
Learning coordinate gradients with multi-task kernels

Yiming Ying and Colin Campbell*

Department of Engineering Mathematics, University of Bristol
Queens Building, University Walk, Bristol, BS8 1TR, UK.
{enxyy, C.Campbell}@bris.ac.uk

Abstract

Coordinate gradient learning is motivated by the problem of variable selection and determining variable covariation. In this paper we propose a novel unifying framework for coordinate gradient learning (MGL) from the perspective of multi-task learning. Our approach relies on multi-task kernels to simulate the structure of gradient learning. This has several appealing properties. Firstly, it allows us to introduce a novel algorithm which appropriately captures the inherent structure of coordinate gradient learning. Secondly, this approach gives rise to a clear algorithmic process: a computational optimization algorithm which is memory and time efficient. Finally, a statistical error analysis ensures convergence of the estimated function and its gradient to the true function and true gradient. We report some preliminary experiments to validate MGL for variable selection as well as determining variable covariation.

1 Introduction

Let $X \subseteq \mathbb{R}^d$ be compact, $Y \subseteq \mathbb{R}$, $Z = X \times Y$, and $\mathbb{N}_n = \{1, 2, \dots, n\}$ for any $n \in \mathbb{N}$. A common theme in machine learning is to learn a target function $f_* : X \rightarrow Y$ from a finite set of input/output samples $\mathbf{z} = \{(x_i, y_i) : i \in \mathbb{N}_m\} \subseteq Z$. However, in many applications, we not only wish to learn the target function, but also want to find which variables are salient and how these variables interact with each other. This problem has practical motivations: to facilitate data visualization and dimensionality reduction, for example. Such a motivation is important when there are many redundant variables and we wish to find the salient features among these. These problems can occur in many contexts. For example, with gene expression array datasets, the vast majority of features may be redundant to a classification task and we need to find a small set of genuinely distinguishing features. These motivations have driven the design of various statistical and machine learning models [8, 11, 21, 22] for variable (feature) selection.

Here, we build on previous contributions [15, 16, 17] by addressing coordinate gradient learning and its use for variable

selection and covariation learning, the interaction between variables. Specifically, for any $x \in X$, we denote x by (x^1, x^2, \dots, x^d) . The target is to learn the gradient of f_* (if it exists) denoted by a vector-valued function $\nabla f_*(x) = (\frac{\partial f_*}{\partial x^1}, \dots, \frac{\partial f_*}{\partial x^d})$. The intuition behind gradient learning for variable selection and coordinate covariation is the following. The inner product between components of ∇f_* indicates the interaction between coordinate variables. Specific norms of $\frac{\partial f_*}{\partial x^p}$ can indicate the salience of the p -th variable: the smaller the norm is, the less important this variable will be.

In this paper we propose a novel unifying formulation of coordinate gradient learning from the perspective of multi-task learning. Learning multiple tasks together has been extensively studied both theoretically and practically in several papers [1, 2, 5, 7, 13, 14]. One way to frame this problem is to learn a vector-valued function where each of its components is a real-valued function and corresponds to a particular task. A key objective in this formulation is to capture an appropriate structure among tasks so that common information is shared across tasks. Here we follow this methodology and employ a vector-valued function $\vec{f} = (f_1, \vec{f}_2) = (f_1, f_2, \dots, f_{d+1})$, where f_1 is used to simulate f_* , and \vec{f}_2 is used to simulate its gradient ∇f_* . We assume that \vec{f} comes from a vector-valued reproducing kernel Hilbert space (RKHS) associated with multi-task (matrix-valued) kernels, see [14]. The rich structure of RKHS space reflects the latent structure of multi-task gradient learning, i.e. the pooling of information across components (tasks) of ∇f_* using multi-task kernels.

The paper is organized as follows. In Section 2, we first review the definition of multi-task kernels and vector-valued RKHS. Then, we propose a unifying formulation of coordinate gradient learning from the perspective of multi-task learning which is referred to as multi-task gradient learning (MGL). The choices of multi-task kernels motivate different learning models [11, 15, 16, 17]. This allows us to introduce a novel choice of multi-task kernel which reveals the inherent structure of gradient learning. Kernel methods [19, 20] usually enjoy the representer theorem which paves the way for designing efficient optimization algorithms. In Section 3 we explore a representer theorem for MGL algorithms. Subsequently, in Section 4 we discuss computational optimization approaches for MGL algorithms, mainly focusing on least square loss and the SVM counterpart for gradient learning.

*We acknowledge support from EPSRC grant EP/E027296/1.

A statistical error analysis in Section 5 ensures the convergence of the estimated function and its gradient to the true function and true gradient. Finally, in Section 6 preliminary numerical experiments are reported to validate our proposed approach.

1.1 Related work

A number of machine learning and statistical models have been proposed for variable (feature) selection. Least absolute shrinkage and selection operator (LASSO) [21] and basis pursuit denoising [8] suggest use of ℓ^1 regularization to remove redundant features. Weston *et al* [22] introduced a method for selecting features by minimizing bounds on the leave-one-out error.

Guyon *et al* [11] proposed recursive feature elimination (RFE) which used a linear kernel SVM: variables with least influence on the weights $\frac{1}{2}\|w\|^2$ are considered least important. Although these algorithms are promising, there remain unresolved issues. For example, they do not indicate variable covariation and the extension of these algorithms to the non-linear case was marginally discussed. Our method outlined here covers variable covariation and nonlinear feature selection. As such, in Section 2, we show that RFE-SVM is a special case of our multi-task formulation.

Motivated by the Taylor expansion of a function at samples $\{x_i : i \in \mathbb{N}_m\}$, Mukherjee *et al* [15, 16, 17] proposed an algorithm for learning the gradient function. They used the norm of its components for variable (feature) selection and spectral decomposition of the covariance of the learned gradient function for dimension reduction [16]. Specifically, let \mathcal{H}_G be a scalar RKHS (see e.g. [3]) and use $f_1 \in \mathcal{H}_G$ to simulate f_* . For any $p \in \mathbb{N}_d$, a function $f_{p+1} \in \mathcal{H}_G$ is used to learn $\partial f_*/\partial x^p$. The results presented by Mukherjee *et al* are quite promising both theoretically and practically, but there is no pooling information shared across the components of the gradient. This may lead to less accurate approximation to the true gradient. We will address all these issues in our unifying framework.

2 Multi-task kernels and learning gradients

In this section we formulate the gradient learning problem from the perspective of multi-task learning. Specifically, we employ a vector-valued RKHS to simulate the target function and its gradient. The abundant structure of vector-valued RKHS enables us to couple information across components of the gradient in terms of multi-task kernels.

2.1 Multi-task model for gradient learning

We begin with a review of the definition of multi-task kernels and introduce vector-valued RKHS (see [14] and the reference therein). Throughout this paper, we use the notation $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ to denote the standard Euclidean inner product and norm respectively.

Definition 1 We say that a function $\mathcal{K} : X \times X \rightarrow \mathbb{R}^{d+1} \times \mathbb{R}^{d+1}$ is a multi-task (matrix-valued) kernel on X if, for any $x, t \in X$, $\mathcal{K}(x, t)^T = \mathcal{K}(t, x)$, and it is positive semi-definite, i.e., for any $m \in \mathbb{N}$, $\{x_j \in X : j \in \mathbb{N}_m\}$ and $\{y_j \in \mathbb{R}^{d+1} :$

$j \in \mathbb{N}_m\}$ there holds

$$\sum_{i,j \in \mathbb{N}_m} \langle y_i, \mathcal{K}(x_i, x_j) y_j \rangle \geq 0. \quad (1)$$

In the spirit of Moore-Aronszajn's theorem, there exists a one-to-one correspondence between the multi-task kernel \mathcal{K} with property (1) and a vector-valued RKHS of functions $\vec{f} : X \rightarrow \mathbb{R}^{d+1}$ with norm $\langle \cdot, \cdot \rangle_{\mathcal{K}}$ denoted by $\mathcal{H}_{\mathcal{K}}$, see e.g. [14]. Moreover, for any $x \in X$, $y \in \mathbb{R}^{d+1}$ and $\vec{f} \in \mathcal{H}_{\mathcal{K}}$, we have the reproducing property

$$\langle \vec{f}(x), y \rangle = \langle \vec{f}, \mathcal{K}_{x,y} \rangle_{\mathcal{K}} \quad (2)$$

where $\mathcal{K}_{x,y} : X \rightarrow \mathbb{R}^{d+1}$ is defined, for any $t \in X$, by $\mathcal{K}_{x,y}(t) := \mathcal{K}(t, x)y$.

In the following we describe our multi-task kernel-based framework for gradient learning. Following Mukherjee *et al* [15, 17], the derivation of gradient learning can be motivated by the Taylor expansion of $f_* : f_*(x_i) \approx f_*(x_j) + \nabla f_*(x_j)(x_i - x_j)^T$. Since we wish to learn f_* with f_1 and ∇f_* with \vec{f}_2 , replacing $f_*(x_i)$ by y_i , the error¹

$$y_i \approx f_1(x_j) + \vec{f}_2(x_j)(x_i - x_j)^T$$

is expected to be small whenever x_i is close to x_j . To enforce the constraint that x_i is close to x_j , we introduce a weight function produced by a Gaussian with deviation s defined by $w_{ij} = \frac{1}{s^{d+2}} e^{-\frac{\|x_i - x_j\|^2}{2s^2}}$. This implies that $w_{ij} \approx 0$ if x_i is far away from x_j .

We now propose the following multi-task formulation for gradient learning (MGL):

$$\vec{f}_{\mathbf{z}} = \arg \min_{\vec{f} \in \mathcal{H}_{\mathcal{K}}} \left\{ \frac{1}{m^2} \sum_{i,j} w_{ij} L(y_i, f_1(x_j) + \vec{f}_2(x_j)(x_i - x_j)^T) + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \right\}. \quad (3)$$

where $L : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$ is a prescribed loss function and λ is usually called the *regularization parameter*. The minimum is taken over a vector-valued RKHS with multi-task kernel \mathcal{K} . The first component $f_{1,\mathbf{z}}$ of the minimizer $\vec{f}_{\mathbf{z}}$ of the above algorithm is used to simulate the target function and the other components $\vec{f}_{2,\mathbf{z}} := (f_{2,\mathbf{z}}, \dots, f_{d+1,\mathbf{z}})$ to learn its gradient function. In Section 6, we will discuss how to use the solution $\vec{f}_{\mathbf{z}}$ for variable selection as well as covariation measurement.

Different choice of loss functions yield different gradient learning algorithms. For instance, if the loss function $L(y, t) = (y - t)^2$ then algorithm (3) leads to the least-square multi-task gradient learning (LSMGL):

$$\arg \min_{\vec{f} \in \mathcal{H}_{\mathcal{K}}} \left\{ \frac{1}{m^2} \sum_{i,j \in \mathbb{N}_m} w_{ij} [y_i - f_1(x_j) - \vec{f}_2(x_j)(x_i - x_j)^T]^2 + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \right\}. \quad (4)$$

In classification, the choice of loss function $L(y, t) = (1 - yt)_+$ in algorithm (3) yields the support vector machine for multi-task gradient learning (SVMMGL):

¹Our form of Taylor expansion is slightly different from that used in [15, 17]. However, the essential idea is the same.

$$\arg \min_{\substack{\vec{f} \in \mathcal{H}_{\mathcal{K}} \\ b \in \mathbb{R}}} \left\{ \frac{1}{m^2} \sum_{i,j \in \mathbb{N}_m} w_{ij} [1 - y_i (f_1(x_j) + b + \vec{f}_2(x_j)(x_i - x_j)^T)]_+ + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \right\}. \quad (5)$$

Here, $f_1(x) + b$ is used to learn the target function and \vec{f}_2 , simulating the gradient of the target function. Hence b plays the same role of offset as in the standard SVM formulation. In this case, at each point x_i the error between the output y_i and $f(x_i)$ is now replaced by the error between y_i and the first order Taylor expansion of $f(x_i)$ at x_j , i.e., $f_1(x_j) + \vec{f}_2(x_j)(x_i - x_j)^T$.

