
Supplementary Material of ‘‘Sparse Shrunk Additive Models’’

Guodong Liu ^{*1} Hong Chen ^{*2} Heng Huang ¹

A. Proofs of Theorems 1 and 2

A.0. Error decomposition

The proofs of Theorems 1 and 2 involve a integration of techniques for error analysis with integral operator approximation (Smale & Zhou, 2007; Sun & Wu, 2011; Shi, 2013; Nie & Wang, 2015) and the empirical process theory for analyzing kernel methods (Pinelis, 1994; Wu et al., 2007; Christmann & Zhou, 2016). The proof of Theorem 3 follows the analysis technique for sparse characterization (Shi et al., 2011).

The key to bound $\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_{\rho})$ is a novel error decomposition, where some intermediate functions are constructed as the stepping stone functions. Then, we bound the decomposed terms respectively in terms of operator approximation and concentration equalities for empirical processes.

From Proposition 1 in (Shi, 2013), we know that $L_{K^{(j)}}^T = UL_{\tilde{K}^{(j)}}^{\frac{1}{2}}$ and $L_{K^{(j)}} = L_{\tilde{K}^{(j)}}^{\frac{1}{2}}U^T$ for each $j \in \{1, 2, \dots, d\}$, where U is a partial isometry on $L_{\rho, \mathcal{X}^{(j)}}^2$ with U^TU being the orthogonal prediction onto the RKHS $\mathcal{H}_{\tilde{K}^{(j)}}$.

For any $j \in \{1, 2, \dots, d\}$, define the intermediate function $f_{\lambda}^{(j)}$ by

$$f_{\lambda}^{(j)} = \arg \min_{f \in L_{\rho, \mathcal{X}^{(j)}}^2} \left\{ \|L_{K^{(j)}} f^{(j)} - f_{\rho}^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}^2 + \lambda \|U^T f^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}^2 \right\}. \quad (1)$$

Denote $f_{\lambda} = \sum_{j=1}^d f_{\lambda}^{(j)}$ and $g_{\lambda} = \sum_{j=1}^d g_{\lambda}^{(j)}$ with $g_{\lambda}^{(j)} = L_{K^{(j)}} f_{\lambda}^{(j)}$.

^{*}Equal contribution ¹Department of Electrical and Computer Engineering, University of Pittsburgh, Pittsburgh, United States ²Department of Mathematics and Statistics, College of Science, Huazhong Agricultural University, Wuhan, China. Correspondence to: Hong Chen <chenh@mail.hzau.edu.cn>, Guodong Liu <guodong.liu.e@pitt.edu>.

Define the empirical version of g_{λ} as

$$\hat{g}_{\lambda}(x) = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^d f_{\lambda}^{(j)}(x_i^{(j)}) K^{(j)}(x_i^{(j)}, x^{(j)}), x \in \mathcal{X}. \quad (2)$$

Now we give the following error decomposition.

Proposition 1. For $f_{\mathbf{z}}, \hat{g}_{\lambda}$, there holds

$$\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_{\rho}) \leq E_1 + E_2 + E_3,$$

where

$$\begin{aligned} E_1 &= \mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}_{\mathbf{z}}(\pi(f_{\mathbf{z}})) + \mathcal{E}_{\mathbf{z}}(\hat{g}_{\lambda}) - \mathcal{E}(\hat{g}_{\lambda}), \\ E_2 &= \mathcal{E}(\hat{g}_{\lambda}) - \mathcal{E}(g_{\lambda}) + \lambda \|\hat{g}_{\lambda}\|_{\ell_1}, \end{aligned}$$

and

$$E_3 = \mathcal{E}(g_{\lambda}) - \mathcal{E}(f_{\rho}).$$

Proof. According the definition of $f_{\mathbf{z}}$, we have

$$\begin{aligned} &\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_{\rho}) \\ &\leq \mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}_{\mathbf{z}}(\pi(f_{\mathbf{z}})) + \mathcal{E}_{\mathbf{z}}(\hat{g}_{\lambda}) - \mathcal{E}(f_{\rho}) + \lambda \|\hat{g}_{\lambda}\|_{\ell_1} \\ &\quad + \left\{ \mathcal{E}_{\mathbf{z}}(f_{\mathbf{z}}) + \lambda \|f_{\mathbf{z}}\|_{\ell_1} - (\mathcal{E}_{\mathbf{z}}(\hat{g}_{\lambda}) + \lambda \|\hat{g}_{\lambda}\|_{\ell_1}) \right\} \\ &\leq \mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}_{\mathbf{z}}(\pi(f_{\mathbf{z}})) + \mathcal{E}_{\mathbf{z}}(\hat{g}_{\lambda}) - \mathcal{E}(f_{\rho}) + \lambda \|\hat{g}_{\lambda}\|_{\ell_1} \end{aligned} \quad (3)$$

Note that

$$\begin{aligned} \mathcal{E}_{\mathbf{z}}(\hat{g}_{\lambda}) - \mathcal{E}(f_{\rho}) &= (\mathcal{E}(\hat{g}_{\lambda}) - \mathcal{E}(g_{\lambda})) + \mathcal{E}(g_{\lambda}) - \mathcal{E}(f_{\rho}) \\ &\quad + \mathcal{E}_{\mathbf{z}}(\hat{g}_{\lambda}) - \mathcal{E}(\hat{g}_{\lambda}). \end{aligned} \quad (4)$$

Combining both (3) and (4), we get the desired decomposition.

