



HOKKAIDO UNIVERSITY

Title	Analyses on kernel-specific generalization ability for kernel regressors with training samples
Author(s)	Tanaka, Akira; Miyakoshi, Masaaki
Citation	2010 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), 61-66 https://doi.org/10.1109/ISSPIT.2010.5711725
Issue Date	2010-12-15
Doc URL	https://hdl.handle.net/2115/46942
Rights	© 2011 IEEE. Reprinted, with permission, from Tanaka, A., Miyakoshi, M., Analyses on kernel-specific generalization ability for kernel regressors with training samples, 2010 IEEE International Symposium on Signal Processing and Information Technology (ISSPIT), Dec. 2010. This material is posted here with permission of the IEEE. Such permission of the IEEE does not in any way imply IEEE endorsement of any of Hokkaido University products or services. Internal or personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution must be obtained from the IEEE by writing to pubs-permissions@ieee.org . By choosing to view this document, you agree to all provisions of the copyright laws protecting it.
Type	conference paper
File Information	ISSPIT2010_61-66.pdf



Analyses on Kernel-Specific Generalization Ability for Kernel Regressors with Training Samples

Akira Tanaka and Masaaki Miyakoshi

Division of Computer Science,

Graduate School of Information Science and Technology, Hokkaido University

Kita-14, Nishi-9, Kita-ku, Sapporo, 060-0814, Japan

{takira,miyakosi}@main.ist.hokudai.ac.jp

Abstract—Theoretical analyses on generalization error of a model space for kernel regressors with respect to training samples are given in this paper. In general, the distance between an unknown true function and a model space tends to be small with a larger set of training samples. However, it is not clarified that a larger set of training samples achieves a smaller difference at each point of the unknown true function and the orthogonal projection of it onto the model space, compared with a smaller set of training samples. In this paper, we show that the upper bound of the squared difference at each point of these two functions with a larger set of training samples is not larger than that with a smaller set of training samples. We also give some numerical examples to confirm our theoretical result.

Keywords—kernel regressor; reproducing kernel Hilbert space; generalization error; training samples;

I. INTRODUCTION

Learning based on kernel machines[1], represented by the support vector machine[2] and the kernel ridge regressor[2], is widely known as a powerful tool for various fields of information science such as pattern recognition, regression estimation, and density estimation. In general, an appropriate model selection is required in order to obtain a small generalization error in kernel machines. Although many methods for model selection, such as the leave-one-out cross-validation, were proposed, it is important to analyze generalization error theoretically since it may be useful to improve the performance of the model selection methods.

Generalization error is usually defined as the distance between an unknown true function and an estimated one; and it can be decomposed into two components in kernel machines. One is the distance between the unknown true function and an adopted model space, which is the linear subspace spanned by kernel functions corresponding to points in a set of training samples. The other is the distance between the estimated function and the orthogonal projection of the unknown true function onto the model space. The former one was not sufficiently investigated so far, while the latter one was discussed in many articles (see [3], [4] for instance). In our previous work[5], we investigated the former one and showed that a kernel corresponding to the smallest reproducing kernel Hilbert space including an unknown true function gives the best model space among a class of kernels with an invariant metric.

In this paper, we theoretically analyze the generalization error of the model space with respect to training samples. It is intuitively trivial that a larger set of training samples achieves a better model space compared with a smaller set of training samples. However, it is not trivial that a larger set of training samples achieves a smaller difference at each point of the unknown true function and the orthogonal projection of it onto the model space compared with a smaller set of training samples. We show that the upper bound of the squared difference at each point of these two functions with a large set of training samples is not larger than that with a small set of training samples. We also give some numerical examples to confirm our theoretical result.

II. MATHEMATICAL PRELIMINARIES FOR THE THEORY OF REPRODUCING KERNEL HILBERT SPACES

In this section, we prepare some mathematical tools concerned with the theory of reproducing kernel Hilbert spaces[6], [7].

Definition 1: [6] Let \mathbf{R}^n be an n -dimensional real vector space and let \mathcal{H} be a class of functions defined on $\mathcal{D} \subset \mathbf{R}^n$, forming a Hilbert space of real-valued functions. The function $K(\mathbf{x}, \tilde{\mathbf{x}})$, ($\mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{D}$) is called a reproducing kernel of \mathcal{H} , if

- 1) For every $\tilde{\mathbf{x}} \in \mathcal{D}$, $K(\cdot, \tilde{\mathbf{x}})$ is a function belonging to \mathcal{H} .
- 2) For every $\tilde{\mathbf{x}} \in \mathcal{D}$ and every $f(\cdot) \in \mathcal{H}$,

$$f(\tilde{\mathbf{x}}) = \langle f(\cdot), K(\cdot, \tilde{\mathbf{x}}) \rangle_{\mathcal{H}}, \quad (1)$$

where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ denotes the inner product of the Hilbert space \mathcal{H} .

