estimator in accuracy under the similar computational budget.

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

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

5 REAL-DATA STUDY: YEARPREDICTIONMSD WITH FEATURE SELECTION

We evaluate our proposed estimator in real data prediction by comparing it with the same two methods using the YEARPREDICTIONMSD dataset, which include 90 features and a continuous response.

Data and preprocessing. We use the standard split: the first 463,715 rows form a large training pool and the remaining rows a held–out test pool. For each experiment we sample $\texttt{train}_N \in \{2000, 5000, 10000, 20000\}$ for training and fix testing dataset of size 1,000. All randomization is seeded (rng (42)) for reproducibility.

Feature selection. Before fitting, we perform a filter step on the training data by ranking features by absolute Pearson correlation with the response and keep the top $K \in \{30, 60, 90\}$, with K=90 meaning "all features".

Methods and shared budget. We use an RBF kernel for all three methods. Additionally, the proposed estimator is also tested under Matérn-3/2 and Laplace kernels. We impose a shared subset budget and pair training set of size N and number of landmarks/subsamples.

For FALKON and Nyström KRR, the center/landmark budget equals M=n. For the proposed method, the subsample size equals to n.

- **Proposed.** We compute 500 clustering centers on the selected features using RP-k-means (random projection to 32 dimensions followed by k-means). At estimation time we average B=3 resamples of size n, with bandwidths set by the median distance rule. We fix λ_{tuning} =1 and report results for three kernels: Gaussian RBF (Prop(G)), Matérn-3/2 (Prop(M)), and Laplace (Prop(L)). For runtime, we report the total cost across these three kernels with the center selection time.
- FALKON. We draw M uniform landmarks from the training sets and run the PCG solver for iters = 20 steps with regularization $\lambda = 10^{-6}$. FALKON always uses the Gaussian RBF kernelwith $\sigma = 6$ as decribed in Rudi et al. (2017).
- Nyström KRR. We sample M landmarks uniformly, compute a low–rank decomposition, and solve ridge regression with the same (γ, λ, M) as FALKON, again restricted to Gaussian RBF.

Results. Table 3 summarizes the test MSEs and runtimes. The proposed method is consistently competitive with (and often outperforms) FALKON and Nyström in accuracy, particularly at larger budgets, while remaining efficient. The Laplace and Matérn variants provide additional robustness, with Gaussian RBF performing strongest in some settings. Reported runtimes confirm that the proposed approach scales well even for N=20k with n=2000, where total time cost remains on the order of 20 seconds.

Table 3: YearPredictionMSD data using kernel learning with 1,000 testing data points. Proposed method uses random-projection-k-means centers (L=500), B=3, fixed λ =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 (σ =6).

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

6 Discussion

Kernel machine methods are powerful for nonparametric learning, but face computational limits on large data. We developed a sub-data selection approach to overcome computational difficulties. The proposed method was designed to select an informative small subset from a large-scale sample, which is different from the sketch algorithms designed to sample columns from a large-scale kernel matrix. Although we do not use the entire full data set, we have shown that the proposed method performed better in the out-of-sample prediction than some existing sketch algorithms such as FALKON and Nyström methods, while maintaining the computational cost at a relatively low level. We have also shown that the proposed method has a better IMSE in estimating the unknown nonparametric function than that based on the SRS, and it has a comparable IMSE with the estimator based on full data. The proposed method aims to optimize the sub-data selection for prediction at every given test data points so that the overall prediction is minimized. We have also shown theoretically that the proposed sampling method can effectively uses the information from the full data to improve the rate of convergence of the proposed estimator.

REFERENCES

- Amirhesam Abedsoltan, Mikhail Belkin, and Parthe Pandit. Toward large kernel models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 61–78. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/abedsoltan23a.html.
- Mingyao Ai, Jun Yu, Huiming Zhang, and Haiying Wang. Optimal subsampling algorithms for big data regressions. *Statistica Sinica*, 2018. URL https://api.semanticscholar.org/CorpusID:198455923.
- Ahmed Alaoui and Michael W Mahoney. Fast randomized kernel ridge regression with statistical guarantees. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper_files/paper/2015/file/f3f27a324736617f20abbf2ffd806f6d-Paper.pdf.
- Ming-Chung Chang. Supervised stratified subsampling for predictive analytics. *Journal of Computational and Graphical Statistics*, 33(3):1017–1036, 2024. doi: 10.1080/10618600.2024.2304075. URL https://doi.org/10.1080/10618600.2024.2304075.
- Xueying Chen, Jerry Q. Cheng, and Min-ge Xie. *Divide-and-Conquer Methods for Big Data Analysis*, pp. 1–15. John Wiley Sons, Ltd, 2021. ISBN 9781118445112. doi: https://doi.org/10.1002/9781118445112.stat08298. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118445112.stat08298.
- Qianshun Cheng, HaiYing Wang, and Min Yang. Information-based optimal subdata selection for big data logistic regression. *Journal of Statistical Planning and Inference*, 209:112–122, 2020.
- Petros Drineas, Michael W Mahoney, and Shan Muthukrishnan. Sampling algorithms for 1 2 regression and applications. In *Proceedings of the seventeenth annual ACM-SIAM symposium on Discrete algorithm*, pp. 1127–1136, 2006.
- Chong Gu. Smoothing spline ANOVA models, volume 297. Springer, 2013.
- Runze Li, Dennis K.J. Lin, and Bing Li. Statistical inference in massive data sets. *Applied Stochastic Models in Business and Industry*, 29(5):399–409, 2013. doi: https://doi.org/10.1002/asmb.1927. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/asmb.1927.
- Dawei Liu, Xihong Lin, and Debashis Ghosh. Semiparametric regression of multidimensional genetic pathway data: Least-squares kernel machines and linear mixed models. *Biometrics*, 63(4): 1079–1088, 2007.
- Ping Ma and Xiaoxiao Sun. Leveraging for big data regression. *Wiley Interdisciplinary Reviews: Computational Statistics*, 7(1):70–76, 2015.
- SIYUAN MA and Mikhail Belkin. Diving into the shallows: a computational perspective on large-scale shallow learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/bf424cb7b0dea050a42b9739eb261a3a-Paper.pdf.
- Michael W Mahoney. Randomized algorithms for matrices and data. *arXiv preprint arXiv:1104.5557*, 2011.
 - Cameron Musco and Christopher Musco. Recursive sampling for the nystrom method. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/a03fa30821986dff10fc66647c84c9c3-Paper.pdf.

- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In J. Platt, D. Koller, Y. Singer, and S. Roweis (eds.), *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007. URL https://proceedings.neurips.cc/paper_files/paper/2007/file/013a006f03dbc5392effeb8f18fda755-Paper.pdf.
- Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Less is more: Nyström computational regularization. In *Neural Information Processing Systems*, 2015. URL https://api.semanticscholar.org/CorpusID:384343.
- Alessandro Rudi, Raffaello Camoriano, and Lorenzo Rosasco. Falkon: An optimal large scale kernel method. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- Alessandro Rudi, Daniele Calandriello, Luigi Carratino, and Lorenzo Rosasco. On fast leverage score sampling and optimal learning. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS'18, pp. 5677–5687, Red Hook, NY, USA, 2018. Curran Associates Inc.
- Grace Wahba. Spline models for observational data. SIAM, 1990.
- Haiying Wang and Yanyuan Ma. Optimal subsampling for quantile regression in big data. *Biometrika*, 108(1):99–112, 07 2020. ISSN 0006-3444. doi: 10.1093/biomet/asaa043. URL https://doi.org/10.1093/biomet/asaa043.
- HaiYing Wang, Min Yang, and John Stufken. Information-based optimal subdata selection for big data linear regression. *Journal of the American Statistical Association*, 114(525):393–405, 2019.
- Yuedong Wang. Smoothing splines: methods and applications. CRC press, 2011.
- 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.
- Tianbao Yang, Yu-feng Li, Mehrdad Mahdavi, Rong Jin, and Zhi-Hua Zhou. Nyström method vs random fourier features: A theoretical and empirical comparison. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. URL https://proceedings.neurips.cc/paper_files/paper/2012/file/621bf66ddb7c962aa0d22ac97d69b793-Paper.pdf.
- Yaqiong Yao and HaiYing Wang. A review on optimal subsampling methods for massive datasets. *Journal of Data Science*, 19(1):151–172, 2021.
- Ming Yuan and T Tony Cai. A reproducing kernel hilbert space approach to functional linear regression. *The Annals of Statistics*, 38(6):3412–3444, 2010.
- Yuchen Zhang, John Duchi, and Martin Wainwright. Divide and conquer kernel ridge regression: a distributed algorithm with minimax optimal rates. *J. Mach. Learn. Res.*, 16(1):3299–3340, January 2015. ISSN 1532-4435.

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:

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

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

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

So the objective function is equivalent to

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

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 $\hat{f}_{N\lambda_T}$ and the true function f_0 can be decomposed as

$$\hat{f}_{N\lambda_T} - f_0 = (\hat{f}_{N\lambda_T} - \bar{f}_{\infty\lambda_T}) + (\bar{f}_{\infty\lambda_T} - f_0)$$

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

$$\bar{f}_{\infty \lambda_T} = \arg \min_{f \in \mathscr{H}_K} \{ l_{\infty}(f) + \lambda_T ||f||_{\mathscr{H}_K}^2 \}$$

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

$$l_{\infty}(f) := E(l_N(f)) = E(y - f(x))^2$$

= $\sigma^2 + E[f(x) - f_0(x)]^2$.

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

$$\|\hat{f}_{N\lambda_{T}} - f\|^{2} = \int \langle \hat{f}_{N\lambda_{T}} - f, K(x_{0}, \cdot) \rangle_{\mathcal{H}_{K}}^{2} \pi(x_{0}) dx_{0}$$

$$= \int \langle \hat{f}_{N\lambda_{T}} - f, \sum_{j=1}^{\infty} \lambda_{j} \psi_{j}(\cdot) \psi_{j}(x_{0}) \rangle_{\mathcal{H}_{K}}^{2} \pi(x_{0}) dx_{0}$$

$$= \int (\sum_{j=1}^{\infty} \lambda_{j} \langle \hat{f}_{N\lambda_{T}} - f, \psi_{j}(\cdot) \rangle_{\mathcal{H}_{K}}^{2} \psi_{j}(x_{0}))^{2} \pi(x_{0}) dx_{0}$$

$$= \sum_{j=1}^{\infty} \lambda_{j}^{2} \langle \hat{f}_{N\lambda_{T}} - f, \psi_{j}(\cdot) \rangle_{\mathcal{H}_{K}}^{2}$$

$$= \sum_{j=1}^{\infty} \lambda_{j}^{2} \langle \hat{f}_{N\lambda_{T}} - f, \psi_{j}(\cdot) \rangle_{\mathcal{H}_{K}}^{2}$$

$$= \sum_{j=1}^{\infty} \lambda_{j}^{2} \langle \hat{c}_{j} - c_{j} \rangle^{2} \langle \psi_{j}, \psi_{j} \rangle_{\mathcal{H}_{K}}^{2}$$

$$= \sum_{j=1}^{\infty} (\hat{c}_{j} - c_{j})^{2}.$$

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

A.1 Asymptotic Order for Deterministic Error $\bar{f}_{\infty\lambda_T} - f_0$.