The error term E_1 measures the divergence between the empirical risk and the corresponding expected risk, which usually is called sample error in learning theory. In terms of recent theoretical progress for learning with data dependent hypothesis spaces (Shi et al., 2011; Shi, 2013; Feng et al., 2016), we can bound sample error E_1 via concentration inequality associated with empirical covering numbers (Wu et al., 2007; Christmann & Zhou, 2016). The error term E_2 reflects the drift risk for learning with hypothesis spaces $\mathcal{H}_{\mathbf{z}}$ and \mathcal{H} , and hence is called as the hypothesis error.

By relating $\mathcal{E}(\hat{g}_\lambda) - \mathcal{E}(g_\lambda)$ with $\sum_{j=1}^d \|\hat{g}_\lambda^{(j)} - g_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}$, we can estimate this hypothesis error through the inequality in Hilbert space (Pinelis, 1994; Smale & Zhou, 2007). The error term E_2 is called the approximation error, which describes the approximation ability of regularized scheme. Following the approximation analysis with integral operator in (Smale & Zhou, 2007; Shi, 2013; Nie & Wang, 2015), we derive the upper bound of E_2 based on the properties of $L_{\tilde{K}^{(j)}}, 1 \leq j \leq d$.

A.1. Estimate of Approximation Error E_3

In this paper, we use the analysis techniques in (Smale & Zhou, 2007; Shi, 2013) to bound the approximation error E_3 .

The following lemma is used in our analysis, which is proved in Proposition 2 in (Shi, 2013).

Lemma 1. *From the definition of $f_\lambda^{(j)}$ and $g_\lambda^{(j)} = L_{K^{(j)}} f_\lambda^{(j)}, j \in \{1, 2, \dots, d\}$, there are*

$$f_\lambda^{(j)} = U(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^{\frac{1}{2}} f_\rho^{(j)}$$

and

$$\|f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 = \|U^T f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2.$$

Lemma 2. *Under Assumption 1, there holds*

$$\begin{aligned} & \|L_{\tilde{K}^{(j)}} f_\lambda^{(j)} - f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 + \lambda \|f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & \leq \lambda^{\min\{1, 2r\}} \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 (2 + \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 + \|L_{\tilde{K}^{(j)}}^{r-1}\|^2). \end{aligned}$$

Proof. Recall that $\{\lambda_i^{(j)}, \psi_i^{(j)}\}_{i \geq 1}$ are the normalized eigenpairs of the integral operator $L_{\tilde{K}^{(j)}}$ and $\{\psi_i^{(j)}\}_{i \geq 1}$ form an orthogonal basis of $L^2_{\rho_{\mathcal{X}^{(j)}}}$. Let $g_\rho^{(j)} = L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)} = \sum_{t=1}^\infty a_t \psi_t^{(j)}$. Then $\|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 = \sum_{t=1}^\infty (a_t^{(j)})^2 < \infty$.

If Assumption 1 holds for some $r \in (0, \frac{1}{2})$, then from Lemma 1 we have

$$\begin{aligned} & \|f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 = \|U^T f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \|U^T U(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^{\frac{1}{2}} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^{\frac{1}{2}} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^{\frac{1}{2}+r} L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \end{aligned}$$

Moreover,

$$\begin{aligned} & \lambda \|f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \lambda \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^{\frac{1}{2}+r} L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \lambda \left\| \sum_{t \geq 1} \frac{(\lambda_t^{(j)})^{\frac{1}{2}+r}}{\lambda_t^{(j)} + \lambda} a_t^{(j)} \psi_t^{(j)} \right\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \lambda \sum_{t \geq 1} \frac{\lambda_t^{(j)}}{\lambda_t^{(j)} + \lambda} \cdot \frac{\lambda_t^{2r}}{\lambda_t + \lambda} (a_t^{(j)})^2 \\ & \leq \lambda^{2r} \sum_{t \geq 1} \frac{\lambda_t^{(j)}}{\lambda_t^{(j)} + \lambda} (a_t^{(j)})^2 \\ & \leq \lambda^{2r} \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2, \end{aligned} \tag{5}$$

where the first inequality follows from Lemma 1 in (Nie & Wang, 2015) and the second inequality is obtained based on the definition of $g_\rho^{(j)}$.

If Assumption 1 is true for some $r \geq \frac{1}{2}$,

$$\begin{aligned} & \lambda \|f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \lambda \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}} L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & \leq \lambda \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 \cdot \|L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & \leq \lambda \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 \cdot \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2. \end{aligned} \tag{6}$$

Now turn to bound $\|g_\lambda^{(j)} - f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2$. From Lemma 1, we can deduce that

$$\begin{aligned} g_\lambda^{(j)} & = L_{K^{(j)}} f_\lambda^{(j)} = L_{\tilde{K}^{(j)}} (\lambda I + L_{\tilde{K}^{(j)}})^{-1} f_\rho^{(j)} \\ & = (\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}} f_\rho^{(j)} \end{aligned}$$

and

$$\|g_\lambda^{(j)} - f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 = \lambda^2 \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2.$$

For $r \in (0, 1)$, we have

$$\begin{aligned} & \lambda \|g_\lambda^{(j)} - f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & = \lambda^2 \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^r L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \\ & \leq \lambda^2 \sum_{t \geq 1} (a_t^{(j)})^2 \left(\frac{(\lambda_t^{(j)})^r}{\lambda_t^{(j)} + \lambda} \right)^2 \\ & \leq \lambda^{2r} \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2. \end{aligned} \tag{7}$$

