

Error Analysis on Hérmité Learning with Gradient Data*

Baohuai SHENG¹ Jianli WANG¹ Daohong XIANG²

Abstract This paper deals with Hérmité learning which aims at obtaining the target function from the samples of function values and the gradient values. Error analysis is conducted for these algorithms by means of approaches from convex analysis in the framework of multi-task vector learning and the improved learning rates are derived.

Keywords Hérmité learning, Gradient learning, Learning rate, Convex analysis, Multitask learning, Differentiable reproducing kernel Hilbert space

2000 MR Subject Classification 41A25, 46E22, 68Q32, 68T05, 90C25

1 Introduction

We considered Hérmité learning with gradient vectors (see [1–4]) in this paper. It can produce a smoothness function when the function value samples and gradient data is provided.

Let X be a compact subset of the Euclidean space \mathbb{R}^d on which the learning or function approximation is considered. For each $x = (x^1, x^2, \dots, x^d)^T \in X$, the gradient of function $f : X \rightarrow R$ at x is denoted by the vector

$$\nabla_x f(x) = \left(\frac{\partial f(x)}{\partial x^1}, \dots, \frac{\partial f(x)}{\partial x^d} \right)^T$$

if the partial derivative for each variable exists. It is known (see [2]) that the essence of Hérmité learning is to obtain the regression function $f_\rho(x) = \int_R y d\rho(y|x)$ from samples $z = \{(x_i, y_i)\}_{i=1}^m$ with $y_i = (y_i^0, \tilde{y}_i)$, where $\{(x_i, y_i^0)\}_{i=1}^m$ are drawn independently according to $\rho(x, y) = \rho(y|x)\rho_X(x)$ and $\tilde{y}_i \approx \nabla f(x_i)$.

Let $K : X \times X \rightarrow R$ be a Mercer kernel which is continuous, symmetric and positive semi-definite, which means that the matrix $(K(x_i, x_j))_{i,j=1}^l$ is positive semi-definite for any finite set of points $\{x_1, \dots, x_l\} \subset X$. The associated reproducing kernel Hilbert space (RKHS) \mathcal{H}_K is defined (see [5–7]) as the completion of the linear span of the set of the function $\{K_x = K(x, \cdot) : x \in X\}$ with the inner product given by

$$\langle K_x, K_y \rangle_K = K(x, y).$$

Let $L^2(\mu)$ be the class of all square integrable functions with respect to the measure μ with the norm

$$\|f\|_{L^2(\mu)} = \left(\int_X |f(x)|^2 d\mu \right)^{\frac{1}{2}} < +\infty.$$

Manuscript received February 9, 2015. Revised May 28, 2016.

¹Department of Mathematics, Shaoxing University, Shaoxing 312000, Zhejiang, China.

E-mail: bhsheng@usx.edu.cn wjl@usx.edu.cn

²Department of Mathematics, Zhejiang Normal University, Jinhua 321004, Zhejiang, China.

E-mail: daohongxiang@zjnu.cn

*This work was supported by the National Natural Science Foundation of China (No. 11471292).

Define the integral operator \mathcal{L}_K as

$$\mathcal{L}_K(f, x) = \int_X f(u) K(x, u) \, d\mu, \quad f \in L^2(\mu).$$

Since K is a positive semi-definite, \mathcal{L}_K is a compact positive operator. Let λ_k be the k -th positive eigenvalue of $\mathcal{L}_K(f)$ and $\phi_k(x)$ be the corresponding continuous orthonormal eigenfunction. Then, by Mercer's theorem, for all $x, y \in X$, there holds

$$K(x, y) = \sum_{k=0}^{+\infty} \lambda_k \phi_k(x) \phi_k(y), \quad x, y \in X, \quad (1.1)$$

where the convergence is absolute (for each $x, y \in X$) and uniform on $X \times X$. Then we know (see [8–9])

$$\mathcal{H}_K = \mathcal{L}_K^{\frac{1}{2}}(L_2(\mu)) = \left\{ f(x) = \sum_{k=0}^{\infty} a_k(f) \phi_k(x) : \sum_{k=0}^{\infty} \frac{|a_k(f)|^2}{\lambda_k} < +\infty \right\}$$

with inner product

$$\langle f, g \rangle_{\mathcal{H}_K} = \sum_{k=0}^{\infty} \frac{a_k(f) a_k(g)}{\lambda_k}$$

and $a_k(f) = \int_X f(y) \phi_k(y) \, d\mu(y)$.

When X is a compact set, there is a constant $k > 0$ such that

$$\|f\|_{\infty, X} \leq k \|f\|_{\mathcal{H}_K}, \quad f \in \mathcal{H}_K. \quad (1.2)$$

Let $M > 0$ be a given positive real number and $Y = [-M, M]^{d+1}$. We denote by $y \in Y$ as $y = (y^0, \tilde{y})$ with $\tilde{y} = (y^1, y^2, \dots, y^d)^T$ and

$$\vec{f}_\rho(x) = (f_\rho(x), \tilde{f}_\rho(x)),$$

where

$$f_\rho(x) = \int_{[-M, M]} y^0 \, d\rho(y|x)$$

and

$$\begin{aligned} \tilde{f}_\rho(x) &= \int_{[-M, M]^d} \tilde{y} \, d\rho(y|x) \\ &= \left(\int_{[-M, M]^d} y^1 \, d\rho(y|x), \dots, \int_{[-M, M]^d} y^d \, d\rho(y|x) \right)^T. \end{aligned}$$

Let $y_i^0 \in [-M, M]$ and $\tilde{y}_i = (y_i^1, \dots, y_i^d)^T \in [-M, M]^d$. Then, the Hérmité learning algorithm corresponding to the samples z is (see [1])

$$f_{z, \lambda} := \arg \min_{f \in \mathcal{H}_K} \frac{1}{m} \sum_{i=1}^m ((y_i^0 - f(x_i))^2 + \|\tilde{y}_i - \nabla f(x_i)\|_{\mathbb{R}^d}^2) + \lambda \|f\|_{\mathcal{H}_K}^2, \quad (1.3)$$

where $\lambda > 0$ is the regularization parameter, $\|\cdot\|_{\mathbb{R}^d}$ is the usual norm of the Euclidean space \mathbb{R}^d .

The representer theorem for scheme (1.3) was provided in [1], which showed that the minimization over the possibly infinite dimensional space \mathcal{H}_K can be achieved in a finite dimensional subspace generated by $\{K_{x_i}(\cdot)\}$ and their partial derivatives. The explicit solution to (1.3) and an upper bound for the learning rate with the integral operator approach were given in [2].

We notice that the problem of learning multiple tasks with kernel methods has been a developing topic (see [10–14]). The representer theorem for these frameworks have been given, but the results of convergence quantitative analysis are relatively few, and very effective methods for bounding the convergence rates does not appear. An aim of the present paper is to show that model (1.3) is among the group of the multiple tasks. For this purpose, we rewrite model (1.3) from the view of vector valued functions.