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

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

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

$$\begin{split} l_{\infty}(f) &= \sigma^2 + E[f(X) - f_0(X)]^2 \\ &= \sigma^2 + E[\sum_{j=1}^{\infty} c_j \psi_j(X) - \sum_{l=1}^{\infty} a_l \psi_l(X)]^2 \\ &= \sigma^2 + E[\sum_{j=1}^{\infty} \sum_{l=1}^{\infty} (c_j - a_j)(c_l - a_l) \psi_j(x) \psi_l(x)] \\ &= \sigma^2 + \sum_{j=l}^{\infty} (c_j - a_j)^2 E[\psi_j(x)^2] + \sum_{j \neq l}^{\infty} (c_j - a_j)(c_l - a_l) E[\psi_j(x) \psi_l(x)] \\ &= \sigma^2 + \sum_{j=l}^{\infty} (c_j - a_j)^2 \int_{R} \psi_j(x)^2 \pi(x) dx + \sum_{j \neq l}^{\infty} (c_j - a_j)(c_l - a_l) \int_{R} \psi_j(x) \psi_l(x) \pi(x) dx \\ &= \sigma^2 + \sum_{j=1}^{\infty} (c_j - a_j)^2. \end{split}$$

Then the corresponding objective function can be expressed as:

$$\bar{f}_{\infty\lambda_T}(c_j) = \arg\min\{Q_{\infty}(c_j)\} = \arg\min\{\sigma^2 + \sum_{j=1}^{\infty} (c_j - a_j)^2 + \lambda_T \sum_{j=1}^{\infty} \frac{c_j^2}{\lambda_j}\}.$$

We then take derivative w.r.t. c_i and obtain

$$Q_{\infty}'(c_j) = 2c_j - 2a_j + 2\lambda_T \lambda_j^{-1} c_j.$$

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

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

where $\bar{c_j} = a_j/(1 + \lambda_T \lambda_i^{-1})$.

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 $a_j^2 \asymp j^{-(a+k)}$::

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

Let J solve $\lambda_J \simeq \lambda_T$ so $J \simeq \lambda_T^{-1/k}$. Split the sum at J:

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

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

Therefore,

$$\|\bar{f}_{\infty\lambda_T} - f_0\|^2 \approx \lambda_T^{\frac{a+k-1}{k}} \qquad (a > 1, k > 1).$$

A.2 Asymptotic Order for Stochastic Error $\hat{f}_{N\lambda_T} - \bar{f}_{\infty\lambda_T}$.

Recall that

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

and the objective function for $\hat{f}_{N\lambda_T}$ can be written as:

$$\hat{f}_{N\lambda_T} = \arg\min_{f} \{Q_N(f)\} = \arg\min_{f} \{l_N(f) + \lambda_T ||f||_{\mathscr{H}_K}^2\}.$$

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

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

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

$$\hat{f}_{N\lambda_T} - \bar{f}_{\infty\lambda_T} = (\hat{f}_{N\lambda_T} - \tilde{f}) + (\tilde{f} - \bar{f}_{\infty\lambda_T}). \tag{3}$$

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

$$\begin{split} Dl_N(f) \cdot \eta &= \frac{d}{dh} \frac{1}{N} \sum_{i=1}^N (y_i - \langle f + h \eta, K(x_i, \cdot) \rangle_{\mathscr{H}_K})^2 \big|_{h=0} \\ &= \frac{d}{dh} \frac{1}{N} \sum_{i=1}^N (y_i - \langle f, K(x_i, \cdot) \rangle_{\mathscr{H}_K} - h \langle \eta, K(x_i, \cdot) \rangle_{\mathscr{H}_K})^2 \big|_{h=0} \\ &= -\frac{2}{N} \sum_{i=1}^N \langle \eta, K(x_i, \cdot) \rangle_{\mathscr{H}_K} \left(y_i - \langle f, K(x_i, \cdot) \rangle_{\mathscr{H}_K} \right) \\ &= -\frac{2}{N} \sum_{i=1}^N \eta(x_i) (y_i - f(x_i)). \\ Dl_\infty(f) \cdot \eta &= -2 \int \eta(x) (f_0(x) - f(x)) \pi(x) dx \\ &= -2 \int \langle \eta, K(x, \cdot) \rangle_{\mathscr{H}_K} \langle f_0 - f, K(x, \cdot) \rangle_{\mathscr{H}_K} \pi(x) dx \\ D^2 l_N(f) \cdot \eta \cdot g &= \frac{2}{N} \sum_{i=1}^N \langle \eta, K(x_i, \cdot) \rangle_{\mathscr{H}_K} \langle g, K(x_i, \cdot) \rangle_{\mathscr{H}_K} \\ &= \frac{2}{N} \sum_{i=1}^N \langle \langle \eta, K(x_i, \cdot) \rangle_{\mathscr{H}_K} K(x_i, \cdot), g \rangle_{\mathscr{H}_K}. \\ D^2 l_\infty(f) \cdot \eta \cdot g &= 2 \langle \int \langle \eta, K(x, \cdot) \rangle_{\mathscr{H}_K} K(x, \cdot) \pi(x) dx, g \rangle_{\mathscr{H}_K} \\ &= 2 \sum_{j=1}^\infty \lambda_j^2 \langle \eta, \psi_j \rangle_{\mathscr{H}_K} \langle g, \psi_j \rangle_{\mathscr{H}_K}. \\ D \|f\|_{\mathscr{H}_K}^2 \cdot \eta &= 2 \sum_{j=1}^\infty \lambda_j \langle f, \psi_j \rangle_{\mathscr{H}_K} \langle \eta, \psi_j \rangle_{\mathscr{H}_K} \quad \text{and} \\ D^2 \|f\|_{\mathscr{H}_K}^2 \cdot \eta \cdot g &= 2 \sum_{j=1}^\infty \lambda_j \langle g, \psi_j \rangle_{\mathscr{H}_K} \langle \eta, \psi_j \rangle_{\mathscr{H}_K}. \end{split}$$

A.2.1 Evaluate the order of $\tilde{f}_{N\lambda_T} - \bar{f}_{\infty\lambda_T}$.

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

$$Dl_{N\lambda_T}(\bar{f}) = Dl_{N\lambda_T}(\bar{f}) - Dl_{\infty\lambda_T}(\bar{f}) = Dl_N(\bar{f}) - Dl_{\infty}(\bar{f}).$$

For any function η , we have

$$\begin{split} E[Dl_{N\lambda_T}(\bar{f}) \cdot \eta]^2 &= E[Dl_N(\bar{f}) \cdot \eta - Dl_{\infty}(\bar{f}) \cdot \eta]^2 \\ &= E\Big[-\frac{2}{N} \sum_{i=1}^N (y_i - \bar{f}(x_i)) \eta(x_i) + 2 \int (\bar{f}(x) - f_0(x)) \eta(x) \pi(x) dx \Big]^2 \\ &= \frac{4}{N^2} E\Big\{ \sum_{i=1}^N \Big[(y_i - \bar{f}(x_i)) \eta(x_i) - E[(\bar{f}(x) - f_0(x)) \eta(x)] \Big] \Big\}^2 \\ &= \frac{4}{N} Var \big[(y - \bar{f}(x)) \eta(x) \big] \le \frac{4}{N} E\big[(y - \bar{f}(x)) \eta(x) \big]^2 \approx \frac{4}{N}. \end{split}$$

By the definition of G_{λ_T} , we have

$$\begin{split} G_{\lambda_T} \cdot \eta \cdot g &= \frac{1}{2} D_{\infty \lambda_T}^2 \cdot \eta \cdot g = \frac{1}{2} D l_{\infty}^2(f) \cdot \eta \cdot g + \frac{1}{2} D^2 \|f\|_{\mathcal{H}_{\mathcal{K}}}^2 \cdot \eta \cdot g \\ &= \sum_{j=1}^{\infty} \lambda_j^2 \langle \eta, \psi_j \rangle_{\mathcal{H}_{\mathcal{K}}} \langle g, \psi_j \rangle_{\mathcal{H}_{\mathcal{K}}} + \sum_{j=1}^{\infty} \lambda_T^{'} \lambda_j < \eta, \psi_j >_{\mathcal{H}_{\mathcal{K}}} < g, \psi_j >_{\mathcal{H}_{\mathcal{K}}} \\ &= \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda_j^{-1} \lambda_T) \langle g, \psi_j \rangle_{\mathcal{H}_{\mathcal{K}}} \langle \eta, \psi_j \rangle_{\mathcal{H}_{\mathcal{K}}}. \end{split}$$

Let $\eta = \psi_m$ gives

$$\langle G_{\lambda_T} g, \psi_m \rangle_{\mathscr{H}_K} = (\lambda_m + \lambda_T) \langle g, \psi_m \rangle_{\mathscr{H}_K}$$

which leads to

$$\langle G_{\lambda_T}^{-1} g, \psi_m \rangle_{\mathscr{H}_K} = (\lambda_m + \lambda_T)^{-1} \langle g, \psi_m \rangle_{\mathscr{H}_K}.$$

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

$$\begin{split} E \| \tilde{f}_{N\lambda_T} - \bar{f}_{\infty\lambda_T} \|^2 &= E \Big\| \frac{1}{2} G_{\lambda_T}^{-1} D l_{N\lambda_T} (\bar{f}_{\infty\lambda_T}) \Big\|^2 \\ &= \frac{1}{4} E \Big[\sum_{j=1}^{\infty} \lambda_j^2 \left\langle G_{\lambda_T}^{-1} D l_{N\lambda_T} (\bar{f}_{\infty\lambda_T}), \psi_j \right\rangle_{\mathcal{H}_K}^2 \Big] \\ &= \frac{1}{4} E \Big[\sum_{j=1}^{\infty} \lambda_j^2 \frac{1}{(\lambda_j + \lambda_T)^2} \left\langle D l_{N\lambda_T} (\bar{f}_{\infty\lambda_T}), \psi_j \right\rangle_{\mathcal{H}_K}^2 \Big] \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \frac{\lambda_j^2}{(\lambda_j + \lambda_T)^2} E \Big[\left\langle D l_{N\lambda_T} (\bar{f}_{\infty\lambda_T}), \psi_j \right\rangle_{\mathcal{H}_K}^2 \Big] \\ &\lesssim \frac{1}{N} \sum_{j=1}^{\infty} \frac{\lambda_j^2}{(\lambda_j + \lambda_T)^2} \quad \text{(using } E[\langle D l_{N\lambda_T} (\bar{f}), \psi_j \rangle^2] \times N^{-1}) \\ &\lesssim \frac{1}{N} \int_1^{\infty} \frac{1}{(1 + \lambda_T j^k)^2} \, dj \quad \text{(since } \lambda_j \times j^{-k}, k > 1) \\ &\lesssim \frac{1}{N} \lambda_T^{-1/k}. \end{split}$$

A.2.2 EVALUATE THE ORDER OF $\hat{f}_{N\lambda_T} - \tilde{f}_{N\lambda_T}$.