2.2 Choice of multi-task kernels

We note that if \mathcal{K} is a diagonal matrix-valued kernel, then each component of a vector-valued function in the associated RKHS of \mathcal{K} can be represented, independently of the other components, as a function in the RKHS of a scalar kernel. Consequently, for a scalar kernel G if we choose the multi-task kernel \mathcal{K} given, for any $x, t \in X$, by $\mathcal{K}(x, t) = G(x, t)I_{d+1}$ then the MGL algorithm (3) is reduced to the gradient learning algorithm proposed in [15, 16, 17] using $(d+1)$ -folds of scalar RKHS. There, under some conditions on the underlying distribution ρ , it has been proven that $f_{1,\mathbf{z}} \rightarrow f_*$ and $\vec{f}_{2\mathbf{z}} \rightarrow \nabla f_*$ when the number of samples tends to infinity. Although their results are promising both theoretically and practically, a more inherent structure would be $\vec{f}_{2\mathbf{z}} = \nabla f_{1,\mathbf{z}}$. In our MGL framework (3), we can recover this structure by choosing the multi-task kernel appropriately.

Our alternative choice of multi-task kernel is stimulated by the Hessian of Gaussian kernel proposed in [7]. For any scalar kernel G and any $x, t \in X$, we introduce the function

$$\mathcal{K}(x, t) = \begin{pmatrix} G(x, t), & (\nabla_t G(x, t))^T \\ \nabla_x G(x, t), & \nabla_{xt}^2 G(x, t) \end{pmatrix} \quad (6)$$

which we will show to be a multi-task kernel. To see this, let ℓ^2 be the Hilbert space with norm $\|w\|_{\ell^2}^2 = \sum_{j=1}^{\infty} w_j^2$. Suppose that G has a feature representation, i.e., $G(x, t) = \langle \phi(x), \phi(t) \rangle_{\ell^2}$ and, for any $f \in \mathcal{H}_G$, there exists a vector $w \in \ell^2$ such that $f(x) = \langle w, \phi(x) \rangle_{\ell^2}$ and

$$\|f\|_G = \|w\|_{\ell^2}.$$

Indeed, if the input space X is compact and $G : X \times X \rightarrow \mathbb{R}$ is a Mercer kernel, i.e., it is continuous, symmetric and positive semi-definite, then, according to Mercer theorem, G always has the above feature representation (see e.g. [9]).

Now we have the following proposition about \mathcal{K} defined by equation (6). Let \tilde{e}_p be the p -th coordinate basis in \mathbb{R}^{d+1} .

Theorem 2 *For any smooth scalar Mercer kernel G , define function \mathcal{K} by equation (6). Then, \mathcal{K} is a multi-task kernel and, for any $\vec{f} = (f_1, \vec{f}_2) \in \mathcal{H}_{\mathcal{K}}$ there holds*

$$\vec{f}_2 = \nabla f_1. \quad (7)$$

Proof: Since G is a scalar kernel, for any $x, t \in X$ we have that $G(x, t) = G(t, x)$. Therefore, $\mathcal{K}(x, t)^T = \mathcal{K}(t, x)$.

Moreover, G is assumed to be a Mercer kernel which implies that it has a feature representation $G(x, t) = \langle \phi(x), \phi(t) \rangle_{\ell^2}$. Consequently, $\nabla_t G(x, t) = \langle \phi(x), \nabla \phi(t) \rangle_{\ell^2}$ and $\nabla_x G(x, t) = \langle \nabla \phi(x), \phi(t) \rangle_{\ell^2}$, and $\nabla_{xt}^2 G(x, t) = \langle \nabla \phi(x), \nabla \phi(t) \rangle_{\ell^2}$. Then, we introduce, for any $w \in \ell^2, x \in X, \mathbf{y} \in \mathbb{R}^{d+1}$, the feature map $\Phi(x) : \ell^2 \rightarrow \mathbb{R}^{d+1}$ defined by

$$\Phi(x)w := (\langle \phi(x), w \rangle_{\ell^2}, \langle \partial_1 \phi(x), w \rangle_{\ell^2}, \dots, \langle \partial_d \phi(x), w \rangle_{\ell^2})^T.$$

Its adjoint map Φ^* is given, for any $t \in X$ and $\mathbf{y} \in \mathbb{R}^{d+1}$, by $\Phi^*(t)\mathbf{y} := \phi(x)\mathbf{y}^1 + \sum_{p \in \mathbb{N}_d} \partial_p \phi(x)\mathbf{y}^{p+1}$. Hence, $\mathcal{K}(x, t)\mathbf{y} = \Phi(x)\Phi^*(t)\mathbf{y}$. Consequently, for any $m \in \mathbb{N}$, any $i, j \in \mathbb{N}_m$ and $\mathbf{y}_i, \mathbf{y}_j \in \mathbb{R}^{d+1}$, it follows $\sum_{i,j \in \mathbb{N}_m} \langle \mathbf{y}_i, K(x_i, x_j)\mathbf{y}_j \rangle = \|\sum_{i \in \mathbb{N}_m} \Phi^*(x_i)\mathbf{y}_i\|_{\ell^2}^2$ is nonnegative which tells us that \mathcal{K} is a multi-task kernel.

We turn to the second assertion. When \vec{f} is in the form of a finite combination of kernel section $\{\mathcal{K}_{x,\mathbf{y}} : \mathbf{y} \in \mathbb{R}^{d+1}, x \in X\}$, the second assertion follows directly from the definition of \mathcal{K} . For the general case, we use the fact that the vector-valued RKHS is the closure of the span of kernel sections, see [14]. To this end, assume that there exists a sequence $\{\vec{f}_j = (f_1^j, f_2^j, \dots, f_{d+1}^j)\}$ of finite combination of kernel sections such that $\vec{f}_j \rightarrow \vec{f} \in \mathcal{H}_{\mathcal{K}}$ w.r.t. the RKHS norm. Hence, by the reproducing property (2), for any $x \in X$ and $p \in \mathbb{N}_d$, $|f_{p+1}^j(x) - f_{p+1}(x)| = |\langle \tilde{e}_{p+1}, f_j^j(x) - \vec{f}(x) \rangle| = |\langle \vec{f}_j - \vec{f}, \mathcal{K}_x \tilde{e}_{p+1} \rangle_{\mathcal{K}}| \leq \|\vec{f}_j - \vec{f}\|_{\mathcal{K}} \sqrt{\tilde{e}_{p+1}^T \mathcal{K}(x, x) \tilde{e}_{p+1}}$ which tends to zeros as j tends to infinity. Consequently, it follows, for any $x \in X$,

$$f_{p+1}^j(x) \rightarrow f_{p+1}(x), \quad \text{as } j \rightarrow \infty. \quad (8)$$

Let $\delta_p \in \mathbb{R}^d$ be a vector with its p -th component $\delta > 0$ and others equal zero. Applying the reproducing property (2) yields that

$$\begin{aligned} & \left| \frac{[f_1^j(x+\delta_p) - f_1^j(x)] - [f_1(x+\delta_p) - f_1(x)]}{\delta} \right| \\ &= \left| \frac{1}{\delta} \langle \vec{f}_j - \vec{f}, \mathcal{K}_{x+\delta_p} \tilde{e}_1 - \mathcal{K}_x \tilde{e}_1 \rangle_{\mathcal{K}} \right| \\ &\leq \|\vec{f}_j - \vec{f}\|_{\mathcal{K}} \langle \tilde{e}_1, \frac{1}{\delta^2} [\mathcal{K}(x+\delta_p, x+\delta_p) + \mathcal{K}(x, x) - \mathcal{K}(x, x+\delta_p) - \mathcal{K}(x+\delta_p, x)] \tilde{e}_1 \rangle^{\frac{1}{2}} \\ &= \|\vec{f}_j - \vec{f}\|_{\mathcal{K}} \left(\frac{1}{\delta^2} [G(x+\delta_p, x+\delta_p) - G(x, x) - G(x, x+\delta_p) - G(x+\delta_p, x)] \right)^{\frac{1}{2}}. \end{aligned}$$

Since G is smooth and X is compact there exists an absolute constant $\tilde{c} > 0$ such that, for any $\delta > 0$, the above equation is furthermore bounded by

$$\left| \frac{[f_1^j(x+\delta_p) - f_1^j(x)] - [f_1(x+\delta_p) - f_1(x)]}{\delta} \right| \leq \tilde{c} \|\vec{f}_j - \vec{f}\|_{\mathcal{K}}.$$

Consequently, letting $\delta \rightarrow 0$ in the above equation it follows $|\partial_p f_1^j(x) - \partial_p f_1(x)| \rightarrow 0$ as j tends to infinity. Combining this with equation (8) and the fact that $f_{p+1}^j(x) = \partial_p f_1^j(x)$ implies that $\partial_p f_1(x) = f_{p+1}(x)$ which completes the theorem. \blacksquare

The scalar kernel G plays the role of a *hyper-parameter* to produce the multi-task kernel \mathcal{K} given by equation (6). By the above theorem, if we choose \mathcal{K} to be defined by equation (6) then any solution $f_{\mathbf{z}} = (f_{1,\mathbf{z}}, \vec{f}_{2\mathbf{z}})$ of algorithm (3) enjoys the structure $\vec{f}_{2\mathbf{z}} = \nabla f_{1,\mathbf{z}}$.

Further specifying the kernel G in the definition (6) of multi-task kernel \mathcal{K} , we can recover the RFE feature ranking algorithm for a linear SVM [11]. To see this, let G be a linear kernel. In the next section, we will see that, for any solution $\vec{f}_{\mathbf{z}}$ of MGL algorithm (3), there exists $\{c_{j,\mathbf{z}} \in \mathbb{R}^{d+1} : j \in \mathbb{N}_m\}$ such that $\vec{f}_{\mathbf{z}} = \sum_{j \in \mathbb{N}_m} \mathcal{K}_{x_j} c_{j,\mathbf{z}}$. Since G is linear, combining this with Theorem 2 we know that $f_{1,\mathbf{z}}(x) = W_{\mathbf{z}}^T x$ with $W_{\mathbf{z}} = \sum_j (x_j^T, 1) c_{j,\mathbf{z}} \in \mathbb{R}$ and $\vec{f}_{2\mathbf{z}} = \nabla f_{1,\mathbf{z}} = W_{\mathbf{z}}^T$. Consequently, in the case we have that

$$f_{1,\mathbf{z}}(x_j) + \vec{f}_{2\mathbf{z}}(x_j)(x_i - x_j) = W_{\mathbf{z}}^T x_i = f_{1,\mathbf{z}}(x_i).$$

Moreover, by the reproducing property (2) we can check that

$$\|\vec{f}_{\mathbf{z}}\|_{\mathcal{K}}^2 = \|f_{1,\mathbf{z}}\|_G^2 = \|W_{\mathbf{z}}\|^2.$$

Putting the above equations together, in this special case we know that the SVMMLG algorithm (5) is reduced, with the choice of $w_{ij} = 1$, to the classical learning algorithm:

$$\min_{W \in \mathbb{R}^d} \left\{ \frac{1}{m} \sum_{i \in \mathbb{N}_m} (1 - y_i(W^T x_i + b))_+ + \lambda \|W\|^2 \right\}.$$

Hence, our formulation of gradient learning (3) can be regarded as a generalization of RFE-SVM [11] to the nonlinear case.

In the subsequent sections we discuss a general representation theorem and computational optimization problems motivated by MGL algorithms.

3 Representer theorem

In this section we investigate the representer theorem for the MGL algorithm (3). This forms a foundation for the derivation of a computationally efficient algorithm for MGL in Section 4.

Recall that \tilde{e}_p is the p -th coordinate basis in \mathbb{R}^{d+1} and, for any $x \in \mathbb{R}^d$, denote the vector \tilde{x}^T by $(0, x^T)$. By the reproducing property (2), we have that $f_1(x_j) = \langle \vec{f}(x_j), \tilde{e}_1 \rangle = \langle \vec{f}, \mathcal{K}_{x_j} \tilde{e}_1 \rangle_{\mathcal{K}}$ and likewise, $\vec{f}_2(x_j)(x_i - x_j)^T = \langle \vec{f}(x_j), \tilde{x}_i - \tilde{x}_j \rangle = \langle \vec{f}, \mathcal{K}_{x_j}(\tilde{x}_i - \tilde{x}_j) \rangle_{\mathcal{K}}$. Then, the algorithm (3) can be rewritten by

$$\arg \min_{\vec{f} \in \mathcal{H}_{\mathcal{K}}} \left\{ \frac{1}{m^2} \sum_{i,j \in \mathbb{N}_m} w_{ij} L(y_i, \langle \vec{f}, \mathcal{K}_{x_j}(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j) \rangle_{\mathcal{K}}) + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \right\}. \quad (9)$$

In analogy with standard kernel methods [19, 20], we have the following representer theorem for MGL by using the properties of multi-task kernels.