Let $C^{(1)}$ be the class of all the real functions $f(x)$ defined on X with continuous partial derivatives and $C^{(2)}(X)$ denote the class of all the real functions $f(x)$ defined on X with continuous partial derivatives $\frac{\partial^2 f}{\partial x^\alpha \partial y^\beta}$ for $\alpha = 1, 2, \dots, d$. By $K \in C^{(2)}(X \times X)$, we mean all the partial derivatives $\frac{\partial^2 \partial^2 K(x, y)}{\partial (x^\alpha)^2 \partial (x^\beta)^2}$ that are continuous on $X \times X$. By [15–17] we know that if the kernels $K(x, y)$ have the form of (1.1) and $K(x, y) \in C^{(2)}(X \times X)$, then, \mathcal{H}_K can be embedded into both $C^{(1)}(X)$ and $C^{(2)}(X)$, and for any given $x \in X$ and all $\alpha = 1, 2, \dots, d$, the following relations hold:

$$\frac{\partial K_x(\cdot)}{\partial x^\alpha} \in \mathcal{H}_K, \quad \frac{\partial^2 K_x(\cdot)}{\partial (x^\alpha)^2} \in \mathcal{H}_K, \quad (1.4)$$

$$\frac{\partial f(x)}{\partial x^\alpha} = \left\langle f, \frac{\partial}{\partial x^\alpha} K_x(\cdot) \right\rangle_{\mathcal{H}_K}, \quad f \in \mathcal{H}_K \quad (1.5)$$

and

$$\left| \frac{\partial f(x)}{\partial x^\alpha} \right| \leq k_1 \|f\|_{\mathcal{H}_K}, \quad \forall x \in X, \quad \forall \alpha = 0, 1, 2, \dots, d, \quad (1.6)$$

where $\frac{\partial f}{\partial x^0} = f(x)$, $k_1 = \sup_{\substack{x, y \in X \\ 0 \leq \alpha, \beta \leq d}} \sqrt{\left| \frac{\partial^2 K(x, y)}{\partial x^\alpha \partial y^\beta} \right|}$.

Define

$$\mathcal{H}_K^\alpha = \left\{ f_\alpha(x) = \frac{\partial f(x)}{\partial x^\alpha} : f(x) \in \mathcal{H}_K \right\}, \quad \alpha = 1, 2, \dots, d$$

and

$$\begin{aligned} \overrightarrow{\mathcal{H}_K} &= \mathcal{H}_K \times \mathcal{H}_K^1 \times \dots \times \mathcal{H}_K^d \\ &= \{ \vec{f} = (f, f_1, \dots, f_d)^\top : f \in \mathcal{H}_K, f_\alpha \in \mathcal{H}_K^\alpha, \alpha = 1, 2, \dots, d \}. \end{aligned}$$

Then, (1.3) can be rewritten as

$$\vec{f}_{z, \lambda} := \arg \min_{\vec{f} = (f, f_1, \dots, f_d)^\top \in \overrightarrow{\mathcal{H}_K}} \mathcal{E}_z(\vec{f}) + \lambda \|f\|_{\mathcal{H}_K}^2, \quad (1.7)$$

where

$$\begin{aligned} \mathcal{E}_z(\vec{f}) &= \frac{1}{m} \sum_{i=1}^m \|\vec{y}_i - \vec{f}(x_i)\|_{\mathbb{R}^{d+1}}^2, \\ \vec{y}_i &= (y_i^0, y_i^1, \dots, y_i^d)^\top, \quad i = 1, 2, \dots, m \end{aligned}$$

and

$$\overrightarrow{f_{z,\lambda}} = (f^{(z,\lambda)}, f_1^{(z,\lambda)}, \dots, f_d^{(z,\lambda)})^T.$$

Algorithm (1.7) is neither the same as the usual least square regression (see [18–21]), nor the current multitask learning model since it uses the penalty $\|f\|_{\mathcal{H}_K}^2$ not $\|\vec{f}\|_{\mathcal{H}_K}^2$. However, it is really a concrete example of multitask learning models, the study method of its performance can provide useful reference for the research of other multitask models. This is the first motivation for writing this paper. On the other hand, such approaches may provide a way of thinking for dealing with other gradient learning models. The structure of model (1.7) is close to the usual least square learning models, this fact creates an opportunity of studying the gradient learning with existing methods, e.g., the convex analysis method (see [20–24]). This is the second motivation for writing this paper.

We form an improved convex method with the help of Gâteaux derivative and the optimality conditions of convex functions and use it to bound the convergence rates of model (1.7). To show the main results of the present paper, we restate the notion of covering number.

For a distance space S and a real number $\eta > 0$. The covering number $\mathcal{N}(S, \eta)$ is defined to be the minimal positive integer number l such that there exists l disks in S with radius η covering S .

We call a compact subset E of a distance space $(\mathcal{B}, \|\cdot\|_{\mathcal{B}})$ logarithmic complexity exponent $s \geq 0$ if there is a constant $c_s > 0$ such that the closed ball of radius R centered at origin, i.e.,

$$\mathcal{B}_R = \{f \in E : \|f\|_{\mathcal{B}} \leq R\}$$

satisfies

$$\log \mathcal{N}(\mathcal{B}_R, \eta) \leq c_s \left(\frac{R}{\eta}\right)^s, \quad \forall \eta > 0. \quad (1.8)$$

We now give the following Theorem 1.1.

Theorem 1.1 *Let $K(x, y)$ be a Mercer-like kernel satisfying $K(x, y) \in C^{(2)}(X \times X)$, $\overrightarrow{f_{z,\lambda}}$ be the solution of (1.7) and $C = 512 k_1^4 (d+1)^3 M$. If \mathcal{H}_k has logarithmic complexity exponent $s \geq 0$ in the uniform continuity norm $\|\cdot\|_{C(X)}$ and $\lambda \leq \frac{k_1^2 d}{M^2} \times D(\overrightarrow{f_\rho}, \lambda)$, then, for any $0 < \delta < 1$ and $0 < \delta < \frac{4}{e^{c_s}}$, with confidence $1 - \delta$, we have*

$$\|\overrightarrow{f_{z,\lambda}} - \overrightarrow{f_\rho}\|_{L^2(\rho_X)}^2 \leq \frac{C \sqrt{D(\overrightarrow{f_\rho}, \lambda)}}{\lambda^2 \sqrt{m}} \times \left(\frac{1}{\sqrt{m}} + \frac{1}{2+\sqrt{s}m} \right) \left(\log \frac{4}{\delta} \right)^2 + D(\overrightarrow{f_\rho}, \lambda), \quad (1.9)$$

where

$$D(\overrightarrow{f_\rho}, \lambda) = \inf_{\vec{h} \in \mathcal{H}_K} (\mathcal{E}_\rho(\vec{h}) - \mathcal{E}_\rho(\overrightarrow{f_\rho}) + \lambda \|h\|_{\mathcal{H}_K}^2)$$

and

$$\mathcal{E}_\rho(\vec{f}) = \int_Z \|\vec{y} - \vec{f}(x)\|_{\mathbb{R}^{d+1}}^2 d\rho.$$

We now give some comments on (1.9).