Next we need to find the stochastic order of the second term $\hat{f}_{N\lambda_T} - \tilde{f}_{N\lambda_T}$ in the expression (3). Note that, $\hat{f} = \hat{f}_{N\lambda_T}$ is the solution of the following first order equation

$$Dl_{N\lambda_T}(\hat{f}) = Dl_{N\lambda_T}(\bar{f}) + D^2 l_{N\lambda_T}(\bar{f}) \cdot (\hat{f} - \bar{f}) = 0$$

and by the definition of $\tilde{f}_{N\lambda_T}$, we have

$$Dl_{N\lambda_T}(\bar{f}) + D^2 l_{\infty\lambda_T}(\bar{f}) \cdot (\tilde{f}_{N\lambda_T} - \bar{f}) = 0$$

From the above two equations, we can find:

$$D^{2}l_{N\lambda_{T}}(\bar{f})\cdot(\hat{f}-\bar{f})=D^{2}l_{\infty\lambda_{T}}(\bar{f})\cdot(\hat{f}-\bar{f}).$$

Then we can write:

$$\begin{split} D^2 l_{\infty \lambda_T}(\bar{f}) \cdot (\hat{f}_{N \lambda_T} - \tilde{f}) &= D^2 l_{\infty \lambda_T}(\bar{f}) \cdot (\hat{f}_{N \lambda_T} - \bar{f}) + D^2 l_{\infty \lambda_T}(\bar{f}) \cdot (\bar{f} - \tilde{f}) \\ &= D^2 l_{\infty \lambda_T}(\bar{f}) \cdot (\hat{f}_{N \lambda_T} - \bar{f}) - D^2 l_{N \lambda_T}(\bar{f}) \cdot (\hat{f}_{N \lambda_T} - \bar{f}) \\ &= D^2 l_{\infty}(\bar{f}) \cdot (\hat{f}_{N \lambda_T} - \bar{f}) - D^2 l_{N}(\bar{f}) \cdot (\hat{f}_{N \lambda_T} - \bar{f}). \end{split}$$

Using the definition of G_{λ} , we have

$$\hat{f}_{N\lambda_T} - \tilde{f} = \frac{1}{2} G_{\lambda}^{-1} D^2 l_{\infty\lambda_T}(\bar{f}) \cdot (\hat{f}_{N\lambda_T} - \tilde{f}) = \frac{1}{2} G_{\lambda}^{-1} \{ D^2 l_{\infty}(\bar{f}) \cdot (\hat{f}_{N\lambda_T} - \bar{f}) - D^2 l_N(\bar{f}) \cdot (\hat{f}_{N\lambda_T} - \bar{f}) \}.$$

Using the similar steps in evaluating $E\|\tilde{f} - \bar{f}_{\infty\lambda_T}\|_{\mathscr{H}_K}^2$, we have

$$\begin{split} &\|\hat{f}_{N\lambda_{T}} - \tilde{f}\|^{2} \\ &= \left\| \frac{1}{2} G_{\lambda_{T}}^{-1} \left(D^{2} l_{\infty}(\bar{f}) (\hat{f} - \bar{f}) - D^{2} l_{N}(\bar{f}) (\hat{f} - \bar{f}) \right) \right\|^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} \left\langle G_{\lambda_{T}}^{-1} \left(D^{2} l_{\infty}(\bar{f}) (\hat{f} - \bar{f}) - D^{2} l_{N}(\bar{f}) (\hat{f} - \bar{f}) \right), \psi_{j} \right\rangle_{\mathcal{H}_{K}}^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} \frac{1}{(\lambda_{j} + \lambda_{T})^{2}} \left(\left\langle D^{2} l_{\infty}(\bar{f}) (\hat{f} - \bar{f}), \psi_{j} \right\rangle_{\mathcal{H}_{K}} - \left\langle D^{2} l_{N}(\bar{f}) (\hat{f} - \bar{f}), \psi_{j} \right\rangle_{\mathcal{H}_{K}}^{2} \right)^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(\left\langle D^{2} l_{\infty}(\bar{f}) (\hat{f} - \bar{f}), \psi_{j} \right\rangle_{\mathcal{H}_{K}} - \left\langle D^{2} l_{N}(\bar{f}) (\hat{f} - \bar{f}), \psi_{j} \right\rangle_{\mathcal{H}_{K}}^{2} \right)^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{\infty} \lambda_{l}^{2} < \hat{f} - \bar{f}, \psi_{l} \right\rangle_{\mathcal{H}_{K}} < \psi_{j}, \psi_{l} \right\rangle_{\mathcal{H}_{K}}^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{N} \sum_{k=1}^{\infty} \lambda_{l} < \hat{f} - \bar{f}, \psi_{l} \right\rangle_{\mathcal{H}_{K}}^{2} < \psi_{j}, \psi_{k} \right\rangle_{\mathcal{H}_{K}}^{2} \psi_{l}(x_{i}) \psi_{k}(x_{i}) \right)^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{N} \sum_{l=1}^{\infty} \lambda_{l} < \hat{f} - \bar{f}, \psi_{l} \right\rangle_{\mathcal{H}_{K}}^{2} \psi_{l}(x_{i}) \psi_{k}(x_{i}) \right)^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{N} \sum_{l=1}^{\infty} \lambda_{l} < \hat{f} - \bar{f}, \psi_{l} \right\rangle_{\mathcal{H}_{K}}^{2} \psi_{l}(x_{i}) \psi_{k}(x_{i}) \right)^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{N} \sum_{l=1}^{\infty} \lambda_{l} < \hat{f} - \bar{f}, \psi_{l} \right\rangle_{\mathcal{H}_{K}}^{2} \psi_{l}(x_{i}) \psi_{k}(x_{i}) \right)^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{N} \sum_{l=1}^{\infty} \lambda_{l} < \hat{f} - \bar{f}, \psi_{l} \right\rangle_{\mathcal{H}_{K}}^{2} \psi_{l}(x_{i}) \psi_{l}(x_{i}) \right)^{2} \\ &= \frac{1}{4} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{N} \sum_{l=1}^{\infty} \lambda_{l} < \hat{f} - \bar{f}, \psi_{l} \right\rangle_{\mathcal{H}_{K}}^{2} \psi_{l}(x_{i}) \psi_{l}(x_{i}) \right)^{2} \\ &= \frac{1}{4} \sum_{l=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \left(2 \sum_{l=1}^{\infty} \lambda_{l} < \hat{f} - \bar{f}, \psi_{l} \right)_{\mathcal{H}_{K}}^{2} \psi_{l}(x_{i}) \psi_{l}(x_{i}) \right)^{2} \\ &=$$

$$= N^{-2} \sum_{j=1}^{\infty} \lambda_j^2 (1 + \lambda_T \lambda_j^{-1})^{-2} \|\hat{f} - \bar{f}\|_{\mathcal{H}_{\mathcal{K}}}^2 \underbrace{\left[\sum_{l=1}^{\infty} \lambda_l (\sum_{i=1}^N E[\psi_j(x_i) \psi_l(x_i)] - \sum_{i=1}^N \psi_j(x_i) \psi_l(x_i))^2 \right]}_{\text{Term } \Lambda}.$$

For Term A, we can expand it as:

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

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:

$$\begin{split} &E[\lambda_{j}\sum_{i=1}^{N}[E(\psi_{j}^{2}(x_{i}))-\psi_{j}^{2}(x_{i})]^{2}]\\ &=E\left[\lambda_{j}\sum_{i=1}^{N}[\psi_{j}^{2}(x_{i})-1]^{2}+\sum_{l\neq j}^{\infty}\lambda_{l}\sum_{i=1}^{N}(\psi_{j}(x_{i})\psi_{l}(x_{i}))^{2}\right]\\ &=\lambda_{j}NE[\psi_{j}^{2}(x_{i})-1]^{2}+NE\left[\sum_{l\neq j}^{\infty}\lambda_{l}\psi_{j}^{2}(x_{i})\psi_{l}^{2}(x_{i})\right]\\ &\leq\lambda_{j}NVar[\psi_{j}^{2}(x)]+N\left\{E\left(\sum_{l=1}^{\infty}\lambda_{l}\psi_{l}^{2}(x)\right)^{2}+Var\left(\sum_{l=1}^{\infty}\lambda_{l}\psi_{l}^{2}(x)\right)\right\}^{\frac{1}{2}}(E(\psi_{j}^{4}(x)))^{\frac{1}{2}}\\ &=\lambda_{j}NVar[\psi_{j}^{2}(x)]+N\left\{E[K(x,x)]^{2}+Var(K(x,x))\right\}^{\frac{1}{2}}(E(\psi_{j}^{4}(x)))^{\frac{1}{2}}\\ &\asymp N. \end{split}$$

$$Var\Big[\sum_{l=1}^{\infty} \lambda_{l} \sum_{i \neq k}^{N} \left(E[\psi_{j}(x_{i})\psi_{l}(x_{i})] - \psi_{j}(x_{i})\psi_{l}(x_{i})\right) \left(E[\psi_{j}(x_{k})\psi_{l}(x_{k})] - \psi_{j}(x_{k})\psi_{l}(x_{k})\right)\Big]$$