Theorem 3 *For any multi-task kernel \mathcal{K} , consider the gradient learning algorithm (3). Then, there exists representer coefficients $\{c_{j,\mathbf{z}} \in \mathbb{R}^{d+1} : j \in \mathbb{N}_m\}$ such that*

$$\vec{f}_{\mathbf{z}} = \sum_{j \in \mathbb{N}_m} \mathcal{K}_{x_j} c_{j,\mathbf{z}}$$

and, for every $j \in \mathbb{N}_m$, the representer coefficient $c_{j,\mathbf{z}} \in \text{span}\{\tilde{e}_1, \tilde{x}_i : i \in \mathbb{N}_m\}$.

Proof: We can write any minimizer $\vec{f}_{\mathbf{z}} \in \mathcal{H}_{\mathcal{K}}$ as $\vec{f}_{\mathbf{z}} = \vec{f}_{\parallel} + \vec{f}_{\perp}$ where \vec{f}_{\parallel} is in the span $\{\mathcal{K}_{x_j} \tilde{e}_1, \mathcal{K}_{x_j} \tilde{x}_i, i, j \in \mathbb{N}_m\}$ and \vec{f}_{\perp} is perpendicular to this span space. By the reproducing property (2), we have that $\langle \vec{f}(x_j), \tilde{e}_1 + \tilde{x}_i - \tilde{x}_j \rangle = \langle \vec{f}, \mathcal{K}_{x_j}(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j) \rangle_{\mathcal{K}} = \langle \vec{f}_{\parallel}, \mathcal{K}_{x_j}(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j) \rangle_{\mathcal{K}}$. Hence, \vec{f}_{\perp} makes no contribution to the loss function in the MGL algorithm (9) (i.e. algorithm (3)). However, the norm $\|\vec{f}\|_{\mathcal{K}}^2 = \|\vec{f}_{\parallel}\|_{\mathcal{K}}^2 + \|\vec{f}_{\perp}\|_{\mathcal{K}}^2 > \|\vec{f}_{\parallel}\|_{\mathcal{K}}^2$ unless $f_{\perp} = 0$. This implies, any solution $\vec{f}_{\mathbf{z}}$ belongs to the span space $\{\mathcal{K}_{x_j} \tilde{e}_1, \mathcal{K}_{x_j} \tilde{x}_i, i, j \in \mathbb{N}_m\}$ and the corresponding representer coefficients belong to the span of $\{\tilde{e}_1, \tilde{x}_i : i \in \mathbb{N}_m\}$. ■

The representer theorem above tells us that the optimal solution $\vec{f}_{\mathbf{z}}$ of algorithm (3) lives in the finite span of training samples which paves the way for designing efficient optimization algorithms for multi-task gradient learning.

4 Optimization and solution

In this section, by the above representer theorem, we explore efficient algorithms for computing the representer coefficients. For clarity, we mainly focus on least-square multi-task gradient learning algorithms (LSMGL). At the end of this section, the support vector machine for gradient learning (SVMMLG) in classification will be briefly discussed. One can apply the subsequent procedures to other loss functions.

4.1 Computation of representer coefficients

To specify the solution of LSMGL, we denote the column vector $C_{\mathbf{z}} \in \mathbb{R}^{m(d+1)}$ by consecutively catenating all column vectors $\{c_{j,\mathbf{z}} \in \mathbb{R}^{d+1} : j \in \mathbb{N}_m\}$ and, likewise we define a column vector $\mathbb{Y} \in \mathbb{R}^{m(d+1)}$ by catenating column vectors $\{y_i \in \mathbb{R}^{d+1} : i \in \mathbb{N}_m\}$. Moreover, we introduce an $m(d+1) \times m(d+1)$ matrix by catenating all $(d+1) \times (d+1)$ matrix $\mathcal{K}(x_i, x_j)$ denoted by

$$\mathcal{K}_{\mathbf{x}} = (\mathcal{K}(x_i, x_j))_{i,j \in \mathbb{N}_m}.$$

Finally, we introduce a system of equations

$$m^2 \lambda c_j + B_j \sum_{l \in \mathbb{N}_m} \mathcal{K}(x_j, x_l) c_l = y_j, \quad \forall j \in \mathbb{N}_m \quad (10)$$

where $B_j = \sum_{i \in \mathbb{N}_m} w_{ij} (\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j)(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j)^T$, $y_j = \sum_{i \in \mathbb{N}_m} w_{ij} y_i (\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j)$.

We now can solve the LSMGL algorithm by the following theorem.

Theorem 4 *For any $j \in \mathbb{N}_m$, the vectors B_j, y_j be defined by equation (10). Then, the representer coefficients $C_{\mathbf{z}}$ for the solution of the LSMGL algorithm are given by the following equation*

$$\mathbb{Y} = \left(m^2 \lambda I_{m(d+1)} + \text{diag}(B_1, \dots, B_m) \mathcal{K}_{\mathbf{x}} \right)_{i,j=1}^m C_{\mathbf{z}}. \quad (11)$$

Proof: By Theorem 3, there exists $\{c_{j,\mathbf{z}} \in \mathbb{R}^{d+1} : j \in \mathbb{N}_m\}$ such that $\vec{f}_{\mathbf{z}} = \sum_{j \in \mathbb{N}_m} \mathcal{K}_{x_j} c_{j,\mathbf{z}}$. However, taking the

functional derivative of algorithm (3) with respect to f yields that $\frac{1}{m^2} \sum_{i,j \in \mathbb{N}_m} w_{ij} (\langle \vec{f}_z, \mathcal{K}_{x_j}(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j) \rangle_{\mathcal{K}} - y_i) \mathcal{K}_{x_j}(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j) + \lambda \vec{f}_z = 0$ which means that $c_{j,z} = \frac{1}{m^2 \lambda} \sum_{i \in \mathbb{N}_m} w_{ij} (y_i - \langle \vec{f}_z, \mathcal{K}_{x_j}(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j) \rangle_{\mathcal{K}}) (\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j)$. Equivalently, equation (10) holds true, and hence completes the assertion. \blacksquare

Solving equation (11) involves the inversion of an $m(d+1) \times m(d+1)$ matrix whose time complexity is usually $O((md)^3)$. However, it is computationally prohibitive since the coordinate (feature) dimension d is very large in many applications. Fortunately, as suggested in Theorem 3, the representer coefficients $\{c_{j,z} : j \in \mathbb{N}_m\}$ can be represented by the span of column vectors of matrix

$$\widetilde{M}_x = \{\tilde{e}_1, \tilde{x}_1, \dots, \tilde{x}_{m-1}, \tilde{x}_m\}.$$

This observation suggests the possibility of reduction of the original high dimensional problem in \mathbb{R}^{d+1} to the low dimensional space spanned by \widetilde{M}_x . This low dimensional space can naturally be introduced by singular vectors of \widetilde{M}_x .

To this end, we consider the representation of the matrix \widetilde{M}_x by its singular vectors. It will be proven to be useful to represent matrix \widetilde{M}_x from the singular value decomposition (SVD) of the data matrix defined by

$$M_x = [x_1, x_2, \dots, x_{m-1}, x_m].$$

Apparently, the rank s of M_x is at most $\min(m, d)$. The SVD of M_x tells us that there exists orthogonal matrices $V_{d \times d}$ and $U_{m \times m}$ such that

$$\begin{aligned} M_x &= [V_1, \dots, V_d] \Sigma \begin{bmatrix} U_1 \\ \vdots \\ U_m \end{bmatrix}^T \\ &= [V_1, \dots, V_s] (\beta_1, \dots, \beta_m) \end{aligned} \quad (12)$$

Here, the $d \times m$ matrix $\Sigma = \begin{bmatrix} \text{diag}\{\sigma_1, \dots, \sigma_s\} & 0 \\ 0 & 0 \end{bmatrix}$. For any $j \in \mathbb{N}_m$, we use the notation $U_j = (U_{1j}, \dots, U_{mj})$ and $\beta_j^T = (\sigma_1 U_{1j}, \sigma_2 U_{2j}, \dots, \sigma_s U_{sj}) \in \mathbb{R}^s$. From now on we also denote

$$\mathcal{V} = [V_1, \dots, V_s] \quad (13)$$

Hence, we have, for any $j \in \mathbb{N}_m$, that $x_j = \mathcal{V} \beta_j$.

We are now ready to specify the representation of \widetilde{M}_x from the above SVD of M_x . To see this, for any $l \in \mathbb{N}_s$ and $j \in \mathbb{N}_m$, let $\tilde{V}_l^T = (0, V_l^T)$, $\tilde{\beta}_j^T = (0, \beta_j^T)$. In addition, we introduce the $(d+1) \times (s+1)$ matrix

$$\tilde{\mathcal{V}} = \begin{pmatrix} 1 & 0 \\ 0 & \mathcal{V} \end{pmatrix} = [\tilde{e}_1, \tilde{V}_1, \dots, \tilde{V}_s] \quad (14)$$

which induces a one-to-one mapping $\tilde{\mathcal{V}} : \mathbb{R}^{s+1} \rightarrow \mathbb{R}^{d+1}$ defined, for any $\beta \in \mathbb{R}^{s+1}$, by $x = \tilde{\mathcal{V}} \beta \in \mathbb{R}^{d+1}$ since column vectors in $\tilde{\mathcal{V}}$ are orthogonal to each other. Consequently, it follows that

$$\widetilde{M}_x = \tilde{\mathcal{V}} [e_1, \tilde{\beta}_1, \dots, \tilde{\beta}_m],$$

where e_1 is the standard first coordinate basis in \mathbb{R}^{s+1} . Equivalently, for any $i, j \in \mathbb{N}_m$,

$$\tilde{e}_1 = \tilde{\mathcal{V}} e_1, \quad \tilde{x}_j = \tilde{\mathcal{V}} \tilde{\beta}_j. \quad (15)$$

We now assemble all material to state the reduced system associated with equation (11). For this purpose, firstly we define the kernel $\tilde{\mathcal{K}}$, for any $x, t \in X$, by

$$\tilde{\mathcal{K}}(x, t) := \tilde{\mathcal{V}}^T \mathcal{K}(x, t) \tilde{\mathcal{V}},$$

and introduce the $m(s+1) \times m(s+1)$ matrix by concatenating all $(s+1) \times (s+1)$ matrices $\tilde{\mathcal{K}}(x_i, x_j)$:

$$\tilde{\mathcal{K}}_x = (\tilde{\mathcal{K}}(x_i, x_j))_{i,j \in \mathbb{N}_m}. \quad (16)$$

Secondly, for any $j \in \mathbb{N}_m$, set $\mathcal{B}_j = \sum_{i \in \mathbb{N}_m} w_{ij} (e_1 + \tilde{\beta}_i - \tilde{\beta}_j)(e_1 + \tilde{\beta}_i - \tilde{\beta}_j)^T$ and $\mathcal{Y}_j = \sum_{i \in \mathbb{N}_m} w_{ij} y_i (e_1 + \tilde{\beta}_i - \tilde{\beta}_j)$. Thirdly, associated with the system (10), for any $j \in \mathbb{N}_m$ and $\gamma_j \in \mathbb{R}^{s+1}$, we define the system in reduced low dimensional space \mathbb{R}^{s+1}

$$m^2 \lambda \gamma_j + \mathcal{B}_j \sum_{l \in \mathbb{N}_m} \tilde{\mathcal{K}}(x_j, x_l) \gamma_l = \mathcal{Y}_j. \quad (17)$$

Finally, in analogy with the notation \mathbb{Y} , the column vector $\mathcal{Y} \in \mathbb{R}^{m(s+1)}$ is defined by successively concatenating column vectors $\{\mathcal{Y}_i \in \mathbb{R}^{s+1} : i \in \mathbb{N}_m\}$. Likewise we can define γ_z by concatenating column vectors $\{\gamma_{j,z} \in \mathbb{R}^{s+1} : j \in \mathbb{N}_m\}$.

With the above preparation we have the following result.

Theorem 5 *If the $\{\gamma_{j,z} \in \mathbb{R}^{s+1} : j \in \mathbb{N}_m\}$ is the solution of system (17), i.e.,*

$$\mathcal{Y} = \left(m^2 \lambda I_{m(s+1)} + \text{diag}(\mathcal{B}_1, \dots, \mathcal{B}_m) \tilde{\mathcal{K}}_x \right) \gamma, \quad (18)$$

then the coefficient C_z defined, for any $j \in \mathbb{N}_m$, by $c_{j,z} = \tilde{\mathcal{V}} \gamma_{j,z}$ is one of the solution of system (11), and thus yields representation coefficients of the solution \vec{f}_z for the LSMGL algorithm (3).