(1) Let $H_{\rho_X}^1$ be the Sobolev space consisting of all the functions $f \in L^2(\rho_X)$ with all partial derivatives belonging to $L^2(\rho_X)$, whose norm $\|f\|_{H_{\rho_X}^1}$ is induced by the inner product

$$\langle f, g \rangle_{H_{\rho_X}^1} = \int_X (f(x) g(x) + \nabla_x f(x) \nabla_x g(x)) d\rho_X.$$

If ρ is perfect (see [2]), i.e., $\tilde{f}_\rho(x) = \nabla f_\rho(x)$, then, $\vec{f}_\rho(x) = (f_\rho(x), \nabla f_\rho(x))$ and

$$\|\vec{f}_{z,\lambda} - \vec{f}_\rho\|_{L^2(\rho_X)} = \|f_{z,\lambda} - f_\rho\|_{H_{\rho_X}^1}.$$

(2) By Lemma 2.2 afterward we have

$$D(\vec{f}_\rho, \lambda) = \inf_{\vec{h} \in \vec{\mathcal{H}}_K} (\|\vec{f}_\rho - \vec{h}\|_{L^2(\rho_X)}^2 + \lambda \|h\|_{\mathcal{H}_K}^2).$$

As usual, we assume that there are given constants $c > 0$ and $0 < \beta < 4$ such that $D(\vec{f}_\rho, \lambda) \leq c\lambda^\beta$ if $\vec{\mathcal{H}}_K$ is density in $L^2(\rho_X)$. In this case, if $s = 0$, then we have by (1.9) that

$$\|f_{z,\lambda} - f_\rho\|_{H_{\rho_X}^1}^2 = O\left(\frac{1}{m \lambda^{2-\frac{\beta}{2}}} \times \left(\log \frac{4}{\delta}\right)^2 + \lambda^\beta\right). \quad (1.10)$$

Further, if $\lambda = m^{-\theta}$ and $0 < \theta < \frac{2}{4-\beta}$, then, we have by (1.10) the following estimate

$$\|f_{z,\lambda} - f_\rho\|_{H_{\rho_X}^1}^2 = O\left(\frac{(\log \frac{4}{\delta})^2}{m^{1-(2-\frac{\beta}{2})\theta}} + m^{-\beta\theta}\right). \quad (1.11)$$

(3) Define a new integral operator

$$L = L_K : H_{\rho_X}^1 \rightarrow H_{\rho_X}^1$$

associated with K and ρ_X by

$$L_K(f, x) = \int_X (K(x, y) f(y) + \nabla_y K(x, y) \nabla_y f(y)) d\rho_X(y), \quad x \in X, f \in H_{\rho_X}^1.$$

Let L_K^r be defined as

$$L_K^r \left(\sum c_k \phi_k(x) \right) = \sum c_k \lambda_k^r \phi_k(x)$$

for $r > 0$. Then, by taking $\lambda = 8(d+1) k^2 \log(\frac{4}{\delta}) m^{-\beta}$, with confidence $1 - \delta$, we have (see [2])

$$\|f_{z,\lambda} - f_\rho\|_{H_{\rho_X}^1} \leq 8 \log\left(\frac{4}{\delta}\right) \{M + (d+1) \|K\|_{C^2}^r \times \|L_K^{-r} f_\rho\|_{H_{\rho_X}^1}\} m^{-r\beta}, \quad (1.12)$$

where

$$\beta = \begin{cases} \frac{1}{2r+1}, & \text{if } r > \frac{1}{2}, \\ \frac{1}{2}, & \text{if } 0 < r \leq \frac{1}{2}. \end{cases}$$

It is easy to see that (1.11) sharps (1.12).

2 Proofs

To show Theorem 1.1, we need the notation of the Gâteaux derivative and some related lemmas.

Let $(\mathcal{H}, \|\cdot\|_{\mathcal{H}})$ be a Hilbert space, $F(f) : \mathcal{H} \rightarrow \mathbb{R} \cup \{\mp\infty\}$ be a real function. We say F is Gâteaux differentiable at $f \in \mathcal{H}$, if there is a $\xi \in \mathcal{H}$ such that for any $g \in \mathcal{H}$, there holds

$$\lim_{t \rightarrow 0} \frac{F(f+tg) - F(f)}{t} = \langle g, \xi \rangle_{\mathcal{H}} \quad (2.1)$$

and write $\nabla_f F(f) = \xi$ as the Gâteaux derivative of $F(f)$ at f .

Lemma 2.1 Let $F(f) : \mathcal{H} \rightarrow R \cup \{\mp\infty\}$ be a function defined on Hilbert space \mathcal{H} . Then, we have following results:

(i) If $F(f)$ is a convex function, then, $F(f)$ attains minimal value at f_0 if and only if $\nabla_f F(f_0) = 0$.

(ii) If $F(f) : \mathcal{H} \rightarrow R \cup \{\mp\infty\}$ is a Gâteaux differentiable function, then, $F(f)$ is a convex on \mathcal{H} if and only if for any $f, g \in \mathcal{H}$ we have

$$F(g + f) - F(f) \geq \langle g, \nabla_f F(f) \rangle_{\mathcal{H}}. \quad (2.2)$$

Proof We have (i) from Proposition 17.4 of [25] and we have (ii) from Proposition 17.10 and Proposition 17.12 of [25].

Lemma 2.2 \vec{f}_ρ satisfies the relation

$$\vec{f}_\rho = \arg \min_{\vec{f}=(f, f_1, \dots, f_d)^T} \mathcal{E}_\rho(\vec{f}) \quad (2.3)$$

and the equation

$$\mathcal{E}_\rho(\vec{f}) - \mathcal{E}_\rho(\vec{f}_\rho) = \int_Z \|\vec{f}_\rho(x) - \vec{f}(x)\|_{\mathbb{R}^{d+1}}^2 d\rho_X. \quad (2.4)$$

Proof Since equality

$$\|a + b\|_{\mathbb{R}^{d+1}}^2 = \|a\|_{\mathbb{R}^{d+1}}^2 + 2\langle a, b \rangle_{\mathbb{R}^{d+1}} + \|b\|_{\mathbb{R}^{d+1}}^2$$

holds for any $a, b \in \mathbb{R}^{d+1}$, we have

$$\begin{aligned} \mathcal{E}_\rho(\vec{f}) &= \int_Z \|\vec{y} - \vec{f}_\rho\|_{\mathbb{R}^{d+1}}^2 d\rho - 2 \int_Z \langle \vec{y} - \vec{f}_\rho(x), \vec{f}(x) - \vec{f}_\rho(x) \rangle_{\mathbb{R}^{d+1}} d\rho \\ &\quad + \int_Z \|\vec{f}_\rho(x) - \vec{f}(x)\|_{\mathbb{R}^{d+1}}^2 d\rho \\ &= \mathcal{E}_\rho(\vec{f}_\rho) - 2 \int_X \left\langle \int_Y (\vec{y} - \vec{f}_\rho(x)) d\rho(y|x), \vec{f}(x) - \vec{f}_\rho(x) \right\rangle_{\mathbb{R}^{d+1}} d\rho_X \\ &\quad + \int_Z \|\vec{f}_\rho(x) - \vec{f}(x)\|_{\mathbb{R}^{d+1}}^2 d\rho_X \\ &= \mathcal{E}_\rho(\vec{f}_\rho) + \int_Z \|\vec{f}_\rho(x) - \vec{f}(x)\|_{\mathbb{R}^{d+1}}^2 d\rho_X. \end{aligned} \quad (2.5)$$

Let $\vec{h}_\lambda = (h^{(\lambda)}, h_1^{(\lambda)}, \dots, h_d^{(\lambda)})^T$ be defined as

$$\vec{h}_\lambda = \arg \min_{\vec{h}=(h, h_1, \dots, h_d)^T \in \vec{\mathcal{H}}_K} (\mathcal{E}_\rho(\vec{h}) + \lambda \|h\|_{\mathcal{H}_K}^2). \quad (2.6)$$

Then, we have the following Lemma 2.3.