The Hilbert space \mathcal{H} that has a reproducing kernel is called a reproducing kernel Hilbert space (RKHS). The reproducing property Eq.(1) enables us to treat a value of a function at a point in \mathcal{D} . Note that reproducing kernels are positive definite [6]:

$$\sum_{i,j=1}^N c_i c_j K(\mathbf{x}_i, \mathbf{x}_j) \geq 0, \quad (2)$$

for any N , $c_1, \dots, c_N \in \mathbf{R}$, and $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{D}$. In addition, $K(\mathbf{x}, \tilde{\mathbf{x}}) = K(\tilde{\mathbf{x}}, \mathbf{x})$ for any $\mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{D}$ is followed[6]. If a reproducing kernel $K(\mathbf{x}, \tilde{\mathbf{x}})$ exists, it is unique[6]. Conversely,

every positive definite function $K(\mathbf{x}, \tilde{\mathbf{x}})$ has the unique corresponding RKHS [6].

Next, we introduce the Schatten product [8] that is a convenient tool to reveal the reproducing property of kernels.

Definition 2: [8] Let \mathcal{H}_1 and \mathcal{H}_2 be Hilbert spaces. The Schatten product of $g \in \mathcal{H}_2$ and $h \in \mathcal{H}_1$ is defined by

$$(g \otimes h)f = \langle f, h \rangle_{\mathcal{H}_1} g, \quad f \in \mathcal{H}_1. \quad (3)$$

Note that $(g \otimes h)$ is a linear operator from \mathcal{H}_1 onto \mathcal{H}_2 . It is easy to show that the following relations hold for $h, v \in \mathcal{H}_1$, $g, u \in \mathcal{H}_2$.

$$(h \otimes g)^* = (g \otimes h), \quad (4)$$

$$(h \otimes g)(u \otimes v) = \langle u, g \rangle_{\mathcal{H}_2} (h \otimes v), \quad (5)$$

where the superscript $*$ denotes the adjoint operator.

III. FORMULATION OF LEARNING PROBLEMS

Let $\{(y_i, \mathbf{x}_i) | i = 1, \dots, \ell\}$ be a given training data set with $y_i \in \mathbf{R}$, $\mathbf{x}_i \in \mathbf{R}^n$, satisfying

$$y_i = f(\mathbf{x}_i) + n_i, \quad (6)$$

where $f(\cdot)$ denotes the unknown true function and n_i denotes a zero-mean additive noise. The aim of machine learning is to estimate the unknown true function $f(\cdot)$ by using the given training data set and statistical properties of noise.

In this paper, we assume that the unknown true function $f(\cdot)$ belongs to the RKHS \mathcal{H}_K corresponding to a certain kernel function K . If $f(\cdot) \in \mathcal{H}_K$, then Eq.(6) is rewritten as

$$y_i = \langle f(\cdot), K(\cdot, \mathbf{x}_i) \rangle_{\mathcal{H}_K} + n_i, \quad (7)$$

on the basis of the reproducing property of kernels. Let $\mathbf{y} = [y_1, \dots, y_\ell]'$ and $\mathbf{n} = [n_1, \dots, n_\ell]'$ with the superscript $'$ denoting the transposed matrix (or vector), then applying the Schatten product to Eq.(7) yields

$$\mathbf{y} = \left(\sum_{k=1}^{\ell} [e_k^{(\ell)} \otimes K(\cdot, \mathbf{x}_k)] \right) f(\cdot) + \mathbf{n}, \quad (8)$$

where $e_k^{(\ell)}$ denotes the k -th vector of the canonical basis of \mathbf{R}^ℓ . For a convenience of description, we write

$$A_{K,X} = \left(\sum_{k=1}^{\ell} [e_k^{(\ell)} \otimes K(\cdot, \mathbf{x}_k)] \right), \quad (9)$$

where $X = \{\mathbf{x}_k | k \in \{1, \dots, \ell\}\}$. The operator $A_{K,X}$ is a linear operator that maps an element in \mathcal{H}_K onto \mathbf{R}^ℓ and Eq.(8) can be written by

$$\mathbf{y} = A_{K,X} f(\cdot) + \mathbf{n}, \quad (10)$$

which represents the relation between the unknown true function $f(\cdot)$ and an output vector \mathbf{y} . Therefore, a machine learning problem can be interpreted as an inversion problem of the linear equation Eq.(10)[9].

IV. KERNEL SPECIFIC GENERALIZATION ABILITY

In general, a learning result by kernel machines is represented by a linear combination of $K(\cdot, \mathbf{x}_i)$, which means that the learning result is an element in $\mathcal{R}(A_{K,X}^*)$ (the range space of the linear operator $A_{K,X}^*$) since

$$\begin{aligned} \hat{f}(\cdot) &= A_{K,X}^* \boldsymbol{\alpha} \\ &= \left(\sum_{k=1}^{\ell} [K(\cdot, \mathbf{x}_k) \otimes e_k^{(\ell)}] \right) \boldsymbol{\alpha} \\ &= \sum_{k=1}^{\ell} \alpha_k K(\cdot, \mathbf{x}_k) \end{aligned} \quad (11)$$

holds, where $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_\ell]'$ denotes an arbitrary vector in \mathbf{R}^ℓ . The point at issue of this paper is selection of a model space, that is, the generalization ability of $\mathcal{R}(A_{K,X}^*)$ which is independent from criteria of learning machines. In order to discuss $\mathcal{R}(A_{K,X}^*)$, the orthogonal projector onto $\mathcal{R}(A_{K,X}^*)$ in \mathcal{H}_K , written as $P_{K,X}$, plays an important role.