Proof: Let γ_z is the solution of system (17). Since $\tilde{\mathcal{V}}$ is orthogonal, the system (17) is equivalent to the following equation

$$m^2 \lambda \tilde{\mathcal{V}} \gamma_{j,z} + \tilde{\mathcal{V}} \mathcal{B}_j \sum_{l \in \mathbb{N}_m} \tilde{\mathcal{K}}(x_j, x_l) \gamma_{l,z} = \tilde{\mathcal{V}} \mathcal{Y}_j.$$

Recall, for any $j \in \mathbb{N}_m$, that $B_j = \tilde{\mathcal{V}} \mathcal{B}_j \tilde{\mathcal{V}}^T$, $\tilde{\mathcal{V}} e_1 e_1^T \tilde{\mathcal{V}}^T = \tilde{e}_1 \tilde{e}_1^T$, $\tilde{\mathcal{V}}^T \tilde{\mathcal{V}} = I_{s+1}$, and $Y_j = \tilde{\mathcal{V}} \mathcal{Y}_j$. Hence, the above system is identical to system (11) (i.e. system (10)) with C_j replaced by $\tilde{\mathcal{V}} \gamma_{j,z}$ which completes the assertion. \blacksquare

We end this subsection with a brief discussion of the solution of the SVMGL algorithm. Since the hinge loss is not differentiable, we cannot use the above techniques to derive an optimization algorithm for SVMGL. Instead, we can consider its dual problem. To this end, we introduce slack variables $\{\xi_{ij} : i, j \in \mathbb{N}_m\}$ and rewrite SVMGL as follows:

$$\begin{cases} \arg \min_{\vec{f}, \xi, b} \left\{ \frac{1}{m^2} \sum_{i,j \in \mathbb{N}_m} w_{ij} \xi_{ij} + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \right\} \\ \text{s.t.} & y_i (\langle \vec{f}, \mathcal{K}_{x_j}(\tilde{e}_1 + \tilde{x}_i - \tilde{x}_j) \rangle_{\mathcal{K}} + b) \geq 1 - \xi_{ij}, \\ & \xi_{ij} \geq 0, \quad \forall i, j \in \mathbb{N}_m. \end{cases} \quad (19)$$

- Given a multi-task kernel \mathcal{K} produced by scalar kernel G , $\lambda > 0$ and inputs $\{(x_i, y_i) : i \in \mathbb{N}_m\}$
1. Compute projection mapping $\tilde{\mathcal{V}}$ and reduced vector $\tilde{\beta}_j$ from equations (12), (13), and (14).
 2. Compute $\tilde{\mathcal{K}}(x_j, x_i) = \tilde{\mathcal{V}}^T \mathcal{K}(x_j, x_i) \tilde{\mathcal{V}}$ (i.e. equation (16))
 3. Solving equation (18) to get coefficient $\gamma_{\mathbf{z}}$, see Theorem 5 (equivalently, equation (25) when G is linear or RBF kernel, see Theorem 6).
 4. Output vector-valued function: $\vec{f}_{\mathbf{z}}(\cdot) = \sum_{j \in \mathbb{N}_m} \mathcal{K}(\cdot, x_j) (\tilde{\mathcal{V}} \gamma_{j, \mathbf{z}})$.
 5. Compute variable covariance and ranking variables using Proposition 1 in Section 6.

Table 1: Pseudo-code for least square multi-task gradient learning

Parallel to the derivation of the dual problem of standard SVM (e.g. [20, 23]), using Lagrangian theory we can obtain the following dual problem of SVMMLG:

$$\begin{cases} \arg \max_{\alpha} \sum_{i, j \in \mathbb{N}_m} \alpha_{ij} - \frac{1}{4m^2 \lambda} \sum_{i, j, i', j' \in \mathbb{N}_m} \alpha_{ij} y_i \alpha_{i' j'} y_{i'} \\ \times [(\tilde{\mathbf{e}}_1 + \tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j)^T \mathcal{K}(x_j, x_{j'}) (\tilde{\mathbf{e}}_1 + \tilde{\mathbf{x}}_{i'} - \tilde{\mathbf{x}}_{j'})] \\ \text{s.t. } \sum_{i, j \in \mathbb{N}_m} y_i \alpha_{ij} = 0, 0 \leq \alpha_{ij} \leq w_{ij}, \forall i, j \in \mathbb{N}_m. \end{cases} \quad (20)$$

Moreover, if the solution of dual problem is $\alpha_{\mathbf{z}} = \{\alpha_{ij, \mathbf{z}} : i, j \in \mathbb{N}_m\}$ then the solution of SVMMLG can be represented by

$$\vec{f}_{\mathbf{z}} = \frac{1}{2m^2 \lambda} \sum_{i, j \in \mathbb{N}_m} y_i \alpha_{ij, \mathbf{z}} \mathcal{K}_{x_j}(\tilde{\mathbf{e}}_1 + \tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j).$$

Note that

$$((\tilde{\mathbf{e}}_1 + \tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j)^T \mathcal{K}(x_j, x_{j'}) (\tilde{\mathbf{e}}_1 + \tilde{\mathbf{x}}_{i'} - \tilde{\mathbf{x}}_{j'}))_{(i, j), (i', j')}$$

is a scalar kernel matrix with double indices (i, j) and (i', j') . Then, when the number of samples is small the dual problem of SVMMLG can be efficiently solved by quadratic programming with $\alpha \in \mathbb{R}^{m^2}$.

4.2 Further low dimensional formulation

Consider the multi-task kernel \mathcal{K} defined by equation (6) with scalar kernel G . In this section, by further specifying G we show that LSMGL algorithm (4) with input/output $\{(x_i, y_i) : i \in \mathbb{N}_m\}$ can be reduced to its low dimensional formulation with input/output $\{(\beta_i, y_i) : i \in \mathbb{N}_m\}$, where β_j is defined by equation (12). This clarification will provide a more computationally efficient algorithm.

To this end, consider the scalar kernel G defined on $\mathbb{R}^d \times \mathbb{R}^d$. By definition of kernels (called restriction theorem in [3]), we can see that G is also a reproducing kernel on $\mathbb{R}^s \times \mathbb{R}^s$. Hence, \mathcal{K} defined by equation (6) on \mathbb{R}^d is also a multi-task kernel on the underlying space \mathbb{R}^s ; we use the same notation \mathcal{K} when no confusion arises. Therefore, associated with the LSMGL algorithm (4) in \mathbb{R}^d we have an LSMGL in low dimensional input space \mathbb{R}^s :

$$\vec{g}_{\mathbf{z}} = \arg \min_{\vec{g} \in \mathcal{H}_{\mathcal{K}}} \left\{ \frac{1}{m^2} \sum_{i, j} w_{ij} [y_i - g_1(\beta_j) - \vec{g}_2(\beta_j) (\beta_i - \beta_j)^T]^2 + \lambda \|\vec{g}\|_{\mathcal{K}}^2 \right\}. \quad (21)$$

In analogy with the derivation of the system (10), for any $j \in \mathbb{N}_m$ and $\gamma_j \in \mathbb{R}^{s+1}$, we know that the representer coefficients of the LSMGL algorithm (21) in reduced low dimensional space \mathbb{R}^{s+1} satisfy that

$$m^2 \lambda \gamma_j + \mathcal{B}_j \sum_{l \in \mathbb{N}_m} \mathcal{K}(\beta_j, \beta_l) \gamma_l = \mathcal{Y}_j. \quad (22)$$

We are now ready to discuss the relation between representer coefficients of the LSMGL algorithm (4) and those of reduced LSMGL algorithm (21). For this purpose, let G to satisfy, for any $d \times s$ matrix V , that $V^T V = I_s$ and $\beta, \beta' \in \mathbb{R}^s$, that

$$G(V\beta, V\beta') = G(\beta, \beta'). \quad (23)$$

There exists abundant functions G satisfying the above property. For instance, linear product kernel $G(x, t) = x^T t$, Gaussian kernel $G(x, t) = e^{-\|x-t\|^2/2\sigma}$, and sigmoid kernel $G(x, t) = \tanh(ax^T t + r)$ with parameters $a, r \in \mathbb{R}$. More generally, kernel G satisfies property (23) if it is produced by a radial basis function (RBF) $h : (0, \infty) \rightarrow \mathbb{R}$ defined, for any $x, t \in X$, by

$$G(x, t) = h(\|x - t\|^2). \quad (24)$$

We say a function $h : (0, \infty) \rightarrow \mathbb{R}$ is *complete monotone* if it is smooth and, for any $r > 0$ and $k \in \mathbb{N}$, $(-1)^k f^{(k)}(r) \geq 0$. Here $h^{(k)}$ denotes the k -th derivative of h . According to the well-known Schoenberg's theorem [18], if h is complete monotone then the function G defined by equation (24) is positive semi-definite, and hence becomes a scalar kernel. For instance, the choice of $h(t) = e^{-\frac{t}{2\sigma}}$ with standard deviation $\sigma > 0$ and $h(t) = (\sigma^2 + \|x - t\|^2)^{-\alpha}$ with parameter $\alpha > 0$ yield Laplacian kernel and inverse polynomial kernel respectively.

Now we are in a position to summarize the reduction theorem for multi-task kernels (6) produced by scalar kernels G satisfying (23). Here we also use the convention that $\mathcal{K}_{\beta} = (\mathcal{K}(\beta_i, \beta_j))_{i, j \in \mathbb{N}_m}$.

Theorem 6 *Let G have the property (23) and \mathcal{K} be defined by equation (6). Suppose $\{\gamma_{j, \mathbf{z}} : j \in \mathbb{N}_m\}$ are the representer coefficients of algorithm (21), i.e., $\gamma_{\mathbf{z}}$ solves the equation*

$$\mathcal{Y} = \left(m^2 \lambda I_{m(s+1)} + \text{diag}(\mathcal{B}_1, \dots, \mathcal{B}_m) \mathcal{K}_{\beta} \right) \gamma_{\mathbf{z}}, \quad (25)$$

Then, the representer coefficients $\{c_{j, \mathbf{z}} : j \in \mathbb{N}_m\}$ of algorithm (4) are given by

$$c_{j, \mathbf{z}} = \tilde{\mathcal{V}} \gamma_{j, \mathbf{z}}.$$

Proof: Suppose the multi-task kernel \mathcal{K} is produced by G with property (23). Recall that $\mathcal{V}^T \mathcal{V} = I_s$ with \mathcal{V} given by (14). Then, kernel \mathcal{K} satisfies, for any $x = \mathcal{V}\beta$ and $t = \mathcal{V}\beta'$ with \mathcal{V} , that

$$\begin{aligned} \mathcal{K}(x, t) &= \mathcal{K}(\mathcal{V}\beta_i, \mathcal{V}\beta_j) \\ &= \begin{pmatrix} G(\beta, \beta'), & (\mathcal{V}\nabla_{\beta_i} G(\beta_i, \beta_j))^T \\ \mathcal{V}\nabla_{\beta_j} G(\beta_i, \beta_j), & \mathcal{V}(\nabla_{\beta_i} \nabla_{\beta_j} G(\beta_i, \beta_j))\mathcal{V}^T \end{pmatrix}. \end{aligned}$$

Hence, it follows, for any $x_i, x_j \in X$ and $i, j \in \mathbb{N}_m$, that

$$\tilde{\mathcal{K}}(x_i, x_j) = \begin{pmatrix} 1 & 0 \\ 0 & \mathcal{V} \end{pmatrix}^T \mathcal{K}(x_i, x_j) \begin{pmatrix} 1 & 0 \\ 0 & \mathcal{V} \end{pmatrix} = \mathcal{K}(\beta_i, \beta_j).$$

Therefore, the system (17) is identical to the system (22). Consequently, the desired assertion follows directly from Theorem 5. \blacksquare

Equipped with Theorem 6, the time complexity and computer memory can be further reduced by directly computing $m(s+1) \times m(s+1)$ matrix \mathcal{K}_β instead of first computing $m(d+1) \times m(d+1)$ matrix \mathcal{K}_x and then $\tilde{\mathcal{K}}_x$ in Theorem 5. Theorem 6 also gives an appealing insight into multi-task gradient learning framework. Roughly speaking, learning gradient in the high dimensional space is equivalent to learning them in the low dimensional projection space spanned by the input data.