$$\leq \sum_{l=1}^{\infty} \lambda_{l}^{2} (N-1)^{2} E\Big[\left(E[\psi_{j}(x_{i})\psi_{l}(x_{i})] - \psi_{j}(x_{i})\psi_{l}(x_{i})\right)^{2}\Big] E\Big[\left(E(\psi_{j}(x_{k})\psi_{l}(x_{k})) - \psi_{j}(x_{k})\psi_{l}(x_{k})\right)^{2}\Big]$$

$$+ \sum_{l \neq v}^{\infty} \lambda_{l} \lambda_{v} (N-1)^{2} Var(\psi_{j}(x)\psi_{l}(x)) Var(\psi_{j}(x)\psi_{v}(x))$$

$$= (N-1)^{2} \sum_{l=1}^{\infty} \lambda_{l}^{2} Var^{2}(\psi_{j}(x)\psi_{l}(x)) + \sum_{l=1}^{\infty} \lambda_{l} \lambda_{v} (N-1)^{2} Var(\psi_{j}(x)\psi_{l}(x)) Var(\psi_{j}(x)\psi_{v}(x))$$

$$\leq (N-1)^{2} \left[\sum_{l=1}^{\infty} \lambda_{l}^{2} Var^{2}(\psi_{j}(x)\psi_{l}(x)) + (\sum_{l=1}^{\infty} \lambda_{l}^{2})^{\frac{1}{2}} (\sum_{v=1}^{\infty} \lambda_{l}^{2})^{\frac{1}{2}} Var(\psi_{j}(x)\psi_{l}(x)) Var(\psi_{j}(x)\psi_{v}(x)) \right] \\ \approx 2(N-1)^{2} \left(\int_{1}^{\infty} l^{-2k} dl \right) \approx N^{2}.$$

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:

$$\|\hat{f}_{N\lambda_{T}} - \tilde{f}\|_{\mathcal{H}_{\mathcal{K}}}^{2} = N^{-1} \sum_{j=1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \|\hat{f} - \bar{f}\|_{\mathcal{H}_{\mathcal{K}}}^{2}$$

$$= N^{-1} \int_{1}^{\infty} \lambda_{j}^{2} (1 + \lambda_{T} \lambda_{j}^{-1})^{-2} \|\hat{f} - \bar{f}\|_{\mathcal{H}_{\mathcal{K}}}^{2}$$

$$\approx N^{-1} \lambda_{T}^{-1/k} \|\hat{f} - \bar{f}\|_{\mathcal{H}_{\mathcal{K}}}^{2}.$$

Now, it follows that

$$\|\hat{f}_{N\lambda_T} - \tilde{f}\|_{\mathscr{H}_{\mathscr{K}}}^2 = O_p(N^{-1}\lambda_T^{-1/k}\|\hat{f} - \bar{f}\|_{\mathscr{H}_{\mathscr{K}}}^2).$$

If
$$N^{-1}\lambda_T^{-1/k} \to 0$$
,

$$\|\hat{f}_{N\lambda_T} - \tilde{f}\|_{\mathcal{H}_{\mathcal{K}}}^2 = o_p(\|\hat{f} - \bar{f}\|_{\mathcal{H}_{\mathcal{K}}}) = o_p(1)\|\hat{f} - \bar{f}\|_{\mathcal{H}_{\mathcal{K}}}.$$

Observed that

$$\|\tilde{f} - \bar{f}\|_{\mathscr{H}_{\mathscr{K}}}^{2} \ge \|\hat{f} - \bar{f}\|_{\mathscr{H}_{\mathscr{K}}}^{2} - \|\hat{f} - \tilde{f}\|_{\mathscr{H}_{\mathscr{K}}}^{2} = (1 - o_{p}(1))\|\hat{f} - \bar{f}\|_{\mathscr{H}_{\mathscr{K}}}^{2},$$

then

$$\|\hat{f} - \bar{f}\|_{\mathscr{H}_{\mathscr{K}}}^2 = O_p(\|\tilde{f} - \bar{f}\|_{\mathscr{H}_{\mathscr{K}}}) = O_p(N^{-1}\lambda_T^{-1/k}).$$

A.3 ASYMPTOTIC ORDER FOR MSE AND TUNING.

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

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

To find the optimal order of λ_T , 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

$$d^{\frac{a+k}{k}} \times N^{-1} \implies \lambda_T \times N^{-\frac{k}{a+k}}.$$

Also check

$$N^{-1}\lambda_T^{-1/k} = N^{-1} \big(N^{-k/(a+k)}\big)^{-1/k} = N^{-\frac{a+k-1}{a+k}} \to 0 \quad (a>1,\; k>1).$$

Therefore,

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

This completes the proof of Theorem 1.

B PROOF OF THE PROPOSED METHOD

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:

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

Denote the resampling subset of size n for the given x_0 is $\mathscr{S}(x_0) = \{x_1^*, \dots, x_n^*\}$, with replacement using probabilities $\{\omega_i(x_0)\}$. Let $I_i(x_0)$ be the count of how many times index i appears in $\mathscr{S}(x_0)$ 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 $\hat{f}_{n\lambda'_T,x_0}(x_0)$, which is given by:

$$\hat{f}_{n\lambda'_{T},x_{0}} = \arg\min_{f \in \mathcal{H}_{K}} \left\{ \frac{1}{n} \sum_{i \in \mathcal{S}(x_{0})} (y_{i} - f(x_{i}))^{2} + \lambda'_{T} ||f||_{\mathcal{H}_{K}}^{2} \right\}.$$

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

$$\|\hat{f}_{n\lambda'_{T}} - f\|^{2} = \int \langle \hat{f}_{n\lambda'_{T},x_{0}} - f, K(x_{0},\cdot) \rangle_{\mathcal{H}_{K}}^{2} \pi(x_{0}) dx_{0}$$

$$= \int \langle \hat{f}_{n\lambda'_{T},x_{0}} - f, \sum_{j=1}^{\infty} \lambda_{j} \psi_{j}(\cdot) \psi_{j}(x_{0}) \rangle_{\mathcal{H}_{K}}^{2} \pi(x_{0}) dx_{0}$$

$$= \int (\sum_{j=1}^{\infty} \lambda_{j} \langle \hat{f}_{n\lambda'_{T},x_{0}} - f, \psi_{j}(\cdot) \rangle_{\mathcal{H}_{K}} \psi_{j}(x_{0}))^{2} \pi(x_{0}) dx_{0}.$$

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

$$\begin{split} l_{\infty}(f) &:= \mathbb{E}\left[\frac{1}{n} \sum_{i \in \mathscr{S}(x_0)} (y_i - f(x_i))^2\right] \\ &= \mathbb{E}\left[\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^N I_i(x_0) \left(y_i - f(x_i)\right)^2 \,\middle|\, \{x_i\}_{i=1}^N\right]\right] \\ &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^N \mathbb{E}[I_i(x_0) \,|\, \{x_i\}_{i=1}^N\right] \left(y_i - f(x_i)\right)^2\right] \\ &= \mathbb{E}\left[\sum_{i=1}^N \omega_i(x_0) \left(y_i - f(x_i)\right)^2\right] \qquad \left(\mathbb{E}[I_i(x_0) \,|\, \{x_i\}] = n \,\omega_i(x_0)\right) \\ &= \mathbb{E}\left[\sum_{i=1}^N \omega_i(x_0) \,\Big\{\sigma^2 + \left(f(x_i) - f_0(x_i)\right)^2\Big\}\right] \\ &= \sigma^2 \,+\, \mathbb{E}\left[\sum_{i=1}^N \omega_i(x_0) \,\Big\langle f - f_0, \, K(x_i, \cdot) \Big\rangle_{\mathscr{H}_K}^2\right]. \end{split}$$

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

$$\hat{f}_{x_0} - f_0 = (\hat{f}_{x_0} - \bar{f}_{x_0}) + (\bar{f}_{x_0} - f_0),$$

where \bar{f}_{x_0} is the solution of the following objective function

$$\bar{f}_{x_0} = \arg\min l_{\infty}(f) + \lambda_T' ||f||_{\mathscr{H}_K}^2.$$

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

$$Dl_{\infty}(f) \cdot \eta = \frac{d}{dh} \mathbb{E} \left[\sum_{i=1}^{N} \omega_{i}(x_{0}) \left(f(x_{i}) + h\eta(x_{i}) - f_{0}(x_{i}) \right)^{2} \right]_{h=0}$$

$$= 2 \mathbb{E} \left[\sum_{i=1}^{N} \omega_{i}(x_{0}) \eta(x_{i}) \left(f(x_{i}) - f_{0}(x_{i}) \right) \right]$$

$$= 2 \mathbb{E} \left[\sum_{i=1}^{N} \omega_{i}(x_{0}) \left\langle \eta, K(x_{i}, \cdot) \right\rangle_{\mathscr{H}_{K}} \left\langle f - f_{0}, K(x_{i}, \cdot) \right\rangle_{\mathscr{H}_{K}} \right].$$

$$\begin{split} D^{2}l_{\infty}(f) \cdot \eta \cdot g &= 2 \, \mathbb{E} \Bigg[\sum_{i=1}^{N} \omega_{i}(x_{0}) \, \big\langle \eta, K(x_{i}, \cdot) \big\rangle_{\mathscr{H}_{K}} \, \big\langle g, K(x_{i}, \cdot) \big\rangle_{\mathscr{H}_{K}} \Bigg] \\ &= 2 \sum_{j=1}^{\infty} \lambda_{j}^{2} \, \big\langle \eta, \psi_{j} \big\rangle_{\mathscr{H}_{K}} \, \big\langle g, \psi_{j} \big\rangle_{\mathscr{H}_{K}} \, \underbrace{\mathbb{E} \Bigg[\sum_{i=1}^{N} \omega_{i}(x_{0}) \, \psi_{j}(x_{i})^{2} \Bigg]}_{=: \, \omega_{x_{0}, j}}. \end{split}$$