Lemma 2.3 Let \vec{h}_λ be a solution of (2.6). Then, we have

$$\nabla_{\vec{f}}(\mathcal{E}_\rho(\vec{f}))(\cdot) = -2 \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \quad (2.7)$$

and

$$\nabla_{\vec{f}}(\mathcal{E}_z(\vec{f}))(\cdot) = -\frac{2}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{f}(x_i) \rangle_{\mathbb{R}^{d+1}}, \quad (2.8)$$

where $\nabla_x^* K_x(\cdot) = (\frac{\partial K_x(\cdot)}{\partial x^0}, \frac{\partial K_x(\cdot)}{\partial x^1}, \dots, \frac{\partial K_x(\cdot)}{\partial x^d})^T$ and $\frac{\partial K_x(\cdot)}{\partial x^0} = K_x(\cdot)$.

Proof By the equality

$$a^2 - b^2 = \langle a - b, \quad 2b \rangle_{\mathbb{R}^{d+1}} + \|a - b\|_{\mathbb{R}^{d+1}}^2, \quad a \in \mathbb{R}^{d+1}, \quad b \in \mathbb{R}^{d+1}. \quad (2.9)$$

we have

$$\begin{aligned} & \lim_{t \rightarrow 0} \frac{\mathcal{E}_\rho(\vec{f} + t\vec{g}) - \mathcal{E}_\rho(\vec{f})}{t} \\ &= \lim_{t \rightarrow 0} \frac{1}{t} \left(\int_Z \|\vec{y} - (\vec{f}(x) + t\vec{g}(x))\|_{\mathbb{R}^{d+1}}^2 d\rho - \int_Z \|\vec{y} - \vec{f}(x)\|_{\mathbb{R}^{d+1}}^2 d\rho \right) \\ &= -2 \int_Z \langle \vec{g}(x), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} d\rho. \end{aligned}$$

Take $g_0(x) = g(x)$, $f_0(x) = f(x)$. Then, by (1.5) we have

$$\begin{aligned} \langle \vec{g}(x), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} &= \sum_{k=0}^d g_k(x) \times (y^k - f_k(x)) \\ &= \sum_{k=0}^d \left(g, \frac{\partial K_x(\cdot)}{\partial x^k} \right)_{\mathcal{H}_K} \times (y^k - f_k(x)) \\ &= \left(g, \sum_{k=0}^d \frac{\partial K_x(\cdot)}{\partial x^k} (y^k - f_k(x)) \right)_{\mathcal{H}_K} \\ &= (g, \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}})_{\mathcal{H}_K}. \end{aligned} \quad (2.10)$$

It follows

$$\lim_{t \rightarrow 0} \frac{\mathcal{E}_\rho(\vec{f} + t\vec{g}) - \mathcal{E}_\rho(\vec{f})}{t} = \left(g, 2 \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right)_{\mathcal{H}_K}.$$

(2.7) thus holds. (2.8) can be proved in the same way.

Lemma 2.4 Let \vec{h}_λ be a solution of (2.6) and let $\vec{f}_{z,\lambda}$ be a solution of (1.7). Then, we have

$$\lambda h^{(\lambda)}(\cdot) = \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{h}_\lambda(x) \rangle_{\mathbb{R}^{d+1}} d\rho \quad (2.11)$$

and

$$\lambda f^{(z,\lambda)}(\cdot) = \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{f}_{z,\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}}. \quad (2.12)$$

Proof Since \vec{h}_λ is the solution of (2.6), we have

$$\begin{aligned} & \nabla_{\vec{h}} (\mathcal{E}_\rho(\vec{h}) + \lambda \|h\|_{\mathcal{H}_K}^2) |_{\vec{h}=\vec{h}_\lambda} \\ &= \nabla_{\vec{h}} \mathcal{E}_\rho(\vec{h}) |_{\vec{h}=\vec{h}_\lambda} + \lambda \nabla_{\vec{h}} (\|h\|_{\mathcal{H}_K}^2) |_{\vec{h}=\vec{h}_\lambda} \\ &= -2 \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{h}_\lambda(x) \rangle_{\mathbb{R}^{d+1}} d\rho(x, y) + 2\lambda h^{(\lambda)}(\cdot) = 0, \end{aligned} \quad (2.13)$$

where we have used the fact that

$$\nabla_{\vec{h}} (\|h\|_{\mathcal{H}_K}^2) = \nabla_h (\|h\|_{\mathcal{H}_K}^2) = 2h.$$

(2.11) thus holds. (2.12) can be proved in the same way.

Lemma 2.5 (2.6) has a unique solution \vec{h}_λ and (1.7) has a unique solution $\vec{f}_{z,\lambda}$. They satisfy the following inequalities:

$$\|h^{(\lambda)}\|_{\mathcal{H}_K} \leq \sqrt{\frac{D(\vec{f}_\rho, \lambda)}{\lambda}} \quad (2.14)$$

and

$$\|f^{(z,\lambda)}\|_{\mathcal{H}_K} \leq \frac{M}{\sqrt{\lambda}}. \quad (2.15)$$

Proof Let $\vec{f} = (f, f_1, \dots, f_d)^T \in \overrightarrow{\mathcal{H}_K}$ and $\vec{g} = (g, g_1, \dots, g_d)^T \in \overrightarrow{\mathcal{H}_K}$. Then, by the equality (2.9) we have

$$\begin{aligned} \mathcal{E}_\rho(\vec{g}) - \mathcal{E}_\rho(\vec{f}) &= \int_Z (\|\vec{y} - \vec{g}(x)\|_{\mathbb{R}^{d+1}}^2 - \|\vec{y} - \vec{f}(x)\|_{\mathbb{R}^{d+1}}^2) d\rho \\ &= -2 \int_Z \langle \vec{g}(x) - \vec{f}(x), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} d\rho + \int_Z \|\vec{f}(x) - \vec{g}(x)\|_{\mathbb{R}^{d+1}}^2 d\rho \\ &\geq -2 \int_Z \langle \vec{g}(x) - \vec{f}(x), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} d\rho. \end{aligned}$$

Since (2.10), we have

$$\langle \vec{g}(x) - \vec{f}(x), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} = (g - f, \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}})_{\mathcal{H}_K}.$$

Therefore, by (2.7) we have

$$\begin{aligned} \mathcal{E}_\rho(\vec{g}) - \mathcal{E}_\rho(\vec{f}) &\geq \left(g - f, -2 \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right)_{\mathcal{H}_K} \\ &= (g - f, \nabla_{\vec{f}}(\mathcal{E}_\rho(\vec{f})))_{\mathcal{H}_K}. \end{aligned}$$

By (ii) of Lemma 2.1 we know $\mathcal{E}_\rho(\vec{f})$ is a convex function on $\overrightarrow{\mathcal{H}_K}$. Since $\|f\|_{\mathcal{H}_K}^2$ is a strictly convex function on $\overrightarrow{\mathcal{H}_K}$ and $\lambda > 0$, we know

$$\mathcal{E}_\rho(\vec{h}) + \lambda \|h\|_{\mathcal{H}_K}^2$$

is strictly convex on $\overrightarrow{\mathcal{H}_K}$ and the optimal solution \vec{h}_λ is unique. By the same way we can show the uniqueness of $\vec{f}_{z,\lambda}$.

By the definition of \vec{h}_λ we have

$$\lambda \|h^{(\lambda)}\|_{\mathcal{H}_K}^2 \leq \mathcal{E}_\rho(\vec{h}_\lambda) - \mathcal{E}_\rho(\vec{f}_\rho) + \lambda \|h_\lambda\|_{\mathcal{H}_K}^2 = D(\vec{f}_\rho, \lambda).$$

(2.14) then holds. By the definition of $\vec{f}_{z,\lambda}$ we have

$$\mathcal{E}_z(\vec{f}_{z,\lambda}) + \lambda \|f_{z,\lambda}\|_{\mathcal{H}_K}^2 \leq \mathcal{E}_z(\vec{0}) = \frac{1}{m} \sum_{i=1}^m \|\vec{y}_i\|_{\mathbb{R}^{d+1}}^2 \leq M^2.$$

(2.15) thus holds.