5 Statistical error analysis

In this section we give an error analysis for least square MGL algorithms. Our target is to show that the learned vector-valued function from our algorithm statistically converges to the true function and true gradient.

For the least square loss, denote by $\rho_X(\cdot)$ the marginal distribution on X and, for any $x \in X$, let $\rho(\cdot|x)$ to be the conditional distribution on Y . Then, the target function is the regression function f_ρ minimizing the generalization error

$$\mathcal{E}(f) = \int_Z (y - f(x))^2 d\rho(x, y).$$

Specifically, the regression function is defined, for any $x \in X$, by

$$f_\rho(x) = \arg \min_{t \in \mathbb{R}} \int_Y (y - t)^2 \rho(y|x) = \int_Y y d\rho(y|x).$$

Hence, in this case the purpose of error analysis is to show that solution \vec{f}_ρ of LSMGL algorithm (4) statistically converges to $\vec{f}_\rho = (f_\rho, \nabla f_\rho)$ as $m \rightarrow \infty$, $s = s(m) \rightarrow 0$ and $\lambda = \lambda(m) \rightarrow 0$.

To this end, we introduce some notations and the following hypotheses are assumed to be true throughout this section. Firstly, we assume that $Y \subseteq [-M, M]$ with $M > 0$. Since X is compact the diameter of X denoted by $D = \sup_{x, u \in X} \|x - u\|^2$ is finite. Secondly, denote by $L_{\rho_X}^2$ the space of square integral functions $\vec{f} : X \rightarrow \mathbb{R}^{d+1}$ with norm $\|\vec{f}\|_{\rho_X}^2 = \int_X \|\vec{f}(x)\|^2 d\rho_X(x)$. Finally, denote the boundary of X by ∂X . We assume, for some constant $c_\rho > 0$, that the marginal distribution satisfies that

$$\rho_X(x \in X : \text{dist}(x, \partial X) < s) \leq c_\rho s, \quad \forall 0 < s < D. \quad (26)$$

and, for some parameter $0 < \theta \leq 1$, the density function $p(x)$ of ρ_X satisfies θ -Hölder continuous condition, i.e., for any $x, u \in X$ there holds

$$|p(x) - p(u)| \leq c_\rho \|x - u\|^\theta, \quad \forall x, u \in X. \quad (27)$$

Of course, p is a bounded function on X since it is continuous and X is compact. For instance, if the boundary of X is piecewise smooth and ρ_X is the uniform distribution over X then the marginal distribution ρ_X satisfies conditions (26) and (27) with parameter $\theta = 1$.

We are ready to present our statistical error analysis of LSMGL algorithms. Recall here we used the notation $\vec{f}_\rho = (f_\rho, \nabla f_\rho)$.

Theorem 7 *Suppose that the marginal distribution ρ_X satisfies (26) and (27). For any multi-task kernel \mathcal{K} , let \vec{f}_ρ be the solution of LSMGL algorithm. If $\vec{f}_\rho \in \mathcal{H}_\mathcal{K}$ then there exists a constant c such that, for any $m \in \mathbb{N}$, with the choice of $\lambda = s^{2\theta}$ and $s = m^{-\frac{1}{3(d+2)+4\theta}}$, there holds*

$$\mathbb{E}[\|\vec{f}_\rho - \vec{f}_\rho\|_{\rho_X}^2] \leq cm^{-\frac{\theta}{3(d+2)+4\theta}}.$$

If moreover ρ_X is a uniform distribution then, choosing $\lambda = s^\theta$ and $s = m^{-\frac{1}{3(d+2)+5\theta}}$, there holds

$$\mathbb{E}[\|\vec{f}_\rho - \vec{f}_\rho\|_{\rho_X}^2] \leq cm^{-\frac{\theta}{3(d+2)+\theta}}.$$

The proof of this theorem needs several steps which are postponed to the appendix. More accurate error rates in terms of probability inequality are possible using techniques in [17, 15]. It would also be interesting to extend this theorem to other loss functions such as the SVMML algorithm.

6 Experimental validation

In this section we will only preliminarily validate the MGL algorithm (3) on the problem of variable selection and covariance measurement.

By the representer Theorem 3 in Section 3, the solution of MGL denoted by $\vec{f}_\rho = (f_{1,\mathbf{z}}, \vec{f}_{2\mathbf{z}}) = (f_{1,\mathbf{z}}, f_{2,\mathbf{z}}, \dots, f_{d+1,\mathbf{z}})$ can be rewritten as $\vec{f}_\rho = \sum_{j \in \mathbb{N}_m} \mathcal{K}_{x_j} c_{j,\mathbf{z}}$. Since it only belongs to a vector-valued RKHS $\mathcal{H}_\mathcal{K}$, we need to find a common *criterion inner product (norm)* $\langle \cdot, \cdot \rangle_r$ to measure each component of the learned gradient $\vec{f}_{2\mathbf{z}} = (f_{2,\mathbf{z}}, \dots, f_{d+1,\mathbf{z}})$. Once we find the criterion inner product $\langle \cdot, \cdot \rangle_r$, we can use the coordinate covariance

$$\text{Cov}(\vec{f}_{2\mathbf{z}}) = \left(\langle f_{p+1,\mathbf{z}}, f_{q+1,\mathbf{z}} \rangle_r \right)_{p,q \in \mathbb{N}_d} \quad (28)$$

to measure how the variables covary. Also, the variable (feature) ranking can be done according to the following relative magnitude of norm of each component of $\vec{f}_{2\mathbf{z}}$:

$$s_p = \frac{\|f_{p+1,\mathbf{z}}\|_r}{(\sum_{q \in \mathbb{N}_d} \|f_{q+1,\mathbf{z}}\|_r^2)^{1/2}}. \quad (29)$$

If the scalar kernel G is a linear kernel then every component of $\vec{f}_{2\mathbf{z}}$ is a constant. In this case, we can choose the standard Euclidean inner product to be the criterion inner product (norm). When the kernel G is an RBF kernel, we show in the following proposition that we can select the criterion inner product $\langle \cdot, \cdot \rangle_r$ to be the RKHS inner product $\langle \cdot, \cdot \rangle_G$ in \mathcal{H}_G . The computation is summarized in the following proposition.

Proposition 1 Suppose the scalar kernel G has a feature representation and the multi-task kernel \mathcal{K} is defined by equation (6). Then, for any solution $\vec{f}_z = \sum_{j \in \mathbb{N}_m} \mathcal{K}_{x_j} c_{j,z} \in \mathcal{H}_{\mathcal{K}}$ of MGL algorithm (3), the following hold true.

1. If G is a linear kernel then the coordinate covariance is defined by

$$\text{Cov}(\vec{f}_{2z}) = \vec{f}_{2z}^T \vec{f}_{2z} = \sum_{i,j \in \mathbb{N}_m} (x_i, I_d) c_{i,z} c_{j,z}^T (x_j, I_d)^T.$$

Moreover, for LSMGL algorithm the above equation can be more efficiently computed by

$$\text{Cov}(\vec{f}_{2z}) = \mathcal{V} \left[\sum_{i,j \in \mathbb{N}_m} (\beta_i, I_s) \gamma_i \gamma_j^T (\beta_j, I_s)^T \right] \mathcal{V}^T.$$

2. If G is a smooth RBF kernel then $f_{p+1,z} \in \mathcal{H}_G$ and the coordinate covariance $\text{Cov}(\vec{f}_{2z}) = ((f_{p+1,z}, f_{q+1,z})_G)_{p,q \in \mathbb{N}_d}$ can be computed by

$$\langle f_{p+1,z}, f_{q+1,z} \rangle_G = C_z^T (\mathbb{K}_{pq}(x_i - x_j))_{i,j=1}^m C_z, \quad (30)$$

where the kernel matrix $\mathbb{K}_{pq}(x_i - x_j)$ defined, for any $i, j \in \mathbb{N}_m$, by

$$\begin{pmatrix} -(\partial_{pq}^2 G)(x_i - x_j), & ((\nabla \partial_{pq}^2 G)(x_i - x_j))^T \\ -(\nabla \partial_{pq}^2 G)(x_i - x_j), & (\nabla^2 \partial_{pq}^2 G)(x_i - x_j) \end{pmatrix}.$$

The proof is postponed to the appendix where the computation of \mathbb{K}_{pq} is also given if G is a Gaussian.

We run our experiment on two artificial datasets and one gene expression dataset following [17]. In the first experiment, the target function $f_\rho : \mathbb{R}^d \rightarrow \mathbb{R}$ with notation $x = (x^1, \dots, x^d) \in \mathbb{R}^d$ and $d = 80$. The output y is contaminated by a Gaussian noise

$$y = f_\rho(x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_y).$$

As depicted in Figure 1 (leftmost), the thirty inputs whose relevant features are $[1, 10] \cup [11, 20] \cup [41, 50]$ are generated as follows:

1. For samples from 1 to 10, $x^p \sim \mathcal{N}(1, 0.05)$, for $p \in [1, 10]$ and $x^p \sim \mathcal{N}(0, 0.1)$, for $p \in [11, 80]$.
2. For samples from 11 to 30, $x^p \sim \mathcal{N}(1, 0.05)$, for $p \in [11, 20]$ and $x^p \sim \mathcal{N}(0, 0.1)$, for $p \in [1, 10] \cup [31, 80]$.
3. For samples from 11 to 30, features are in the form of $x^p \sim \mathcal{N}(1, 0.05)$, for $p \in [41, 50]$, and $x^p \sim \mathcal{N}(0, 0.1)$, for $p \in [1, 40] \cup [51, 80]$.

We let the regression function f_ρ to be a linear function. Specifically, we choose a noise parameter $\sigma_y = 3$ and the regression function is defined by $f_\rho(x_i) = w_1^T x_i$ for $i \in [1, 10]$, $f_\rho(x_i) = w_2^T x_i$ for $i \in [11, 20]$, and $f_\rho(x_i) = w_3^T x_i$ for $i \in [21, 30]$, where, $w_1^k = 2 + 0.5 \sin(2\pi k/10)$ for $k \in [1, 10]$ and otherwise zero, $w_2^k = -2 - 0.5 \sin(2\pi k/10)$ for $k \in [11, 20]$ and zero otherwise. The vector w_3 is defined by $w_3^k = -2 - 0.5 \sin(2\pi k/10)$ for $k \in [41, 50]$ and zero otherwise.

In this linear case, we use the kernel $G(x, t) = x^T t$ as a basic scalar kernel and the multi-task kernel \mathcal{K} defined by (6) in LSMGL algorithm (4). As in [11, 15], the regularization parameter λ is set to be a fixed number such as 0.1 (variation

in this parameter made little difference to feature ranking). The parameter s in the weight coefficients w_{ij} is set to be the median pairwise distance between inputs. In Figure 1, the result of LSMGL is shown in (b) for variable covariation and in (c) for feature selection respectively. We also ran the algorithm (3) with the choice of kernel $\mathcal{K}(x, t) = G(x, t) I_{d+1}$ ([15, 16, 17]). The results are shown in (d) and (e) of Figure 1. We see that both algorithms worked well. The LSMGL algorithm works slightly better: the reason maybe be that it captures the inherent structure of gradient learning as mentioned before. We also ran LSMGL algorithm on this dataset; the result is no essentially different from SVMMLG.

In the second experiment, we use the SVMMLG algorithm for classification. For this dataset, only the first two features are relevant to the classification task. The remaining 78 redundant features are distributed according to a small Gaussian random deviate. The distribution for the first two features is shown in (f). In SVMMLG, the parameter s and λ are the same as those in the first example. The scalar kernel is set to be a Gaussian $G(x, t) = e^{-\|x-t\|^2/2\sigma^2}$ where σ is also the median pairwise distance between inputs. The feature selection results for the SVMMLG algorithm are illustrated respectively in (g) and (h) with different choices of multi-task kernels \mathcal{K} given by equation (6) and $\mathcal{K}(x, t) = G(x, t) I_{d+1}$. Both algorithms picked up the two important features.

