SVM Learning and L^p Approximation by Gaussians on Riemannian Manifolds

Gui-Bo Ye

School of Mathematical Sciences, Fudan University Shanghai 200433, China yeguibo@hotmail.com

Ding-Xuan Zhou

Department of Mathematics, City University of Hong Kong Kowloon, Hong Kong, China mazhou@cityu.edu.hk

April 30, 2007

Abstract

We confirm by the multi-Gaussian support vector machine (SVM) classification that the intrinsic dimension of Riemannian manifolds improves the efficiency (learning rates) of learning algorithms. The essential analysis lies in the study of approximation in L^p (1 p < p) of L^p functions by their convolutions with the Gaussian kernel with variance p0. This covers the SVM case when the approximated function is the Bayes rule and is not continuous in general. The approximation error is estimated by imposing some regularity conditions on the approximated function to lie in some interpolation spaces. Then the learning rates for multi-Gaussian regularized classifiers with general classification loss functions are derived, and the rates depend on the intrinsic dimension of the Riemannian manifold, not the dimension of the underlying Euclidean space. Here the input space is assumed to be connected compact p0. Riemannian submanifold of p1. Riemannian manifold and the radial basis form of Gaussian kernels play an important role.

Key words and phrases: Manifold learning, reproducing kernel Hilbert spaces, Gaussian kernels, approximation, multi-kernel regularized classifier, general loss function.

1 Introduction and Multi-Gaussian SVM

Manifold learning has become a hot topic in machine learning recently. It appears naturally from the observation or belief that high-dimensional data or functions arising from physical or biological systems can be effectively modeled or analyzed as being concentrated on a low-dimensional manifold. There have been many tasks for manifold learning demanded by many applications such as dimensionality reduction [4], feature selection [6, 17, 18], semi-supervised learning [3], and learning topological statistics [14]. Compared with vast practical observations and empirical testing, rigorous mathematical analysis in the topic of manifold learning is rather limited [2, 17, 14, 15, 18].

In [25] we consider the approximation of continuous functions on Riemannian manifolds by functions from reproducing kernel Hilbert spaces associated with Gaussian kernels. The obtained order of approximation is applied to the multi-kernel least-square regularization scheme generated by Gaussians with flexible variances. The derived learning rate is better than that in the setting of Euclidean space domains, which confirms the belief that the low dimensionality of manifolds improves the efficiency of learning algorithms.

Many problems in machine learning are about classification where an essential mathematical problem is the approximation of functions in spaces like $L^p(X)$, not in C(X). So in this paper we study the **approximation in** $L^p(X)$ by reproducing kernel Hilbert spaces associated with Gaussian kernels K_{σ} with variances σ 0. Then we apply the approximation order to get learning rates of **multi-Gaussian regularized classifiers** with general classification loss functions. The obtained learning rates depend on the **intrinsic dimension** of the Riemannian manifold, not the dimension of the underlying Euclidean space.

Let us mention the setting of binary classification and the special example of support vector machines.

1.1 Binary classification

Two classes dealt with by binary classification learning algorithms can be represented by $Y = \{1, -1\}$. The events for which the prediction of classes is desired are points from a

metric space X (called the input space, maybe a subset of \mathbb{R}^n). A model used in learning theory is to assume a probability measure ρ on $Z := X \times Y$, then the conditional distribution of ρ at x - X is a probability distribution $\rho(\cdot | x)$ on Y. For y = 1 or -1 in Y, P(y|x) stands for the probability for x to belong to the class y. The marginal distribution ρ_X of ρ on X measures how the events are distributed in X.

A binary classifier is a function \mathcal{C} from X to Y. It gives a prediction of class $\mathcal{C}(x)$ Y for each event x X. The **misclassification error** for the classifier \mathcal{C} is defined as the probability of wrong prediction

$$R(f) := \operatorname{Prob}\{C(x) = y\} = \sum_{X} P(y = C(x)/x) d\rho_{X}.$$
(1.1)

By discussing for each event x = X, we can easily see that a best classifier minimizing the misclassification error, called the **Bayes rule** (e.g. [11]), can be expressed as

$$f_c(x) = \begin{cases} 1, & \text{if } P(y = 1/x) \quad P(y = -1/x), \\ -1, & \text{if } P(y = 1/x) < P(y = -1/x). \end{cases}$$
 (1.2)

The purpose of classification algorithms is to find good classifier approximations C_Z of the Bayes rule from the random sample $\mathbf{z} = \{(x_i, y_i)\}_{i=1}^m$ drawn according to the probability measure ρ . We hope that the approximating classifier C_Z will approach the Bayes rule when the number of samples increases, in the sense that the **excess misclassification error** $R(C_Z) - R(f_c)$ tends to zero with confidence as m.

Most practical algorithms can be expressed mathematically as minimizers of some functionals over some spaces of continuous functions. A binary classifier can be derived from a continuous function f: X \mathbb{R} as $\operatorname{sgn}(f): X$ Y given by $\operatorname{sgn}(f)(x) = 1$ if f(x) = 0 and -1 otherwise. To measure how well $\operatorname{sgn}(f)$ can be used for binary classification, we often use a loss function $\phi: \mathbb{R}$ \mathbb{R}_+ and the value $\phi(yf(x))$ as the local error incurred in predicting y from f(x). Define the **generalization error** of f with respect to the loss function ϕ as

$$E_{\phi}(f) = \int_{Z} \phi(yf(x))d\rho \tag{1.3}$$

and the empirical error with respect to the loss function ϕ as

$$E_{\phi}^{\mathbf{Z}}(f) = \sum_{i=1}^{\mathbf{X}^n} \phi(y_i f(x_i)). \tag{1.4}$$

Many learning algorithms for classification involve this empirical error. Their error analysis (for the misclassification error) can be done by estimating the excess generalization error

 $E_{\phi}(f) - E_{\phi}(f_{\rho}^{\phi})$ where f_{ρ}^{ϕ} is the target function defined by

$$f_{\rho}^{\phi} = \arg \min \ E_{\phi}(f) : f \text{ is a measurable function from } X \text{ to } \mathbb{R} \ .$$
 (1.5)

1.2 Support vector machines

Support vector machines form an important class of classification algorithms. As a special kernel method, they can be expressed in reproducing kernel Hilbert spaces associated with Mercer kernels.

We say that $K: X \times X$ \mathbb{R} is a **Mercer kernel** if it is continuous, symmetric and positive semidefinite.

The **Reproducing kernel Hilbert space** (RKHS) H_K associated with the kernel K is defined [1] to be the completion of the linear span of the set of functions $\{K_x := K(x, \cdot)\}_{x \ge X}$ with the inner product $\cdot, \cdot K$ given by $K_x, K_y K = K(x, y)$.

For example, the Gaussian kernel with variance σ (0,) defined by

$$K_{\sigma}(x,y) = \exp \left(-\frac{|x-y|^2}{2\sigma^2}\right), \quad x,y \quad X \quad \mathbb{R}^n$$
 (1.6)

is a Mercer kernel (e. g. [9]).

The important role played by Mercer kernels in kernel methods can be seen from the **regularization scheme** for classification problem associated with the RKHS H_K and ϕ :

$$f_{\mathsf{Z},\lambda,K} = \arg\min_{f \ge \mathsf{H}_K} \frac{1}{m} \sum_{i=1}^{\mathsf{X}^n} \phi \ y_i f(x_i) + \lambda \ f \ _K^2 \ ,$$
 (1.7)

where $\lambda > 0$ is a constant called the **regularization parameter**. The classifier is given by the sign function $\operatorname{sgn}(f_{z,\lambda,K})$. The reproducing property of H_K :

$$K_x, f_K = f(x) \qquad x \quad X, f \quad H_K \tag{1.8}$$

together with the orthogonal projection in the Hilbert space H_K tells us that the minimizer in (1.7) has the form $f_{\mathsf{Z},\lambda,K} = \bigcap_{i=1}^m c_i^\mathsf{Z} K_{x_i}$. The coefficients $(c_i^\mathsf{Z})_{i=1}^m$ can be computed by solving an optimization problem which is convex when the loss function ϕ is admissible.

Definition 1. We say that $\phi : \mathbb{R} - \mathbb{R}_+$ is an admissible loss function if it is convex and differentiable at 0 with $\phi^0(0) < 0$.

A special setting is given by the hinge loss $\phi(x) = (1 - x)_+ := \max\{0, 1 - x\}$. The algorithm (1.7) with this special loss function ϕ is the support vector machine in the regularization form [13]. Its special form ensures that the convex optimization problem for finding $(c_i^z)_{i=1}^m$ in $f_{z,\lambda,K}$ is actually a convex quadratic programming one and many efficient interior point methods are available from optimization theory.

The minimizer f_{ρ}^{ϕ} for ϕ being the hinge loss is exactly the Bayes rule f_c . For the error analysis in this special case, Zhang [27] showed that $R(\operatorname{sgn}(f)) - R(f_c) = E_{\phi}(f) - E_{\phi}(f_c)$ which in turn [8] can be bounded by $f - f_c |_{L_{\rho_X}^1}$. Thus to estimate the excess misclassification error $R(\operatorname{sgn}(f_{Z,\lambda,K})) - R(f_c)$ for the efficiency of the SVM algorithm, we need to consider $f_{Z,\lambda,K} - f_c |_{L_{\rho_X}^1}$, the approximation of the generally discontinuous function f_c in the function space $L_{\rho_X}^1$, not in the space C(X) of all continuous functions on X. So the result from [25] cannot be used for SVM in the manifold setting. This is one motivation of our study in this paper.

1.3 Multi-kernel regularized classifiers

It was shown in [8, 19] that when f_{ρ} lies in the Sobolev space $H^{s}(X)$ with s > 0 and X is a domain of Euclidean space with nonempty interior, the learning rate of the algorithm (1.7) with a fixed Gaussian kernel and the least square loss $\phi(t) = (1 - t)^{2}$ is only $O((\log m)^{-s/2})$. If we allow flexible variances of Gaussian kernels, things are totally different and getting polynomial decay for the learning rate is possible [26, 25]. This confirms the usefulness of flexible variances in applying Gaussian kernels in practice.

In this paper we consider the multi-kernel regularized classifier $\operatorname{sgn}(f_{Z,\lambda})$ generated by the regularization scheme associated with the general loss function ϕ and Gaussians (1.6) with flexible variances $\{K_{\sigma}: 0 < \sigma < \}$ defined as

with flexible variances
$$\{K_{\sigma}: 0 < \sigma < \}$$
 defined as
$$()$$

$$f_{\mathsf{Z},\lambda} = \arg\min_{\sigma^{2}(0,+1)} \min_{f^{\mathsf{ZH}}} \frac{1}{m} \sum_{i=1}^{\mathsf{X}^{n}} \phi\left(y_{i}f(x_{i})\right) + \lambda f_{K_{\sigma}}^{2} .$$

$$(1.9)$$

Though multi-kernel algorithms for regression and classification have been applied extensively, their error analysis is well understood only if the input space X is a domain of \mathbb{R}^n with nonempty interior and the learning rate is not as fast as expected when the dimension n of the Euclidean space is large. It was pointed out in [26] that when the input space X is a low-dimensional manifold embedded in the large-dimensional Euclidean space, the learning rates may be greatly improved. In such a manifold setting the Fourier transform technique in [26] can no longer be used and other methods are required.

We took a step toward this problem in [25], where we obtained satisfactory learning rates for the multi-kernel regression algorithm with the least square loss by means of uniformly normal neighborhoods in Riemannian manifolds. But the involved approximation scheme there cannot be directly used for the classification setting as mentioned in §1.2. The minimizing function f_{ρ}^{ϕ} associated with the loss function ϕ in the classification problem is often discontinuous and we need to consider the approximation problem in the space $L^{p}(X)$, not in C(X).

2 Main Results on Riemannian Manifolds

In what follows we assume that X is a connected compact C^1 submanifold of \mathbb{R}^n without boundary which is isometrically embedded and of dimension d. Under this assumption, X is a metric space with a metric d_X and the inclusion map $\Phi: (X, d_X) \subset (\mathbb{R}^n, \cdot)$ is well defined and continuous (actually it is C^1). Here \cdot is the norm in \mathbb{R}^n . Our assumption that the embedding map Φ is the inclusion map is essential. For a general embedding map (which always exits according to the Nash Embedding Theorem), we still do not know how to establish similar results.

For 1 p < R, the space $L^p(X)$ on X consists of all measurable functions on X with the norm $f_{L^p(X)} = \frac{1}{X} / f(x) / p dV(x)$ finite, where Y is the Riemannian volume measure of X.