Lemma 2.6 Let $\overrightarrow{h_\lambda}$ be defined as (2.6) and let $\overrightarrow{f_{z,\lambda}}$ be the solution of (1.7). Then, there holds

$$\begin{aligned} \|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K} &\leq \frac{2}{\lambda} \left\| \int_Z \langle \nabla_x^* K_x(\cdot), \overrightarrow{y} - \overrightarrow{h_\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right. \\ &\quad \left. - \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \overrightarrow{y}_i - \overrightarrow{h_\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}} \right\|_{\mathcal{H}_K}. \end{aligned} \quad (2.16)$$

Proof By (2.9), we have

$$\begin{aligned} \mathcal{E}_z(\overrightarrow{f_{z,\lambda}}) - \mathcal{E}_z(\overrightarrow{h_\lambda}) &= \frac{1}{m} \sum_{i=1}^m \|\overrightarrow{y}_i - \overrightarrow{f_{z,\lambda}}(x_i)\|_{\mathbb{R}^{d+1}}^2 - \frac{1}{m} \sum_{i=1}^m \|\overrightarrow{y}_i - \overrightarrow{h_\lambda}(x_i)\|_{\mathbb{R}^{d+1}}^2 \\ &\geq \frac{1}{m} \sum_{i=1}^m \langle \overrightarrow{f_{z,\lambda}}(x_i) - \overrightarrow{h_\lambda}(x_i), -2(\overrightarrow{y}_i - \overrightarrow{h_\lambda}(x_i)) \rangle_{\mathbb{R}^{d+1}}. \end{aligned}$$

Let $f_0^{(z,\lambda)}(x) = f^{(z,\lambda)}(x)$ and $h^{(\lambda)} = h_0^{(\lambda)}$. Then, by (1.5) we have

$$\begin{aligned} &\langle \overrightarrow{f_{z,\lambda}}(x_i) - \overrightarrow{h_\lambda}(x_i), -2(\overrightarrow{y}_i - \overrightarrow{h_\lambda}(x_i)) \rangle_{\mathbb{R}^{d+1}} \\ &= -2 \sum_{j=0}^d \left(\frac{\partial f^{(z,\lambda)}(x_i)}{\partial x^j} - \frac{h^{(\lambda)}(x_i)}{\partial x^j} \right) \times (y_i^{(j)} - h_j^{(\lambda)}(x_i)) \\ &= -2 \sum_{j=0}^d \left(f^{(z,\lambda)} - h^{(\lambda)}, \frac{\partial K_{x_i}(\cdot)}{\partial x^j} \times (y_i^{(j)} - h_j^{(\lambda)}(x_i)) \right)_{\mathcal{H}_K} \\ &= \left(f^{(z,\lambda)} - h^{(\lambda)}, -2 \sum_{j=0}^d \frac{\partial K_{x_i}(\cdot)}{\partial x^j} \times (y_i^{(j)} - h_j^{(\lambda)}(x_i)) \right)_{\mathcal{H}_K} \\ &= (f^{(z,\lambda)} - h^{(\lambda)}, -2 \langle \nabla_x^* K_{x_i}(\cdot), (\overrightarrow{y}_i - \overrightarrow{h_\lambda}(x_i)) \rangle_{\mathbb{R}^{d+1}})_{\mathcal{H}_K}. \end{aligned}$$

Therefore,

$$\mathcal{E}_z(\overrightarrow{f_{z,\lambda}}) - \mathcal{E}_z(\overrightarrow{h_\lambda}) \geq \left(f^{(z,\lambda)} - h^{(\lambda)}, -\frac{2}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), (\overrightarrow{y}_i - \overrightarrow{h_\lambda}(x_i)) \rangle_{\mathbb{R}^{d+1}} \right)_{\mathcal{H}_K}. \quad (2.17)$$

By the definition of $\overrightarrow{f_{z,\lambda}}$ we have

$$0 \geq (\mathcal{E}_z(\overrightarrow{f_{z,\lambda}}) + \lambda \|f^{(z,\lambda)}\|_{\mathcal{H}_K}^2) - (\mathcal{E}_z(\overrightarrow{h_\lambda}) + \lambda \|h^{(\lambda)}\|_{\mathcal{H}_K}^2). \quad (2.18)$$

Since \mathcal{H}_K is a Hilbert space, the parallelogram law shows

$$\|f^{(z,\lambda)}\|_{\mathcal{H}_K}^2 - \|h^{(\lambda)}\|_{\mathcal{H}_K}^2 = (f^{(z,\lambda)} - h^{(\lambda)}, 2h^{(\lambda)})_K + \|f^{(z,\lambda)} - h^{(\lambda)}\|_{\mathcal{H}_K}^2.$$

It follows by (2.17)–(2.18) that

$$\begin{aligned} 0 &\geq \left(f^{(z,\lambda)} - h^{(\lambda)}, -\frac{2}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), (\overrightarrow{y}_i - \overrightarrow{h_\lambda}(x_i)) \rangle_{\mathbb{R}^{d+1}} \right)_{\mathcal{H}_K} \\ &\quad + (f^{(z,\lambda)} - h^{(\lambda)}, 2\lambda h^{(\lambda)}) + \lambda \|f^{(z,\lambda)} - h^{(\lambda)}\|_{\mathcal{H}_K}^2. \end{aligned}$$

By (2.11) we have

$$\begin{aligned}
0 &\geq \left(f^{(z,\lambda)} - h^{(\lambda)}, -\frac{2}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), (\vec{y}_i - \vec{h}_\lambda(x_i)) \rangle_{\mathbb{R}^{d+1}} + 2\lambda h^{(\lambda)} \right)_{\mathcal{H}_K} \\
&\quad + \lambda \|f^{(z,\lambda)} - h^{(\lambda)}\|_{\mathcal{H}_K}^2 \\
&= 2 \left(f^{(z,\lambda)} - h^{(\lambda)}, \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{h}^{(\lambda)}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right. \\
&\quad \left. - \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), (\vec{y}_i - \vec{h}_\lambda(x_i)) \rangle_{\mathbb{R}^{d+1}} \right)_{\mathcal{H}_K} \\
&\quad + \lambda \|f^{(z,\lambda)} - h^{(\lambda)}\|_{\mathcal{H}_K}^2. \tag{2.19}
\end{aligned}$$

(2.19) gives (2.16).