Finally, we apply our LSMGL algorithm to a well-studied expression dataset. This dataset has two classes: acute myeloid leukemia (AML) and acute lymphoblastic leukemia (ALL), see e.g. [10]. There are a total of 7129 genes (variables) and 72 patients, split into a training set of 38 examples and a test set of 34 examples. In the training set, 27 examples belong to ALL and 11 belong to AML, and the test set is composed of 20 ALL and 14 AML. Various variable selection algorithms have been applied to this dataset by choosing features based on training set, and then performing classification on the test set with the selected features. We ran LSMGL with the choice of multi-task \mathcal{K} given by equation (6) where G is a linear kernel. The solution \vec{f}_z is learned from the training set for ranking the genes according to the values of s_p defined by equation (29). Then, ridge regression is run on the training set with truncated features to build a classifier to predict the labels on the test set. The regularization parameter of LSMGL is fixed to be 0.1 while the regularization parameter in ridge regression is tuned using leave-one-out cross-validation in the training set. The test error with selected top ranked genes is reported in Table 2. The classification accuracy is quite comparable to the gradient learning algorithm using individual RKHSs [15, 17]. However, [11, 15, 17] did the recursive techniques to rank features and employed SVM for classification while our method showed that ridge regression for classification and non-recursive technique for feature ranking also worked well in this data set. It would be interesting to further explore this issue.

The preliminary experiments above validated our proposed MGL algorithms. However further experiments need to be performed to evaluate our multi-task framework for gradient learning.

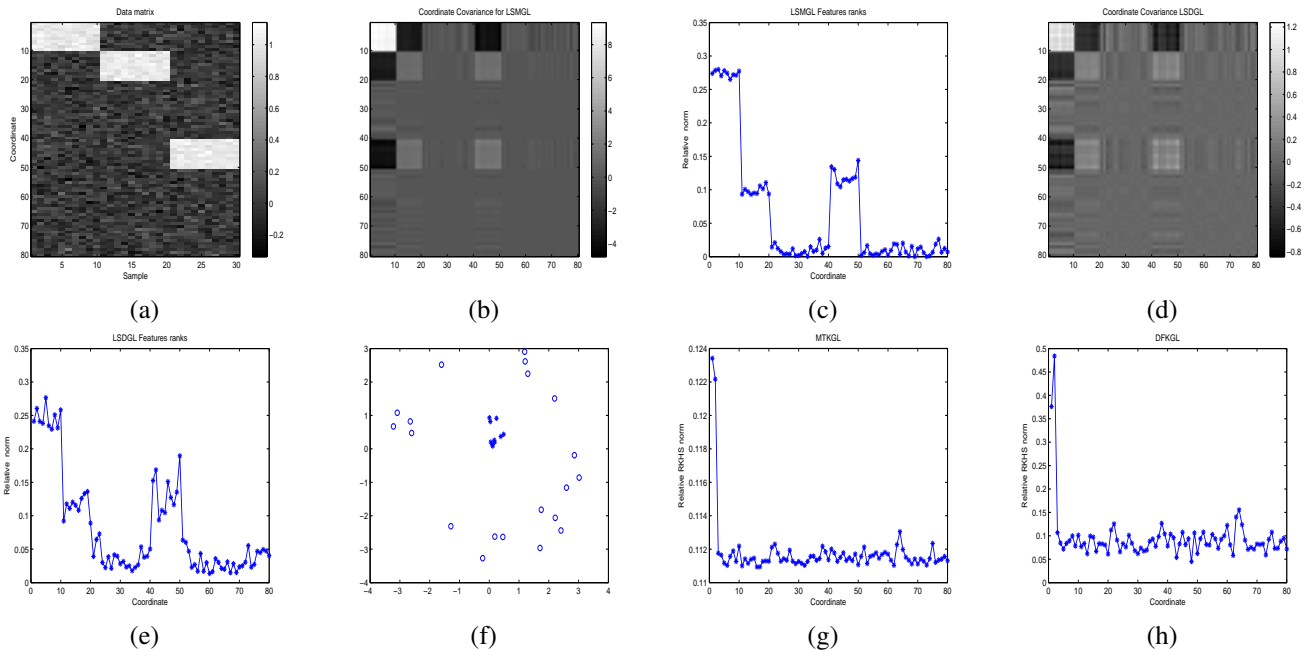


Figure 1: LSMGL and SVMML feature ranking

genes	10	40	80	100	200	500
test error	2	1	0	0	1	1
genes	1000	2000	3000	4000	6000	7129
test error	2	1	1	1	1	1

Table 2: Number of test error using ridge regression algorithm versus the number of top ranked genes selected by LSMGL algorithm.

7 Conclusions

In this paper, our main contribution was to provide a novel unifying framework for gradient learning from the perspective of multi-task learning. Various variable selection methods in the literature can be recovered by the choice of multi-task kernels. More importantly, this framework allows us to introduce a novel choice of multi-task kernel to capture the inherent structure of gradient learning. An appealing representation theorem was presented which facilitates the design of efficient optimization algorithms, especially for datasets with high dimension and few training examples. Finally, a statistical error analysis was provided to ensure the convergence of the learned function to true function and true gradient.

Here we only preliminarily validated the method. A more extensive benchmark study remains to be pursued. In future we will explore more experiments on biomedical datasets and compare our MGL algorithms with previous related methods for feature selection, such as those in [21, 22] etc. It will be interesting to implement different loss functions in the MGL algorithms for regression and classification, apply the spectral decomposition of the gradient outer products to dimension reduction (see e.g. [16]), and possible use for network inference from the covariance of the learned gradient

function.

References

- [1] R. K. Ando & T. Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. *J. Machine Learning Research*, 1817–1853, 2005
- [2] A. Argyriou, T. Evgeniou, & M. Pontil. Multi-task feature learning. *NIPS*, 2006.
- [3] N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.* 68: 337–404, 1950.
- [4] P. L. Bartlett & S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *J. of Machine Learning Research*, 3:463–482, 2002.
- [5] S. Ben-David & R. Schuller. Exploiting task relatedness for multiple task learning. *COLT*, 2003.
- [6] D. R. Chen, Q. Wu, Y. Ying, & D. X. Zhou. Support vector machine soft margin classifiers: Error analysis. *J. of Machine Learning Research*, 5:1143–1175, 2004.
- [7] A. Caponnetto, C.A. Micchelli, M. Pontil, & Y. Ying. Universal multi-task kernels, Preprint, 2007.
- [8] S.S. Chen, D.L. Donoho & M.A. Saunders. Atomic decomposition pursuits. *SIAM J. of Scientific Computing*, 20: 33-61,1999.
- [9] F. Cucker & S.Smale. On the mathematical foundations of learning, *Bull. Amer. Math. Soc.* 39: 149, 2001.
- [10] T. R. Golub et. al. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring, *Science*, 286: 531-537, 1999.
- [11] I. Guyon, J.Weston, S. Barnhill, & V. Vapnik. Gene selection for cancer classification using support vector machines. *Machine learning* 46: 389–422, 2002.
- [12] V.I. Koltchinskii & D. Panchenko. Rademacher processes and bounding the risk of function learning. In *J.*

Wellner E. Gine, D. Mason, editor, High Dimensional Probability II, pages 443-459, 2000.

- [13] T. Evgeniou, C. A. Micchelli & M. Pontil. Learning multiple tasks with kernel methods. *J. Machine Learning Research*, 6: 615–637, 2005.
- [14] C. A. Micchelli & M. Pontil. On learning vector-valued functions. *Neural Computation*, 17: 177-204, 2005.
- [15] S. Mukherjee & Q. Wu. Estimation of gradients and coordinate covariation in classification. *J. of Machine Learning Research* 7: 2481-2514, 2006.
- [16] S. Mukherjee, Q. Wu, & D. X. Zhou. Learning gradients and feature selection on manifolds. Preprint, 2007.
- [17] S. Mukherjee & D. X. Zhou. Learning coordinate covariances via gradients, *J. of Machine Learning Research* 7: 519-549, 2006.
- [18] I. J. Schoenberg. Metric spaces and completely monotone functions, *Ann. of Math.* 39: 811-841, 1938.
- [19] B. Schölkopf & A. J. Smola. *Learning with Kernels*. The MIT Press, Cambridge, MA, USA, 2002.
- [20] J. Shawe-Taylor and N. Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, Cambridge, 2004.
- [21] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. Royal. Statist. Soc B.* 58: 267-288, 1996.
- [22] J. Weston, S. Mukherjee, O. Chapelle, M. Pontil, T. Poggio, & V. Vapnik. Feature selection for SVMs, *NIPS*, 2001.
- [23] V. N. Vapnik. *Statistical learning theory*. Wiley, New York, 1998.

Appendix

Let G be a scalar kernel, we use the convention $\partial_p^{(2)}G$ to denote the p -th partial derivative of G with respect to the second argument, and so is the gradient $\nabla^{(2)}G$.

Proof of Proposition 1.

When G is a linear kernel, by the definition (6) of multi-task kernel \mathcal{K} , we have that $\vec{f}_{2\mathbf{z}} = \sum_{j \in \mathbb{N}_m} (x_j, I_d) c_{j,\mathbf{z}}$ which implies that

$$\text{Cov}(\vec{f}_{2\mathbf{z}}) = \vec{f}_{2\mathbf{z}} \vec{f}_{2\mathbf{z}}^T = \sum_{i,j \in \mathbb{N}_m} (x_i, I_d) c_{i,\mathbf{z}} c_{j,\mathbf{z}}^T (x_j, I_d)^T.$$

For the LSMGL algorithm, in Section 4 we showed that $c_{j,\mathbf{z}} = \tilde{\mathcal{V}} \gamma_{j,\mathbf{z}}$ and, for any $j \in \mathbb{N}_m$, $x_j = \mathcal{V} \beta_j$, the above equation can be further simplified to the following:

$$\text{Cov}(\vec{f}_{2\mathbf{z}}) = \mathcal{V} \left[\sum_{i,j \in \mathbb{N}_m} (\beta_i, I_s) \gamma_{i,\mathbf{z}} \gamma_{j,\mathbf{z}}^T (\beta_j, I_s)^T \right] \mathcal{V}^T.$$

When G is an RBF kernel, for any $x, t \in \mathbb{R}^d$ and $p, q \in \mathbb{N}_d$, $G(x, t) = G(x - t)$ and $\partial_{t_q} \partial_{x_p} G(x, t) = -(\partial_p \partial_q G)(x, t)$. Hence, for any $p \in \mathbb{N}_d$ and $x \in \mathbb{R}^d$, we have that

$$f_{p+1,\mathbf{z}}(x) = \sum_{j \in \mathbb{N}_m} (-\partial_p^{(2)} G(x, x_j), -\nabla^{(2)} \partial_p^{(2)} G(x, x_j)) c_j$$

Since G has a feature representation, i.e., for any $x, t \in X$, there holds that $G(x, t) = \langle \phi(x), \phi(t) \rangle_{\ell^2}$. Also, observe that $\partial_p^{(2)} G(x, x_j) = \langle \phi(x), (\partial_p \phi)(x_j) \rangle_{\ell^2}$ and $\partial_q^{(2)} \partial_p^{(2)} G(x, x_j) =$

$\langle \phi(x), (\partial_p \partial_q \phi)(x_j) \rangle_{\ell^2}$. Denote c_j by $(c_j^1, \dots, c_j^{d+1})^T$. Consequently, for any $p \in \mathbb{N}_d$, $f_{p+1,\mathbf{z}}(\cdot) = \langle w_p, \phi(\cdot) \rangle_{\ell^2}$ with $w_p = -\sum_{j \in \mathbb{N}_m} (\partial_p \phi(x_j) c_j^1 + \sum_{q \in \mathbb{N}_d} \partial_q \partial_p \phi(x_j) c_j^{q+1})$. Therefore, for any $p, q \in \mathbb{N}_d$ we have that $f_{p+1,\mathbf{z}} \in \mathcal{H}_G$ and

$$\langle f_{p+1,\mathbf{z}}, f_{q+1,\mathbf{z}} \rangle_G = \langle w_p, w_q \rangle_{\ell^2} = \sum_{i,j \in \mathbb{N}_m} c_i^T \mathbb{K}_{pq}(x_i - x_j) c_j$$

which completes the assertion. \square

Computation of kernel \mathbb{K} .