To measure the regularity of functions on X, we need Sobolev spaces on the Riemannian manifold X. For an integer k and f $C^1(X)$, kf denotes the kth covariant derivative of f (with the convention ${}^0f = f$). As an example, the components of f in local coordinates are given by ($f_{ij} = \partial_{ij}f - \int_{\ell=1}^{d} \Gamma_{ij}^{\ell}\partial_{\ell}f$ where Γ_{ij}^{ℓ} is Christoffel symbols of f with respect to f_{ij}^{d} , see [16] for more details.

Definition 2. Let p-1 and $k-\mathbb{N}$. The Sobolev space $H_k^p(X)$ is the completion of $\mathcal{C}^1(X)$ with respect to the norm

$$f_{H_k^p(X)} = X^k Z_{j=0} X^{-j} f^p dV^{1/p}.$$

Recall that the Bayes rule f_c is discontinuous in general. Its regularity may not be satisfactorily characterized by Sobolev spaces $H_k^p(X)$ with integer indices k. To get suitable

characterizations, we need interpolation spaces [5] from which one can define Sobolev spaces of arbitrary indices [12].

Definition 3. Let $0 < \theta$ 1. The interpolation space $(L^p(X), H_2^p(X))_{\theta}$ between $L^p(X)$ and $H_2^p(X)$ consists of all functions f $L^p(X)$ such that the norm

$$f_{\theta} := \sup_{t>0} \frac{\mathbb{K}(f,t)}{t^{\theta}} \tag{2.1}$$

is finite. Here $\mathbb{K}(f,t)$ is the \mathbb{K} -functional of the pair $(L^p(X),H_2^p(X))$ defined by

$$\mathbb{K}(f,t) = \inf_{g \ge H_2^p(X)} \quad f - g_{L^p(X)} + t \ g_{H_2^p(X)}, \ t > 0.$$
 (2.2)

The interpolation space $(L^p(X), H_2^p(X))_{\theta}$ is a Banach space between $L^p(X)$ (with $\theta = 0$) and $H_2^p(X)$ (with $\theta = 1$). It can be easily seen that the function $\mathbb{K}(f,t)$ of t is continuous, non-decreasing and bounded by $f_{L^p(X)}$ (take g = 0 in (2.2)). Since $H_2^p(X)$ is dense in $L^p(X)$, $\mathbb{K}(f,t)$ tends to zero as t=0. The requirement that the function f lies in $(L^p(X), H_2^p(X))_{\theta}$ is equivalent to the decay condition $\mathbb{K}(f,t) = O(t^{\theta})$.

Example 1. Let $X = S^1 = \{e^{iu} : u \mid \mathbb{R}\}$ $\mathbb{R}^2, p = 2$ and 0 < a < b < 1. Define a function f $L^2(X)$ by $f(e^{iu}) = 1$ for a u b and zero for u $[0, 2\pi] \setminus [a, b]$. Then f is not in $H_1^2(X)$ since f $H_1^2(X)$ would imply f C(X). But for any $0 < \theta$ $\frac{1}{4}$, we have f $(L^p(X), H_2^p(X))_{\theta}$. In fact, by choosing g_t as

$$g_{t}(e^{iu}) = \begin{cases} \mathbf{8} & \text{if } u = [a,b], \\ \frac{2}{t}(u - (a - \bar{t}))^{2} & \text{if } u = [a - \bar{t}, a - \frac{\mathsf{p}_{\bar{t}}}{2}), \\ 1 - \frac{2}{t}(u - a)^{2} & \text{if } u = [a - \frac{\bar{t}}{2}, a), \\ \mathbf{1} - \frac{2}{t}(u - b)^{2} & \text{if } u = (b, b + \frac{\bar{t}}{2}), \\ \frac{2}{t}(u - (b + \bar{t}))^{2} & \text{if } u = (b + \frac{\bar{t}}{2}, b + \bar{t}], \\ 0 & \text{if } u = [0, 2\pi) \sqrt{[a - \bar{t}, b + \bar{t}]}, \end{cases}$$

we have $f - g_{t-L^2(X)} + t g_{t-H_2^p(X)} = O(t^{\frac{1}{4}}).$

Using the regularity condition imposed by interpolation spaces, we can state our first main result concerning the learning rates of the multi-kernel SVM with flexible Gaussians on Riemannian manifolds.

Theorem 1. Let X be a connected compact C^1 submanifold of \mathbb{R}^n without boundary which is isometrically embedded and of dimension d. Let $f_{z,\lambda}$ be defined by (1.9) and $\phi(x) = (1-x)_+$.

If $f_c = (L^1(X), H_2^1(X))_{\theta}$ for some $0 < \theta = 1$, then by taking $\lambda = \frac{\log^2 m}{m}$, we have

$$\mathbb{E}_{\mathsf{Z2}\,\mathsf{Z}^m} \ \ \mathsf{R}(sgn(f_{\mathsf{Z},\lambda}) - \mathsf{R}(f_c)) \qquad \mathbf{\mathfrak{E}} \ \ \frac{\log^2 m}{m} \ \ \frac{\theta}{6\theta + d}, \tag{2.3}$$

where \mathfrak{G} is a constant independent of m.

Theorem 1 is exciting since the learning rate depends only on the intrinsic dimension d of the Riemannian manifold X, not on its extrinsic dimension n. As d is very small and much less than n in many practical problems, our learning rate is satisfactory and convincing theoretical studies in manifold learning. This is another motivation of our investigation in this paper.

Our second main result is about the approximation ability of Gaussians on Riemannian manifolds. This theorem is of importance on its own in approximation theory and it is the key step to prove Theorem 1.

When X has nonempty interior as a subset of $\mathbb{R}^n(d=n)$, the approximation of functions from various function spaces by Gaussians is a classical topic in approximation theory [12] and its application in error analysis has been well understood [19, 10, 22]. Things are totally different when X is a Riemannian submanifold of \mathbb{R}^n and little is known. In [25], we considered the approximation ability of Gaussians on the space C(X). Here we consider the approximation on the space $L^p(X)$.

Let 1
$$p$$
 . Define a family of linear operators $\{I_{\sigma}: L^{p}(X) \quad L^{p}(X)\}_{\sigma>0}$ as
$$I_{\sigma}(f)(x) = \frac{1}{\left(-\frac{1}{2\pi}\sigma\right)^{d}} \underset{X}{\mathsf{Z}^{X}} K_{\sigma}(x,y)f(y)dV(y)$$

$$= \frac{1}{\left(-\frac{1}{2\pi}\sigma\right)^{d}} \underset{X}{\mathsf{exp}} -\frac{x-y^{-2}}{2\sigma^{2}} f(y)dV(y), \qquad x \in X, \qquad (2.4)$$

where V is the Riemannian volume measure of X.

Note that a d-dimensional manifold is, roughly speaking, a topological space which is locally Euclidean of dimension d. That's why we use the scaling factor $\frac{1}{2\pi\sigma^d}$.

Theorem 2. Let X be a connected compact C^1 submanifold of \mathbb{R}^n without boundary which is isometrically embedded and of dimension d. Let p-1. Define $I_{\sigma}: L^p(X) = L^p(X)$ for $\sigma > 0$ by (2.4). If $f = H_2^p(X)$, then we have

$$I_{\sigma}(f) - f_{L^{p}(X)} C_{X} f_{H_{2}^{p}(X)} \sigma^{2} \sigma > 0,$$
 (2.5)

where C_X is a positive constant independent of f or σ .

The main difficulty in the proof of Theorem 2 (given in Section 3) lies in bounding the integrals over uniformly normal neighborhoods of the convolutions with Gaussians. This is different from the approximation in C(X) where only function values need to be bounded [25].

Due to a saturation phenomenon in approximation theory, the order of approximation in (2.5) cannot be further increased for functions in higher order Sobolev spaces.

The methods of deriving learning rates in this paper can be extended to other learning algorithms on Riemannian manifolds such as online learning for classification [24] and other L^p problems on Riemannian manifolds.

3 L^p Approximation on Manifolds by Gaussians

In this section we prove Theorem 2 after some preparation and illustration.

3.1 Some knowledge on Riemannian manifolds

We start with a brief introduction to normal coordinates and uniform normal neighborhoods (see [14] and [25] in detail). These two basic concepts provide good coordinate systems on Riemannian manifolds and make computations easier. Denote the tangent space at \mathbf{p} as $T_{\mathbf{p}}(X)$.

Definition 4. For \mathbf{p} X and v $T_{\mathsf{p}}(X)$, let $\gamma(t,\mathbf{p},v),t>0$, be the geodesic satisfying $\gamma(0,\mathbf{p},v)=\mathbf{p}$ and $\gamma^{\mathsf{0}}(0,\mathbf{p},v)=v$. The **exponential map** $E_{\mathsf{p}}:T_{\mathsf{p}}(X)$ X is defined by $E_{\mathsf{p}}(v)=\gamma(1,\mathbf{p},v)$.

By [7], we know that for each $\mathbf{p} = X$, there exists a strongly convex neighborhood $U_{\mathbf{p}}$ of \mathbf{p} , that is, for any two points $\mathbf{q}_1, \mathbf{q}_2$ in the closure $\overline{U}_{\mathbf{p}}$ of $U_{\mathbf{p}}$, there exists a unique minimizing geodesic γ joining \mathbf{q}_1 and \mathbf{q}_2 whose interior is contained in $U_{\mathbf{p}}$. Denote $B_{\delta}(0) = \{v = T(X) : |v| < \delta\}$ as the ball of T(X) centered at 0 with radius δ .

Definition 5. An open set U = X is called **uniformly normal** if there exists some $\delta > 0$ such that $U = E(B_{\delta}(0))$ for every $\mathbf{q} = U$.

The following proposition (see [25]) tells us the existence of uniform normal neighborhoods having some good properties.

Proposition 1. For every \mathbf{p} X there exist neighborhoods W_{p} and W_{p} and a number $\delta_{\mathsf{p}} > 0$ such that the following conditions hold:

- (a) for every \mathbf{q} $W_{\mathbf{p}}$, the map $\mathbf{E}: B_{\delta_{\mathbf{p}}}(0)$ T(X) X is a diffeomorphism on $B_{\delta_{\mathbf{p}}}(0)$;
- (b) W_{p} is uniformly normal with respect to δ_{p} , that is, $W_{\mathsf{p}} = \mathcal{E}\left(B_{\delta_{\mathsf{p}}}(0)\right)$ for every $\mathbf{q} = W_{\mathsf{p}}$;
- (c) The closure of W_p is contained in W_p and W_p U_p .

Choose an orthonormal basis $\{e_1, e_2, \dots, e_d\}$ of $T_p(X)$, then for each \mathbf{q} U_p , the set of tangent vectors $\{e_1, e_2, \dots, e_d\}$, moved by parallel transport from \mathbf{p} to \mathbf{q} along the unique minimizing geodesic, forms an orthonormal basis of T(X). In addition, this frame depends smoothly on \mathbf{q} . According to (a) of Proposition 1-p the map ϕ from $U = \{u \mid \mathbb{R}^d : u < \delta_p\}$ \mathbb{R}^d to X defined by $\phi(u^1, \dots, u^d) = E(\begin{pmatrix} d \\ i=1 \end{pmatrix} u^i e_i)$ gives a system of coordinates around \mathbf{q} . We call such coordinates \mathbf{q} -normal coordinates. Under these normal coordinates,

$$d_X \mathbf{q}, E \begin{pmatrix} \mathbf{X}^d \\ u^i e_i \end{pmatrix} = u \qquad u < \delta_{\mathbf{p}}, \qquad (3.1)$$

where d_X is the Riemannian metric of X.

In addition, the Riemannian structure g of the isometrically embedded manifold X under the **q**-normal coordinate (U, ϕ) can be expressed as

$$g_{ij}(u^1, \dots, u^d) := d\Phi \quad \frac{\partial}{\partial u^i}(\mathbf{q}) , d\Phi \quad \frac{\partial}{\partial u^j}(\mathbf{q}) \quad \mathbb{R}^n, \ i, j = 1, \dots, d.$$
 (3.2)

Here $\frac{\partial}{\partial u^i}(\mathbf{q}) \stackrel{d}{=1}$ is a basis of T(X) [7, 25] and $d\Phi$ is a map from T(X) to $T(\mathbb{R}^n)$ induced by the inclusion map Φ .

For each $\mathbf{q} = W_{\mathsf{p}}$ and u = U, g_{ij} is well defined and is C^1 as a function on $W_{\mathsf{p}} \times U$. It satisfies $g_{ij}(0) = \delta_{ij}$ and furthermore, we have the following proposition.