For the empirical loss,

$$Dl_{n,x_0}(f) \cdot \eta = \frac{d}{dh} \frac{1}{n} \sum_{i \in \mathscr{S}(x_0)} \left(y_i - \langle f + h\eta, K(x_i, \cdot) \rangle_{\mathscr{H}_K} \right)^2 \Big|_{h=0}$$
$$= -\frac{2}{n} \sum_{i \in \mathscr{S}(x_0)} \left\langle \eta, K(x_i, \cdot) \right\rangle_{\mathscr{H}_K} \left(y_i - \langle f, K(x_i, \cdot) \rangle_{\mathscr{H}_K} \right).$$

$$D^{2}l_{n,x_{0}}(f) \cdot \eta \cdot g = \frac{2}{n} \sum_{i \in \mathscr{S}(x_{0})} \left\langle \eta, K(x_{i}, \cdot) \right\rangle_{\mathscr{H}_{K}} \left\langle g, K(x_{i}, \cdot) \right\rangle_{\mathscr{H}_{K}}$$

$$= \frac{2}{n} \sum_{i \in \mathscr{S}(x_{0})} \left(\sum_{j=1}^{\infty} \lambda_{j} \left\langle \eta, \psi_{j} \right\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{i}) \right) \left(\sum_{k=1}^{\infty} \lambda_{k} \left\langle g, \psi_{k} \right\rangle_{\mathscr{H}_{K}} \psi_{k}(x_{i}) \right)$$

$$= \frac{2}{n} \sum_{i \in \mathscr{S}(x_{0})} \sum_{j,k \geq 1} \lambda_{j} \lambda_{k} \left\langle \eta, \psi_{j} \right\rangle_{\mathscr{H}_{K}} \left\langle g, \psi_{k} \right\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{i}) \psi_{k}(x_{i}).$$

For the RKHS penalty,

$$D\|f\|_{\mathcal{H}_{K}}^{2} \cdot \eta = 2 \langle f, \eta \rangle_{\mathcal{H}_{K}},$$

$$D^{2}\|f\|_{\mathcal{H}_{K}}^{2} \cdot \eta \cdot g = 2 \langle \eta, g \rangle_{\mathcal{H}_{K}} = 2 \sum_{i=1}^{\infty} \lambda_{j} \langle \eta, \psi_{j} \rangle_{\mathcal{H}_{K}} \langle g, \psi_{j} \rangle_{\mathcal{H}_{K}}.$$

To evaluate the stochastic error, we use a similar method in the proof of Theorem 1 by defining an intermediate quantity $\tilde{f}_{x_0} = \bar{f}_{x_0} - \frac{1}{2} G_{\lambda'_n}^{-1} Dl_{n\lambda'_T}(\bar{f})$, with $G_{\lambda'_T} = \frac{1}{2} D_{\infty\lambda'_n}^2(\bar{f}_{x_0})$. Then, we define

$$\tilde{f}_n = \{ \tilde{f}_{x_0} : x_0 \in \mathcal{X} \}.$$

and decompose $\hat{f}_{n\lambda_T'} - \bar{f}_n$ into

$$\hat{f}_{n\lambda'_T} - \bar{f}_n = (\hat{f}_{n\lambda'_T} - \tilde{f}_n) + (\tilde{f}_n - \bar{f}_n).$$

B.1 ASYMPTOTIC ORDER FOR STOCHASTIC ERROR.

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

B.1.1 Step 1: Evaluate the order of $\tilde{f}_n - \bar{f}_n$.

Assume that $\omega_{x_0,j} = E\left[\sum_{i=1}^N \omega_i(x_0)\psi_j^2(x_i)\right]$. To this end, we first obtain the eigenvalues of the operator $G_{\lambda_T'}$.

$$\begin{split} G_{\lambda_T'} \cdot \eta \cdot g &= \frac{1}{2} D_{\infty \lambda_T'}^2 \cdot \eta \cdot g = \frac{1}{2} D l_{\infty}^2(f) \cdot \eta \cdot g + \frac{1}{2} D^2 \|f\|_{\mathcal{H}_{\mathcal{K}}}^2 \cdot \eta \cdot g \\ &= \sum_{j=1}^\infty \lambda_j^2 < \eta, \psi_j >_{\mathcal{H}_{\mathcal{K}}} < g, \psi_j >_{\mathcal{H}_{\mathcal{K}}} E \left[\sum_{i=1}^N \omega_i(x_0) \psi_j^2(x_i) \right] \\ &+ \sum_{j=1}^\infty \lambda_T' \lambda_j < \eta, \psi_j >_{\mathcal{H}_{\mathcal{K}}} < g, \psi_j >_{\mathcal{H}_{\mathcal{K}}} \\ &= \sum_{j=1}^\infty \omega_{x_0,j} \lambda_j^2 < \eta, \psi_j >_{\mathcal{H}_{\mathcal{K}}} < g, \psi_j >_{\mathcal{H}_{\mathcal{K}}} + \sum_{j=1}^\infty \lambda_T' \lambda_j < \eta, \psi_j >_{\mathcal{H}_{\mathcal{K}}} < g, \psi_j >_{\mathcal{H}_{\mathcal{K}}} \\ &= \sum_{j=1}^\infty (\omega_{x_0,j} \lambda_j^2 + \lambda_T' \lambda_j) < \eta, \psi_j >_{\mathcal{H}_{\mathcal{K}}} < g, \psi_j >_{\mathcal{H}_{\mathcal{K}}} \\ &= \sum_{j=1}^\infty \lambda_j^2 (\omega_{x_0,j} + \lambda_T' \lambda_j^{-1}) < \eta, \psi_j >_{\mathcal{H}_{\mathcal{K}}} < g, \psi_j >_{\mathcal{H}_{\mathcal{K}}} . \end{split}$$

Let $\eta = \psi_m$ gives

$$\langle G_{\lambda'_{T}}g, \psi_{m} \rangle_{\mathscr{H}_{K}} = \lambda_{m}(\omega_{x_{0},m} + \lambda'_{T}\lambda_{m}^{-1})\langle g, \psi_{m} \rangle_{\mathscr{H}_{K}}$$

which leads to

$$\langle G_{\lambda_T'}^{-1}g,\psi_m\rangle_{\mathscr{H}_{\mathrm{K}}}=\lambda_m^{-1}(\omega_{x_0,m}+\lambda_T^{'}\lambda_m^{-1})^{-1}\langle g,\psi_m\rangle_{\mathscr{H}_{\mathrm{K}}}.$$

Using the expression of the norm and the expression of the operator $G_{\lambda_{\pi}'}$, we have the following:

1113
1114
1115
$$\|\tilde{f}_{n} - \bar{f}_{n}\|^{2} = \frac{1}{4} \int \left(\sum_{j=1}^{\infty} \lambda_{j} < G_{\lambda_{T}^{\prime}}^{-1} Dl_{n\lambda_{T}^{\prime}}(\bar{f}_{x_{0}}), \psi_{j} >_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \right)^{2} \pi(x_{0}) dx_{0}$$
1116
1117
1118
$$= \frac{1}{4} \int \left[\sum_{j=1}^{\infty} \lambda_{j} \lambda_{j}^{-1} (\omega_{x_{0}, j} + \lambda_{T}^{\prime} \lambda_{j}^{-1})^{-1} < Dl_{n\lambda_{T}^{\prime}}(\bar{f}_{x_{0}}), \psi_{j} >_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \right]^{2} \pi(x_{0}) dx_{0}$$
1121
1122
$$= \frac{1}{4} \int \left[\sum_{j=1}^{\infty} (\omega_{x_{0}, j} + \lambda_{T}^{\prime} \lambda_{j}^{-1})^{-1} < Dl_{n\lambda_{T}^{\prime}}(\bar{f}_{x_{0}}), \psi_{j} >_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \right]^{2} \pi(x_{0}) dx_{0}$$
1125
1126
$$= \frac{1}{4} \int \left[\sum_{j=1}^{\infty} \frac{B_{j}(x_{0})}{\sqrt{\lambda_{j}}} D_{j}(x_{0}) \sqrt{\lambda_{j}} \psi_{j}(x_{0}) \right]^{2} \pi(x_{0}) dx_{0}$$
1128
1129
$$\leq \frac{1}{4} \int \left[\sum_{j=1}^{\infty} \lambda_{j}^{-1} B_{j}^{2}(x_{0}) D_{j}^{2}(x_{0}) \right] \left[\sum_{j=1}^{\infty} \lambda_{j} \psi_{j}^{2}(x_{0}) \right] \pi(x_{0}) dx_{0}$$
1131
1132
$$= \frac{1}{4} \int K(x_{0}, x_{0}) \sum_{j=1}^{\infty} \lambda_{j}^{-1} (\omega_{x_{0}, j} + \lambda_{T}^{\prime} \lambda_{j}^{-1})^{-2} < Dl_{n\lambda_{T}^{\prime}}(\bar{f}_{x_{0}}), \psi_{j} >_{\mathscr{H}_{K}}^{\mathscr{H}_{K}}(x_{0}) \pi(x_{0}) dx_{0}.$$

To find the asymptotic order of $\|\tilde{f}_n - \bar{f}_n\|^2$, we need to derive $\omega_{x_0,j}$, the coefficients $\sum_{j=1}^{\infty} \lambda_j^{-1}(\omega_{x_0,j} + \lambda_T' \lambda_j^{-1})^{-2}$ and $E[< Dl_{n\lambda_T'}(\bar{f}_{x_0}), \psi_j>_{\mathscr{H}_{\mathscr{K}}}^2]$.