For $r \geq 1$, we get

$$\begin{aligned}
 & \lambda \|g_\lambda^{(j)} - f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}^2 \\
 &= \lambda^2 \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}} L_{\tilde{K}^{(j)}}^{r-1} L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}^2 \\
 &\leq \lambda^2 \|L_{\tilde{K}^{(j)}}^{r-1}\|^2 \|L_{\tilde{K}^{(j)}}^{-r} f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}^2 \\
 &\leq \lambda^2 \|L_{\tilde{K}^{(j)}}^{r-1}\|^2 \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}^2. \tag{8}
 \end{aligned}$$

Combining (5)-(8), we get the desired result. \square

Lemma 3. For $j \in \{1, 2, \dots, d\}$ and $g_\lambda^{(j)} = L_{\tilde{K}^{(j)}} f_\lambda^{(j)}$ with $f_\lambda^{(j)}$ defined in Section 4, there hold

$$\begin{aligned}
 \|f_\lambda^{(j)}\|_\infty &\leq \sqrt{2 + \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 + \|L_{\tilde{K}^{(j)}}^{r-1}\|^2} \\
 &\quad \cdot \lambda^{\min\{-\frac{1}{2}, r-1\}} \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}
 \end{aligned}$$

and

$$\begin{aligned}
 \|f_\lambda^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2} &\leq \sqrt{2 + \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 + \|L_{\tilde{K}^{(j)}}^{r-1}\|^2} \\
 &\quad \cdot \lambda^{\min\{0, r-\frac{1}{2}\}} \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}.
 \end{aligned}$$

Proof. Note that

$$\begin{aligned}
 f_\lambda^{(j)} &= U(\lambda I + L_{\tilde{K}^{(j)}})^{-1} L_{\tilde{K}^{(j)}}^{\frac{1}{2}} f_\rho^{(j)} \\
 &= L_{\tilde{K}^{(j)}}^T (\lambda I + L_{\tilde{K}^{(j)}})^{-1} f_\rho^{(j)}
 \end{aligned}$$

and

$$\|g_\lambda^{(j)} - f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2} = \lambda \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}.$$

Therefore,

$$\begin{aligned}
 \|f_\lambda^{(j)}\|_\infty &\leq \|(\lambda I + L_{\tilde{K}^{(j)}})^{-1} f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2} \\
 &= \lambda^{-1} \|L_{\tilde{K}^{(j)}} f_\lambda^{(j)} - f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2} \\
 &\leq \sqrt{2 + \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 + \|L_{\tilde{K}^{(j)}}^{r-1}\|^2} \\
 &\quad \cdot \lambda^{\min\{-\frac{1}{2}, r-1\}} \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}.
 \end{aligned}$$

The second statement follows directly from the result of Lemma 2. \square

Proposition 2. For $g_\lambda = \sum_{j=1}^d g_\lambda^{(j)} = \sum_{j=1}^d L_{\tilde{K}^{(j)}} f_\lambda^{(j)}$, there holds

$$\begin{aligned}
 E_3 &\leq \lambda^{\min\{1, 2r\}} \left(2 + \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 + \|L_{\tilde{K}^{(j)}}^{r-1}\|^2\right) \\
 &\quad \cdot \left(\sum_{j=1}^d \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}\right)^2.
 \end{aligned}$$

Proof. Based on Cauchy-Schwarz inequality, we can observe that

$$\begin{aligned}
 \sqrt{E_3} &= \left(\int_{\mathcal{Z}} (g_\lambda(x) - f_\rho(x))^2 d\rho(x, y)\right)^{\frac{1}{2}} \\
 &= \left(\int_{\mathcal{Z}} \left(\sum_{j=1}^d (g_\lambda^{(j)}(x^{(j)}) - f_\rho^{(j)}(x^{(j)}))^2 d\rho(x, y)\right)^{\frac{1}{2}} \\
 &\leq \sum_{j=1}^d \|L_{\tilde{K}^{(j)}} f_\lambda^{(j)} - f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2} \tag{9}
 \end{aligned}$$

Lemma 2 tells us that $\forall j \in \{1, 2, \dots, d\}$

$$\begin{aligned}
 & \|L_{\tilde{K}^{(j)}} f_\lambda^{(j)} - f_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2} \\
 &\leq \sqrt{2 + \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 + \|L_{\tilde{K}^{(j)}}^{r-1}\|^2} \lambda^{\min\{\frac{1}{2}, r\}} \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}.
 \end{aligned}$$

Combining this estimate with (9), we get the desired upper bound on E_3 . \square

A.2. Estimate of Hypothesis Error E_2

The hypothesis error reflects the divergence between \hat{g}_λ and g_λ on the expected risk and regularization. The following inequality from (Pinelis, 1994; Smale & Zhou, 2007) is used to bound the divergence.

Lemma 4. Let \mathcal{H} be a Hilbert space. For an independent random variable ξ on \mathcal{Z} with values in \mathcal{H} , assume that $\|\xi\|_{\mathcal{H}} \leq M < \infty$ almost surely. For any given independent identical distributed samples $\{z_i\}_{i=1}^m \subset \mathcal{Z}$ and any $\delta \in (0, 1)$, there holds

$$\begin{aligned}
 & \left\| \frac{1}{m} \sum_{i=1}^m \xi(z_i) - E\xi \right\|_{\mathcal{H}} \\
 &\leq \frac{2M \log(2/\delta)}{m} + \sqrt{\frac{2E\|\xi\|_{\mathcal{H}}^2 \log(2/\delta)}{m}}
 \end{aligned}$$

with confidence at least $1 - \delta/2$.