If G is a Gaussian kernel with standard variation σ , that is, for any $x, t \in \mathbb{R}^d$, $G(x, t) = G(x - t) = e^{-\frac{\|x-t\|^2}{2\sigma^2}}$, the computation of \mathbb{K} is listed as follows.

$$1. (\partial_{pq}^2 G)(x) = \left[\frac{x_p x_q}{\sigma^2} - \frac{\delta_{pq}}{\sigma} \right] G(x)$$

$$2. \text{For any } q' \in \mathbb{N}_d, \partial_{q'} \partial_p \partial_q G(x) = \left[\frac{x_q \sigma_{pq'} + x_p \delta_{qq'} + x_{q'} \delta_{pq} - x_p x_q x_{q'}}{\sigma^3} \right] G(x). \text{ Hence,}$$

$$\nabla \partial_{pq}^2 G(x) = \left[\frac{e_p x_q + e_q x_p}{\sigma^2} + \frac{x \delta_{pq}}{\sigma^2} - \frac{x_p x_q x}{\sigma^3} \right] G(x)$$

$$3. \text{For any } p', q' \in \mathbb{N}_d, \partial_{p'} \partial_{q'} \partial_{pq}^2 G(x) = G(x) \left[\frac{\delta_{p'q'} \delta_{qp'}}{\sigma^2} + \frac{\delta_{pp'} \delta_{qq'} + \delta_{p'q'} \delta_{pq}}{\sigma^2} - \frac{1}{\sigma^3} (x_q x_{q'} \sigma_{pp'} + x_p x_q \delta_{p'q'} + x_p x_{q'} \delta_{p'q}) + \frac{x_{p'}}{\sigma} \left(\frac{x_p x_q x_{q'}}{\sigma^3} - \frac{x_q \delta_{pq'} + x_{q'} \delta_{pq} + x_p \delta_{qq'}}{\sigma^2} \right) \right]. \text{ Hence,}$$

$$\nabla^2 \partial_{pq}^2 G(x) = G(x) \left[\frac{e_q^T e_p + e_p^T e_q + \delta_{pq} I_d}{\sigma^2} - \frac{1}{\sigma^3} (x_q (e_p x^T + x e_p^T) + x_p (x e_q^T + e_q x^T)) - \frac{x_p x_q I_d}{\sigma^3} + \frac{x x^T}{\sigma^3} \left(\frac{x_p x_q}{\sigma} - \sigma_{pq} \right) \right].$$

Proof of Theorem 7

We turn our attention to the proof of Theorem 7. We begin with some notations and background materials. First denote by $\text{Er}_{\mathbf{z}}$ the empirical loss in LSMGL algorithm, i.e.,

$$\text{Er}_{\mathbf{z}}(\vec{\mathcal{F}}) = \frac{1}{m^2} \sum_{i,j} w(x_i - x_j) \times (y_i - f_1(x_j) - \vec{\mathcal{F}}_2(x_j)(x_i - x_j)^T)^2,$$

and the modified form of its expectation

$$\text{Er}(\vec{\mathcal{F}}) = \int_Z \int_X w(x - u) \left[y - f_1(u) - \vec{\mathcal{F}}_2(u)(x - u) \right]^2 d\rho(x, y) d\rho_X(u).$$

Since the Gaussian weight $w(x - u) = w_s(x - u)$ is dependent on s , the above definition of $\text{Er}(\vec{\mathcal{F}})$ is depending on the parameter s . In addition, define the Lipschitz constant $|\nabla f_\rho|_{\text{Lip}}$ to be the minimum constant c such that $\|\nabla f_\rho(x + u) - \nabla f_\rho(x)\| \leq c\|u\|$, $\forall x, u \in X$. We say that ∇f_ρ is Lipschitz continuous if $|\nabla f_\rho|_{\text{Lip}}$ is finite.

The error analysis here is divided into two main steps motivated by the techniques in [15]. The first step is to bound the square error $\|\vec{\mathcal{F}}_{\mathbf{z}} - \vec{\mathcal{F}}_\rho\|_{\rho_X}^2$ by the excess error $\text{Er}(\vec{\mathcal{F}}_{\mathbf{z}}) - \text{Er}(\vec{\mathcal{F}}_\rho)$. In the second step, we employ standard error decomposition [6] and Rademacher complexities [4, 12] to estimate the excess error. These two steps will be respectively stated in the

following two propositions. Before we do that, we introduce an auxiliary functional \mathcal{Q}_s defined by

$$\begin{aligned} \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) &= \int \int w(x-u) [f_\rho(u) - f_1(u) \\ &+ (\vec{f}_2(u) - \nabla f_\rho(u))(u-x)^T]^2 d\rho_X(x) d\rho_X(u). \end{aligned}$$

We are ready to present the first step of the error analysis: bounding the square error $\|\vec{f}_z - \vec{f}_\rho\|_{\rho_X}^2$ by the excess error $\text{Er}(\vec{f}_z) - \text{Er}(\vec{f}_\rho)$ which is stated as the following proposition.

Proposition 2 *If $0 < s, \lambda < 1$ then there exists a constant c'_ρ such that*

$$\begin{aligned} \mathbb{E} \left[\|\vec{f}_z - \vec{f}_\rho\|_{\rho_X}^2 \right] &\leq c'_\rho \left(\min [s^{-\theta}, \max_{x \in X} p^{-1}(x)] \right. \\ &\quad \times \mathbb{E} [\text{Er}(\vec{f}_z) - \text{Er}(\vec{f}_\rho)] \\ &\quad \left. + s^\theta (\mathbb{E} [\|\vec{f}_z\|_{\mathcal{K}}^2] + \|\vec{f}_\rho\|_\infty + |\nabla f_\rho|_{\text{Lip}}^2) \right). \end{aligned}$$

The proof of this proposition follows directly from the following Lemmas 8 and 9. For this purpose, let the subset X_s of X be

$$X_s = \{u \in X : \text{dist}(u, \partial X) > s, |p(u)| \geq (1 + c_\rho)s^\theta\} \quad (31)$$

and

$$c_\rho(s) := \min\{p(x) : \|u - x\| \leq s, u \in X_s\}.$$

Recall that \tilde{e}_1 is the first coordinate basis in \mathbb{R}^{d+1} and, for any $x \in \mathbb{R}^d$, $\tilde{x}^T = (0, x)^T \in \mathbb{R}^{d+1}$.

Lemma 8 *If $0 < s < 1$ then there exists a constant c'_ρ such that*

$$\begin{aligned} \mathbb{E} \left[\|\vec{f}_z - \vec{f}_\rho\|_{\rho_X}^2 \right] &\leq c'_\rho \left(s^\theta (\mathbb{E} [\|\vec{f}_z\|_{\mathcal{K}}^2] + \|\vec{f}_\rho\|_\infty) \right. \\ &\quad \left. + \min [s^{-\theta}, \max_{x \in X} p^{-1}(x)] \mathbb{E} \left[\mathcal{Q}_s(\vec{f}_z, \vec{f}_\rho) \right] \right). \end{aligned}$$

Proof: Write $\|\vec{f}_z - \vec{f}_\rho\|_{\rho_X}^2$ by

$$\begin{aligned} \|\vec{f}_z - \vec{f}_\rho\|_{\rho_X}^2 &= \int_{X \setminus X_s} \|\vec{f}_z(u) - \vec{f}_\rho(u)\|^2 d\rho_X(u) \\ &\quad + \int_{X_s} \|\vec{f}_z(u) - \vec{f}_\rho(u)\|^2 d\rho_X(u) \end{aligned} \quad (32)$$

By the definition of X_s , we have that $\rho_X(X \setminus X_s) \leq c_\rho s + c_\rho(1 + c_\rho)|X|s^\theta \leq c'_\rho s^\theta$ where $|X|$ is the Lebesgue measure of X . Hence, the first term of the above equation is bounded by

$$2c'_\rho (\|\vec{f}_z\|_\infty^2 + \|\vec{f}_\rho\|_\infty^2) s^\theta \leq 2c'_\rho (\|\vec{f}_z\|_{\mathcal{K}}^2 + \|\vec{f}_\rho\|_\infty^2) s^\theta.$$

For the second term on the right-hand side of equation (32), observe that, for any $u \in X_s$, $\text{dist}(u, \partial X) > s$ and $\{u : \|u - x\| \leq s, x \in X_s\} \subseteq X$. Moreover, for any $x \in X$ such that $\|u - x\| \leq s$, by the definition of X_s there holds

$$p(x) = p(u) - (p(u) - p(x)) \geq (1 + c_\rho)s^\theta - c_\rho \|u - x\|^\theta \geq s^\theta.$$

Consequently, it follows that

$$c_\rho(s) \geq \max(s^\theta, \min_{x \in X} p(x)), \quad (33)$$

and

$$\begin{aligned} \mathcal{Q}_s(\vec{f}_z, \vec{f}_\rho) &= \int_X \int_X w(x-u) \\ &\quad \times [(\vec{f}_\rho(u) - \vec{f}_z(u))(\tilde{e}_1 + \tilde{u} - \tilde{x})]^2 d\rho_X(x) d\rho_X(u) \\ &\geq \int_{X_s} \left[\int_{\|u-x\| \leq s} ((\vec{f}_\rho(u) - \vec{f}_z(u))(\tilde{e}_1 + \tilde{u} - \tilde{x}))^2 \right. \\ &\quad \left. \times d\rho_X(x) \right] d\rho_X(u) \\ &\geq c_\rho(s) \int_{X_s} \left[\int_{\|u-x\| \leq s} ((\vec{f}_\rho(u) - \vec{f}_z(u))(\tilde{e}_1 + \tilde{u} - \tilde{x}))^2 \right. \\ &\quad \left. \times dx \right] d\rho_X(u). \end{aligned} \quad (34)$$

The integral w.r.t. x on the right-hand side of the above inequality can be written as $(\vec{f}_\rho(u) - \vec{f}_z(u))W(s)(\vec{f}_\rho(u) - \vec{f}_z(u))^T$ with $(d+1) \times (d+1)$ matrix $W(s)$ defined by

$$W(s) = \int_{\|u-x\| \leq s} [(\tilde{e}_1 + \tilde{u} - \tilde{x})(\tilde{e}_1 + \tilde{u} - \tilde{x})^T] dx.$$

Here, $(\tilde{e}_1 + \tilde{u} - \tilde{x})(\tilde{e}_1 + \tilde{u} - \tilde{x})^T$ equals that

$$\begin{pmatrix} 1 & (u-x)^T \\ u-x & (u-x)(u-x)^T \end{pmatrix}$$

Observe that $\int_{\|u-x\| \leq s} w(x-u) dx = s^{-2} \int_{\|t\| \leq 1} e^{-\frac{\|t\|^2}{2}} dt$ and $\int_{\|u-x\| \leq s} w(x-u)(u-x) dx = 0$. In addition, for any $p \neq q \in \mathbb{N}_d$, $\int_{\|u-x\| \leq s} w(x-u)(u^p - x^p)(u^q - x^q) dx = 0$ and $\int_{\|u-x\| \leq s} w(x-u)(x^p - u^p)^2 dx = \int_{\|t\| \leq 1} e^{-\frac{\|t\|^2}{2}} (t^p)^2 dt$. From the above observations, there exists a constant c such that

$$(\vec{f}_\rho(u) - \vec{f}_z(u))W(s)(\vec{f}_\rho(u) - \vec{f}_z(u))^T \geq c \|\vec{f}_\rho(u) - \vec{f}_z(u)\|^2.$$

Recalling the definition of $W(s)$ and substituting this back into equation (34) implies, for any $0 < s < 1$, that

$$c c_\rho(s) \int_{X_s} \|\vec{f}_z(u) - \vec{f}_\rho(u)\|^2 d\rho_X(u) \leq \mathcal{Q}_s(\vec{f}_z, \vec{f}_\rho).$$

Plugging this into equation (34), the desired estimate follows from the estimation of $c_\rho(s)$, i.e., equation (33). \blacksquare

Now we can bound \mathcal{Q}_s by the following lemma.

Lemma 9 *If $0 < s < 1$ then there exists a constant c such that, for any $\vec{f} \in \mathcal{H}_{\mathcal{K}}$, the following equations hold true.*

1. $\mathcal{Q}_s(\vec{f}, \vec{f}_\rho) \leq c \left(s^2 |\nabla f_\rho|_{\text{Lip}}^2 + [\text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho)] \right)$.
2. $\text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) \leq c \left(s^2 |\nabla f_\rho|_{\text{Lip}}^2 + \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) \right)$.