Derivation for
$$E[< Dl_{n\lambda_T^{'}}(\bar{f}_{x_0}), \psi_j>^2_{\mathscr{H}_{\mathscr{K}}}]$$
.

$$\begin{split} & \text{1139} \\ & \text{1140} \\ & & \mathbb{E}\Big[\Big\langle Dl_{n,\lambda_{T}'}(\bar{f}_{x_{0}}),\psi_{j}\Big\rangle_{\mathscr{H}_{K}}^{2}\Big] \\ & \text{1141} \\ & = \mathbb{E}\Big[\Big\langle Dl_{n}(\bar{f}_{x_{0}}) - Dl_{\infty}(\bar{f}_{x_{0}}),\,\psi_{j}\Big\rangle_{\mathscr{H}_{K}}^{2}\Big]^{2} \\ & \text{1142} \\ & = \mathbb{E}\left\{-\frac{2}{n}\sum_{i\in\mathscr{S}(x_{0})}^{N}\left(y_{i} - \bar{f}_{x_{0}}(x_{i})\right)\psi_{j}(x_{i}) + 2\sum_{i=1}^{N}\omega_{i}(x_{0})\,\mathbb{E}\Big[\big(f_{0}(x_{i}) - \bar{f}_{x_{0}}(x_{i})\big)\psi_{j}(x_{i})\,\big|\,x_{i}\big]\Big\}^{2} \\ & \text{1147} \\ & \text{1148} \\ & = \mathbb{E}\left\{-\frac{2}{n}\sum_{i=1}^{N}I_{i}(x_{0})\,(y_{i} - \bar{f}_{x_{0}}(x_{i}))\psi_{j}(x_{i}) + \frac{2}{n}\,\mathbb{E}\left[\sum_{i=1}^{N}I_{i}(x_{0})\,(f_{0}(x_{i}) - \bar{f}_{x_{0}}(x_{i}))\psi_{j}(x_{i})\,\big|\,x_{i},y_{i}\big]\Big\}^{2} \\ & \text{1150} \\ & \text{1151} \\ & = \frac{4}{n^{2}}\,\mathbb{E}\Big\{\sum_{i=1}^{N}\big(I_{i}(x_{0}) - n\,\omega_{i}(x_{0})\big)\,(y_{i} - \bar{f}_{x_{0}}(x_{i}))\psi_{j}(x_{i})\Big\}^{2} \quad (\text{since}\,\,\mathbb{E}\big[I_{i}(x_{0})\,|\,x_{i},y_{i}\big] = n\omega_{i}(x_{0})) \\ & \text{1153} \\ & = \frac{4}{n^{2}}\,\mathbb{E}\Big\{\sum_{i=1}^{N}\big(I_{i}(x_{0}) - n\,\omega_{i}(x_{0})\big)\,(y_{i} - \bar{f}_{x_{0}}(x_{i}))\psi_{j}(x_{i})\Big\}^{2}\,\big|\,x_{i},y_{i}\big\}\Big] \\ & = \frac{4}{n^{2}}\,\mathbb{E}\Big\{\text{Var}\Big(\sum_{i=1}^{N}I_{i}(x_{0})\,a_{i}\,\big|\,x_{i},y_{i}\big) + \Big(\mathbb{E}\Big[\sum_{i=1}^{N}I_{i}(x_{0})\,a_{i}\,|\,x_{i},y_{i}\big] - n\sum_{i=1}^{N}\omega_{i}(x_{0})a_{i}\Big)^{2}\Big] \\ & = \frac{4}{n^{2}}\,\mathbb{E}\Big[\text{Var}\Big(\sum_{i=1}^{N}I_{i}(x_{0})\,a_{i}\,\big|\,x_{i},y_{i}\big)\Big] \end{aligned}$$

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

$$\operatorname{Var}\left(\sum_{i=1}^{N} I_{i}(x_{0}) a_{i} \mid x_{i}, y_{i}\right) = \sum_{i=1}^{N} a_{i}^{2} \operatorname{Var}(I_{i} \mid x_{i}) + 2 \sum_{1 \leq i < k \leq N} a_{i} a_{k} \operatorname{Cov}(I_{i}, I_{k} \mid x)$$

$$= n \sum_{i=1}^{N} \omega_{i} a_{i}^{2} - n \sum_{i=1}^{N} \sum_{k=1}^{N} \omega_{i} \omega_{k} a_{i} a_{k}$$

$$= n \left(\sum_{i=1}^{N} \omega_{i} a_{i}^{2} - \left(\sum_{i=1}^{N} \omega_{i} a_{i}\right)^{2}\right) \leq n \sum_{i=1}^{N} \omega_{i} a_{i}^{2}.$$

Therefore,

1174
1175
1176
$$\mathbb{E}\Big[\Big\langle Dl_{n,\lambda_{T}'}(\bar{f}_{x_{0}}),\psi_{j}\Big\rangle_{\mathcal{H}_{K}}^{2}\Big] \leq \frac{4}{n}\,\mathbb{E}\Big[\sum_{i=1}^{N}\omega_{i}(x_{0})\,a_{i}^{2}\Big]$$

$$= \frac{4}{n}\,\sum_{i=1}^{N}\omega_{i}(x_{0})\,\mathbb{E}\Big[(y_{i}-\bar{f}_{x_{0}}(x_{i}))^{2}\,\psi_{j}(x_{i})^{2}\Big]$$
1180
$$= \frac{4}{n}\,\sum_{i=1}^{N}\omega_{i}(x_{0})\,\mathbb{E}\Big[\big((f_{0}(x_{i})-\bar{f}_{x_{0}}(x_{i}))+\epsilon_{i}\big)^{2}\,\psi_{j}(x_{i})^{2}\Big]$$
1182
$$= \frac{4}{n}\,\sum_{i=1}^{N}\omega_{i}(x_{0})\,\mathbb{E}\Big[\big((f_{0}(x_{i})-\bar{f}_{x_{0}}(x_{i}))^{2}+\sigma^{2}\big)\,\psi_{j}(x_{i})^{2}\Big]$$
1184
1185
1186
1187
$$\leq \frac{4}{n}\,\left\{\sigma^{2}\,\sum_{i=1}^{N}\omega_{i}(x_{0})\,\mathbb{E}\big[\psi_{j}(x_{i})^{2}\big]+\sum_{i=1}^{N}\omega_{i}(x_{0})\,\mathbb{E}\big[(f_{0}(x_{i})-\bar{f}_{x_{0}}(x_{i}))^{2}\,\psi_{j}(x_{i})^{2}\big]\right\} \approx \frac{1}{n}.$$

Derivation for the coefficients $\sum_{j=1}^{\infty} \lambda_j^{-1} (\omega_{x_0,j} + \lambda_T' \lambda_j^{-1})^{-2}$. Assume $\lambda_j \approx j^{-k}$ with k > 1, $a_j^2 \lambda_j^{-1} = j^{-a}$ with a > 1, and there exists $\beta \in [0, k/2)$ such that, $\omega_{x_0,j}$ uniformly in x_0 , then we assume

$$\omega_{x_0,j} \asymp j^{2\beta}$$
, equivalently, $\lambda_j \omega_{x_0,j} \asymp j^{2\beta-k}$.

Define

$$S(x_0) := \sum_{j=1}^{\infty} \lambda_j^{-1} (\omega_{x_0,j} + \lambda_T' \lambda_j^{-1})^{-2}.$$

Then

$$S(x_0) \approx \sum_{j=1}^{\infty} \frac{j^k}{\left(j^{2\beta} + \lambda'_T j^k\right)^2} = \sum_{j=1}^{\infty} \frac{j^{-k}}{\left(j^{-(k-2\beta)} + \lambda'_T\right)^2}.$$

Let J solve $j^{-(k-2\beta)} \simeq \lambda_T'$, so we have $J \simeq (\lambda_T')^{-1/(k-2\beta)}$. Then

$$\sum_{j=1}^{J} \frac{j^{-k}}{\left(j^{-(k-2\beta)} + \lambda_T'\right)^2} \, \asymp \, \sum_{j=1}^{J} j^{k-4\beta} \, \asymp \, J^{k-4\beta+1} \, \asymp \, \lambda_T'^{-\frac{k+1-4\beta}{k-2\beta}},$$

$$\sum_{j=J+1}^{\infty} \frac{j^{-k}}{\left(j^{-(k-2\beta)} + \lambda_T'\right)^2} \; \asymp \; \frac{1}{(\lambda_T')^2} \sum_{j=J+1}^{\infty} j^{-k} \; \asymp \; \lambda_T'^{-2} \, J^{-(k-1)} \; \asymp \; \lambda_T'^{-\frac{k+1-4\beta}{k-2\beta}}.$$

Therefore,

$$S(x_0) = O\left(\lambda_T^{\prime - \frac{k+1-4\beta}{k-2\beta}}\right), \qquad 0 \le 2\beta < k.$$

Using this bound and $\mathbb{E}[\langle Dl_{n,\lambda_T'}(\bar{f}_{x_0}),\psi_j\rangle_{\mathscr{H}_K}^2] \approx n^{-1}$, we obtain

$$\|\tilde{f}_{n} - \bar{f}_{n}\|^{2} \leq \frac{1}{4} \int K(x_{0}, x_{0}) \sum_{j=1}^{\infty} \lambda_{j}^{-1} (\omega_{x_{0}, j} + \lambda_{T}' \lambda_{j}^{-1})^{-2} \mathbb{E}[\langle Dl_{n, \lambda_{T}'}(\bar{f}_{x_{0}}), \psi_{j} \rangle_{\mathscr{H}_{K}}^{2}] \pi(x_{0}) dx_{0}$$

$$\approx \left(\int K(x_{0}, x_{0}) \pi(x_{0}) dx_{0} \right) \lambda_{T}'^{-\frac{k+1-4\beta}{k-2\beta}} \frac{1}{n}$$

$$= O_{p} \left(n^{-1} \lambda_{T}'^{-\frac{k+1-4\beta}{k-2\beta}} \right).$$

B.1.2 Step 2: Evaluate the order of $\hat{f}_{n\lambda_T'} - \tilde{f}_n$.