Lemma 2.7 Let \vec{h}_λ be defined as (2.6) and let $\vec{f}_{z,\lambda}$ be the solution of (1.7). Then, there holds

$$\|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 \leq \frac{4}{\lambda} A(z) B(z), \tag{2.20}$$

where

$$A(z) = \left\| \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{h}_\lambda(x) \rangle_{\mathbb{R}^{d+1}} d\rho - \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{h}_\lambda(x_i) \rangle_{\mathbb{R}^{d+1}} \right\|_{\mathcal{H}_K}$$

and

$$B(z) = \left\| \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho - \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{f}_{z,\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}} \right\|_{\mathcal{H}_K}.$$

Proof By (2.9) we have

$$\begin{aligned}
&\mathcal{E}_\rho(\vec{h}_\lambda) - \mathcal{E}_\rho(\vec{f}_{z,\lambda}) \\
&= \left(\vec{h}_\lambda - \vec{f}_{z,\lambda}, -2 \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right)_{\mathcal{H}_K} + \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2. \tag{2.21}
\end{aligned}$$

The definitions of \vec{h}_λ yields

$$0 \geq (\mathcal{E}_\rho(\vec{h}_\lambda) + \|h^{(\lambda)}\|_{\mathcal{H}_K}^2) - (\mathcal{E}_\rho(\vec{f}_{z,\lambda}) + \|f^{(z,\lambda)}\|_{\mathcal{H}_K}^2). \tag{2.22}$$

Further, by (2.21) and the equality

$$\|h^{(\lambda)}\|_{\mathcal{H}_K}^2 - \|f^{(z,\lambda)}\|_{\mathcal{H}_K}^2 = (h^{(\lambda)} - f^{(z,\lambda)}, 2f^{(z,\lambda)})_K + \|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K}^2,$$

we have

$$\begin{aligned}
0 &\geq \left(h^{(\lambda)} - f^{(z,\lambda)}, -2 \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right)_{\mathcal{H}_K} \\
&\quad + \lambda (h^{(\lambda)} - f^{(z,\lambda)}, 2f^{(z,\lambda)})_{\mathcal{H}_K} + \lambda \|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K}^2 + \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 \\
&= 2 \left(h^{(\lambda)} - f^{(z,\lambda)}, - \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right. \\
&\quad \left. + \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{f}_{z,\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}} \right)_{\mathcal{H}_K} + \lambda \|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K}^2 \\
&\quad + \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2, \tag{2.23}
\end{aligned}$$

where we have used (2.12). By Cauchy's inequality we have

$$\begin{aligned} & \lambda \|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K}^2 + \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 \\ & \leq 2 \|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K} \times \left\| \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right. \\ & \quad \left. - \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{f}_{z,\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}} \right\|_{\mathcal{H}_K}. \end{aligned}$$

Above inequality leads to

$$\lambda \|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K} + \frac{\|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2}{\|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K}} \leq 2 B(z).$$

It follows that

$$\frac{\|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2}{\|h^{(\lambda)} - f^{(z,\lambda)}\|_{\mathcal{H}_K}} \leq 2 B(z). \quad (2.24)$$

By (2.16) and (2.24) we have (2.20).

Lemma 2.8 (see [19]) *Let $(H, \|\cdot\|)$ be a Hilbert space and ξ be a random variable on (Z, ρ) with values in H . Assume that $\|\xi\|_H \leq \tilde{M} < +\infty$ almost surely. Let $\{z_i\}_{i=1}^m$ be independent samples drawers of ρ . For any $0 < \delta < 1$, with confidence $1 - \delta$,*

$$\left\| \frac{1}{m} \sum_{i=1}^m \xi(z_i) - E(\xi) \right\|_H \leq \frac{2\tilde{M} \log(\frac{2}{\delta})}{\sqrt{m}}. \quad (2.25)$$

Lemma 2.9 (see [6]) *Let \mathcal{F} be a family of functions from a probability space Z to \mathbb{R} and $d(\cdot, \cdot)$ be a distance on \mathcal{F} . Let $\mathcal{U} \subset Z$ be of full measure and constants $H, p > 0$ such that*

- (i) $|\xi(z)| \leq H$ for all $\xi \in \mathcal{F}$ and all $z \in \mathcal{U}$, and
- (ii) $|L_z(\xi_1) - L_z(\xi_2)| \leq p d(\xi_1, \xi_2)$ for all $\xi_1, \xi_2 \in \mathcal{F}$ and all $z \in \mathcal{U}^m$, where

$$L_z(\xi) = \int_Z \xi(z) - \frac{1}{m} \sum_{i=1}^m \xi(z_i).$$

Then, for all $\epsilon > 0$,

$$\text{Prob}_{z \in Z^m} \left\{ \sup_{\xi \in \mathcal{F}} |L_z(\xi)| \leq \epsilon \right\} \geq 1 - \mathcal{N}\left(\mathcal{F}, \frac{\epsilon}{2p}\right) \times 2 \exp\left(-\frac{m \epsilon^2}{8 H^2}\right). \quad (2.26)$$

Lemma 2.10 *Under the conditions of Theorem 1.1, we have the following two estimates:*

- (1) *For any $\delta \in (0, 1)$, with confidence $1 - \frac{\delta}{2}$, we have*

$$A(z) \leq \frac{16 k_1^2 d \sqrt{D(\vec{f}_\rho, \lambda)}}{\sqrt{m\lambda}} \times \log \frac{4}{\delta}. \quad (2.27)$$

- (2) *For any $\delta \in (0, 1)$ and $0 < \delta < \frac{4}{e^{c_s}}$, with confidence $1 - \frac{\delta}{2}$, we have*

$$B(z) \leq \frac{4k_1^2 (d+1)^2 M}{\sqrt{\lambda}} \times \left(\frac{1}{\sqrt{m}} + \frac{1}{2^{+s}\sqrt{m}} \right) \log \frac{4}{\delta}. \quad (2.28)$$

Proof of (2.27) Take

$$\xi(x, \vec{y}, \cdot) = \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{h}_\lambda(x) \rangle_{\mathbb{R}^{d+1}}, \quad (x, \vec{y}) \in X \times Y.$$

Since

$$\nabla_x^* K_x(\cdot) = \left(\frac{\partial K_x(\cdot)}{\partial x^0}, \frac{\partial K_x(\cdot)}{\partial x^1}, \dots, \frac{\partial K_x(\cdot)}{\partial x^d} \right)^T, \quad \frac{\partial K_x(\cdot)}{\partial x^0} = K_x(\cdot),$$

we have

$$\langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{h}_\lambda(x) \rangle_{\mathbb{R}^{d+1}} = \sum_{i=0}^d \frac{\partial K_x(\cdot)}{\partial x^i} \times (y^i - h_i^{(\lambda)}(x)). \quad (2.29)$$

By (1.5), we have

$$\begin{aligned} \|\xi(x, \vec{y}, \cdot)\|_{\mathcal{H}_K}^2 &= (\xi(x, \vec{y}, \cdot), \xi(x, \vec{y}, \cdot))_{\mathcal{H}_K} \\ &= \sum_{i,j=0}^d \left(\frac{\partial K_x(\cdot)}{\partial x^i}, \frac{\partial K_x(\cdot)}{\partial x^j} \right)_{\mathcal{H}_{K,\rho_X}} \times (y^i - h_i^{(\lambda)}(x)) \times (y^j - h_j^{(\lambda)}(x)) \\ &= \sum_{i,j=0}^d \frac{\partial^2 K_x(x)}{\partial x^i \partial x^j} \times (y^i - h_i^{(\lambda)}(x)) \times (y^j - h_j^{(\lambda)}(x)) \\ &\leq k_1^2 d^2 \sum_{i=0}^d |y^i - h_i^{(\lambda)}(x)|^2 \\ &= k_1^2 d^2 \|\vec{y} - \vec{h}_\lambda(x)\|_{\mathbb{R}^{d+1}}^2 \\ &\leq 4 k_1^2 d^2 (M^2 + \|\vec{h}_\lambda(x)\|_{\mathbb{R}^{d+1}}^2). \end{aligned} \quad (2.30)$$

Furthermore, by (2.14) we have

$$\begin{aligned} \|\vec{h}_\lambda(x)\|_{\mathbb{R}^{d+1}}^2 &= \sum_{i=0}^d \left| \frac{\partial h^{(\lambda)}(x)}{\partial x^i} \right|^2 \\ &\leq k_1^2 d \|h^{(\lambda)}\|_{\mathcal{H}_K}^2 \leq \frac{k_1^2 d D(\vec{f}_\rho, \lambda)}{\lambda}. \end{aligned} \quad (2.31)$$

Therefore, if $\lambda \leq \frac{k_1^2 d D(\vec{f}_\rho, \lambda)}{M^2}$, then

$$\begin{aligned} \|\xi(x, \vec{y}, \cdot)\|_{\mathcal{H}_K}^2 &\leq 4k_1^2 \left(M^2 + \frac{k_1^2 d D(\vec{f}_\rho, \lambda)}{\lambda} \right) \\ &\leq \frac{8k_1^4 d D(\vec{f}_\rho, \lambda)}{\lambda}. \end{aligned} \quad (2.32)$$

By (2.32) and (2.25), we have (2.27).