Proposition 2. For \mathbf{p} X, choose W_{p} and δ_{p} as in Proposition 1. For each \mathbf{q} W_{p} , choose the \mathbf{q} -normal coordinate (U, ϕ) and the corresponding local representation g_{ij} of the Riemannian structure as above. Then the following two bounds hold with a constant C_{p} independent of \mathbf{q} W_{p} :

$$q \frac{1}{\det(g_{ij})}(u^1, \cdots, u^d) - 1 \qquad C_p \quad u^2 \qquad u \qquad \delta_p, \qquad (3.3)$$

$$d_X(\mathbf{q}, x)^2 - \mathbf{q} - x^2 C_p d_X(\mathbf{q}, x)^4 x E(B_{\delta_p}(0)).$$
 (3.4)

This proposition is a slight variation of Proposition 2.2 in [14] and it is easy to give a self-contained proof as in [25]. So we omit the proof here.

In this paper, we have assumed that X is a Riemannian submanifold of \mathbb{R}^n . For each pair (x,y) of points on X, we have the distance $d_X(x,y)$ under the Riemannian metric and the distance x-y under the Euclidean metric. The following lemma concerning a relationship between these two metrics was proved in detail in [25].

Lemma 1. There exists a positive constant C_0 1 such that

$$d_X(x,y) C_0 x - y x, y X. (3.5)$$

This lemma will be used frequently in the following since in learning processes we do not see the Riemannian metric d_X . We can only use the Euclidean norm \cdot . But in analysis, we can assume the existence of d_X and make good use of it.

3.2 An illustration of computing integrals on manifolds

In order to get some ideas of using \mathbf{q} -normal coordinates system to compute some integrals on the Riemannian manifold X, we prove the following lemma.

Lemma 2. For the Gaussian kernel K_{σ} defined by (1.6), we have

$$Z = K_{\sigma}(\mathbf{q}, y)dV(y) \quad \mathcal{C}_{X}\sigma^{d} \quad \mathbf{q} \quad X, \sigma > 0,$$
(3.6)

where \mathfrak{S}_X is a constant independent of σ or \mathbf{q} .

Proof. Let W_{p} , δ_{p} and C_{p} as in Proposition 1 and Proposition 2. Denote $W_{\mathsf{p}}^{\mathsf{0}} := W_{\mathsf{p}}$ $\mathcal{E}_{\mathsf{p}}(B_{\delta_{\mathsf{p}}/2}(0))$. Since $X_{\mathsf{p2P}}W_{\mathsf{p}}^{\mathsf{0}}$ and X is compact, there exists a finite subset P of X such that $X_{\mathsf{p2P}}W_{\mathsf{p}}^{\mathsf{0}}$.

Let $\delta = \min_{\mathsf{p} \, \mathsf{2P}} \min \frac{1}{2C_{\mathsf{p}}}, \delta_{\mathsf{p}} > 0$. Let the constant C_0 as in (3.5). Let $\mathbf{q} \in \mathbb{R}$ belongs to some W^0_{p} with $\mathbf{p} \in P$. Choose $B_{\sigma} := \{x \in \mathbb{R} \mid X : d_X(\mathbf{q}, x) < C_0 \in \mathbb{R} \mid \overline{2d + 6\sigma} \mid \overline{\log \frac{1}{\sigma}} \}$. Choose a constant $0 < \sigma_0 = 1$ such that $C_0 = \overline{2d + 6\sigma_0} \mid \overline{\log \sigma_0} \mid \overline{\log$

We first consider the case when $0 < \sigma < \sigma_0$. Since E is a differmorphism on $B_{\delta^*}(0)$, using the equality (3.1), we have $B_{\sigma} = E(B_{\delta^*}(0))$ and $B_{\sigma} = E(B_{\delta^*}(0)) = E(B_{\delta^*}(0))$, where

$$\boldsymbol{\beta}_{\sigma} := u \quad \mathbb{R}^{d} : u < C_{0} \quad \overline{2d + 6\sigma} \quad \log \frac{1}{\sigma} . \tag{3.7}$$

For $u=(u^1,\ldots,u^d)$ \mathbb{R}^d , denote ϕ $(u)={\it E}$ $({\stackrel{\mbox{\sf P}}{=}}_{i=1}u^ie_i)$, then $B_\sigma=\phi$ $(\it B\hskip-1.5mu\sigma)$.

Using (3.1) and the inequality (3.4) in Proposition 2, we have

$$d_X^2(\mathbf{q}, \phi(u)) - \mathbf{q} - \phi(u)^2 = u^2 - \mathbf{q} - \phi(u)^2 \qquad C_p d_X^4(\mathbf{q}, \phi(u)) = C_p u^4.$$

By the definition of δ , $d_X(\mathbf{q}, \phi\ (u)) = u < \delta$ min $\frac{1}{2C_{\mathbf{p}}}, \delta_{\mathbf{p}}$. Hence,

$$\frac{1}{2} u^2 \qquad \phi (u) - \mathbf{q}^2 \quad \frac{3}{2} u^2 \qquad u \quad \mathcal{B}_{\sigma}.$$
 (3.8)

In addition, by the inequality (3.3) in Proposition 2,

$$\frac{1}{2} \quad {\mathsf{q}} \frac{1}{\det(g_{ij})} (u^1, u^2, \cdots, u^d) \quad \frac{3}{2} \quad u \quad \mathcal{B}_{\sigma}. \tag{3.9}$$

Decompose the domain X into two parts B_{σ} and $X \setminus B_{\sigma}$. We have

$$\sum_{X} K_{\sigma}(\mathbf{q}, y) dV(y) = \sum_{B_{\sigma}^{\mathbf{q}}} K_{\sigma}(\mathbf{q}, y) dV(y) + \sum_{X \cap B_{\sigma}^{\mathbf{q}}} K_{\sigma}(\mathbf{q}, y) dV(y) := J_{1}(\mathbf{q}) + J_{2}(\mathbf{q}).$$

Using the local representation of the Riemannian volume measure under the **q**-normal coordinates involving a measurable function h on B_{σ} :

the first term $J_1(\mathbf{q})$ is

$$J_{1}(\mathbf{q}) = \sum_{\widetilde{B}_{\sigma}}^{\mathbf{Z}} \exp \left(-\frac{\mathbf{q} - \phi(u)^{2}}{2\sigma^{2}}\right)^{\mathbf{q}} \frac{\mathbf{q}}{\det(g_{ij})(u)du.$$

By the inequalities (3.8) and (3.9), we have

$$J_1(\mathbf{q})$$
 $\frac{3}{2} \sum_{\tilde{B}_{\sigma}}^{\mathbf{Z}} \exp \left(-\frac{u^2}{4\sigma^2}\right) du = \frac{3}{2} \sum_{\mathbb{R}^d}^{\mathbf{Z}} \exp \left(-\frac{u^2}{4\sigma^2}\right) du.$

Using the radial coordinates in \mathbb{R}^d , for any univariate function $\psi(r) : \mathbb{R}_+$ \mathbb{R} , we have the following equality for the radial function $\psi(y)$:

$$\frac{Z}{W^{d}} \psi(y) dy = \frac{2\pi^{d/2}}{\Gamma(d/2)} \sum_{0}^{Z} \psi(r) r^{d-1} dr, \qquad (3.11)$$

where Γ is the Gamma function given for α (0,) by $\Gamma(\alpha) = \frac{\mathsf{R_1}}{\mathsf{0}} r^{\alpha-1} e^{-r} dr$.

Applying (3.11) to the function $\psi(r) = e^{-\frac{r^2}{4\sigma^2}}$, we have $J_1(\mathbf{q}) = 3 \cdot 2^{d-1} \pi^{\frac{d}{2}} \sigma^d$.

As for the second term $J_2(\mathbf{q})$, we notice that for y $X \setminus B_{\sigma}$, the restriction $d_X(\mathbf{q}, y)$ $C_0 = \overline{2d + 6}\sigma = \overline{\log \frac{1}{\sigma}}$ together with (3.5) implies $\mathbf{q} - y = \overline{2d + 6}\sigma = \overline{\log \frac{1}{\sigma}}$. Thus

$$J_{2}(\mathbf{q}) = Z^{X \cap B_{\sigma}^{\mathbf{q}}} \exp \left[-\frac{(2d+6)\sigma^{2} \log \frac{1}{\sigma}}{2\sigma^{2}} dV(y)\right]$$

$$= Z^{X \cap B_{\sigma}^{\mathbf{q}}} \sigma^{d+3} dV(y) \quad \text{Vol}(X)\sigma^{d}.$$

Combining the estimates for $J_1(\mathbf{q})$ and $J_2(\mathbf{q})$, for $0 < \sigma < \sigma_0$, we have

$$K_{\sigma}(\mathbf{q}, y)dV(y) \qquad 3 \cdot 2^{d-1}\pi^{\frac{d}{2}} + \text{Vol}(X) \quad \sigma^{d}.$$

For the case $\sigma = \sigma_0$, it is easy to see from $e^{-r} = 1$ for r = 0 that

$$\frac{\mathbf{Z}}{\mathbf{X}} K_{\sigma}(\mathbf{q}, y) dV(y) \quad \text{Vol}(\mathbf{X}) \quad \frac{\text{Vol}(\mathbf{X})}{\sigma_0^d} \sigma^d.$$

This proves the desired result with the constant $\mathcal{C}_X = 3 \cdot 2^{d-1} \pi^{\frac{d}{2}} + \operatorname{Vol}(X) + \frac{\operatorname{Vol}(X)}{\sigma_0^d}$.

Remark 1. In the proof of Lemma 2, we only need W_p instead of W_p^0 . To be consistent with later discussion, we use W_p^0 here.

3.3 Some ideas for proving Theorem 2

Since $I_{\sigma}(f) - f_{L^p(X)}^p = \frac{\mathsf{R}}{X} \frac{\mathsf{R}}{X} \exp \left(-\frac{\mathsf{k} x \ y \mathsf{k}^2}{2\sigma^2}\right) f(y) dV(y) - f(x)^p dV(x)$ involves two layers of integrals, we need to decompose it twice to make the integral computable in local coordinates.

P Let W_p^0 , P and σ_0 be as in §3.2, we know that $X_{p2P} W_p^0$. Thus $I_{\sigma}(f) - f_{L^p(X)} = \int_{p2P} W_p^0 / I_{\sigma}(f)(x) - f(x) / p dV(x)^{-\frac{1}{p}}$. It will be seen in the following Proposition 4 that the operator $I_{\sigma}: L^p(X) = L^p(X)$ is uniformly bounded (the bound is independent of R_p^0). Furthermore, $C^1(X)$ is dense in $H_2^p(X)$. Thus the problem becomes to estimate $R_{W_p^0} / I_{\sigma}(f)(x) - f(x) / p dV(x)$ for each R_p^0 and R_p^0 are R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 and R_p^0 and R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 and R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 and R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 and R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 are R_p^0 and R_p^0 are R_p^0 and R_p^0 and R_p^0 are R_p^0 are R_p^0 and R_p^0 are R_p^0 and R_p^0 are R_p^0 are R_p^0 are R_p^0 are R_p^0 are R

Note that $I_{\sigma}(f)$ in the expression $\underset{W'_{\mathbf{p}}}{\mathsf{R}}/I_{\sigma}(f)(x) - f(x)/\!\!/^{2}dV(x)$ still contains an integral over the whole manifold X. We need to decompose it further. Let \mathbf{q} $W_{\mathsf{p}}^{\mathsf{0}}$. Choose B_{σ} and \mathcal{B}_{σ} as in §3.2.

Separating the domain
$$X$$
 into two parts B_{σ} and $X \setminus B_{\sigma}$, we have
$$Z = \frac{1}{(-\overline{2\pi}\sigma)^d} \sum_{B_{\sigma}^{\mathbf{q}}} K_{\sigma}(\mathbf{q},y) f(y) dV(y) + \frac{1}{(-\overline{2\pi}\sigma)^d} \sum_{X \cap B_{\sigma}^{\mathbf{q}}} K_{\sigma}(\mathbf{q},y) f(y) dV(y).$$

The second term of the above equation can be easily bounded due to the fast decay of $K_{\sigma}(\mathbf{q},y)$. Using (3.10), the first term equals

$$\frac{1}{(\overline{2\pi}\sigma)^d} \sum_{\widetilde{B}_{\sigma}} \exp -\frac{\mathbf{q} - \phi(u)}{2\sigma^2} f(\phi(u)) \mathbf{q} \frac{\mathbf{d}}{\det(g_{ij})}(u) du. \tag{3.12}$$

For approximation in C(X), the quantity $f \phi(u)$ can be easily bounded by the uniform norm $f_{C(X)}$.