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

$$Dl_{n\lambda_T}(\hat{f}_{x_0}) = Dl_{n\lambda_T}(\bar{f}_{x_0}) + D^2l_{n\lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) = 0.$$
(4)

Using the functional derivatives derived above, we have

$$Dl_{n\lambda_T} = -\frac{2}{n} \sum_{i \in \mathscr{S}(x_0)} K(x_i, \cdot) \left(y_i - \langle \bar{f}_{x_0}, K(x_i, \cdot) \rangle_{\mathscr{H}_K} \right) + 2\lambda_T \langle \bar{f}_{x_0}, \cdot \rangle_{\mathscr{H}_K}$$
(5)

 $D^{2}l_{n\lambda_{T}}(\bar{f}_{x_{0}})(\hat{f}_{x_{0}} - \bar{f}_{x_{0}}) = \frac{2}{n} \sum_{i \in \mathscr{S}(x_{0})} \langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, K(x_{i}, \cdot) \rangle_{\mathscr{H}_{K}} K(x_{i}, \cdot) + 2\lambda_{T} \langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, \cdot \rangle_{\mathscr{H}_{K}}$ (6)

Adding (5) and (6) together, we obtain:

$$(5) + (6) = -\frac{2}{n} \sum_{i \in \mathscr{S}(x_0)} K(x_i, \cdot) (y_i - \langle \hat{f}_{x_0}, K(x_i, \cdot) \rangle_{\mathscr{H}_K}) + 2\lambda_T \langle \hat{f}_{x_0}, \cdot \rangle_{\mathscr{H}_K}$$
$$= Dl_{n\lambda_T}(\hat{f}_{x_0}) = 0.$$

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

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

Combining with the equation: $Dl_{n\lambda_T}(\hat{f}_{x_0}) = Dl_{n\lambda_T}(\bar{f}_{x_0}) + D^2l_{n\lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) = 0$, we get $D^2l_{\infty\lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) = D^2l_{n\lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0})$. Then, we can derive:

$$\begin{split} D^2 l_{\infty \lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \tilde{f}_{x_0}) &= D^2 l_{\infty \lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) + D^2 l_{\infty \lambda_T}(\bar{f}_{x_0})(\bar{f}_{x_0} - \tilde{f}_{x_0}) \\ &= D^2 l_{\infty \lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) - D^2 l_{n\lambda_T}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) \\ &= D^2 l_{\infty}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}) - D^2 l_{n}(\bar{f}_{x_0})(\hat{f}_{x_0} - \bar{f}_{x_0}). \end{split}$$

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

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

So we write:

$$\begin{aligned} &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n,\lambda_{T}'}\|^{2} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n,\lambda_{T}'}\|^{2} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n,\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}}, K(x_{0},\cdot)\rangle_{\mathscr{H}_{K}}^{2} \pi(x_{0}) \, dx_{0} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}}, \psi_{j}\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \Big)^{2} \pi(x_{0}) \, dx_{0} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}}, \psi_{j}\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \Big)^{2} \pi(x_{0}) \, dx_{0} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}}, \psi_{j}\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \Big)^{2} \pi(x_{0}) \, dx_{0} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}}, \psi_{j}\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \Big)^{2} \pi(x_{0}) \, dx_{0} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}}, \psi_{j}\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \Big)^{2} \pi(x_{0}) \, dx_{0} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}_{n\lambda_{T}',x_{0}}, \psi_{j}\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \Big)^{2} \pi(x_{0}) \, dx_{0} \\ &\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n\lambda_{T}',x_{0}} - \tilde{f}$$

where

$$A_{j}(x_{0}) = \left\langle D^{2}l_{\infty,\lambda'_{T}}(\bar{f}_{x_{0}})(\hat{f}_{x_{0}} - \bar{f}_{x_{0}}), \psi_{j} \right\rangle_{\mathscr{H}_{K}} - \left\langle D^{2}l_{n,\lambda'_{T}}(\bar{f}_{x_{0}})(\hat{f}_{x_{0}} - \bar{f}_{x_{0}}), \psi_{j} \right\rangle_{\mathscr{H}_{K}}$$

$$= 2 \sum_{\ell=1}^{\infty} \lambda_{\ell}^{2} \left\langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, \psi_{\ell} \right\rangle_{\mathscr{H}_{K}} \left\langle \psi_{j}, \psi_{\ell} \right\rangle_{\mathscr{H}_{K}} \mathbb{E} \left[\sum_{i=1}^{N} \omega_{i}(x_{0}) \psi_{\ell}^{2}(x_{i}) \right]$$

$$- \frac{2}{n} \sum_{i \in \mathscr{S}(x_{0})} \sum_{k=1}^{\infty} \lambda_{k}^{2} \left\langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, \psi_{k} \right\rangle_{\mathscr{H}_{K}} \left\langle \psi_{j}, \psi_{k} \right\rangle_{\mathscr{H}_{K}} \psi_{k}^{2}(x_{i})$$

$$= 2 \lambda_{j}^{2} \frac{1}{\lambda_{j}} \left\langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, \psi_{j} \right\rangle_{\mathscr{H}_{K}} \mathbb{E} \left[\sum_{i=1}^{N} \omega_{i}(x_{0}) \psi_{j}^{2}(x_{i}) \right]$$

$$- \frac{2}{n} \sum_{i \in \mathscr{S}(x_{0})} \lambda_{j}^{2} \frac{1}{\lambda_{j}} \left\langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, \psi_{j} \right\rangle_{\mathscr{H}_{K}} \psi_{j}^{2}(x_{i})$$

$$= 2 \lambda_{j} \left\langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, \psi_{j} \right\rangle_{\mathscr{H}_{K}} \left\{ \mathbb{E} \left[\sum_{i=1}^{N} \omega_{i}(x_{0}) \psi_{j}^{2}(x_{i}) \right] - \frac{1}{n} \sum_{i \in \mathscr{S}(x_{0})} \psi_{j}^{2}(x_{i}) \right\}.$$

plug it back into 2.1

$$\frac{1298}{1299} = \frac{1}{4} \int \left[\sum_{j=1}^{\infty} \left(\omega_{x_0,j} + \lambda_T' \lambda_j^{-1} \right)^{-1} \left(\left(2 \lambda_j \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K} \right) \left(\omega_{x_0,j} - \frac{1}{n} \sum_{i \in \mathscr{S}(x_0)} \psi_j^2(x_i) \right) \psi_j(x_0) \right) \right]^2 \pi(x_0) \, dx_0$$

$$\frac{1301}{1302} = \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} \lambda_j \left(\omega_{x_0,j} + \lambda_T' \lambda_j^{-1} \right)^{-1} \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K} \psi_j(x_0) \left(n \omega_{x_0,j} - \sum_{i \in \mathscr{S}(x_0)} \psi_j^2(x_i) \right) \right]^2 \pi(x_0) \, dx_0$$

$$\frac{1}{1305} = \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} \lambda_j \left(\omega_{x_0,j} + \lambda_T' \lambda_j^{-1} \right)^{-1} \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K} \psi_j(x_0) \right] \\
\frac{1}{1308} = \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} \lambda_j \left(\omega_{x_0,j} + \lambda_T' \lambda_j^{-1} \right)^{-1} \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K} \psi_j(x_0) \right] \\
\times \left(n \mathbb{E} \left[\sum_{i=1}^{N} \omega_i(x_0) \psi_j^2(x_i) \right] - \sum_{i=1}^{N} I_i(x_0) \psi_j^2(x_i) \right]^2 \pi(x_0) \, dx_0. \tag{2.2}$$

Set

$$\Delta_j(x_0) := n \, \mathbb{E} \Big[\sum_{i=1}^N \omega_i(x_0) \, \psi_j^2(x_i) \Big] - \sum_{i=1}^N I_i(x_0) \, \psi_j^2(x_i).$$

Since $\mathbb{E}[I_i(x_0) \mid x_1, ..., x_N] = n \, \omega_i(x_0)$ 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

$$\mathbb{E}[\Delta_j(x_0)^2] = \mathbb{E}\Big[\operatorname{Var}\Big(\sum_{i=1}^N I_i(x_0)\,\psi_j^2(x_i)\,\Big|\,x_1,\dots,x_N\Big)\Big]$$

$$= \mathbb{E}\Big[n\Big(\sum_{i=1}^N \omega_i(x_0)\,\psi_j^4(x_i) - \Big(\sum_{i=1}^N \omega_i(x_0)\,a_i\Big)^2\Big)\Big]$$

$$\leq \mathbb{E}\Big[n\sum_{i=1}^N \omega_i(x_0)\,\psi_j^4(x_i)\Big]$$

$$= n\,\mathbb{E}\Big[\sum_{i=1}^N \omega_i(x_0)\,\psi_j^4(x_i)\Big] = n\,\mathbb{E}[\psi_j^4(x)].$$

Therefore,

$$\Delta_j(x_0) \simeq O_p \sqrt{n \, \mathbb{E}[\psi_j^4(x)]}.$$

Plug it back to (2.2):

$$\approx \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} \lambda_j \left(\omega_{x_0,j} + \lambda_T' \lambda_j^{-1} \right)^{-1} \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K} \psi_j(x_0) \right]^2 \pi(x_0) \, dx_0$$

$$= \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} (\omega_{x_0,j} + \lambda_T' \lambda_j^{-1})^{-1} \left(\lambda_j \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K} \psi_j(x_0) \right) \right]^2 \pi(x_0) \, dx_0$$

$$\leq \frac{1}{n^2} \int \left[\sum_{j=1}^{\infty} (\omega_{x_0,j} + \lambda_T' \lambda_j^{-1})^{-2} \right] \left[\sum_{j=1}^{\infty} \lambda_j^2 \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K}^2 \psi_j(x_0)^2 \right] \pi(x_0) \, dx_0$$

$$\leq \frac{1}{n^2} \int \underbrace{\left[\sum_{j=1}^{\infty} (\omega_{x_0,j} + \lambda_T' \lambda_j^{-1})^{-2} \right]}_{S_1(x_0)} \left[\sum_{j=1}^{\infty} \lambda_j \left\langle \hat{f}_{x_0} - \bar{f}_{x_0}, \psi_j \right\rangle_{\mathscr{H}_K} \psi_j(x_0) \right]^2 \pi(x_0) \, dx_0.$$

Assume $\lambda_j \asymp j^{-k}$ with k>1 and, uniformly in $x_0, \ \omega_{x_0,j} \asymp j^{2\beta}$ for some $\beta \in [0,k/2)$. Recall

$$S_1(x_0) := \sum_{j=1}^{\infty} \left(\omega_{x_0,j} + \lambda_T' \lambda_j^{-1} \right)^{-2} \, symp \, \sum_{j=1}^{\infty} \frac{1}{\left(j^{2\beta} + \lambda_T' j^k \right)^2}.$$

1356 Let J solve $j^{2\beta}=\lambda_T'j^k$, then $J\asymp (\lambda_T')^{-1/(k-2\beta)}$. Then

$$\sum_{j=1}^{J} \frac{1}{(j^{2\beta})^2} \; = \; \sum_{j=1}^{J} j^{-4\beta} \; \asymp J^{1-4\beta} = \lambda_T^{\prime \frac{4\beta-1}{k-2\beta}}, \qquad \sum_{j>J} \frac{1}{(\lambda_T^{\prime} j^k)^2} \; = \; \lambda_T^{\prime -2} \sum_{j>J} j^{-2k} \; \asymp \; \lambda_T^{\prime \frac{4\beta-1}{k-2\beta}}.$$