Proof of (2.28) By the definition of $\|\cdot\|_{\mathcal{H}_K}$, we have

$$\begin{aligned} B &= \sup_{\|g\|_{\mathcal{H}_K} \leq 1} \left| \left(g, \int_Z \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right. \right. \\ &\quad \left. \left. - \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{f}_{z,\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}} \right)_{\mathcal{H}_K} \right| \\ &= \sup_{\|g\|_{\mathcal{H}_K} \leq 1} \left| \int_Z \left(g, \langle \nabla_x^* K_x(\cdot), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho \right)_{\mathcal{H}_K} \right. \\ &\quad \left. - \left(g, \frac{1}{m} \sum_{i=1}^m \langle \nabla_x^* K_{x_i}(\cdot), \vec{y}_i - \vec{f}_{z,\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}} \right)_{\mathcal{H}_K} \right|. \end{aligned}$$

By (2.10), we have

$$B = \sup_{\|g\|_{\mathcal{H}_K} \leq 1} \left| \int_Z \langle \vec{g}(x), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}} d\rho - \frac{1}{m} \sum_{i=1}^m \langle \vec{g}(x_i), \vec{y}_i - \vec{f}_{z,\lambda}(x_i) \rangle_{\mathbb{R}^{d+1}} \right|.$$

Define

$$\eta(x, \vec{y}) = \langle \vec{g}(x), \vec{y} - \vec{f}_{z,\lambda}(x) \rangle_{\mathbb{R}^{d+1}}.$$

Then

$$|\eta(x, \vec{y})| \leq \|\vec{g}(x)\|_{\mathbb{R}^{d+1}} \times \|\vec{y} - \vec{f}_{z,\lambda}(x)\|_{\mathbb{R}^{d+1}}. \quad (2.33)$$

By (1.6), we have

$$\begin{aligned} \|\vec{g}(x)\|_{\mathbb{R}^{d+1}} &= \left(\sum_{i=0}^d |g_i(x)|^2 \right)^{\frac{1}{2}} \\ &= \left(\sum_{i=0}^d \left| \frac{\partial g(x)}{\partial x^i} \right|^2 \right)^{\frac{1}{2}} \\ &\leq (d+1)k_1 \|g\|_{\mathcal{H}_K} \\ &\leq (d+1)k_1, \end{aligned} \quad (2.34)$$

and by (2.15), we have

$$\begin{aligned} \|\vec{y} - \vec{f}_{z,\lambda}(x)\|_{\mathbb{R}^{d+1}} &\leq \|\vec{y}\|_{\mathbb{R}^{d+1}} + \|\vec{f}_{z,\lambda}(x)\|_{\mathbb{R}^{d+1}} \\ &\leq (d+1)(M + k_1 \|f^{(z,\lambda)}\|_{\mathcal{H}_K}) \\ &\leq (d+1)M \left(1 + \frac{k_1}{\sqrt{\lambda}} \right). \end{aligned}$$

Therefore,

$$\begin{aligned} |\eta(x, \vec{y})| &\leq (d+1)^2 k_1 M \left(1 + \frac{k_1}{\sqrt{\lambda}} \right) \\ &\leq H = \frac{2(d+1)^2 k_1^2 M}{\sqrt{\lambda}} \end{aligned} \quad (2.35)$$

if $\lambda \leq k_1^2$.

Take

$$\mathcal{F} = \{\eta(x, \vec{y}) \in C(Z) : |\eta(x, \vec{y})| \leq H\}.$$

Then,

$$\begin{aligned} B &= \sup_{\eta \in \mathcal{F}} \left| \int_Z \eta(x, \vec{y}) d\rho - \sum_{i=1}^m \eta(x_i, \vec{y}_i) \right| \\ &= \sup_{\eta \in \mathcal{F}} |\mathcal{L}_z(\eta)|, \end{aligned} \quad (2.36)$$

where

$$\mathcal{L}_z(\eta) = \int_Z \eta(x, \vec{y}) d\rho - \frac{1}{m} \sum_{i=1}^m \eta(x_i, \vec{y}_i), \quad z = (x, \vec{y})$$

satisfies the inequality

$$|\mathcal{L}_z(\eta_1) - \mathcal{L}_z(\eta_2)| \leq 2\|\eta_1 - \eta_2\|_{C(Z)}, \quad \eta_1 \in \mathcal{F}, \eta_2 \in \mathcal{F}.$$

By Lemma 2.9, we have for any $\varepsilon > 0$

$$\begin{aligned} \text{Prob}_{z \in Z^m} (|\mathcal{L}_z(\eta)| \leq \varepsilon) &\geq 1 - \mathcal{N}\left(\mathcal{F}, \frac{\varepsilon}{4}\right) 2 \exp\left(-\frac{m\varepsilon^2}{8H^2}\right) \\ &\geq 1 - 2 \exp\left(c_s \left(\frac{4H}{\varepsilon}\right)^s - \frac{m\varepsilon^2}{8H^2}\right). \end{aligned}$$

Take

$$2 \exp\left(c_s \left(\frac{4H}{\varepsilon}\right)^s - \frac{m\varepsilon^2}{8H^2}\right) = \delta.$$

Then, we have

$$\varepsilon^{2+s} - \frac{8H^2\varepsilon^2}{m} \log \frac{4}{\delta} - \frac{8 \times 4^s c_s H^{2+s}}{m} = 0. \quad (2.37)$$

Notice that there is the famous lemma (see [26]).

Let $c_1 > 0$, $c_2 > 0$ and $u > t > 0$. Then the equation

$$x^u - c_1 x^t - c_2 = 0$$

has a unique positive zero x^* . In addition,

$$x^* \leq \max\{(2c_1)^{\frac{1}{u-t}}, (2c_2)^{\frac{1}{u}}\}. \quad (2.38)$$

By (2.37)–(2.38), we have

$$\begin{aligned} \varepsilon &\leq \sqrt{\frac{16H^2}{m} \log \frac{4}{\delta}} + \sqrt[2+s]{\frac{16 \times 4^s c_s H^{2+s}}{m}} \\ &= \frac{8(d+1)^2 k_1^2}{\sqrt{\lambda}} \left(\sqrt{\frac{1}{m} \log \frac{4}{\delta}} + \sqrt[2+s]{\frac{c_s}{m}} \right). \end{aligned} \quad (2.39)$$

By (2.36) and (2.39), we have (2.28).