For L^p approximation, we need to tackle the following problem.

Question 1: How can one bound the expression (3.12) in terms of $f_{H_2^p(X)}$ by treating $f(\phi(u))$ properly?

In the further decompositions, the term $f(\phi(u)) - f(\mathbf{q})$ naturally appears. Since f $H_2^p(X)$, it reminds us of the Taylor expansion of $f(\phi(u))$ in its integral form. We denote by D^kh the kth derivative of a function h on the Euclidean space \mathbb{R}^d . That is, the components of $D^k h$ are given by $(D^k h)_{i_1 i_d} = \partial_{i_1 i_d} h$, where (i_1, \cdots, i_d) $i_1 + \cdots + i_d = k$. Then

$$f(\phi(u)) - f(\mathbf{q}) = f(\mathbf{q}), \sum_{i=1}^{\mathbf{X}^d} e_i u^i + \sum_{0}^{1} (1-y)D^2(f(\phi(yu)))(u,u)dy.$$
(3.13)

Question 2: How can one bound (3.13) in terms of $f_{H_2^p(X)}$ by handling $D^2(f(\phi(yu)))$ properly?

The above two problems will essentially be solved by the following Proposition 3. It gives us bounds of $\frac{1}{W_{\mathbf{p}}}/f \phi(u) / dV(\mathbf{q})$ and $\frac{1}{W_{\mathbf{p}}} D^2(f - \phi)(u) / dV(\mathbf{q})$ in terms of $f_{H_2^p(X)}$.

Proposition 3. Let \mathbf{p} X and ϕ (u) = E $\bigcap_{i=1}^{\mathsf{P}} u^i e_i$ for $u = (u^1, u^2, \dots u^d)$ \mathbb{R}^d and $\mathbf{q} \quad W_{\mathsf{p}}^{\mathsf{0}} = W_{\mathsf{p}} \quad E_{\mathsf{p}}(B_{\delta_{\mathsf{p}}/2}(0)). \quad \text{Then there exists a constant } \delta_{\mathsf{p}}^{\mathsf{0}} \text{ satisfying } 0 < \delta_{\mathsf{p}}^{\mathsf{0}} \quad \frac{\delta_{\mathsf{p}}}{2} \text{ such that for all } u \quad \delta_{\mathsf{p}}^{\mathsf{0}}, \text{ we have}$

$$D^{2}(f \phi)(u)^{p}dV(\mathbf{q}) C_{p}^{00} f_{H_{2}^{p}(X)}^{p} f C_{2}^{p}(X), \qquad (3.15)$$

where $C_{\rm p}^{\rm 0}$ and $C_{\rm p}^{\rm 00}$ are two constants independent of f or u.

Proof. Let $W_{\mathbf{R}}$ be as in Proposition 1. Let $\mathbf{q} = \phi^{\mathsf{p}}(x)$ W_{p} and $\Re(x, u) = h(\phi^{\mathsf{p}}(x), u) = \phi(u)$. Write $W_{\mathbf{p}}$ $f(\phi(u))^p dV(\mathbf{q})$ in the **p**-normal coordinate, we know that

$$\frac{\mathsf{Z}}{W_{\mathbf{p}}'} / f(\phi(u)) / dV(\mathbf{q}) = \frac{\mathsf{Z}}{\int_{\mathsf{R}^d: \phi^{\mathbf{p}}(x)} 2W_{\mathbf{p}}' g} \int_{\mathsf{R}^d: \phi^{\mathbf{p}}(x)} f(\mathbf{R}(x, u)) d\mathbf{q} \frac{\mathsf{Z}}{\det g_{ij}^{\mathbf{p}}(x) dx}. \tag{3.16}$$

From the definition of $g_{ij}^{\mathsf{p}}(x)$, we know that $\gcd(x, y) > 0$ for all $x < \delta_{\mathsf{p}}$. On the other hand, $d_X(\mathbf{p}, \Re(x, u)) = d_X(\mathbf{p}, \phi^{\mathsf{p}}(x)) + d_X(\mathbf{q}, \phi^{\mathsf{p}}(u)) = \frac{\delta_{\mathsf{p}}}{2} + u = \delta_{\mathsf{p}}$ for all $(x, u) = B_1 := \{(x, u) : \phi^{\mathsf{p}}(x) = \overline{W_{\mathsf{p}}^{\mathsf{p}}}, \quad u = \frac{\delta_{\mathsf{p}}}{2}\}$. Thus $\gcd(x, u) = (\phi^{\mathsf{p}})^{-1} = \Re(x, u) > 0$ for all $(x, u) = B_1$. Let $g(x, u) = (\phi^{\mathsf{p}})^{-1} = \Re(x, u)$ and J(x, u) be the Jacobian of g(x, u) with respect to the variable x, i.e.,

$$J(x,u) = \begin{array}{cccc} \frac{\partial g_1(x,u)}{\partial x^1} & \frac{\partial g_1(x,u)}{\partial x^2} & \cdots & \frac{\partial g_1(x,u)}{\partial x^d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g_d(x,u)}{\partial x^1} & \frac{\partial g_d(x,u)}{\partial x^2} & \cdots & \frac{\partial g_d(x,u)}{\partial x^d} \end{array}$$

It is easy to see that $g(x,0) = (\phi^{\mathsf{p}})^{-1} \Re(x,0) = x$. Hence J(x,0) = 1. Since J(x,u) is continuous on the compact set B_1 , there exists a constant $0 < \delta^{\mathsf{0}}_{\mathsf{p}} - \frac{\delta_{\mathsf{p}}}{2}$ such that J(x,u) > 0 for all $(x,u) - B_2 := \{(x,u) : \phi^{\mathsf{p}}(x) - \overline{W^{\mathsf{0}}_{\mathsf{p}}}, u - \delta^{\mathsf{0}}_{\mathsf{p}}\}$.

Therefore, the function $\frac{\overline{\det g_{ij}^{\mathbf{p}}}(x)}{\overline{\det g_{ij}^{\mathbf{p}}} \ (\phi^{\mathbf{p}^{-1}} \ \widetilde{h}(x,u \ J(x,u))}$ is well defined and continuous on the compact set B_2 . So the constant

$$C_{\mathrm{p}}^{\mathrm{0}} := \max_{(x,u \ \mathsf{2}\,B_2)} \mathsf{q} \frac{\frac{\mathsf{q} \ }{\det g_{ij}^{\mathrm{p}}}(x)}{\det g_{ij}^{\mathrm{p}} \quad (\phi^{\mathrm{p}})^{-1} \quad \Re(x,u)J(x,u)}$$

is positive and finite. It follows that (3.16) can be bounded as

by a change of variables $x - \phi(u)$. This proves the inequality (3.14).

As to the second inequality (3.15), denote $\mathcal{F}(u) := f^{-}\phi^{\mathsf{p}}(u)$ and $g(\mathbf{q}, u) := (\phi^{\mathsf{p}})^{-1} - \phi^{-}(u)$, then

$$f(\phi(u)) = f(\phi^p)^{-1}(\phi^p)^{-1}(u) = \mathcal{F}(\mathbf{q}, u).$$

Let $y = g(\mathbf{q}, u)$, by the chain rule,

$$D(f \ \phi)(u)_{i} = \frac{\mathbf{X}^{d}}{k=1} \frac{\partial \mathbf{P}(y)}{\partial y^{k}} \cdot \frac{\partial g_{k}(\mathbf{q}, u)}{\partial u^{i}}$$

and

$$D^{2}(f \quad \phi \)(u)_{ij} = \frac{\mathsf{X}^{d}}{\partial y^{k} \partial y^{\ell}} \cdot \frac{\partial^{2} \mathbf{f}(y)}{\partial u^{i}} \cdot \frac{\partial g_{k}(\mathbf{q}, u)}{\partial u^{i}} \frac{\partial g_{\ell}(\mathbf{q}, u)}{\partial u^{j}} + \frac{\mathsf{X}^{d}}{k=1} \frac{\partial \mathbf{f}(y)}{\partial y^{k}} \frac{\partial^{2} g_{k}(\mathbf{q}, u)}{\partial u^{i} \partial u^{j}}.$$

Since the function $g(\mathbf{q}, u)$ is C^1 on $\{(\mathbf{q}, u) : \mathbf{q} \quad \overline{W_p^0}, u \quad \delta_p \}$, we have

$$C_1 := \sup_{\mathbf{2}\overline{W_{\mathbf{p}}'}, \mathsf{kuk} \ \delta_{\mathbf{p}}, 1 \ k, \ell, i, j \ d} \frac{\partial g_k(\mathbf{q}, u)}{\partial u^i} \cdot \frac{\partial g_\ell(\mathbf{q}, u)}{\partial u^j} < ,$$

$$C_2 := \sup_{2\overline{W_{\mathbf{p}}'}, \mathsf{k}u\mathsf{k} \ \delta_{\mathbf{p}}, \mathsf{1} \ k, i, j \ d} \frac{\partial^2 g_k(\mathbf{q}, u)}{\partial u^i \partial u^j} < .$$

Hence

$$D^{2}(f \quad \phi \)(u)_{ij} \qquad C_{1} \underset{k,\ell=1}{\overset{\mathsf{X}^{d}}{\longrightarrow}} \frac{\partial^{2} \mathcal{F}(y)}{\partial y^{k} \partial y^{\ell}} + C_{2} \underset{k=1}{\overset{\mathsf{X}^{d}}{\longrightarrow}} \frac{\partial \mathcal{F}(y)}{\partial y^{k}} \ . \tag{3.18}$$

Applying the local representation of f and 2f under the **p**-normal coordinate, we know that

$$\frac{\partial^2 \mathbf{P}(y)}{\partial y^k \partial y^\ell} = {}^2 f \phi (u) + \sum_{k\ell=1}^{\mathbf{X}^d} \Gamma_{k\ell}^m g(\mathbf{q}, u) \qquad f \phi (u) \quad , \quad \frac{\partial \mathbf{P}(y)}{\partial y^k} = \int \phi (u) \cdot \int_{k\ell} \phi$$

Denote $C_3 := \sup_{\mathbf{2}\overline{W'_{\mathbf{p}}}, \mathbf{k}u\mathbf{k}} \delta_{\mathbf{p},1} k_{\ell,m} d \Gamma^m_{\ell\ell} g(\mathbf{q}, u)$, we can bound (3.18) further as

$$D^{2}(f \phi)(u)_{ij} C_{1} \sum_{k,\ell=1}^{\mathbf{X}^{d}} f \phi(u) + (C_{2} + d^{2}C_{1}C_{3})_{k=1}^{\mathbf{X}^{d}} f \phi(u)$$

This together with the elementary inequality $(P_{i=1}^{m}/a_{i}/)^{p}$ $m^{p}P_{i=1}^{m}/a_{i}/p$ implies that

$$D^{2}(f \phi)(u)^{p} 3^{p}d^{p+2}(C_{1}+C_{2}+d^{2}C_{1}C_{3})^{p} 2^{p}f \phi(u)^{p} + f \phi(u)^{p}. (3.19)$$

Integrating over W_p^0 with respect to \mathbf{q} and using the inequality (3.14), we get the desired result.

Proposition 3 yields the following lemma that will be frequently used in proving Theorem 2.

Lemma 3. Let $W_{\mathsf{p}}^{\mathsf{0}} = W_{\mathsf{p}}$ E_{p} $B_{\delta_{\mathsf{p}}/2}(0)$ and $\hat{\mathcal{B}}_{\sigma}$ be given by (3.7). Let $\delta_{\mathsf{p}}^{\mathsf{0}}$ and $C_{\mathsf{p}}^{\mathsf{0}}$ be the constants in Proposition 3. Set $\delta = \min_{\mathsf{p} \, \mathsf{2P}} \min \{\delta_{\mathsf{p}}^{\mathsf{0}}, \frac{1}{2C_{\mathsf{p}}}\}$ and σ_{0} (0,1] satisfying C_{0} $\overline{2d + 6}\sigma_{\mathsf{0}}^{\mathsf{p}}$ $\overline{\log \sigma_{\mathsf{0}}^{\mathsf{1}}} < \delta$. Assume h is a measurable function on \mathbb{R}^d and f $L^p(X)$ with p 1. Then

Proof. By the definition of \mathcal{B}_{σ} , we know that for $u = \mathcal{B}_{\sigma}$, $u < \delta = \delta_{\rho}^{0}$. Then the case p = 1 follows from (3.14) and a change of order of integrals.