Hence

$$S_1(x_0) = O\left(\lambda_T'^{-\frac{1-4\beta}{k-2\beta}}\right).$$

Consequently,

$$\|\hat{f}_{n,\lambda_{T}'} - \tilde{f}_{n,\lambda_{T}'}\|^{2} \leq \frac{1}{n^{2}} \int \underbrace{\sum_{j=1}^{\infty} \left(\omega_{x_{0},j} + \lambda_{T}' \lambda_{j}^{-1}\right)^{-2}}_{=S_{1}(x_{0})} \left[\sum_{j=1}^{\infty} \lambda_{j} \left\langle \hat{f}_{x_{0}} - \bar{f}_{x_{0}}, \psi_{j} \right\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \right]^{2} \pi(x_{0}) dx_{0}$$

$$= \frac{S_{1}(x_{0})}{n^{2}} \sum_{j=1}^{\infty} \lambda_{j}^{2} \left\langle \hat{f}_{n,\lambda_{T}'} - \bar{f}_{n,\lambda_{T}'}, \psi_{j} \right\rangle_{\mathscr{H}_{K}}^{2}$$

$$= \frac{S_{1}(x_{0})}{n^{2}} \|\hat{f}_{n,\lambda_{T}'} - \bar{f}_{n,\lambda_{T}'}\|^{2}.$$

Therefore,

$$\|\hat{f}_{n,\lambda_T'} - \tilde{f}_{n,\lambda_T'}\|^2 = O_p \left(n^{-2} \lambda_T'^{-\frac{1-4\beta}{k-2\beta}} \|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2\right).$$

By the triangle inequality,

$$\|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2 \le \|\tilde{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2 + \|\hat{f}_{n,\lambda_T'} - \tilde{f}_{n,\lambda_T'}\|^2.$$

From Step 1 we have

$$\|\tilde{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2 = O_p\left(n^{-1}\lambda_T'^{-\frac{k+1-4\beta}{k-2\beta}}\right).$$

From Step 2 we obtained

$$\|\hat{f}_{n,\lambda_T'} - \tilde{f}_{n,\lambda_T'}\|^2 = O_p\left(n^{-2}\lambda_T'^{-\frac{1-4\beta}{k-2\beta}}\|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2\right).$$

If

$$n^{-2}\lambda_T^{\prime - \frac{1-4\beta}{k-2\beta}} \longrightarrow 0,$$

then

$$\|\hat{f}_{n,\lambda_T'} - \tilde{f}_{n,\lambda_T'}\|^2 = o_p(\|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2) = o_p(1) \|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2.$$

Observed that

$$\|\tilde{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2 \ge \|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2 - \|\hat{f}_{n,\lambda_T'} - \tilde{f}_{n,\lambda_T'}\|^2 = (1 - o_p(1)) \|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2,$$

we conclude

$$\|\hat{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2 = O_p(\|\tilde{f}_{n,\lambda_T'} - \bar{f}_{n,\lambda_T'}\|^2) = O_p(n^{-1}\lambda_T'^{-\frac{k+1-4\beta}{k-2\beta}}).$$

B.2 ASYMPTOTIC ORDER FOR DETERMINISTIC ERROR.

Recall

$$\begin{split} l_{\infty}(f) &:= \mathbb{E}\big[l_{n}(f)\big] = \sigma^{2} \, + \, \mathbb{E}\left[\sum_{i=1}^{N} \omega_{i}(x_{0}) \left\langle f - f_{0}, \, K(x_{i}, \cdot) \right\rangle_{\mathcal{H}_{K}}^{2}\right] \\ &= \sigma^{2} + \sum_{i=1}^{N} \mathbb{E}\left[\omega_{i}(x_{0}) \left(f(x_{i}) - f_{0}(x_{i})\right)^{2}\right] \\ &= \sigma^{2} + \sum_{i=1}^{N} \mathbb{E}\left[\omega_{i}(x_{0}) \left(\sum_{j=1}^{\infty} (c_{j} - a_{j}) \psi_{j}(x_{i})\right)^{2}\right] \\ &= \sigma^{2} + \sum_{j,\ell \geq 1} (c_{j} - a_{j}) (c_{\ell} - a_{\ell}) \sum_{i=1}^{N} \mathbb{E}\left[\omega_{i}(x_{0}) \psi_{j}(x_{i}) \psi_{\ell}(x_{i})\right] \\ &= \sigma^{2} + \sum_{j=1}^{\infty} \omega_{x_{0}, j} (c_{j} - a_{j})^{2}, \end{split}$$

where $\omega_{x_0,j} := \mathbb{E}\left[\sum_{i=1}^N \omega_i(x_0) \psi_j^2(x_i)\right]$.

The objective function about c_i at x_0 becomes

$$Q_{x_0}(c) = \sum_{j>1} \left\{ \omega_{x_0,j} (c_j - a_j)^2 + \lambda_T' \frac{c_j^2}{\lambda_j} \right\}.$$

Taking derivative,

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

hence

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

Recall

$$\omega_{x_0,j} := \mathbb{E}\Big[\sum_{i=1}^N \omega_i(x_0)\psi_j^2(x_i)\Big] = \psi_j^2(x_0) + \Delta_{h,L}(x_0).$$

Assume $\lambda_j \asymp j^{-k}$ with k>1, $a_j^2\lambda_j^{-1}=j^{-a}$ with a>1, and there exists $\beta\in[0,k/2)$ such that, uniformly in x_0 ,

$$\omega_{x_0,j} \asymp j^{2\beta}$$
 (equivalently, $\lambda_j \, \omega_{x_0,j} \asymp j^{2\beta-k}$).

Now we evaluate $\|\bar{f}_{x_0} - f_0\|^2$:

$$\|\bar{f}_{n\lambda'_{T}} - f_{0}\|^{2} := \int \left\langle \bar{f}_{n\lambda'_{T},x_{0}} - f_{0}, K(x_{0},\cdot) \right\rangle_{\mathscr{H}_{K}}^{2} \pi(x_{0}) dx_{0}$$

$$= \int \left(\sum_{j=1}^{\infty} \lambda_{j} \left\langle \bar{f}_{n\lambda'_{T},x_{0}} - f_{0}, \psi_{j} \right\rangle_{\mathscr{H}_{K}} \psi_{j}(x_{0}) \right)^{2} \pi(x_{0}) dx_{0}$$

$$= \int \left(\sum_{j=1}^{\infty} \left(\bar{c}_{j}(x_{0}) - a_{j} \right) \psi_{j}(x_{0}) \right)^{2} \pi(x_{0}) dx_{0}$$

$$= \int \left(\sum_{j=1}^{\infty} \frac{\lambda'_{T}}{\omega_{x_{0},j}\lambda_{j} + \lambda'_{T}} a_{j} \psi_{j}(x_{0}) \right)^{2} \pi(x_{0}) dx_{0}$$

$$= \int \left(\sum_{j=1}^{\infty} \frac{\lambda'_{T}}{j^{2\beta-k} + \lambda'_{T}} a_{j} \psi_{j}(x_{0}) \right)^{2} \pi(x_{0}) dx_{0}$$

1458
1459
$$\asymp \left(\sum_{j=1}^{\infty} \frac{\lambda_T'}{j^{2\beta-k} + \lambda_T'} a_j\right)^2$$

1462 Set $J \simeq \lambda_T'^{-1/(k-2\beta)}$ so that $J^{2\beta-k} \simeq \lambda_T'$. Then

$$\sum_{j \leq J} \frac{\lambda_T'^2 \, j^{-(a+k)}}{(j^{\, 2\beta - k})^2} \, \asymp \, \lambda_T'^2 \sum_{j \leq J} j^{-(a+4\beta - k)} \, \asymp \, \lambda_T'^2 \, J^{\, 1 - (a+4\beta - k)},$$

$$\sum_{j>J} \frac{\lambda_T'^2 \, j^{-(a+k)}}{(\lambda_T')^2} = \sum_{j>J} j^{-(a+k)} \; \asymp \; J^{-(a+k-1)}.$$

With $J \asymp \lambda_T'^{-1/(k-2\beta)}$ both parts scale as $\lambda_T'^{\frac{a+k-1}{k-2\beta}}$.

$$\|\bar{f}_{n\lambda_T} - f_0\|^2 \simeq \lambda_T'^{\frac{a+k-1}{k-2\beta}}, \qquad 0 < 2\beta < k.$$

B.3 ASYMPTOTIC ORDER FOR MSE AND TUNING.

$$\|\hat{f}_{n,\lambda_T'} - f_0\|^2 = O_p \left(\lambda_T'^{\frac{a+k-1}{k-2\beta}} + n^{-1} \lambda_T'^{-\frac{k+1-4\beta}{k-2\beta}} \right),$$

under the condition

$$n^{-2}\lambda_T^{\prime-\frac{1-4\beta}{k-2\beta}}\longrightarrow 0,$$

with
$$a>1, \; k>1, \; 0 \leq \beta < k/2$$
 .

Let

$$M_{\beta}(\lambda_T') = \lambda_T'^p + n^{-1} \lambda_T'^{-q}, \qquad p := \frac{a+k-1}{k-2\beta}, \ q := \frac{k+1-4\beta}{k-2\beta}.$$

Differentiate and set to zero:

$$\frac{dM_{\beta}}{d\lambda_{T}'} \; = \; p \, \lambda_{T}'^{p-1} \; - \; q \, n^{-1} \, \lambda_{T}'^{-q-1} \; = \; 0 \; \Longrightarrow \; \lambda_{T}'^{p+q} \asymp n^{-1}.$$

Hence the optimal order of tuning λ_T' is 1493

$$\lambda_T' \approx n^{-\frac{1}{p+q}} = n^{-\frac{k-2\beta}{a+2k-4\beta}}.$$

1496 At this choice,

$$\lambda_T'^p \ \asymp \ n^{-\frac{p}{p+q}} = n^{-\frac{a+k-1}{a+2k-4\beta}}, \qquad n^{-1}\lambda_T'^{-q} \ \asymp \ n^{-\frac{p}{p+q}} = n^{-\frac{a+k-1}{a+2k-4\beta}},$$

so

$$\|\hat{f}_{n,\lambda_T'} - f_0\|^2 = O_p\left(n^{-\frac{a+k-1}{a+2k-4\beta}}\right).$$

Additionally, the small-o condition from Step 2 holds at the optimized tuning:

$$n^{-2}\lambda_T^{\prime - \frac{1-4\beta}{k-2\beta}} \approx n^{-2+\frac{1-4\beta}{a+2k-4\beta}} \to 0 \quad \text{since } a > 1, \ k > 1, \ 0 \le \beta < \frac{k}{2}.$$

This completes the proof of Theorem 2.