Lemma 1: [5]

$$P_{K,X} = \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} (G_{K,X}^+)_{i,j} [K(\cdot, \mathbf{x}_i) \otimes K(\cdot, \mathbf{x}_j)], \quad (12)$$

where $G_{K,X}$ denotes the Gramian matrix of K with X and the superscript $+$ denotes the Moore-Penrose generalized inverse matrix[10].

Let $\|\cdot\|_{\mathcal{H}_K}$ be the induced norm in \mathcal{H}_K , then a generalization error of $\mathcal{R}(A_{K,X}^*)$ for $f(\cdot) \in \mathcal{H}_K$ can be defined as

$$J_{K,X}^{(N)} = \|f(\cdot) - P_{K,X} f(\cdot)\|_{\mathcal{H}_K}^2, \quad (13)$$

which is called 'norm-based generalization error.' In [5], we discussed the norm-based generalization error of a class of kernels with an invariant metric; and obtained the following result concerned with a selection of RKHS itself.

Theorem 1: [5] Let K_1 and K_2 be kernels satisfying the following properties:

- (1) $\mathcal{H}_{K_1} \subset \mathcal{H}_{K_2}$ (nested).
- (2) $\langle f(\cdot), g(\cdot) \rangle_{\mathcal{H}_{K_1}} = \langle f(\cdot), g(\cdot) \rangle_{\mathcal{H}_{K_2}}$ for any $f(\cdot), g(\cdot) \in \mathcal{H}_{K_1}$ (invariant metric).

Then, for any $f(\cdot) \in \mathcal{H}_{K_1}$ and any set of input vectors X ,

$$J_{K_1,X}^{(N)} \leq J_{K_2,X}^{(N)} \quad (14)$$

holds.

This theorem claims that the kernel corresponding to the smallest RKHS gives the best model among a class of kernels with an invariant metric whose corresponding RKHS's include the unknown true function.

On the other hand in this paper, we discuss the properties of $\mathcal{R}(A_{K,X}^*)$ with respect to a set of training samples. We show

the following lemmas as preparations.

Lemma 2: [5] For any $f(\cdot) \in \mathcal{H}_K$ and X ,

$$\|P_{K,X}f(\cdot)\|_{\mathcal{H}_K}^2 = \mathbf{f}'G_{K,X}^+\mathbf{f} \quad (15)$$

holds, where $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_\ell)]'$.

Lemma 3: Let

$$G = \begin{bmatrix} A & B \\ B' & C \end{bmatrix} \in \mathbf{R}^{(n+m) \times (n+m)} \quad (16)$$

be an *n.n.d.* symmetric matrix with $A \in \mathbf{R}^{n \times n}$, $C \in \mathbf{R}^{m \times m}$, and $B \in \mathbf{R}^{n \times m}$ and let $\mathbf{v} \in \mathcal{R}(G)$. Then,

$$\mathbf{v}' \left(G^+ - \begin{bmatrix} A^+ & O_{n,m} \\ O_{m,n} & O_{m,m} \end{bmatrix} \right) \mathbf{v} \geq 0 \quad (17)$$

holds, where $O_{m,n} \in \mathbf{R}^{m \times n}$ denotes the zero matrix.

Proof: Since G is *n.n.d.*, it is trivial that ZGZ' is also *n.n.d.* for any matrix $Z \in \mathbf{R}^{p \times (m+n)}$. Let

$$Z = \begin{bmatrix} I_n - AA^+ & O_{n,m} \\ -B'A^+ & I_m \end{bmatrix},$$

where I_n denotes the identity matrix of degree n . Then, using the facts that $I_n - A^+A$ is a symmetric orthogonal projector and $A^+A = AA^+$ holds for any symmetric matrix A yield that

$$\begin{aligned} ZG'Z &= \begin{bmatrix} I_n - AA^+ & O_{n,m} \\ -B'A^+ & I_m \end{bmatrix} \begin{bmatrix} A & B \\ B' & C \end{bmatrix} \\ &\times \begin{bmatrix} I_n - AA^+ & O_{n,m} \\ -B'A^+ & I_m \end{bmatrix}' \\ &= \begin{bmatrix} I_n - AA^+ & O_{n,m} \\ -B'A^+ & I_m \end{bmatrix} \begin{bmatrix} A & B \\ B' & C \end{bmatrix} \\ &\times \begin{bmatrix} I_n - AA^+ & -A^+B \\ O_{m,n} & I_m \end{bmatrix} \\ &= \begin{bmatrix} O_{n,n} & (I_n - AA^+)B \\ B'(I_n - A^+A) & C - B'A^+B \end{bmatrix} \end{aligned}$$