When p > 1, let $q = \frac{p}{p-1}$ and write $h(u) = h(u)^{\frac{1}{p} + \frac{1}{q}}$. Then using the Hölder inequality and (3.14), we have

This proves (3.20) in the case p > 1.

3.4 Uniform boundedness of linear operators

We give the uniform boundedness of the operator $I_{\sigma}: L^{p}(X) = L^{p}(X)$ defined by (2.4). It will be used not only for the proof of Theorem 2, but also for deriving Theorem 1.

Proposition 4. Let $I_{\sigma}: L^{p}(X)$ $L^{p}(X)$ be defined by (2.4). Then

$$I_{\sigma}(f) \ _{L^{p}(X)} C_{X}^{0} f \ _{L^{p}(X)} \sigma > 0, f \ L^{p}(X),$$
 (3.21)

where C_X^0 is a constant independent of σ or f.

Proof. Let $W_p^0, \delta_p^0, \delta$ and σ_0 be given in Lemma 3 and P be a finite subset of X such that

$$X \qquad _{\mathsf{p2P}} W_{\mathsf{p}}^{\mathsf{0}}. \text{ For } \sigma \quad \sigma_{\mathsf{0}},$$

$$I_{\sigma}(f) \ _{L^{p}(X)} \ = \ \frac{1}{(\overline{2\pi}\sigma)^{d}} \sum_{\substack{X \ Z \ Z \ \\ \overline{(2\pi}\sigma)^{d} \ \\ }} \exp \ - \frac{x-y^{2}}{2\sigma^{2}} \ f(y)dV(y)^{\frac{p}{d}}dV(x)^{\frac{1}{p}}$$

$$= \frac{1}{(\overline{2\pi}\sigma)^{d}} \sum_{\substack{X \ X \ X \ X \ \\ \hline{(Vol(X))^{\frac{1}{p}}} \ Z \ \\ \hline{(\frac{Vol(X))^{\frac{1}{p}}}{(\overline{2\pi}\sigma)^{d}}} \sum_{\substack{X \ X \ X \ X \ \\ \hline{(f(y)/dV(y)} \ \frac{Vol(X)}{(\overline{2\pi}\sigma_{0})^{d}} \ f_{L^{p}(X)}.$$

For $0 < \sigma < \sigma_0$, we know from the inequality

$$I_{\sigma}(f)$$
 $L^{p}(X)$ X $I_{\sigma}(f)(\mathbf{q})/p dV(\mathbf{q})$

that it is sufficient to prove for each p - P,

$$Z = \int_{W_{\mathbf{p}}^{\prime}} |I_{\sigma}(f)(\mathbf{q})|^{p} dV(\mathbf{q}) \qquad C_{X,p}^{0} \quad f_{L^{p}(X)} \qquad 0 < \sigma < \sigma_{0}.$$

$$(3.22)$$

In fact we get (3.21) by setting $C_X^0 = \max \frac{\operatorname{Vol}(X)}{\sqrt[n]{2\pi}\sigma_0}, P_{p2P} C_{X,p}^0$

Now we prove (3.22) for each **p** P. Let $0 < \sigma < \sigma_0$.

For any $\mathbf{q} = W_{\mathfrak{p}}^{0}$, choose B_{σ} , \mathcal{B}_{σ} as in §3.2. Decompose the domain X into two parts B_{σ} and $X \setminus B_{\sigma}$. We have from (3.10)

$$I_{\sigma}(f)(\mathbf{q}) = \frac{1}{(2\pi\sigma)^{d}} \sum_{\tilde{P}_{\mathbf{Z}}}^{\mathbf{Z}} \exp \left[-\frac{\mathbf{q} - \phi(u)^{2}}{2\sigma^{2}} f(\phi(u))\right]^{\mathbf{q}} \frac{\mathbf{q}}{\det(g_{ij})}(u)du$$

$$+ \frac{1}{(2\pi\sigma)^{d}} \sum_{X \cap B_{\sigma}^{\mathbf{q}}}^{\mathbf{q}} K_{\sigma}(\mathbf{q}, y)f(y)dV(y)$$

$$:= J_{1}(\mathbf{q}) + J_{2}(\mathbf{q}).$$

For the first term $J_1(\mathbf{q})$, we see from (3.8) and (3.9) that

$$J_{1}(\mathbf{q}) = \frac{3}{2(\overline{2\pi}\sigma)^{d}} \sum_{\tilde{B}_{\sigma}}^{\mathbf{Z}} \exp \left(-\frac{u^{2}}{4\sigma^{2}}\right) / f(\phi(u)) / du.$$
(3.23)

By inequality (3.20) in Lemma 3 with $h(u)=\frac{3}{2(\sqrt[p]{2\pi}\sigma^d}\exp^--\frac{kuk^2}{4\sigma^2}$, we have

By a change of variables $\frac{u}{\sigma}$ and the equation (3.11), we have

$$Z \qquad Z \qquad Z \qquad Z \qquad \qquad Z \qquad \qquad \lambda du = 3 \cdot 2^{\frac{d}{2}} \quad \Delta du = 3 \cdot 2^{\frac{d}{2} - 1}.$$

Therefore,

$$Z = \int_{W_{\mathbf{p}}'} J_{1}(\mathbf{q}) p^{p} dV(\mathbf{q}) = 3 \cdot 2^{\frac{d}{2} - 1} (C_{\mathbf{p}}^{\mathbf{0}})^{\frac{1}{p}} f_{L^{p}(X)}.$$
 (3.24)

As for the second term $J_2(\mathbf{q})$, we notice that for $y = X \setminus B_{\sigma}$, the restriction $d_X(\mathbf{q}, y)$ $C_0 = \frac{1}{2d+6\sigma}$ yields $\mathbf{q} - y = \frac{1}{2d+6\sigma}$ Thus

$$|J_{2}(\mathbf{q})| = \frac{1}{(\overline{2\pi}\sigma)^{d}} \mathbf{Z}_{X \cap B_{\sigma}^{\mathbf{q}}} K_{\sigma}(\mathbf{q}, y) f(y) dV(y)$$

$$= \frac{1}{(\overline{2\pi}\sigma)^{d}} \mathbf{Z}_{X \cap B_{\sigma}^{\mathbf{q}}} \exp -\frac{(2d+6)\sigma^{2} \log \frac{1}{\sigma}}{2\sigma^{2}} |f(y)| dV(y)$$

$$= (2\pi)^{d/2} \sigma^{3} /f(y) /dV(y).$$

Hence

Combining this with the inequality (3.24), we get the desired bound (3.22).

3.5 Proof of Theorem 2

Now we can prove Theorem 2.

Proof of Theorem 2: Let $W_{\mathsf{p}}^0, \delta_{\mathsf{p}}^0, \delta$ and σ_0 be given in Lemma 3 and P be a finite subset of X such that $X = \Pr(X) = \Pr$

$$Z = \int_{W_{\mathbf{p}}'} |I_{\sigma}(f)(x) - f(x)|^{p} dV(x) \qquad C_{X,\mathbf{p}} f_{H_{2}^{p}(X)} \sigma^{2} \qquad 0 < \sigma < \sigma_{0}.$$
 (3.26)

We prove (3.26) in three steps. Let $0 < \sigma < \sigma_0$.

Step 1: Decomposition. Let $\mathbf{q} = W_{\mathsf{p}}^{\mathsf{0}}$. Choose B_{σ} , \mathcal{B}_{σ} as in §3.2. By the identity $\frac{1}{(\mathsf{p} \frac{1}{2\pi\sigma^{-d}} - \mathsf{R}^d)} \exp\{-\frac{\mathsf{k} u \mathsf{k}^2}{2\sigma^2}\} du = 1$, we can decompose $f(\mathbf{q})$ as

$$f(\mathbf{q}) = \frac{1}{(\overline{2\pi}\sigma)^d} \sum_{\widetilde{B}_{\sigma}}^{\mathbf{Z}} \exp \left(-\frac{u^2}{2\sigma^2} f(\mathbf{q})du + \frac{1}{(\overline{2\pi}\sigma)^d} \sum_{\mathbb{R}^d \mathsf{n}\widetilde{B}_{\sigma}}^{\mathbf{Z}} \exp \left(-\frac{u^2}{2\sigma^2} f(\mathbf{q})du\right)\right)$$

Separate the integral on X for $I_{\sigma}(f)(\mathbf{q})$ to two parts on B_{σ} and $X \setminus B_{\sigma}$, we have

$$I_{\sigma}(f)(\mathbf{q}) - f(\mathbf{q}) = J_1(\mathbf{q}) + J_2(\mathbf{q}) \tag{3.27}$$

where

$$J_{1}(\mathbf{q}) = \frac{1}{(\overline{2\pi}\sigma)^{d}} \sum_{\tilde{B}_{\sigma}}^{\mathbf{Z}} \exp \left[-\frac{\mathbf{q} - \phi(u)^{2}}{2\sigma^{2}} f(\phi(u))^{\mathbf{q}} \frac{\mathbf{d}}{\det(g_{ij})}(u)\right]$$

$$-\exp \left[-\frac{u^{2}}{2\sigma^{2}} f(\mathbf{q}) du,\right]$$

$$J_{2}(\mathbf{q}) = \frac{1}{(\overline{2\pi}\sigma)^{d}} \sum_{X \cap B_{\sigma}^{\mathbf{q}}}^{\mathbf{q}} K_{\sigma}(\mathbf{q}, y) f(y) dV(y) - \frac{1}{(\overline{2\pi}\sigma)^{d}} \sum_{\mathbb{R}^{d} \cap \tilde{B}_{\sigma}}^{\mathbf{q}} \exp \left[-\frac{u^{2}}{2\sigma^{2}} f(\mathbf{q}) du.\right]$$

Step 2: Estimation of $R_{W_{\mathbf{p}}'}/J_1(\mathbf{q})/p^dV(\mathbf{q})^{-\frac{1}{p}}$. We separate the error further as

$$J_{1}(\mathbf{q}) = \frac{1}{(2\pi\sigma)^{d}} \sum_{\tilde{\mathbf{Z}}}^{\tilde{\mathbf{P}}} \exp \left[-\frac{\mathbf{q} - \phi(u))^{2}}{2\sigma^{2}} - \exp \left[-\frac{u^{2}}{2\sigma^{2}}\right] f(\phi(u))du \right]$$

$$+ \frac{1}{(2\pi\sigma)^{d}} \sum_{\tilde{\mathbf{Z}}}^{\tilde{\mathbf{B}}_{\sigma}} \exp \left[-\frac{\mathbf{q} - \phi(u)^{2}}{2\sigma^{2}}\right] f(\phi(u)) + \frac{1}{\det g_{ij}}(u) - 1 du$$

$$+ \frac{1}{(2\pi\sigma)^{d}} \sum_{\tilde{\mathbf{B}}_{\sigma}}^{\tilde{\mathbf{B}}_{\sigma}} \exp \left[-\frac{u^{2}}{2\sigma^{2}}\right] f(\phi(u)) - f(\mathbf{q}) du$$

$$:= J_{11}(\mathbf{q}) + J_{12}(\mathbf{q}) + J_{13}(\mathbf{q}).$$

For $J_{11}(\mathbf{q})$, we use (3.4) and the elementary inequality $|e^{-a} - e^{-b}| - |a - b| \max\{e^{-a}, e^{-b}\}$ (valid for any a, b > 0) and find that

$$\frac{1}{(2\pi\sigma)^{d}} Z \exp -\frac{\mathbf{q} - \phi (u)^{2}}{2\sigma^{2}} - \exp -\frac{u^{2}}{2\sigma^{2}} /f(\phi (u))/du
\frac{1}{(2\pi\sigma)^{d}} Z^{\widetilde{B}_{\sigma}} \max \exp -\frac{\mathbf{q} - \phi (u)^{2}}{2\sigma^{2}} , \exp -\frac{u^{2}}{2\sigma^{2}} \frac{f(\phi (u))}{2\sigma^{2}} /f(\phi (u))/du.$$

So by (3.8),

$$/J_{11}(\mathbf{q})/$$
 $\frac{\mathsf{Z}}{\tilde{B}_{\sigma}} \frac{1}{(\overline{2\pi}\sigma)^d} \exp -\frac{u^2}{4\sigma^2} \frac{C_{\mathsf{p}} u^4}{2\sigma^2} /f(\phi(u))/du.$

It follows from (3.20) in Lemma 3 with $h(u) := \frac{1}{(P \overline{2\pi}\sigma^d)} \exp -\frac{kuk^2}{4\sigma^2} \frac{C_{\mathbf{p}}kuk^4}{2\sigma^2}$ that