Proposition 3. For any $0 < \delta < 1$, with confidence $1 - \delta$, we have

$$\begin{aligned}
 & \mathcal{E}(\hat{g}_\lambda) - \mathcal{E}(g_\lambda) \\
 &\leq 16\sqrt{c} \lambda^{\min\{0, r-\frac{1}{2}\}} \left(\frac{\log(2/\delta)}{m} + \sqrt{\frac{\log(2/\delta)}{m}}\right) \\
 &\quad \cdot \left(\sum_{j=1}^d \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2} + \left(\sum_{j=1}^d \|g_\rho^{(j)}\|_{L_{\rho, \mathcal{X}^{(j)}}^2}\right)^2\right)
 \end{aligned}$$

and

$$E_2 \leq c_2 \left(\lambda^{\min\{0, r-\frac{1}{2}\}} \sqrt{\frac{\log(2/\delta)}{m}} + \lambda^{\min\{\frac{1}{2}, r\}}\right),$$

where $c = 2 + \|L_{\tilde{K}^{(j)}}^{r-\frac{1}{2}}\|^2 + \|L_{\tilde{K}^{(j)}}^{r-1}\|^2$ and c_2 is a positive constant independent of m, δ .

Proof. From Cauchy-Schwarz inequality, we can see that

$$\begin{aligned}
 & \mathcal{E}(\hat{g}_\lambda) - \mathcal{E}(g_\lambda) \\
 & \leq \left(\int_{\mathcal{Z}} (2y - \hat{g}_\lambda(x) - g_\lambda(x))^2 d\rho(x, y) \right)^{\frac{1}{2}} \\
 & \quad \cdot \left(\int_{\mathcal{Z}} (\hat{g}_\lambda(x) - g_\lambda(x))^2 d\rho(x, y) \right)^{\frac{1}{2}} \\
 & \leq \left(8 + 2 \int_{\mathcal{Z}} (\hat{g}_\lambda(x) - g_\lambda(x))^2 d\rho(x, y) \right)^{\frac{1}{2}} \\
 & \quad \cdot \left(\int_{\mathcal{Z}} (\hat{g}_\lambda(x) - g_\lambda(x))^2 d\rho(x, y) \right)^{\frac{1}{2}} \\
 & \leq \left(\sqrt{8} + \sqrt{2} \sum_{j=1}^d \|\hat{g}_\lambda^{(j)} - g_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \right) \\
 & \quad \cdot \sum_{j=1}^d \|\hat{g}_\lambda^{(j)} - g_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}. \tag{10}
 \end{aligned}$$

Denote $\xi^{(j)} = f_\lambda^{(j)}(x^{(j)})K(x^{(j)}, u)$ for any $j \in \{1, 2, \dots, d\}$. Then, from Lemma 3, we can deduce that

$$\|\xi^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \leq \|f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \leq \sqrt{c}\lambda^{\min\{0, r-\frac{1}{2}\}} \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}$$

and

$$E\|\xi^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \leq \|f_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2 \leq c\lambda^{\min\{0, 2r-1\}} \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}^2.$$

Moreover, for any $j \in \{1, \dots, d\}$ and $u \in \mathcal{X}^{(j)}$,

$$\begin{aligned}
 & \|\hat{g}_\lambda^{(j)} - g_\lambda^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \\
 & = \left\| \frac{1}{m} \sum_{i=1}^m \xi_i^{(j)} - E\xi^{(j)} \right\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \\
 & \leq \frac{2\sqrt{c}\lambda^{\min\{0, r-\frac{1}{2}\}} \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \log(2/\delta)}{m} \\
 & \quad + \lambda^{\min\{0, r-\frac{1}{2}\}} \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \sqrt{\frac{2c \log(2/\delta)}{m}}, \tag{11}
 \end{aligned}$$

where the last inequality is derived from Lemma 4. Then, we obtain the first statement by combining the estimates (10) and (11).

Now consider the upper bound of $\lambda\|\hat{g}_\lambda\|_{\ell_1}$. From the definition of \hat{g}_λ , we have

$$\begin{aligned}
 \lambda\|\hat{g}_\lambda\|_{\ell_1} & \leq \lambda \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty \leq \sum_{j=1}^d \|L_{K^{(j)}} f_\lambda^{(j)} - f_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}} \\
 & \leq \sqrt{c}\lambda^{\min\{0, r-\frac{1}{2}\}} \sum_{j=1}^d \|g_\rho^{(j)}\|_{L^2_{\rho_{\mathcal{X}^{(j)}}}}.
 \end{aligned}$$

Combining this estimate with the first statement, we derive the desired upper bound of E_2 . \square

A.3. Estimate of Sample Error E_1

In this paper, the sample error is estimated by the analysis technique associated with the empirical covering numbers. The empirical covering numbers with ℓ_2 -metric is denoted by $\mathcal{N}_2(\mathcal{F}, \varepsilon)$ and its detail definition can be founded in (Van der Vaart & Wellner, 1996; Shi et al., 2011).