is also *n.n.d.*

Since $\mathbf{v} \in \mathcal{R}(G)$, there exist \mathbf{z} such that $\mathbf{v} = G\mathbf{z}$. Thus,

$$\begin{aligned} \mathbf{v}' \left(G^+ - \begin{bmatrix} A^+ & O_{n,m} \\ O_{m,n} & O_{m,m} \end{bmatrix} \right) \mathbf{v} &= \mathbf{z}'G \left(G^+ - \begin{bmatrix} A^+ & O_{n,m} \\ O_{m,n} & O_{m,m} \end{bmatrix} \right) G\mathbf{z} \\ &= \mathbf{z}' \left(\begin{bmatrix} A & B \\ B' & C \end{bmatrix} - \begin{bmatrix} A & AA^+B \\ B'A^+A & B'A^+B \end{bmatrix} \right) \mathbf{z} \\ &= \mathbf{z}' \begin{bmatrix} O_{n,n} & (I_n - AA^+)B \\ B'(I_n - A^+A) & C - B'A^+B \end{bmatrix} \mathbf{z} \geq 0 \end{aligned}$$

is obtained, which concludes the proof. \square

Let $X_S = \{\mathbf{x}_i \mid i \in \{1, \dots, \ell_S\}\}$ and let $X_L = \{\mathbf{x}_i \mid i \in \{1, \dots, \ell_L\}\}$ with $\ell_S < \ell_L$. Note that $X_S \subset X_L$ holds. It is intuitively trivial that $\mathcal{R}(A_{K,X_L}^*)$ gives a better model space than $\mathcal{R}(A_{K,X_S}^*)$ since $\mathcal{R}(A_{K,X_S}^*) \subset \mathcal{R}(A_{K,X_L}^*)$. Its theoretical ground is given by the following theorem.

Theorem 2: Let $f(\cdot) \in \mathcal{H}_K$, then

$$J_{K,X_L}^{(N)} \leq J_{K,X_S}^{(N)} \quad (18)$$

holds.

Proof: From Lemma 2 and the Pythagorean theorem, we have

$$\begin{aligned} J_{K,X}^{(N)} &= \|f(\cdot) - P_{K,X}f(\cdot)\|_{\mathcal{H}_K}^2 \\ &= \|f(\cdot)\|_{\mathcal{H}_K}^2 - \|P_{K,X}f(\cdot)\|_{\mathcal{H}_K}^2 \\ &= \|f(\cdot)\|_{\mathcal{H}_K}^2 - \mathbf{f}'G_{K,X}^+\mathbf{f}, \end{aligned} \quad (19)$$

where $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_\ell)]'$. Thus,

$$\begin{aligned} J_{K,X_S}^{(N)} - J_{K,X_L}^{(N)} &= \|f(\cdot) - P_{K,X_S}f(\cdot)\|_{\mathcal{H}_K}^2 \\ &\quad - \|f(\cdot) - P_{K,X_L}f(\cdot)\|_{\mathcal{H}_K}^2 \\ &= (\|f(\cdot)\|_{\mathcal{H}_K}^2 - \mathbf{f}'_S G_{K,X_S}^+ \mathbf{f}_S) \\ &\quad - (\|f(\cdot)\|_{\mathcal{H}_K}^2 - \mathbf{f}'_L G_{K,X_L}^+ \mathbf{f}_L) \\ &= \mathbf{f}'_L G_{K,X_L}^+ \mathbf{f}_L - \mathbf{f}'_S G_{K,X_S}^+ \mathbf{f}_S \end{aligned} \quad (20)$$

is obtained, where

$$\begin{aligned} \mathbf{f}_S &= [f(\mathbf{x}_1), \dots, f(\mathbf{x}_{\ell_S})]', \\ \mathbf{f}_L &= [f(\mathbf{x}_1), \dots, f(\mathbf{x}_{\ell_L})]'. \end{aligned}$$

Note that

$$\mathbf{f}_L \in \mathcal{R}(A_{K,X_L}) = \mathcal{R}(A_{K,X_L} A_{K,X_L}^*) = \mathcal{R}(G_{K,X_L})$$

holds since $f(\cdot) \in \mathcal{H}_K$, $\mathbf{f}_L = A_{K,X_L} f(\cdot)$, and