Lemma 2.11 For any $\delta \in (0, 1)$, with confidence $1 - \delta$, we have

$$\begin{aligned} \|\vec{f_{z,\lambda}} - \vec{h_\lambda}\|_{L^2(\rho_X)}^2 &\leq \frac{512 k_1^4 (d+1)^3 M \sqrt{D(\vec{f_\rho}, \lambda)}}{\lambda^2 \sqrt{m}} \\ &\quad \times \left(\frac{1}{\sqrt{m}} + \frac{1}{\sqrt[2+s]{m}} \right) \times \left(\log \frac{4}{\delta} \right)^2. \end{aligned} \quad (2.40)$$

Proof By (2.20), we have

$$\frac{\lambda}{4} \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 \leq A(z) B(z).$$

Then,

$$\begin{aligned} & \left\{ z \in Z^m : \frac{\lambda}{4} \|\vec{f}_{z,\lambda} - \vec{h}_\lambda\|_{L^2(\rho_X)}^2 \leq \frac{128 k_1^4 (d+1)^3 M \sqrt{D(\vec{f}_\rho, \lambda)}}{m\lambda} \times \left(\log \frac{4}{\delta} \right)^2 \right\} \\ & \supset \left\{ z \in Z^m : A(z) \leq \frac{16 k_1^2 d \sqrt{D(\vec{f}_\rho, \lambda)}}{\sqrt{m\lambda}} \times \log \frac{4}{\delta} \right\} \\ & \cap \left\{ z \in Z^m : B(z) \leq \frac{8 k_1^2 (d+1)^2 M}{\sqrt{\lambda}} \times \left(\frac{1}{\sqrt{m}} + \frac{1}{\sqrt[2+\varepsilon]{m}} \right) \times \log \frac{4}{\delta} \right\}. \end{aligned} \quad (2.41)$$

By (2.27)–(2.28) and (2.41), we have (2.40).

Proof of (1.9) By the definition of \vec{h}_λ , (2.3)–(2.4) we have

$$\begin{aligned} \|\vec{f}_\rho - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 & \leq 4 \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 + 4 \|\vec{f}_\rho - \vec{h}_\lambda\|_{L^2(\rho_X)}^2 \\ & \leq 4 \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 + 4 (\|\vec{f}_\rho - \vec{h}_\lambda\|_{L^2(\rho_X)}^2 + \lambda \|h_\lambda\|_{\mathcal{H}_K}^2) \\ & = 4 \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 + 4 (\mathcal{E}_\rho(\vec{h}_\lambda) - \mathcal{E}(\vec{f}_\rho) + \lambda \|h_\lambda\|_{\mathcal{H}_K}^2) \\ & = 4 \|\vec{h}_\lambda - \vec{f}_{z,\lambda}\|_{L^2(\rho_X)}^2 + 4 D(\vec{f}_\rho, \lambda). \end{aligned} \quad (2.42)$$

(2.40) and (2.42) give (1.9).

References

- [1] Zhou, D. X., Derivative reproducing properties for kernel methods in learning theory, *Journal of Computational and Applied Mathematics*, **220**, 2008, 456–463.
- [2] Shi, L., Guo, X. and Zhou, D. X., Hermite learning with gradient data, *Journal of Computational and Applied Mathematics*, **233**, 2010, 3046–3056.
- [3] Mukherjee, S., Wu, Q. and Zhou, D. X., Learning gradients on manifolds, *Bernoulli*, **16**(1), 2010, 181–207.
- [4] Ying, Y. M., Wu, Q. and Campbell, C., Learning the coordinate gradients, *Advances in Computational Mathematics*, **37**(3), 2012, 355–378.
- [5] Cucker, F. and Smale, F., On the mathematical foundations of learning theory, *Bull. Amer. Math.*, **39**, 2001, 1–49.
- [6] Cucker, F. and Zhou, D. X., *Learning Theory: An Approximation Theory Viewpoint*, Cambridge University Press, New York, 2007.
- [7] Aronszajn, N., Theory of reproducing kernels, *Transactions of the American Mathematical Society*, **68**, 1950, 337–404.
- [8] Belkin, M., Niyogi, P. and Sindhvani, V., Manifold regularization: A geometric framework for learning from labeled and unlabeled examples, *Journal of Machine Learning Research*, **7**, 2006, 2399–2434.
- [9] Smale, S. and Zhou, D. X., Estimating the approximation error in learning theory, *Analysis and Applications*, **1**(1), 2003, 17–41.
- [10] Solnon, M., Arlot, S. and Bach, F., Multi-task regression using minimal penalties, *Journal of Machine Learning Research*, **13**, 2012, 2773–2812.
- [11] Caponnetto, A., Micchelli, C. A., Pontil, M. and Ying, Y. M., Universal multi-task kernels, *Journal of Machine Learning Research*, **9**, 2008, 1615–1646.
- [12] Wang, J. Y., Bensmail, H. and Gao, X., Feature selection and multi-kernel learning for sparse representation on a manifold, *Neural Networks*, **51**, 2014, 9–16.

- [13] Evgeniou, T., Micchelli, C. A. and Pontil, M., Learning multiple tasks with kernel methods, *Journal of Machine Learning Research*, **6**, 2005, 615–637.
- [14] Zhang, H. Z. and Zhang, J., Vector-valued reproducing kernel Banach spaces with applications to multi-task learning, *Journal of Complexity*, **29**(2), 2013, 195–215.
- [15] Jordão, J. and Menegatto, V. A., Reproducing properties of differentiable Mercer-like kernels on the sphere, *Numerical Functional Analysis and Optimization*, **33**(10), 2012, 1221–1243.
- [16] Ferreira, J. C. and Menegatto, V. A., Reproducing properties of differentiable Mercer-like kernels, *Mathematische Nachrichten*, **285**(8–9), 2012, 959–973.
- [17] Jordão, T. and Menegatto, V. A., Weighted Fourier-Laplace transforms in reproducing kernel Hilbert spaces on the sphere, *Journal of Mathematical Analysis and Applications*, **411**, 2014, 732–741.
- [18] Sun, H. W. and Wu, Q., A note on application of integral operator in learning, *Applied and Computational Harmonic Analysis*, **26**, 2009, 416–421.
- [19] Smale, S. and Zhou, D. X., Learning theory estimates via integral operators and their applications, *Constructive Approximation*, **26**, 2007, 153–172.
- [20] Sheng, B. H. and Ye, P. X., The learning rates of regularized regression based on reproducing kernel Banach spaces, *Abstract and Applied Analysis*, 2013, Article ID 694181, 10 pages, [http://dx. doi .org/10.1155/2013/694181](http://dx.doi.org/10.1155/2013/694181).
- [21] Sheng, B. H., The convergence rates of Shannon sampling learning algorithms, *Sciences in China: Mathematics*, **55**(6), 2012, 1243–1256.
- [22] Christmann, A. and Steinwart, I., Consistency and robustness of kernel-based regression in convex risk minimization, *Bernoulli*, **13**(3), 2007, 799–819.
- [23] Steinwart, I., Sparseness of support vector machines, *Journal of Machine Learning Research*, **4**, 2003, 1071–1105.
- [24] Sheng, B. H., Xiang, D. H. and Ye, P. X., Convergence rate of semi-supervised gradient learning algorithms, *International Journal of Wavelets, Multiresolution and Information Processing*, **13**(4), 2015, 1550021 (26 pages).
- [25] Bauschke, H. H. and Combettes, P. L., *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Springer-Verlag, New York, 2010.
- [26] Cucker, F. and Smale, S., Best choices for regularization parameters in learning theory: On the bias-variance problem, *Foundation of Computational Mathematics*, **2**, 2002, 413–428.