By a change of variables $\frac{u}{\sigma}$ and the equation (3.11), we have

$$\mathsf{Z} \underset{\widetilde{B}_{\sigma}}{h(u)du} \quad \frac{C_{\mathsf{p}}}{2}(2\pi)^{-\frac{d}{2}}\sigma^2 \underset{\mathbb{R}^d}{\mathsf{exp}} \quad -\frac{u^{-2}}{4} \quad u^{-4}du = \frac{2^{\frac{d}{2}+3}C_{\mathsf{p}}\,\Gamma(\frac{d+4}{2})}{\Gamma(\frac{d}{2})}\sigma^2.$$

Hence

$$Z \atop W_{\mathbf{p}} / J_{11}(\mathbf{q}) / dV(\mathbf{q}) \stackrel{\frac{1}{p}}{=} \frac{2^{\frac{d}{2} + 3} C_{\mathbf{p}} (C_{\mathbf{p}}^{\mathbf{0}})^{\frac{1}{p}} \Gamma(\frac{d+4}{2})}{\Gamma(\frac{d}{2})} f_{L^{p}(X)} \sigma^{2}.$$
(3.28)

For $J_{12}(\mathbf{q})$, we use (3.3) and (3.8) and obtain

$$/J_{12}(\mathbf{q})/\frac{C_{\mathsf{p}}}{(\overline{2\pi}\sigma)^d} \sum_{\widetilde{B}_{\sigma}}^{\mathsf{Z}} \exp -\frac{u^2}{4\sigma^2} u^2/f(\phi(u))/du.$$

Thus using (3.20) in Lemma 3 with $h(u) = \frac{C_p}{\sqrt{2\pi}\sigma^d} \exp \left(-\frac{kuk^2}{4\sigma^2}\right) u^2$, we obtain

$$\frac{Z}{J_{12}(\mathbf{q})} \int dV(\mathbf{q}) \frac{1}{p} \frac{2^{\frac{d}{2}+2} C_{p} (C_{p}^{0})^{\frac{1}{p}} \Gamma(\frac{d+2}{2})}{\Gamma(\frac{d}{2})} f_{L^{p}(X)} \sigma^{2}.$$
(3.29)

For the last term $J_{13}(\mathbf{q})$ of $J_1(\mathbf{q})$, we apply the Taylor expansion (3.13) and get the following further decomposition:

$$J_{13}(\mathbf{q}) = \frac{1}{(\overline{2\pi}\sigma)^{d}} \sum_{\widetilde{B}_{\sigma}}^{\mathbf{Z}} \exp \left[-\frac{u^{2}}{2\sigma^{2}} + f(\mathbf{q}), \sum_{i=1}^{d} e_{i} u_{i} du \right] + \frac{1}{(\overline{2\pi}\sigma)^{d}} \exp \left[-\frac{u^{2}}{2\sigma^{2}} + \frac{Z_{1}}{(\overline{2\sigma^{2}})^{d}} + \frac{U_{13}(\mathbf{q})}{(1-y)D^{2}(f(\phi(yu)))(u,u)dydu} \right] = J_{13}^{0}(\mathbf{q}) + J_{13}^{00}(\mathbf{q}).$$
(3.30)

Since

$$\frac{1}{(\overline{2\pi}\sigma)^d} \mathop{\rm R}^d \exp -\frac{u^2}{2\sigma^2} \qquad f(\mathbf{q}), \quad f(\mathbf{q}) = 0,$$

we have

$$/J_{13}^{\mathbf{0}}(\mathbf{q})/\frac{1}{(\overline{2\pi}\sigma)^d} \underset{\mathbb{R}^d \cap \tilde{B}_{\sigma}}{\overset{\mathbf{Z}}{=}} \exp -\frac{u^2}{2\sigma^2} / f(\mathbf{q})/u \ du.$$
 (3.31)

Let $q = \frac{p}{p-1}$ when p > 1 and denote $h_{\sigma}(u) := \frac{1}{(P \frac{1}{2\pi\sigma} d)} \exp(-\frac{kuk^2}{2\sigma^2} u)$. Then using the Hölder inequality, we get

$$Z = \frac{|J_{13}^{0}(\mathbf{q})|^{p}dV(\mathbf{q})}{|Z|^{N_{p}^{\prime}}} Z = Z \qquad \frac{p}{q} Z \qquad \frac{1}{p} Z \qquad \frac{h_{\sigma}(u)du}{|R^{d}n\tilde{B}_{\sigma}|} h_{\sigma}(u)/|f(\mathbf{q})|^{p}dudV(\mathbf{q}) \qquad \frac{1}{p} Z \qquad \frac{h_{\sigma}(u)du}{|R^{d}n\tilde{B}_{\sigma}|} exp - \frac{u^{2}}{2\sigma^{2}} = u du$$

$$= \int_{H_{2}^{p}(X)} \frac{2^{1-d}\sigma}{\Gamma(\frac{d}{2})} \frac{1}{r^{2}} c_{0}^{p} \frac{exp}{2d+6} = \frac{r^{2}}{\log \frac{1}{\sigma}} exp - \frac{r^{2}}{2} r^{d}dr$$

$$= \int_{H_{2}^{p}(X)} \frac{2^{1-d}\sigma}{\Gamma(\frac{d}{2})} \frac{1}{r^{2}} c_{0}^{p} \frac{exp}{2d+6} = \frac{-r^{2}}{\log \frac{1}{\sigma}} exp - \frac{C_{0}^{2}(2d+6)(\log \sigma^{-1})}{4} = exp - \frac{r^{2}}{4} r^{d}dr$$

$$= \frac{2^{1-d}}{\Gamma(\frac{d}{2})} \frac{1}{r^{2}} \frac{exp}{r^{2}} \frac{1}{r^{2}} exp - \frac{r^{2}}{r^{2}} r^{d}dr = \frac{2\Gamma(\frac{d+1}{2})}{\Gamma(\frac{d}{2})} f_{H_{2}^{p}(X)} \sigma^{2}. \qquad (3.32)$$

The proof for (3.32) in the case p = 1 follows directly from a change of order of integrals.

As for the term $J_{13}^{00}(\mathbf{q})$, it is easy to see that

$$/J_{13}^{00}(\mathbf{q})/$$
 $\frac{1}{(\overline{2\pi}\sigma)^d}\sum_{\widetilde{B}_{\sigma}}^{\mathbf{Z}}\exp -\frac{u^2}{2\sigma^2}\sum_{0}^{2}(1-y)D^2(f-\phi)(yu)u^2dydu.$

For $u = \mathcal{B}_{\sigma}$, by (3.15) and the Hölder inequality when p > 1, we get

Denote $\Re_{\sigma}(u) := \frac{1}{(-2\pi\sigma)^d} \exp(-\frac{kuk^2}{2\sigma^2}) u^2$. Then using the Hölder inequality, we know that $\frac{1}{2} (C_{\mathsf{p}}^{\mathsf{oq}})^{\frac{1}{p}} f_{H_{2}^{p}(X)} - \frac{1}{\tilde{B}_{\sigma}} \frac{1}{(\overline{2\pi}\sigma)^{d}} \exp - \frac{u^{2}}{2\sigma^{2}} u^{2} du$ $(C_{\mathsf{p}}^{00})^{\frac{1}{p}} \frac{\Gamma(\frac{d+2}{2})}{\Gamma(\frac{d}{2})} f_{H_2^p(X)} \sigma^2.$

This together with (3.32) yields

$$Z = \int_{W_{\mathbf{p}}^{\prime}} J_{13}(\mathbf{q}) / p dV(\mathbf{q})^{\frac{1}{p}} = ((C_{\mathbf{p}}^{00})^{\frac{1}{p}} + 2) \frac{\Gamma(\frac{d+2}{2})}{\Gamma(\frac{d}{2})} f_{H_{2}^{p}(X)} \sigma^{2}.$$

Combining the estimates for $J_{11}(\mathbf{q}), J_{12}(\mathbf{q})$ and $J_{13}(\mathbf{q})$, we have

$$(Z_{W_{\mathbf{p}}^{\prime}})^{-1} / J_{1}(\mathbf{q}) / dV(\mathbf{q})$$

$$2^{\frac{d}{2}+4} C_{\mathbf{p}} (C_{\mathbf{p}}^{\mathbf{0}})^{\frac{1}{p}} + (C_{\mathbf{p}}^{\mathbf{0}})^{\frac{1}{p}} + 2 \frac{\Gamma(\frac{d+4}{2})}{\Gamma(\frac{d}{2})} f_{H_{2}^{p}(X)} \sigma^{2}.$$

$$(3.33)$$

Step 3: Estimation of $R_{W_{\mathbf{p}}'}/J_2(\mathbf{q})/p^dV(\mathbf{q})^{-\frac{1}{p}}$. Denote the two terms in the expression of $J_2(\mathbf{q})$ as $J_2^0(\mathbf{q})$ and $J_2^{00}(\mathbf{q})$. The first term $J_2^0(\mathbf{q})$ of $J_2(\mathbf{q})$ has been estimated in the proof of Proposition 4 as (3.25). Hence

$$Z = \int_{W_{\mathbf{p}}'} J_2^{\mathbf{0}}(\mathbf{q}) / p dV(\mathbf{q}) = (2\pi)^{-d/2} \operatorname{Vol}(X) f_{L^p(X)} \sigma^2.$$

$$(3.34)$$

Now we bound the second term of
$$J_2(\mathbf{q})$$
. Using (3.11) again, we have
$$\begin{split} |J_2^{00}(\mathbf{q})| &= \frac{1}{(\frac{1}{2\pi}\sigma)^d} \sum_{\mathbf{R}^d \cap \tilde{B}_\sigma} \exp{-\frac{u^2}{2\sigma^2}} f(\mathbf{q}) du \\ &= \frac{|f(\mathbf{q})|}{(\frac{1}{2\pi}\sigma)^d} \sup_{\mathbf{k}u\mathbf{k}} \frac{|f(\mathbf{q})|}{|\mathbf{q}|} \exp{-\frac{u^2}{2\sigma^2}} du \\ &= \frac{2^{1-\frac{d}{2}}}{\Gamma(\frac{d}{2})} |f(\mathbf{q})| \exp{-\frac{r^2}{2d+6}(\log \sigma^{-1-1/2})} \exp{-\frac{r^2}{2}} r^{d-1} dr \\ &= \frac{2^{1-\frac{d}{2}}}{\Gamma(\frac{d}{2})} |f(\mathbf{q})| \exp{-\frac{r^2}{2d+6}(\log \sigma^{-1-1/2})} \exp{-\frac{C_0^2(2d+6)(\log \sigma^{-1})}{4}} \exp{-\frac{r^2}{4}} r^{d-1} dr \\ &= \frac{2^{1-\frac{d}{2}}}{\Gamma(\frac{d}{2})} |f(\mathbf{q})| exp - \frac{r^2}{2d+6}(\log \sigma^{-1-1/2}) \exp{-\frac{r^2}{4}} r^{d-1} dr = 2^{\frac{d}{2}} |f(\mathbf{q})| exp - \frac{r^2}{4} r^{d-1} dr = 2^{\frac{d}{2}} |f(\mathbf{q})| exp - \frac$$

But C_0 1 and d 1, so there holds

$${\sf Z} \underset{W_{\bf p}'}{/J_2^{\bf 00}({\bf q})/\!dV({\bf q})} \ ^{\frac{1}{p}} \ \ 2^{\frac{d}{2}}\sigma^2 \ \ \underset{W_{\bf p}'}{/J({\bf q})/\!\!/^{\!\!p}} dV({\bf q}) \ \ ^{\frac{1}{p}} \ \ 2^{\frac{d}{2}} \ f_{L^p(X)}\sigma^2.$$

Combining this with (3.34), we get

Z
$$/J_2(\mathbf{q})/dV(\mathbf{q}) \stackrel{\frac{1}{p}}{\qquad} (2\pi)^{\frac{d}{2}} \operatorname{Vol}(X) + 2^{\frac{d}{2}} \quad f_{L^p(X)} \sigma^2.$$

This together with (3.33) yields the desired estimate (3.26).

4 Learning Rates

In this section, we derive learning rates for the multi-kernel classification algorithm (1.9) in the manifold setting, especially for the case of SVM. This is done by balancing the sample error and the regularization error [SZ2, WYZ, WZ].

We need the following result from [26] where we have changed some notation in order to make it consistent with this paper.