$$\begin{aligned} A_{K,X_L} A_{K,X_L}^* &= \left(\sum_{k=1}^{\ell_L} [\mathbf{e}_k^{(\ell_L)} \otimes K(\cdot, \mathbf{x}_k)] \right) \left(\sum_{j=1}^{\ell_L} [\mathbf{e}_j^{(\ell_L)} \otimes K(\cdot, \mathbf{x}_j)] \right)^* \\ &= \left(\sum_{k=1}^{\ell_L} [\mathbf{e}_k^{(\ell_L)} \otimes K(\cdot, \mathbf{x}_k)] \right) \left(\sum_{j=1}^{\ell_L} [K(\cdot, \mathbf{x}_j) \otimes \mathbf{e}_j^{(\ell_L)}] \right) \\ &= \sum_{k=1}^{\ell_L} \sum_{j=1}^{\ell_L} K(\mathbf{x}_k, \mathbf{x}_j) (\mathbf{e}_k^{(\ell_L)} \otimes \mathbf{e}_j^{(\ell_L)}) \\ &= \sum_{k=1}^{\ell_L} \sum_{j=1}^{\ell_L} K(\mathbf{x}_k, \mathbf{x}_j) \mathbf{e}_k^{(\ell_L)} (\mathbf{e}_j^{(\ell_L)})' = G_{K,X_L} \end{aligned}$$

hold. Let

$$E = \begin{bmatrix} I_{\ell_S} \\ O_{\ell_D, \ell_S} \end{bmatrix} \in \mathbf{R}^{\ell_L \times \ell_S}, \quad (21)$$

where $\ell_D = \ell_L - \ell_S$, then

$$E'G_{K,X_L}E = G_{K,X_S}.$$

Thus, applying Lemma 3 to Eq.(20) yields

$$\begin{aligned} & \mathbf{f}'_L G_{K,X_L}^+ \mathbf{f}_L - \mathbf{f}'_S G_{K,X_S}^+ \mathbf{f}_S \\ &= \mathbf{f}'_L G_{K,X_L}^+ \mathbf{f}_L - \mathbf{f}'_L E G_{K,X_S}^+ E' \mathbf{f}_L \\ &= \mathbf{f}'_L G_{K,X_L}^+ \mathbf{f}_L - \mathbf{f}'_L \begin{bmatrix} G_{K,X_S}^+ & O_{\ell_S, \ell_D} \\ O_{\ell_D, \ell_S} & O_{\ell_D, \ell_D} \end{bmatrix} \mathbf{f}_L \\ &= \mathbf{f}'_L \left(G_{K,X_L}^+ - \begin{bmatrix} G_{K,X_S}^+ & O_{\ell_S, \ell_D} \\ O_{\ell_D, \ell_S} & O_{\ell_D, \ell_D} \end{bmatrix} \right) \mathbf{f}_L \\ &\geq 0 \end{aligned}$$

is obtained, which concludes the proof. \square

As mentioned above, the statement in Theorem 2 is rather trivial since $\mathcal{R}(A_{K,X_S}^*) \subset \mathcal{R}(A_{K,X_L}^*)$. However, it is not trivial that a larger set of training samples achieves smaller squared difference at each point $\mathbf{x} \in \mathcal{D}$ of $f(\cdot)$ and $P_{K,X}f(\cdot)$, that is, Eq.(18) does not always mean

$$(f(\mathbf{x}) - P_{K,X_L}f(\mathbf{x}))^2 \leq (f(\mathbf{x}) - P_{K,X_S}f(\mathbf{x}))^2 \quad (22)$$

at an arbitrarily fixed point $\mathbf{x} \in \mathcal{D}$.

In [11], we analyzed the difference of $f(\cdot)$ and $P_{K,X}f(\cdot)$ at a point $\mathbf{x} \in \mathcal{D}$.

Lemma 4: [11] Let $f(\cdot) \in \mathcal{H}_K$, then

$$(f(\mathbf{x}) - P_{K,X}f(\mathbf{x}))^2 \leq \|f\|_{\mathcal{H}_K}^2 E_{K,X}(\mathbf{x}) \quad (23)$$

holds for any $\mathbf{x} \in \mathcal{D}$, where

$$\begin{aligned} E_{K,X}(\mathbf{x}) &= K(\mathbf{x}, \mathbf{x}) \\ &\quad - \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} K(\mathbf{x}, \mathbf{x}_j)(G_{K,X}^+)_{i,j} K(\mathbf{x}, \mathbf{x}_i). \end{aligned} \quad (24)$$

Proof: [11] Using the property $P_{K,X} = P_{K,X}^*$, the reproducing property Eq.(1), the Schwarz's inequality, and the Pythagorean theorem, we have

$$\begin{aligned} & (f(\mathbf{x}) - P_{K,X}f(\mathbf{x}))^2 \\ &= (\langle f(\cdot), K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K} - \langle P_{K,X}f(\cdot), K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K})^2 \\ &= (\langle f(\cdot), K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K} - \langle f(\cdot), P_{K,X}K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K})^2 \\ &= (\langle f(\cdot), K(\cdot, \mathbf{x}) - P_{K,X}K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K})^2 \\ &\leq \|f(\cdot)\|_{\mathcal{H}_K}^2 \|K(\cdot, \mathbf{x}) - P_{K,X}K(\cdot, \mathbf{x})\|_{\mathcal{H}_K}^2 \\ &= \|f(\cdot)\|_{\mathcal{H}_K}^2 (\|K(\cdot, \mathbf{x})\|_{\mathcal{H}_K}^2 - \|P_{K,X}K(\cdot, \mathbf{x})\|_{\mathcal{H}_K}^2) \\ &= \|f(\cdot)\|_{\mathcal{H}_K}^2 E_{K,X}(\mathbf{x}), \end{aligned}$$