Proof: Observe that $[y - f_1(u) - \vec{f}_2(u)(x-u)^T]^2 = [y - f_\rho(u) - \nabla f_\rho(u)(x-u)^T]^2 + 2[y - f_\rho(u) - \nabla f_\rho(u)(x-u)^T][f_\rho(u) - f_1(u) + (\vec{f}_2(u) - \nabla f_\rho(u))(u-x)^T] + [f_\rho(u) - f_1(u) + (\vec{f}_2(u) - \nabla f_\rho(u))(u-x)^T]^2$. Then, taking the integral of both sides of the above equality and using the fact that $f_\rho(x) = \int_Y y d\rho_X(x)$ we have that

$$\begin{aligned} \text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) &= \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) + 2 \int_X \int_X w(x-u) [f_\rho(x) \\ &\quad - f_\rho(u) - \nabla f_\rho(u)(x-u)^T] [f_\rho(u) - f_1(u) + \\ &\quad (\vec{f}_2(u) - \nabla f_\rho(u))(u-x)^T] d\rho_X(x) d\rho_X(u) \\ &\geq \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) - 2 \left(\int_X \int_X w(x-u) [f_\rho(x) - f_\rho(u) \right. \\ &\quad \left. - \nabla f_\rho(u)(x-u)^T]^2 d\rho_X(x) d\rho_X(u) \right)^{\frac{1}{2}} \left(\mathcal{Q}_s(\vec{f}, \vec{f}_\rho) \right)^{\frac{1}{2}}. \end{aligned}$$

Applying the inequality, for any $a, b > 0$, that $-2a^2 - \frac{1}{2}b^2 \leq -2ab$, from the above equality we further have that

$$\begin{aligned} \text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) &\geq \frac{1}{2} \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) - 2 \int_X \int_X w(x-u) \\ &\quad [f_\rho(x) - f_\rho(u) - \nabla f_\rho(u)(x-u)^T]^2 d\rho_X(x) d\rho_X(u). \end{aligned} \quad (35)$$

Likewise,

$$\begin{aligned} \text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) &= \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) + 2 \int_X \int_X w(x-u) [f_\rho(x) \\ &\quad - f_\rho(u) - \nabla f_\rho(u)(x-u)^T] [f_\rho(u) - f_1(u) + \\ &\quad (\vec{f}_2(u) - \nabla f_\rho(u))(u-x)^T] d\rho_X(x) d\rho_X(u) \\ &\leq \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) + 2 \left(\int_X \int_X w(x-u) [f_\rho(x) - f_\rho(u) \right. \\ &\quad \left. - \nabla f_\rho(u)(x-u)^T]^2 d\rho_X(x) d\rho_X(u) \right)^{\frac{1}{2}} \left(\mathcal{Q}_s(\vec{f}, \vec{f}_\rho) \right)^{\frac{1}{2}}. \end{aligned}$$

Applying the inequality $2ab \leq a^2 + b^2$ to the above inequality yields that

$$\begin{aligned} \text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) &\leq 2\mathcal{Q}_s(\vec{f}, \vec{f}_\rho) + \int_X \int_X w(x-u) \\ &\quad [f_\rho(x) - f_\rho(u) - \nabla f_\rho(u)(x-u)^T]^2 d\rho_X(x) d\rho_X(u) \end{aligned} \quad (36)$$

However, $|f_\rho(x) - f_\rho(u) - \nabla f_\rho(u)(x-u)^T| = \left| \int_0^1 (\nabla f_\rho(tx + (1-t)u) - \nabla f_\rho(u))(x-u)^T dt \right| \leq |\nabla f_\rho|_{\text{Lip}} \|x-u\|^2$ and the density $p(x)$ of ρ_X is a bounded function since we assume it is θ -Hölder continuous and X is compact. Therefore,

$$\begin{aligned} &\int_X \int_X w(x-u) [f_\rho(x) - f_\rho(u) - \nabla f_\rho(u)(x-u)^T]^2 \\ &\quad \times d\rho_X(x) d\rho_X(u) \\ &\leq \|p\|_\infty |\nabla f_\rho|_{\text{Lip}}^2 \left[\int_{\mathbb{R}^d} \frac{1}{s^{d+2}} e^{-\frac{\|x\|^2}{2s^2}} \|x\|^4 dx \right] \\ &\leq c \|p\|_\infty |\nabla f_\rho|_{\text{Lip}}^2 s^2. \end{aligned}$$

Putting this into Equations (35) and (36) and arranging the terms involved yields the desired result. \blacksquare

From Property (1) of Lemma 9, for any $\vec{f} \in \mathcal{H}_\mathcal{K}$ we have that

$$\text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) \geq -cs^2 |\nabla f_\rho|_{\text{Lip}}^2. \quad (37)$$

We now turn our attention to the second step of the error analysis: the estimation of the excess error $\text{Er}(\vec{f}_\mathbf{z}) - \text{Er}(\vec{f}_\rho) + \lambda \|\vec{f}_\mathbf{z}\|_{\mathcal{K}}^2$. To do this, let

$$\vec{f}_\lambda = \arg \inf_{\vec{f} \in \mathcal{H}_\mathcal{K}} \left\{ \text{Er}(\vec{f}) + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \right\}$$

By the *error decomposition* technique in [6], we get the following estimation.

Proposition 3 *If \vec{f}_λ is defined above then there exists a constant c such that*

$$\begin{aligned} \text{Er}(\vec{f}_\mathbf{z}) - \text{Er}(\vec{f}_\rho) + \lambda \|\vec{f}_\mathbf{z}\|_{\mathcal{K}}^2 &\leq \mathcal{S}(\mathbf{z}) \\ &\quad + c(s^2 |\nabla f_\rho|_{\text{Lip}}^2 + \mathcal{A}(\lambda, s)), \end{aligned}$$

where

$$\mathcal{S}(\mathbf{z}) = \text{Er}(\vec{f}_\mathbf{z}) - \text{Er}_\mathbf{z}(\vec{f}_\mathbf{z}) + \text{Er}_\mathbf{z}(\vec{f}_\lambda) - \text{Er}(\vec{f}_\lambda)$$

is referred to the *sample error* and

$$\mathcal{A}(\lambda, s) = \inf_{\vec{f} \in \mathcal{H}_\mathcal{K}} \left\{ \mathcal{Q}_s(\vec{f}, \vec{f}_\rho) + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \right\}$$

is called the *approximation error*.

Proof: Note that $\text{Er}(\vec{f}_\mathbf{z}) - \text{Er}(\vec{f}_\rho) + \lambda \|\vec{f}_\mathbf{z}\|_{\mathcal{K}}^2 = [\text{Er}(\vec{f}_\mathbf{z}) - \text{Er}_\mathbf{z}(\vec{f}_\mathbf{z}) + \text{Er}_\mathbf{z}(\vec{f}_\lambda) - \text{Er}(\vec{f}_\lambda)] + [(\text{Er}_\mathbf{z}(\vec{f}_\mathbf{z}) + \lambda \|\vec{f}_\mathbf{z}\|_{\mathcal{K}}^2) - (\text{Er}_\mathbf{z}(\vec{f}_\lambda) + \lambda \|\vec{f}_\lambda\|_{\mathcal{K}}^2)] + [\text{Er}(\vec{f}_\lambda) - \text{Er}(\vec{f}_\rho) + \lambda \|\vec{f}_\lambda\|_{\mathcal{K}}^2]$. By the definition of $\vec{f}_\mathbf{z}$, we know that the second term in parenthesis on the right-hand side of the above equation is negative. Hence, by the definition of f_λ , we get that $\text{Er}(\vec{f}_\mathbf{z}) - \text{Er}(\vec{f}_\rho) + \lambda \|\vec{f}_\mathbf{z}\|_{\mathcal{K}}^2 \leq \mathcal{S}(\mathbf{z}) + \inf_{f \in \mathcal{H}_\mathcal{K}} \{ \text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \}$. By the property (2) in Lemma 9, we also have, for any $\vec{f} \in \mathcal{H}_\mathcal{K}$, that $\text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) \leq c \left((s^2 |\nabla f_\rho|_{\text{Lip}}^2 + \mathcal{Q}_s(\vec{f}, \vec{f}_\rho)) \right)$ which implies that

$$\inf \{ \text{Er}(\vec{f}) - \text{Er}(\vec{f}_\rho) + \lambda \|\vec{f}\|_{\mathcal{K}}^2 \} \leq c \left(s^2 |\nabla f_\rho|_{\text{Lip}}^2 + \mathcal{A}(\lambda, s) \right).$$

This completes the proposition. \blacksquare

Now it suffice to estimate the sample error $\mathcal{S}(\mathbf{z})$. To this end, observe that $\text{Er}_\mathbf{z}(\vec{f}_\mathbf{z}) + \lambda \|\vec{f}_\mathbf{z}\|_{\mathcal{K}}^2 \leq \text{Er}_\mathbf{z}(0) + \lambda \|0\|_{\mathcal{K}}^2 \leq \frac{M^2}{s^{d+2}}$ which implies that $\|\vec{f}_\mathbf{z}\|_{\mathcal{K}} \leq Ms^{-(d+2)/2}$. Likewise, $\text{Er}(\vec{f}_\lambda) + \lambda \|\vec{f}_\lambda\|_{\mathcal{K}}^2 \leq \text{Er}(0) + \lambda \|0\|_{\mathcal{K}}^2 \leq \frac{M^2}{s^{d+2}}$ which tells us that $\|\vec{f}_\mathbf{z}\|_{\mathcal{K}} \leq Ms^{-(d+2)/2}$. Using these bounds on $\|\vec{f}_\mathbf{z}\|_{\mathcal{K}} + \|\vec{f}_\lambda\|_{\mathcal{K}}$, we can use the Rademacher averages (see e.g. [4, 12]) for its definition and properties) to get the following estimation for the sample error.

Lemma 10 *For any $0 < \lambda < 1$, there exists a constant c such that, for any $m \in \mathbb{N}$, there holds*

$$\mathbb{E}[\mathcal{S}(\mathbf{z})] \leq c \left(\frac{1}{s^{2(d+2)} \lambda m} + \frac{1}{s^{3(d+2)/2} \sqrt{\lambda m}} \right).$$

Since the proof of the above lemma is rather a standard approach and indeed parallel to the proof of Lemma 26 (replacing $r = Ms^{-(d+2)/2}$ there) in the appendix of [15], for the simplicity we omit the details here.

We have assembled all the materials to prove Theorem 7.

Proof of Theorem 7

Since we assume that $\vec{f}_\rho \in \mathcal{H}_\mathcal{K}$, for any $|\partial_p f_\rho(x+u) - \partial_p \vec{f}_\rho(u)| = |\langle \tilde{e}_{p+1}, \vec{f}_\rho(x+u) - \vec{f}_\rho(u) \rangle| = |\langle \vec{f}_\rho, \mathcal{K}_{x+u} \tilde{e}_{p+1} - \mathcal{K}_u \tilde{e}_{p+1} \rangle| \leq \|\vec{f}_\rho\|_{\mathcal{K}} \left[\tilde{e}_{p+1}^T (\mathcal{K}(x+u, x+u) + \mathcal{K}(u, u) - \mathcal{K}(x+u, u) - \mathcal{K}(u, x+u)) \tilde{e}_{p+1} \right]^{\frac{1}{2}} \leq c \|\vec{f}_\rho\|_{\mathcal{K}} \|u\|$, and hence $|\nabla f_\rho|_{\text{Lip}} \leq c \|\vec{f}_\rho\|_{\mathcal{K}}$. Moreover,

$$\mathcal{A}(\lambda, s) \leq \mathcal{Q}(\vec{f}_\rho, \vec{f}_\rho) + \lambda \|\vec{f}_\rho\|_{\mathcal{K}}^2 = \lambda \|\vec{f}_\rho\|_{\mathcal{K}}^2.$$

Hence, we know from Proposition 3 and equation (37) that $\lambda \|\vec{f}_\mathbf{z}\|_{\mathcal{K}}^2 \leq \mathcal{S}(\mathbf{z}) + c'(s^2 + \lambda)$.

Combining the above equations with Propositions 2 and 3, there exists a constant c such that

$$\begin{aligned} \mathbb{E}[\|\vec{f}_\mathbf{z} - \vec{f}_\rho\|_{\rho_X}^2] &\leq c \left(\min(s^{-\theta}, \max_{x \in X} p^{-1}(x)) + \frac{s^\theta}{\lambda} \right) \\ &\quad \times [\mathbb{E}[\mathcal{S}(\mathbf{z})] + s^2 + \lambda] + s^\theta. \end{aligned}$$

If we choose $\lambda = s^{2\theta}$ and $s = m^{-\frac{1}{3(d+2)+4\theta}}$ yields the first assertion.

If ρ_X is the uniform distribution over X , then we have that $\min(s^{-\theta}, \min_{x \in X} p(x)) = 1$. Hence, choosing $\lambda = s^\theta$ and $s = m^{-\frac{1}{3(d+2)+5\theta}}$ we have the desired second assertion. This completes the theorem.