The regularization error of the algorithm (1.9) is defined as

$$D(\lambda) = \min_{\sigma \in (0,1]} \min_{f \ni H_{K_{\sigma}}} E_{\phi}(f) - E_{\phi}(f_{\rho}^{\phi}) + \lambda f_{K_{\sigma}}^{2}, \quad \lambda > 0.$$
 (4.1)

Denote

$$C_{\lambda} = \sup \max\{/\phi^{0}(t)/, /\phi^{0}_{+}(t)/\} : /t/$$
 $r = \frac{|\phi(0)/|}{\lambda}$ (4.2)

Proposition 5. Let X be a subset of \mathbb{R}^n and ϕ be admissible with $C_{\lambda} <$. Define $f_{\mathsf{z},\lambda}$ by (1.9). Then we have

$$\mathbb{E} E_{\phi}(f_{\mathsf{z},\lambda}) - E_{\phi}(f_{\rho}^{\phi}) \qquad C_{\lambda} \frac{\overline{C^{0}/\phi(0)/}}{\lambda} \frac{\log^{2} m}{m} + \frac{2/\phi(0)/\overline{m}}{\overline{m}} + D(\lambda), \tag{4.3}$$

where C^0 is a constant independent of m or λ .

For the SVM case, we have a simplified version of Proposition 5.

Proposition 6. Let $\phi(x) = (1 - x)_+ = \max\{0, 1 - x\}$ and $X = \mathbb{R}^n$. Define $f_{z,\lambda}$ by (1.9). If $0 < \lambda < 1$, then there exists a constant C^0 independent of m or λ such that

$$\mathbb{E} E_{\phi}(f_{\mathsf{z},\lambda}) - E_{\phi}(f_c) \qquad \frac{\Gamma}{\lambda} \frac{\overline{C^0}}{m} \stackrel{1/4}{-} + D(\lambda). \tag{4.4}$$

Thus if we can estimate the regularization error $D(\lambda)$ in the manifold setting, the excess misclassification error can be easily derived using the comparison theorem between the excess misclassification error and excess generalization error.

4.1 Learning rates for the multi-Gaussian SVM classifier

We begin with the special case of multi-Gaussian SVM classifier. Since $f_{\rho}^{\phi} = f_c$ and $\phi(x) - \phi(y) / (x - y)$ for the hinge loss, we can bound the regularization error $D(\lambda)$ as

$$D(\lambda) \quad \inf_{\sigma 2 \, (0,1)} \inf_{f \, 2H} \inf_{K_\sigma} \quad f - f_c \, _{L^1_{\rho_X}(X)} + \lambda \, f \, _{K_\sigma}^2 \ .$$

Thus we only need to bound $\inf_{\sigma^2(0,1)} \inf_{f^2 H_{K_{\sigma}}} f - f_c L^1_{\rho_X}(X + \lambda f)^2_{K_{\sigma}}$ for the SVM case.

We shall choose $f = I_{\sigma}(f_c)$ to bound $D(\lambda)$. So we need to estimate $I_{\sigma}(f_c)$ K_{σ} .

Lemma 4. Let $1 p < and f L^p(X)$. Then the function $I_{\sigma}(f)$ is in $\mathcal{H}_{K_{\sigma}}(X)$ and

$$I_{\sigma}(f) K_{\sigma} = \begin{cases} 8 \\ \frac{1}{(\frac{p-1}{2\pi}d)} (\mathfrak{S}_{X})^{\frac{p-1}{p}} f_{L^{p}(X)} \sigma^{\frac{d}{p}}, & \text{if } 1 = p = 2, \\ \vdots & \frac{1}{(\frac{p-1}{2\pi}d)} (\mathfrak{S}_{X})^{\frac{1}{2}} (Vol(X))^{\frac{1}{2}-\frac{1}{p}} f_{L^{p}(X)} \sigma^{\frac{d}{2}}, & \text{if } p > 2 \end{cases}$$

$$= \frac{1}{(\frac{1}{2\pi})^{d}} (\mathfrak{S}_{X})^{1-\frac{1}{\min\{p,2\}}} (Vol(X))^{\frac{1}{2}-\frac{1}{\max\{p,2\}}} f_{L^{p}(X)} \sigma^{\frac{d}{\min\{p,2\}}}$$
(4.5)

where \mathfrak{S}_X is the constant given by Lemma 2.

Proof. By the definition of $I_{\sigma}(f)$ and the equation

$$K_{\sigma}(\cdot, y), K_{\sigma}(\cdot, z) |_{K_{\sigma}} = K_{\sigma}(y, z),$$
 (4.6)

we have

$$I_{\sigma}(f) \ _{K_{\sigma}}^{2} = \frac{1}{(-\overline{2\pi}\sigma)^{2d}} \underset{X=X}{\mathsf{Z}} K_{\sigma}(y,z) f(y) f(z) dV(y) dV(z).$$

When p = 1, we get from $K_{\sigma}(y,z) \mathbf{Z} \mathbf{Z} \mathbf{Z}$ 1 that

$$I_{\sigma}(f) \stackrel{2}{}_{K_{\sigma}} = \frac{1}{(\overline{2\pi}\sigma)^{2d}} \sum_{X=X} |f(y)|/|f(z)|dV(y)dV(z) = \frac{1}{(2\pi)^{d}} \sigma^{-2d} f \stackrel{2}{}_{L^{1}(X)},$$

which proves (4.5) in this case.

When $1 , we set <math>q = \frac{p}{p-1}$ and apply the Hölder inequality to the function

$$(K_{\sigma}(y,z))^{\frac{1}{q}}/f(y)/^{\frac{p}{q}} \cdot (K_{\sigma}(y,z))^{\frac{1}{p}}/f(y)/^{1-\frac{p}{q}}/f(z)/. \text{ Then}$$

$$Z Z$$

$$I_{\sigma}(f)^{2}_{K_{\sigma}} \frac{1}{(\overline{2\pi}\sigma)^{2d}} K_{\sigma}(y,z)/f(y)/^{p}dV(y)dV(z)$$

$$Z Z$$

$$K_{\sigma}(y,z)/f(y)/^{p(1-\frac{p}{q})}/f(z)/^{p}dV(y)dV(z)$$

$$K_{\sigma}(y,z)/f(y)/^{p}dV(y)dV(z)$$

Lemma 2 tells us that \nearrow $K_{\sigma}(y,z)dV(z)$ $\mathscr{E}_{X}\sigma^{d}$ for each y X. So

$$K_{\sigma}(y,z)/f(y)/p^{d}V(y)dV(z) \qquad \mathcal{E}_{X}\sigma^{d} \quad f_{L^{p}(X)}^{p}.$$
 On the other hand, for z X , we apply the Hölder inequality and find Z

It follows that

$$I_{\sigma}(f) \stackrel{2}{\underset{K_{\sigma}}{\stackrel{}{=}}} \qquad \frac{1}{(\overline{2\pi}\sigma)^{2d}} \mathscr{E}_{X}\sigma^{d} f \stackrel{p}{\underset{L^{p}(X)}{\stackrel{}{=}}} \qquad f \stackrel{\frac{1}{q}}{\underset{L^{p}(X)}{\stackrel{}{=}}} (\mathscr{E}_{X}\sigma^{d})^{\frac{p}{q}} \stackrel{\mathbf{0}}{\underset{p}{\stackrel{}{=}}}$$

$$= \frac{\mathscr{E}_{X}^{2} \stackrel{2}{\underset{p}{\stackrel{}{=}}}}{(2\pi)^{d}} f \stackrel{2}{\underset{L^{p}(X)}{\stackrel{}{=}}} \sigma \stackrel{\frac{2d}{p}}{\stackrel{}{=}}.$$

That is,

$$I_{\sigma}(f)_{K_{\sigma}} = \frac{\left(\mathbf{\mathfrak{G}}_{X}\right)^{\frac{p-1}{p}}}{\left(\overline{2\pi}\right)^{d}} f_{L^{p}(X)} \sigma^{\frac{d}{p}}.$$

When p > 2, each function in $L^p(X)$ lies in $L^2(X)$:

$$f_{L^{2}(X)} = \begin{cases} Z & \frac{1}{2} \\ /f(x) / {p} dV(x) & \\ Z & \frac{2}{p} & Z \\ /f(x) / {p} dV(x) & dV(x) \\ X & X & X \end{cases}$$

$$= (Vol(X))^{\frac{1}{2} - \frac{1}{p}} f_{L^{p}(X)}.$$

So the desired bound follows from the case p = 2:

$$I_{\sigma}(f)_{K_{\sigma}} = \frac{(\mathcal{C}_{X})^{\frac{1}{2}}}{(\overline{2\pi})^{d}} f_{L^{2}(X)} \sigma^{-\frac{d}{2}} = \frac{(\mathcal{C}_{X})^{\frac{1}{2}} (\operatorname{Vol}(X))^{\frac{1}{2} - \frac{1}{p}}}{(\overline{2\pi})^{d}} f_{L^{p}(X)} \sigma^{-\frac{d}{2}}.$$

This proves the desired inequality (4.5).

Proposition 7. Let X be a connected compact C^1 submanifold of \mathbb{R}^n without boundary which is isometrically embedded and of dimension d. If $f_c = (L^1(X), H^1_2(X))_{\theta}$ for some $0 < \theta = 1$, then

$$D(\lambda) \qquad (C_X^0 + 1 + C_X) f_c \theta + \frac{\mathfrak{E}_X (\operatorname{Vol}(X))^2}{(2\pi)^d} \lambda^{\frac{2\theta}{2\theta + d}}.$$

Proof. Let σ (0,). For any $g = H_2^1(X)$, we get from Theorem 2 that

$$I_{\sigma}(g) - g_{L^{1}(X)} C_{X} g_{H_{2}^{1}(X)} \sigma^{2}.$$

By Proposition 4, we also have

$$I_{\sigma}(f_c) - I_{\sigma}(g) _{L^1(X)} = I_{\sigma}(f_c - g) _{L^1(X)} C_X^{\mathbf{0}} g - f_c _{L^1(X)}.$$

Since $f_c = (L^1(X), H_2^1(X))_{\theta}$, we know that

$$I_{\sigma}(f_c) - f_{c \ L^{1}(X)} \qquad I_{\sigma}(f_c) - I_{\sigma}(g) \ _{L^{1}(X)} + \ I_{\sigma}(g) - g \ _{L^{1}(X)} + \ g - f_{c \ L^{1}(X)}$$

$$(C_X^{0} + 1) \ g - f_{c \ L^{1}(X)} + C_X \ g \ _{H^{1}_{\sigma}(X)} \sigma^{2}.$$

Taking infimum over $g = H_2^1(X)$, we get

$$I_{\sigma}(f_{c}) - f_{c} \ _{L^{1}(X)} \qquad (C_{X}^{0} + 1 + C_{X}) \inf_{g2 H_{2}^{1}(X)} g - f_{c} \ _{L^{1}(X)} + \sigma^{2} g \ _{H_{2}^{1}(X)}$$

$$= (C_{X}^{0} + 1 + C_{X}) \mathbb{K}(f_{c}, \sigma^{2}) \quad (C_{X}^{0} + 1 + C_{X}) f_{c} \ _{\theta}\sigma^{2\theta}.$$

Since f_c is the Bayes rule, $|f_c(x)| = 1$ for all x = X. Thus $|f_c|_{L^2(X)} = \frac{\mathsf{p}}{\mathsf{Vol}(X)}$ and the regularization error can be bounded as

$$D(\lambda) \qquad \inf_{\substack{\sigma^{2}(0,1) \\ \sigma^{2}(0,1)}} I_{\sigma}(f_{c}) - f_{c-L^{1}(X)} + \lambda I_{\sigma}(f_{c}) \frac{2}{K_{\sigma}}$$

$$\inf_{\substack{\sigma^{2}(0,1) \\ \sigma^{2}(0,1) \\ (C_{X}^{0} + 1 + C_{X})}} (C_{X}^{0} + 1 + C_{X}) f_{c-\theta} \sigma^{2\theta} + \lambda \cdot \frac{\mathfrak{S}_{X} Vol(X)}{(2\pi)^{d}} \sigma^{-d}$$

$$(C_{X}^{0} + 1 + C_{X}) f_{c-\theta} + \frac{\mathfrak{S}_{X} Vol(X)}{(2\pi)^{d}} \lambda^{\frac{2\theta}{2\theta + d}},$$

where we have chosen $\sigma = \lambda^{\frac{1}{2\theta+d}}$.

Proof of Theorem 1: An important relation between the excess misclassification error and the excess generalization error for the hinge loss asserts that [27] for any measurable function f: X \mathbb{R}

$$R \operatorname{sgn}(f) - R(f_c) \quad E_{\phi}(f) - E_{\phi}(f_c). \tag{4.8}$$

This together with Proposition 6 and Proposition 7 yields the desired results.