Definition 1. For a function set \mathcal{F} and $\mathbf{u} = (u_i)_{i=1}^k \in \mathcal{X}$, the metric $d_{2, \mathbf{u}}$ is defined by

$$d_{2, \mathbf{u}}(f, g) = \sqrt{\frac{1}{k} \sum_{i=1}^k (f(u_i) - g(u_i))^2}, \forall f, g \in \mathcal{F}.$$

For every $\varepsilon > 0$, the empirical covering number is defined as $\mathcal{N}_2(\mathcal{F}, \varepsilon) = \sup_{k \in \mathbb{N}} \sup_{\mathbf{u} \in \mathcal{X}^k} \mathcal{N}_{2, \mathbf{u}}(\mathcal{F}, \varepsilon)$, where

$$\begin{aligned}
 \mathcal{N}_{2, \mathbf{u}}(\mathcal{F}, \varepsilon) & = \inf \left\{ l \in \mathbb{N} : \exists \{f_i\}_{i=1}^l \text{ such that} \right. \\
 & \quad \left. \mathcal{F} \subset \cup_{i=1}^l \{f \in \mathcal{F} : d_{2, \mathbf{u}}(f, f_i) \leq \varepsilon\} \right\}.
 \end{aligned}$$

The following concentration inequality is established in (Wu et al., 2007).

Lemma 5. Let \mathcal{F} be a measurable function set on \mathcal{Z} . Assume that, for any $f \in \mathcal{F}$, $\|f\|_\infty \leq B$ and $E(f^2) \leq cEf$ for some positive constants B, c . If for some $a > 0$ and $s \in (0, 2)$, $\log \mathcal{N}_2(\mathcal{F}, \varepsilon) \leq a\varepsilon^{-p}$ for any $\varepsilon > 0$, then there exists a constant c'_p such that for any $\delta \in (0, 1)$,

$$\begin{aligned}
 & \left| Ef - \frac{1}{m} \sum_{i=1}^m f(z_i) \right| \\
 & \leq \frac{1}{2} Ef + c'_p \max\{c^{\frac{2-p}{2+p}}, B^{\frac{2-p}{2+p}}\} \left(\frac{a}{m}\right)^{\frac{2}{2+p}} \\
 & \quad + \frac{(2c + 18B) \log(1/\delta)}{m}
 \end{aligned}$$

with confidence at least $1 - 2\delta$.

For any $R > 0$, denote

$$\begin{aligned}
 \mathcal{B}_R^{(j)} & = \left\{ f^{(j)} = \sum_{i=1}^m \alpha_i^{(j)} K^{(j)}(u_i^{(j)}, \cdot) \in \mathcal{H}^{(j)} : \right. \\
 & \quad \left. \|f^{(j)}\|_{\ell_1} \leq R \right\}
 \end{aligned}$$

and

$$\mathcal{B}_R = \left\{ f = \sum_{j=1}^d f^{(j)} : \|f\|_{\ell_1} \leq R \right\},$$

where

$$\|f\|_{\ell_1} = \inf \left\{ \sum_{j=1}^d \|f^{(j)}\|_{\ell_1} : f = \sum_{j=1}^d f^{(j)}, f^{(j)} \in \mathcal{H}^{(j)} \right\}.$$

Now we state the estimate on the empirical covering numbers of \mathcal{B}_1 . Similar analysis can be found in (Christmann & Zhou, 2016) for \mathcal{B}_1 in reproducing kernel Hilbert spaces.

Lemma 6. For any $j \in \{1, 2, \dots, d\}$, assume that $K^{(j)} \in C^s$ for some $s > 0$. Then,

$$\log \mathcal{N}_2(\mathcal{B}_1, \varepsilon) \leq d^{1+p} c_p \varepsilon^{-p},$$

where p is defined in Section 3 and c_p is a constant independent of ε .

Proof. For every $j \in \{1, 2, \dots, d\}$ and $\mathbf{x}^{(j)} \in (\mathcal{X}^{(j)})^S$, there exists a set $\{f_i^{(j)}\}_{i=1}^{N_j}$ with $N_j = \mathcal{N}_2(\mathcal{B}_1^{(j)}, \varepsilon)$ such that

$$\begin{aligned} \forall f^{(j)} \in \mathcal{B}_1^{(j)}, \quad \exists \quad i_j \in \{1, 2, \dots, N_j\}, \text{ s.t.}, \\ d_{2, \mathbf{x}^{(j)}}(f^{(j)} - f_{i_j}^{(j)}) \leq \varepsilon. \end{aligned}$$

For $f = \sum_{j=1}^d f^{(j)} \in \mathcal{B}_1$, we know $f^{(j)} \in \mathcal{B}_1^{(j)}$. For every $\mathbf{x} = (x_\ell)_{\ell=1}^S \in \mathcal{X}^S$, we have $\mathbf{x}^{(j)} = (x_\ell^{(j)})_{\ell=1}^S \in (\mathcal{X}^{(j)})^S, j \in \{1, 2, \dots, d\}$. Let $\tilde{f} = \sum_{j=1}^d \tilde{f}_{i_j}^{(j)}$. Then

$$\begin{aligned} & d_{2, \mathbf{x}}(f, \tilde{f}) \\ &= \left\{ \frac{1}{S} \sum_{\ell=1}^S (f(x_\ell) - \tilde{f}(x_\ell))^2 \right\}^{\frac{1}{2}} \\ &= \left\{ \frac{1}{S} \sum_{\ell=1}^S \left(\sum_{j=1}^d f^{(j)}(x_\ell^{(j)}) - \sum_{j=1}^d \tilde{f}_{i_j}^{(j)}(x_\ell^{(j)}) \right)^2 \right\}^{\frac{1}{2}} \\ &\leq \sum_{j=1}^d \left\{ \frac{1}{S} \sum_{\ell=1}^S (f^{(j)}(x_\ell^{(j)}) - \tilde{f}_{i_j}^{(j)}(x_\ell^{(j)}))^2 \right\}^{\frac{1}{2}} \\ &\leq \sum_{j=1}^d d_{2, \mathbf{x}^{(j)}}(f^{(j)} - \tilde{f}_{i_j}^{(j)}) \\ &\leq d\varepsilon. \end{aligned}$$