which concludes the proof. \square

This lemma claims that the upper bound of the squared difference between $f(\cdot)$ and $P_{K,X}f(\cdot)$ at a point $\mathbf{x} \in \mathcal{D}$ is

proportional to $\|f(\cdot)\|_{\mathcal{H}_K}^2$ and $E_{K,X}(\mathbf{x})$. Since a kernel K is fixed in the context of this paper, we define the relative upper bound of the squared difference of $f(\cdot)$ and $P_{K,X}f(\cdot)$ at a point $\mathbf{x} \in \mathcal{D}$ by $E_{K,X}(\mathbf{x})$. Incorporating this lemma, we give the following theorem which is the main result of this paper.

Theorem 3: For any $\mathbf{x} \in \mathcal{D}$,

$$E_{K,X_L}(\mathbf{x}) \leq E_{K,X_S}(\mathbf{x}) \quad (25)$$

holds.

Proof: Let

$$\mathbf{k}_L = [K(\mathbf{x}, \mathbf{x}_1), \dots, K(\mathbf{x}, \mathbf{x}_{\ell_L})]'$$

then

$$\mathbf{k}_L \in \mathcal{R}(A_{K,X_L}) = \mathcal{R}(A_{K,X_L} A_{K,X_L}^*) = \mathcal{R}(G_{K,X_L})$$

since $K(\mathbf{x}, \cdot) \in \mathcal{H}_K$ and $\mathbf{k}_L = A_{K,X_L} K(\mathbf{x}, \cdot)$ holds for an arbitrarily fixed $\mathbf{x} \in \mathcal{D}$. Let

$$\mathbf{k}_S = [K(\mathbf{x}, \mathbf{x}_1), \dots, K(\mathbf{x}, \mathbf{x}_{\ell_S})]'$$

then \mathbf{k}_S can be written as

$$\mathbf{k}_S = E' \mathbf{k}_L$$

with E in Eq.(21). Thus,

$$\begin{aligned} & E_{K,X_S}(\mathbf{x}) - E_{K,X_L}(\mathbf{x}) \\ &= \left(K(\mathbf{x}, \mathbf{x}) - \sum_{i=1}^{\ell_S} \sum_{j=1}^{\ell_S} K(\mathbf{x}, \mathbf{x}_j)(G_{K,X_S}^+)_{i,j} K(\mathbf{x}, \mathbf{x}_i) \right) \\ &\quad - \left(K(\mathbf{x}, \mathbf{x}) - \sum_{i=1}^{\ell_L} \sum_{j=1}^{\ell_L} K(\mathbf{x}, \mathbf{x}_j)(G_{K,X_L}^+)_{i,j} K(\mathbf{x}, \mathbf{x}_i) \right) \\ &= \mathbf{k}'_L G_{K,X_L}^+ \mathbf{k}_L - \mathbf{k}'_S G_{K,X_S}^+ \mathbf{k}_S \\ &= \mathbf{k}'_L G_{K,X_L}^+ \mathbf{k}_L - \mathbf{k}'_L E G_{K,X_S}^+ E' \mathbf{k}_L \\ &= \mathbf{k}'_L \left(G_{K,X_L}^+ - \begin{bmatrix} G_{K,X_S}^+ & O_{\ell_S, \ell_D} \\ O_{\ell_D, \ell_S} & O_{\ell_D, \ell_D} \end{bmatrix} \right) \mathbf{k}_L \\ &\geq 0 \end{aligned}$$

is obtained from Lemma 3, which concludes the proof. \square

According to Theorem 3, it is guaranteed that the upper bound of the squared difference at each point $\mathbf{x} \in \mathcal{D}$ of $f(\cdot)$ and $P_{K,X_L}f(\cdot)$ is not greater than that of $f(\cdot)$ and $P_{K,X_S}f(\cdot)$.

Since the range of $E_{K,X}(\mathbf{x})$ depends on $\mathbf{x} \in \mathcal{D}$, it is significant to consider the normalized version of it. $E_{K,X}(\mathbf{x})$ is non-negative since it is the squared norm of $K(\cdot, \mathbf{x}) - P_{K,X}K(\cdot, \mathbf{x})$. Moreover,

$$\sum_{i=1}^{\ell} \sum_{j=1}^{\ell} K(\mathbf{x}, \mathbf{x}_j)(G_{K,X}^+)_{i,j} K(\mathbf{x}, \mathbf{x}_i)$$

in Eq.(24) is also non-negative since it is a quadratic form with the *n.n.d.* matrix $G_{K,X}^+$. Therefore, we have

$$K(\mathbf{x}, \mathbf{x}) \geq \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} K(\mathbf{x}, \mathbf{x}_j)(G_{K,X}^+)_{i,j} K(\mathbf{x}, \mathbf{x}_i) \geq 0. \quad (26)$$

Note that when $K(\mathbf{x}, \mathbf{x}) = 0$, $E_{K,X}(\mathbf{x})$ is necessarily reduced to 0. When $K(\mathbf{x}, \mathbf{x}) \neq 0$, we also have

$$1 \geq \frac{1}{K(\mathbf{x}, \mathbf{x})} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} K(\mathbf{x}, \mathbf{x}_j)(G_{K,X}^+)_{i,j} K(\mathbf{x}, \mathbf{x}_i) \geq 0 \quad (27)$$

from Eq.(26).