4.2Learning rates for general loss functions

We need the following relationship between $E_{\phi}(f) - E_{\phi}(f_{\rho}^{\phi})$ and $f - f_{\rho}^{\phi}|_{L^{p}(X)}$ which can be derived as Theorem 17 in [21].

Proposition 8. Let X be a subset of \mathbb{R}^n and ϕ be admissible.

(a) If ϕ is a Lipschitz s classification loss function on \mathbb{R} with Lipschitz constant C, then for any measurable function $f: X_{\mathbf{7}}$

(b) If ϕ is C^1 and its derivative is Lipschitz s on \mathbb{R} with Lipschitz constant C, then

$$E_{\phi}(f) - E_{\phi}(f_{\rho}^{\phi}) \quad C \quad f - f_{\rho}^{\phi} \quad {}^{1+s}_{L_{\rho_X}^{1+s}} \quad C \quad f - f_{\rho}^{\phi} \quad {}^{\frac{1+s}{2}}_{L_{\rho_X}^2}.$$

Motivated by this result, in the following, we assume that for any measurable function IR, the excess generalization error satisfies

$$E_{\phi}(f) - E_{\phi}(f_{\phi}^{\phi}) \qquad C \quad f - f_{\phi}^{\phi} \quad {}_{L_{P}(X)}^{\alpha}, \qquad (4.9)$$

where C, p \mathbb{R}_+ are constants independent of f. 1 and α

Recall C_{λ} defined by (4.2). As in [24], we assume that

$$C_{\lambda} = C_0 \lambda^{-\beta}$$
, for some $\beta = \mathbb{R}_+$, (4.10)

where C_0 is a constant.

Theorem 3. Let ϕ be an admissible loss function with $\phi^{00}(0) > 0$ satisfying (4.9) and (4.10). Define $f_{z,\lambda}$ by (1.9) and f_{ρ}^{ϕ} by (1.5). Assume f_{ρ}^{ϕ} $(L^{p}(X), H_{2}^{p}(X))_{\theta}$ for some $0 < \theta$ and p 1. Then by taking $\lambda = \frac{\log^{2} m}{m}$, we have

and
$$p-1$$
. Then by taking $\lambda = \frac{\log^2 m}{m}$, we have

$$\mathbb{E}_{\mathsf{Z2}\,Z^{m}} \ \ \mathsf{R}(sgn(f_{\mathsf{Z},\lambda}) - \mathsf{R}(f_{c}) = O \quad \frac{\log^{2} m}{m} \quad \frac{\frac{\theta\alpha\min\{p,2\}}{4(3+2\beta)\theta\alpha\min\{p,2\}+4(2\beta+1)d}}{m} \, . \tag{4.11}$$

Proof. As in the proof of Proposition 7, we have

$$I_{\sigma}(f_{\rho}^{\phi}) - f_{\rho}^{\phi}|_{L^{p}(X)} \qquad \{C_{X}^{0} + 1 + C_{X}\} f_{\rho}^{\phi}|_{\theta}\sigma^{2\theta},$$

It follows from inequalities (4.9) and (4.5) that

$$D(\lambda) \qquad \inf_{\substack{\sigma^{2}(0,1)\\ \sigma^{2}(0,1)}} C I_{\sigma}(f_{\rho}^{\phi}) - f_{\rho}^{\phi} {}_{L^{p}(X)}^{\alpha} + \lambda I_{\sigma}(f_{\rho}^{\phi}) {}_{K_{\sigma}}^{2} \\ \inf_{\substack{\sigma^{2}(0,1)\\ \sigma^{2}(0,1)}} C \{C_{X}^{0} + 1 + C_{X}\}^{\alpha} f_{\rho}^{\phi} {}_{\theta}^{\alpha} \sigma^{2\alpha\theta} + \lambda \cdot \mathcal{C}_{X}^{0} \sigma^{\frac{2d}{\min\{p,2\}}} {}^{0},$$

where $\mathcal{C}_{X}^{0} = \frac{1}{(2\pi^{-d})} (\mathcal{C}_{X})^{2} \frac{\frac{2}{\min\{p,2\}}}{\min\{p,2\}} (\operatorname{Vol}(X))^{1} \frac{2}{\max\{p,2\}} f_{\rho}^{\phi} f_{L^{p}(X)}^{2}$. Taking $\sigma = \lambda^{\frac{\min\{p,2\}}{2\alpha\theta \min\{p,2\}+2d}}$, we have

$$D(\lambda) \qquad C\{C_X^0 + 1 + C_X\}^{\alpha} f_{\rho}^{\phi} f_{\theta}^{\alpha} + \mathcal{C}_X^{0} \lambda^{\frac{\alpha\theta \min\{p,2\}}{\alpha\theta \min\{p,2\}+d}}.$$

$$(4.12)$$

Since ϕ is an admissible loss function with $\phi^{00}(0) > 0$, it is shown in [8] that there exists a constant C_{ϕ} depending only on ϕ such that for any measurable function f: X \mathbb{R} .

$$R(\operatorname{sgn}(f)) - R(f_c) \quad c_{\phi} \quad \overline{E_{\phi}(f) - E_{\phi}(f_{\rho}^{\phi})}. \tag{4.13}$$

Then the stated error bound follows from Proposition 5 and (4.10). This proves Theorem 3.

Now we apply Theorem 3 to q-norm hinge loss $\phi(x) = (1-x)_+^q$ with q > 1.

Corollary 1. Let $\phi(x) = (1-x)^q_+$ with q > 1. Define $f_{z,\lambda}$ by (1.9) and f^{ϕ}_{ρ} by (1.5). Suppose that there exists a positive constant C_{ρ} such that $d\rho_X = C_{\rho}dV$. Assume f^{ϕ}_{ρ}

$$(L^q(X), H_2^q(X))_{\theta}$$
 for some $0 < \theta$ 1. If $1 = q$ 2, then by taking $\lambda = \frac{\log^2 m}{m}$ we have

$$\mathbb{E}_{\mathbf{z}\mathbf{2}Z^{m}} R(sgn(f_{\mathbf{z},\lambda}) - R(f_{c}) = O \otimes \frac{\log^{2} m}{m} \frac{\frac{q^{2}\theta}{4(2+q)q^{2}\theta + 4qd}}{\mathbf{A}} \mathbf{A} . \tag{4.14}$$

If q > 2, then by taking $\lambda = \frac{\log^2 m}{m}$, we have

$$\mathbb{E}_{\mathsf{Z}\mathsf{Z}Z^m} \ \ \mathsf{R}(sgn(f_{\mathsf{Z},\lambda}) - \mathsf{R}(f_c) = O \quad \frac{\log^2 m}{m} \quad \frac{\frac{\theta}{4(2+q)\theta + 2qd}}{} \ . \tag{4.15}$$

Proof. For $\phi(x) = (1-x)_+^q (q > 1)$, we have the following estimate for the excess generalization error by Theorem 25 in [8].

$$E_{\phi}(f) - E_{\phi}(f_{\rho}^{\phi}) = \begin{cases} f - f_{\rho}^{\phi} q & \text{if } 1 < q = 2, \\ q2^{q-1} f - f_{\rho}^{\phi} L_{\rho_{X}}^{q}(2^{q-1} + f - f_{\rho}^{\phi} L_{\rho_{X}}^{q}), & \text{if } q > 2. \end{cases}$$

Together with the assumption that $d\rho_X = C_{\rho}dV$, for the case 1 < q = 2, we can get the desired results by using Theorem 3 with $\alpha = q$ and $\beta = \frac{q-1}{2}$. As to the case q > 2, choose $\alpha = 1$ and $\beta = \frac{q-1}{2}$. This is the end of the proof.

Acknowledgments

This work is supported partially by the Research Grants Council of Hong Kong [Project No. CityU 103206], City University of Hong Kong [Project No. 7002126], National Science Fund for Distinguished Young Scholars of China [Project No. 10529101], and National Basic Research Program of China [Project No. 973-2006CB303102]. The corresponding author is Ding-Xuan Zhou.

References

- [1] N. Aronszajn, Theory of reproducing kernels, *Trans. Amer. Math. Soc.* **68** (1950), 337–404.
- [2] M. Belkin and P. Niyogi, Towards a theoretical foundation for Laplacian-Based manifold methods, *COLT* (2005), P. Auer and R. Meir (eds.), pp. 486–500, 2005.
- [3] M. Belkin and P. Niyogi, Semi-supervised learning on Riemannian manifolds, *Mach. Learning* **56** (2004), 209–239.
- [4] M. Belkin and P. Niyogi, Laplacian eigenmaps for dimensionality reduction and data representation, *Neural Comput.* **15** (2003), 1373–1396.
- [5] J. Bergh and J. Löfström, Interpolation Spaces, Springer-Verlag, 1976.
- [6] O. Bousquet, O. Chapelle, and M. Hein, Measure based regularization, NIPS (2003).
- [7] M. do Carmo, Riemannian Geometry, Birkhäuser, Boston, 1992.
- [8] D. R. Chen, Q. Wu, Y. Ying and D. X. Zhou, Support vector machine soft margin classifiers: error analysis, *J. Mach. Learning Res.* 5 (2004), 1143–1175.

- [9] F. Cucker and D. X. Zhou, Learning Theory: An Approximation Theory Viewpoint, Cambridge University Press, 2007.
- [10] E. De Vito, A. Caponnetto, and L. Rosasco, Model selection for regularized least-squares algorithm in learning theory, *Found. Comput. Math.* **5** (2005), 59–85.
- [11] L. Devroye, L. Györfi, and G. Lugosi, A Probabilistic Theory of Pattern Recognition, Springer-Verlag, New York, 1997.
- [12] Z. Ditzian and V. Totik, Moduli of Smoothness, Springer-Verlag, New York, 1987.
- [13] T. Evgeniou, M. Pontil and T. Poggio, Regularization networks and suport vector machines, Adv. Comput. Math. 13(2000),1-50
- [14] E. Gine and V. Koltchinskii, Empirical graph Laplacian approximation of Laplace-Beltrami operators: large sample results, to appear in the Proceedings of the 4th International Conference on High Dimensional Probability.
- [15] D. Hardin, I. Tsamardinos, and C. F. Aliferis, A theoretical characterization of linear SVM-based feature selection, Proc. of the 21st Int. Conf. on Machine Learning, Banff, Canada, 2004.
- [16] E. Hebey, Sobolev Spaces on Riemannian Manifolds, Lecture Notes in Mathematics 1635, Springer-Verlag, 1996.
- [17] U. von Luxburg, M. Belkin, and O. Bousquet, Consistency of spectral clustering, preprint, 2004.
- [18] S. Mukherjee, Q. Wu, and D. X. Zhou, Learning gradients and feature selection on manifolds, preprint, 2006.
- [19] S. Smale and D. X. Zhou, Estimating the approximation error in learning theory, *Anal. Appl.* **1** (2003), 17–41.
- [20] S. Smale and D. X. Zhou, Learning theory estimates via integral operators and their applications, *Constr. Approx.* 2007.
- [21] Q. Wu, Y. Ying and D. X. Zhou, Learning theory: from regression to classification, Topics in Multivariate Approximation and Interpolation edited by K. Jetter et al., 2005.

- [22] Q. Wu, Y. Ying and D. X. Zhou, Multi-kernel regularized classifiers, *J. Complexity*, **23** (2007), 108-134.
- [23] Q. Wu and D. X. Zhou, SVM soft margin classifiers: linear programming versus quadratic programming, *Neural Comp.* **17** (2005), 1160–1187.
- [24] G. B. Ye and D. X. Zhou, Fully online classification by regularization, *Appl. Comput. Harmonic Anal.*, in press. http://dx.doi.org/10.1016/j.acha.2006.12.001.
- [25] G. B. Ye and D. X. Zhou, Learning and approximation by Gaussians on Riemannian manifolds, *Adv. Comput. Math.*, to appear.
- [26] Y. Ying and D. X. Zhou, Learnability of Gaussians with flexible variances, *J. Machine Learning Res.* 8 (2007), 249–276.
- [27] T. Zhang, Statistical behavior and consistency of classification methods based on convex risk minimization, *Ann. Stat.* **32** (2004), 56–85.
- [28] D. X. Zhou, The covering number in learning theory, J. Complexity 18 (2002), 739–767.
- [29] D. X. Zhou, Capacity of reproducing kernel spaces in learning theory, *IEEE Trans. Inform. Theory* **49** (2003), 1743–1752.