Therefore,

$$\log \mathcal{N}_2(\mathcal{B}_1, d\varepsilon) \leq \sum_{j=1}^d \log \mathcal{N}_2(\mathcal{B}_1^{(j)}, d\varepsilon).$$

According to Theorem 2 in (Shi et al., 2011) (also see Lemmas 2 and 3 in (Shi, 2013)) and considering $\|f^{(j)}\|_{\ell_1} \leq \sqrt{m} \|f^{(j)}\|_{\ell_2}^2$, we further get

$$\log \mathcal{N}_2(\mathcal{B}_1, d\varepsilon) \leq d c_p \varepsilon^{-p}.$$

Setting $\tilde{\varepsilon} = d\varepsilon$, we get the desired result. \square

Proposition 4. Under Assumptions 1 and 2, for any $\delta \in (0, 1)$, there holds

$$\begin{aligned} E_1 &\leq \frac{1}{2}(\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho)) + \frac{1}{2}(\mathcal{E}(\hat{g}_\lambda) - \mathcal{E}(g_\lambda)) \\ &\quad + E_3 + C_1 \log(2/\delta) (\lambda^{-\frac{2p}{2+p}} m^{-\frac{2}{2+p}} \\ &\quad + \lambda^{\min\{-1, 2r-2\}} m^{-\frac{2}{2+p}} + m^{-1}) \end{aligned}$$

with confidence $1 - \delta$, where C_1 is a positive constant independent m, λ, δ , and p is defined in Section 3.

Proof. The sample error E_1 can be decomposed as

$$E_{11} = \mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho) - (\mathcal{E}_{\mathbf{z}}(\pi(f_{\mathbf{z}})) - \mathcal{E}_{\mathbf{z}}(f_\rho))$$

and

$$E_{12} = \mathcal{E}_{\mathbf{z}}(\hat{g}_\lambda) - \mathcal{E}_{\mathbf{z}}(f_\rho) - (\mathcal{E}(\hat{g}_\lambda) - \mathcal{E}(f_\rho)).$$

In the sequel, we will bound E_{11} and E_{12} respectively.

Denote

$$\mathcal{G}_R = \{g(z) = (y - \pi(f)(x))^2 - (y - f_\rho(x))^2 : f \in \mathcal{B}_R\}.$$

For any $g \in \mathcal{G}_R$, we can deduce that $|g(z)| \leq 8$ and $Eg^2 \leq 16Eg$. Let $g_1, g_2 \in \mathcal{G}_R$ associated with f_1, f_2 respectively. It can be seen that

$$\begin{aligned} |g_1(z) - g_2(z)| &\leq 4|\pi(f_1)(x) - \pi(f_2)(x)| \\ &\leq 4|f_1(x) - f_2(x)|. \end{aligned}$$

This means

$$\begin{aligned} \log \mathcal{N}_2(\mathcal{G}_R, \varepsilon) &\leq \log \mathcal{N}_2(\mathcal{B}_R, \frac{\varepsilon}{4}) \leq \log \mathcal{N}_2(\mathcal{B}_1, \frac{\varepsilon}{4R}) \\ &\leq c_p d^{1+p} (4R)^p \varepsilon^{-p}, \end{aligned}$$

where the last inequality follows from Lemma 6.

Applying Lemma 5 to \mathcal{G}_R , we have with confidence $1 - \frac{\delta}{2}$

$$\begin{aligned} Eg - \frac{1}{m} \sum_{i=1}^m g(z_i) &\leq \frac{1}{2}(\mathcal{E}(\pi(f)) - \mathcal{E}(f_\rho)) \\ &\quad + \tilde{c}_1 (R^{\frac{2p}{2+p}} m^{-\frac{2}{2+p}} + m^{-1} \log(2/\delta)), \forall g \in \mathcal{G}_R, \end{aligned}$$

where \tilde{c}_1 is a constant independent of m, δ .

From the definition of $f_{\mathbf{z}}$ in Section 2, we know $f_{\mathbf{z}} \in \mathcal{B}_R$ with $R = \lambda^{-1}$. Then

$$\begin{aligned} E_{11} &\leq \frac{1}{2}(\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho)) \\ &\quad + \tilde{c}_1 (\lambda^{-\frac{2p}{2+p}} m^{-\frac{2}{2+p}} + m^{-1} \log(1/\delta)) \quad (12) \end{aligned}$$

with confidence $1 - \frac{\delta}{2}$.