Accordingly, we define the normalized point-wise generalization error of $\mathcal{R}(A_{K,X}^*)$ at a point $\mathbf{x} \in \mathcal{D}$ by

$$J_{K,X}^{(P)}(\mathbf{x}) = \begin{cases} \frac{E_{K,X}(\mathbf{x})}{K(\mathbf{x}, \mathbf{x})} & : (K(\mathbf{x}, \mathbf{x}) \neq 0) \\ 0 & : (K(\mathbf{x}, \mathbf{x}) = 0) \end{cases} \quad (28)$$

and it satisfies

$$0 \leq J_{K,X}^{(P)}(\mathbf{x}) \leq 1 \quad (29)$$

for any $\mathbf{x} \in \mathcal{D}$. When $J_{K,X}^{(P)}(\mathbf{x})$ is close to 0, the generalization error at $\mathbf{x} \in \mathcal{D}$ is small; and when $J_{K,X}^{(P)}(\mathbf{x})$ is close to 1, the generalization error at $\mathbf{x} \in \mathcal{D}$ is close to its upper bound. For $J_{K,X}^{(P)}(\mathbf{x})$, we have the following corollary concerned with a set of training samples.

Corollary 1:

$$J_{K,X_L}^{(P)}(\mathbf{x}) \leq J_{K,X_S}^{(P)}(\mathbf{x}) \quad (30)$$

holds for any $\mathbf{x} \in \mathcal{D}$.

Proof: It is trivial in case of $K(\mathbf{x}, \mathbf{x}) = 0$.

When $K(\mathbf{x}, \mathbf{x}) \neq 0$, for an arbitrarily fixed $\mathbf{x} \in \mathcal{D}$,

$$\begin{aligned} & J_{K,X_S}^{(P)}(\mathbf{x}) - J_{K,X_L}^{(P)}(\mathbf{x}) \\ &= \frac{1}{K(\mathbf{x}, \mathbf{x})} E_{K,X_S}(\mathbf{x}) - \frac{1}{K(\mathbf{x}, \mathbf{x})} E_{K,X_L}(\mathbf{x}) \\ &= \frac{1}{K(\mathbf{x}, \mathbf{x})} (E_{K,X_S}(\mathbf{x}) - E_{K,X_L}(\mathbf{x})) \end{aligned}$$

is obtained. Thus, Theorem 3 and the fact that $K(\mathbf{x}, \mathbf{x}) > 0$ immediately concludes the proof. \square

We can evaluate $J_{K,X}^{(P)}(\mathbf{x})$ at the point which yields the worst generalization error instead of each point $\mathbf{x} \in \mathcal{D}$ by

$$J_{K,X}^{(WP)} = \sup_{\mathbf{x} \in \mathcal{D}} J_{K,X}^{(P)}(\mathbf{x}). \quad (31)$$

We have the following corollary for $J_{K,X}^{(WP)}$ as the same with $J_{K,X}^{(P)}(\mathbf{x})$.

Corollary 2:

$$J_{K,X_L}^{(WP)} \leq J_{K,X_S}^{(WP)}. \quad (32)$$

Proof: Let

$$\mathbf{x}_L = \arg \sup_{\mathbf{x} \in \mathcal{D}} J_{K,X_L}^{(P)}(\mathbf{x}).$$

then Corollary 1 yields

$$J_{K,X_L}^{(P)}(\mathbf{x}_L) \leq J_{K,X_S}^{(P)}(\mathbf{x}_L).$$

On the other hand, it is trivial that

$$J_{K,X_S}^{(P)}(\mathbf{x}_L) \leq \sup_{\mathbf{x} \in \mathcal{D}} J_{K,X_S}^{(P)}(\mathbf{x}),$$

which concludes the proof. \square

In the norm-based generalization error, the norm of the unknown true function plays a crucial role since

$$J_{K,X}^{(N)} = \|f\|_{\mathcal{H}_K}^2 - \mathbf{f}' G_{K,X}^+ \mathbf{f}$$

as shown in the previous section. Thus, $J_{K,X}^{(N)}$ is rather meaningless for evaluation of the generalization ability of $\mathcal{R}(A_{K,X}^*)$ itself.