Now we turn to bound E_{12} . Denote

$$\hat{\mathcal{G}} = \left\{ \hat{g} = \sum_{j=1}^d \hat{g}_\lambda^{(j)} : \hat{g}_\lambda^{(j)} = \frac{1}{m} \sum_{i=1}^m f_\lambda^{(j)}(v_i^{(j)}) K(v_i^{(j)}, \cdot) \right\}$$

and

$$\hat{\mathcal{H}} = \left\{ h : h(z) = (y - \hat{g}(x))^2 - (y - f_\rho(x))^2, \hat{g} \in \hat{\mathcal{G}} \right\}.$$

We can verify that

$$\begin{aligned} \|h\|_\infty &= \sup |2y - \hat{g}(x) - f_\rho(x)| \cdot |\hat{g}(x) - f_\rho(x)| \\ &\leq (3 + \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty)^2 \end{aligned}$$

and

$$Eh^2 \leq (3 + \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty)^2 Eh.$$

For any given $\hat{g}_1, \hat{g}_2 \in \hat{\mathcal{H}}$, the corresponding $h_1, h_2 \in \hat{\mathcal{H}}$ satisfy

$$|h_1(z) - h_2(z)| \leq 2(1 + \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty) |\hat{g}_1(x) - \hat{g}_2(x)|.$$

Then, from Lemma 6 and $\hat{g} \in \mathcal{B}_R$ with $R = \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty$, we have

$$\begin{aligned} &\log \mathcal{N}_2(\hat{\mathcal{H}}, \varepsilon) \\ &\leq \log \mathcal{N}_2\left(\hat{\mathcal{G}}, \frac{\varepsilon}{2(1 + \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty)}\right) \\ &\leq \log \mathcal{N}_2\left(\mathcal{B}_1, \frac{\varepsilon}{2 \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty (1 + \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty)}\right) \\ &\leq c_p d^{1+p} 2^p \left(\sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty + \sum_{j=1}^d \|f_\lambda^{(j)}\|_\infty^2 \right)^p \varepsilon^{-p}. \end{aligned}$$

Applying Lemma 5 to $\hat{\mathcal{H}}$, with confidence $1 - \frac{\delta}{2}$ we have

$$\begin{aligned} E_{12} &= \sum_{i=1}^m h(z_i) - Eh \leq \frac{1}{2} (\mathcal{E}(\hat{g}_\lambda) - \mathcal{E}(f_\rho)) \\ &\quad + \tilde{c}_2 \|f_\lambda\|_\infty^2 (m^{-\frac{2}{2+p}} + m^{-1} \log(2/\delta)) \\ &\leq \frac{1}{2} (\mathcal{E}(\hat{g}_\lambda) - \mathcal{E}(g_\lambda) + E_3) \\ &\quad + d \tilde{c}_2' \lambda^{\min\{-1, 2r-2\}} (m^{-\frac{2}{2+p}} + m^{-1} \log(2/\delta)), \end{aligned}$$

where the last inequality follows from Lemma 3 and $\tilde{c}_2, \tilde{c}_2'$ are some positive constants.

Combining this with the estimates of E_{11} in (12), we get the upper bound on E_1 . \square

A.4. Proof of Theorem 1

Proof. Combining Propositions 1-4, we have with confidence $1 - 4\delta$

$$\begin{aligned} &\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho) \\ &\leq C \log(2/\delta) (\lambda^{\min\{1, 2r\}} + \lambda^{\min\{0, r-\frac{1}{2}\}} m^{-\frac{1}{2}} \\ &\quad + \lambda^{\min\{-1, 2r-2\}} m^{-\frac{2}{2+p}} + \lambda^{-\frac{2p}{2+p}} m^{-\frac{2}{2+p}}). \end{aligned}$$

When $r \in (0, \frac{1}{2})$, by setting $\lambda = m^{-\theta_1}$ with $0 < \theta_1 < \min\{\frac{1}{p}, \frac{1}{(2+p)(1-r)}\}$, we get with confidence $1 - 4\delta$

$$\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho) \leq 4C \log(2/\delta) m^{-\gamma_1},$$

where $\gamma_1 = \min\{2r\theta_1, \frac{1}{2} + (r - \frac{1}{2})\theta_1, \frac{2}{2+p} - (2 - 2r)\theta_1, \frac{2}{2+p} - \frac{2p\theta_1}{2+p}\}$.

When $r \geq \frac{1}{2}$, taking $\lambda = m^{-\theta_2}$ with some $0 < \theta_2 < \min\{\frac{1}{p}, \frac{2}{2+p}\}$, we have with confidence $1 - 4\delta$

$$\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho) \leq 4C \log(2/\delta) m^{-\gamma_2},$$

where

$$\gamma_2 = \min\left\{\theta_2, \frac{1}{2}, \frac{2}{2+p} - \theta_2, \frac{2}{2+p} - \frac{2p\theta_2}{2+p}\right\},$$

This completes the proof. \square

A.5. Proof of Theorem 2

Theorem 2 is dependent on much stronger conditions on f_ρ than Theorem 1. The proof can be obtained directly by the estimate of E_{11} in Proposition 4.

Proof. Since $f_\rho^{(j)} \in \mathcal{H}^{(j)}$ for each $j \in \{1, 2, \dots, d\}$, we know that $f_\rho \in \mathcal{H}$. Then,

$$\begin{aligned} &\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho) \\ &\leq \mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}_{\mathbf{z}}(\pi(f_{\mathbf{z}})) \\ &\quad + \{\mathcal{E}_{\mathbf{z}}(f_{\mathbf{z}}) + \lambda \|f_{\mathbf{z}}\|_{\ell_1} - (\mathcal{E}_{\mathbf{z}}(f_\rho) + \lambda \|f_\rho\|_{\ell_1})\} \\ &\leq \mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho) - (\mathcal{E}_{\mathbf{z}}(\pi(f_{\mathbf{z}})) - \mathcal{E}_{\mathbf{z}}(f_\rho)) + \lambda \|f_\rho\|_{\ell_1} \\ &= E_{11} + \lambda \|f_\rho\|_{\ell_1}. \end{aligned}$$

From the estimate of E_{11} in (12), with confidence $1 - \delta$ we have

$$\mathcal{E}(\pi(f_{\mathbf{z}})) - \mathcal{E}(f_\rho) \leq \bar{c} \log(1/\delta) (\lambda^{-\frac{2p}{2+p}} m^{-\frac{2}{2+p}} + \lambda),$$

where \bar{c} is a positive constant independent of m, λ .

Taking λ such that $\lambda^{-\frac{2p}{2+p}} m^{-\frac{2}{2+p}} = \lambda$, we get the desired result. \square

B. Proof of Theorem 3

Proof. Denote $\alpha = (\alpha_t^{(j)})_{t,j} \in \mathbb{R}^{md}$, where $t \in \{1, 2, \dots, m\}$ and $j \in \{1, 2, \dots, d\}$. Define

$$G(\alpha) = \frac{1}{m} \sum_{i=1}^m (y_i - \sum_{j=1}^d \sum_{t=1}^m \alpha_t^{(j)} K^{(j)}(x_t^{(j)}, x_i^{(j)}))^2 + \lambda \sum_{j=1}^d \sum_{t=1}^m |\alpha_t^{(j)}|.$$

Recall that $f_{\mathbf{z}} = \sum_{j=1}^d \sum_{t=1}^m \hat{\alpha}_t^{(j)} K^{(j)}(x_t^{(j)}, \cdot)$ and $\hat{\alpha} = (\hat{\alpha}_t^{(j)})_{t,j}$ is the maximizer of $G(\alpha)$.

Let $I_+ = \{(t, j) : \hat{\alpha}_t^{(j)} > 0\}$, $I_- = \{(t, j) : \hat{\alpha}_t^{(j)} < 0\}$, and $I_0 = \{(t, j) : \hat{\alpha}_t^{(j)} = 0\}$.

For $(t, j) \in I_+$, we get

$$\begin{aligned} & \left. \frac{\partial G(\alpha)}{\partial \alpha_t^{(j)}} \right|_{\alpha=\hat{\alpha}} \\ &= -\frac{2}{m} \sum_{i=1}^m (y_i - f_{\mathbf{z}}(x_i)) K^{(j)}(x_t^{(j)}, x_i^{(j)}) + \lambda \\ &= 0 \end{aligned}$$

This means

$$\frac{1}{m} \sum_{i=1}^m (y_i - f_{\mathbf{z}}(x_i)) K^{(j)}(x_t^{(j)}, x_i^{(j)}) = \frac{\lambda}{2}.$$

Similarly, for $(t, j) \in I_-$, there exists

$$\frac{1}{m} \sum_{i=1}^m (y_i - f_{\mathbf{z}}(x_i)) K^{(j)}(x_t^{(j)}, x_i^{(j)}) = -\frac{\lambda}{2}.$$

For $(t, j) \in I_0$, there holds

$$\begin{aligned} -\frac{2}{m} \sum_{i=1}^m (y_i - f_{\mathbf{z}}(x_i)) K^{(j)}(x_t^{(j)}, x_i^{(j)}) - \lambda &\leq \left. \frac{\partial G(\alpha)}{\partial \alpha_t^{(j)}} \right|_{\alpha=\hat{\alpha}} \\ &\leq -\frac{2}{m} \sum_{i=1}^m (y_i - f_{\mathbf{z}}(x_i)) K^{(j)}(x_t^{(j)}, x_i^{(j)}) + \lambda. \end{aligned}$$

This means, for any $(t, j) \in I_0$,

$$\left| \frac{1}{m} \sum_{i=1}^m (y_i - f_{\mathbf{z}}(x_i)) K^{(j)}(x_t^{(j)}, x_i^{(j)}) \right| < \frac{\lambda}{2}.$$

This completes the proof. \square

References

- Christmann, A. and Zhou, D. X. Learning rates for the risk of kernel-based quantile regression estimators in additive models. *Analysis and Applications*, 14(3):449–477, 2016.
- Feng, Y., Lv, S., Hang, H., and Suykens, J. A. Kernelized elastic net regularization: Generalization bounds, and sparse recovery. *Neural Comput.*, 28(3):525–562, 2016.
- Nie, W. and Wang, C. Constructive analysis for coefficient regularization regression algorithms. *Journal of Mathematical Analysis and Applications*, 431(2):1153–1171, 2015.
- Pinelis, I. Optimum bounds for the distributions of martingales in banach spaces. *The Annals of Probability*, pp. 1679–1706, 1994.
- Shi, L. Learning theory estimates for coefficient-based regularized regression. *Appl. Comput. Harmon. Anal.*, 34(2):252–265, 2013.
- Shi, L., Feng, Y., and Zhou, D. X. Concentration estimates for learning with ℓ_1 -regularizer and data dependent hypothesis spaces. *Appl. Comput. Harmon. Anal.*, 31(2):286–302, 2011.
- Smale, S. and Zhou, D. X. Learning theory estimates via integral operators and their approximations. *Constr. Approx.*, 26(2):153–172, 2007.
- Sun, H. and Wu, Q. Least square regression with indefinite kernels and coefficient regularization. *Appl. Comput. Harmon. Anal.*, 30(1):96–109, 2011.
- Van der Vaart, A. and Wellner, J. *Weak Convergence and Empirical Processes*. Springer, 1996.
- Wu, Q., Ying, Y., and Zhou, D.-X. Multi-kernel regularized classifiers. *J. Complexity*, 23(1):108–134, 2007.