On the other hand in $E_{K,X}(\mathbf{x})$, the norm of the unknown true function appears as a coefficient. Thus, its influence to the generalization error can be vanished by normalization as in $J_{K,X}^{(P)}(\mathbf{x})$ (and in $J_{K,X}^{(WP)}$), which enables us to have a quantitative goodness of $\mathcal{R}(A_{K,X}^*)$ which is in the interval between the worst case (upper bound) and the best case. This is an essential value of considering $J_{K,X}^{(P)}(\mathbf{x})$ or $J_{K,X}^{(WP)}$.

V. NUMERICAL EXAMPLES

In this section, we confirm the theoretical results in the previous section by numerical examples for 1-dimensional function estimation. We adopt the Gaussian kernel given by

$$K(x, y) = \exp\left(-\frac{(x-y)^2}{4}\right), \quad x, y \in [0, 16] \in \mathbf{R} \quad (33)$$

as a kernel and

$$f(\cdot) = \sum_{k=1}^{100} c_k K(\cdot, x_k) \in \mathcal{H}_K \quad (34)$$

as an unknown true function, where x_i is randomly generated from the uniform distribution on $[0, 16]$ written as $U(0, 16)$ and c_i is randomly generated from the standard normal distribution $N(0, 1)$. Coefficient c_i is normalized so that $\|f(\cdot)\|_{\mathcal{H}_K}^2 = 1$. As a smaller and a larger sets of training samples, we adopt

$$\begin{aligned} X_S &= \{x_i \sim U(0, 16) \mid i \in \{1, \dots, 10\}\}, \\ X_L &= X_S \cup X_E \end{aligned}$$

where $X_E = \{x_i \sim U(0, 16) \mid i \in \{1, \dots, 5\}\}$. Figure 1 shows $f(\cdot)$, $P_{K,X_S} f(\cdot)$, $P_{K,X_L} f(\cdot)$, and training samples in which points in X_S are denoted by 'x' and those in X_E are denoted by '+'.
Figure 2 shows $J_{K,X_S}^{(P)}$ and $J_{K,X_L}^{(P)}$ with respect to $x \in [0, 16]$.

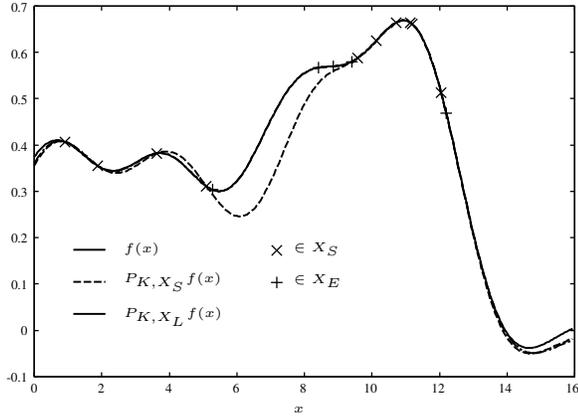


Fig. 1. An unknown true function, sampling points, and its orthogonal projections by P_{K, X_S} and P_{K, X_L} .

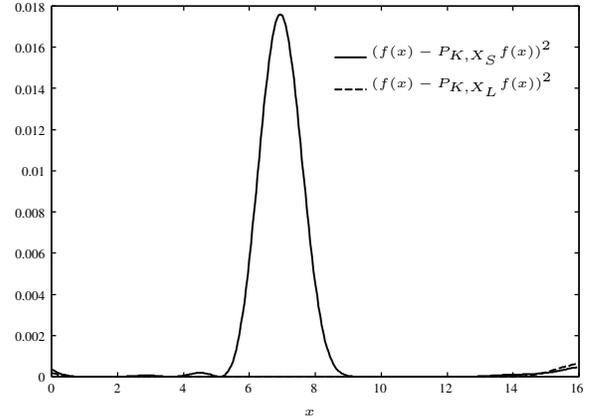


Fig. 3. Squared difference of $f(x)$ and $P_{K, X} f(x)$ with X_S and X_L .

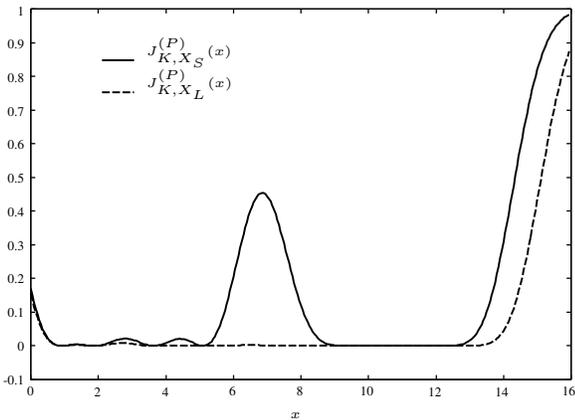


Fig. 2. $J_{K, X_L}^{(P)}$ and $J_{K, X_S}^{(P)}$ with respect to x .

According to this result, it is confirmed that the upper bound of the squared difference at each point in \mathcal{D} with a larger set of training samples is smaller than that with a smaller set of training samples, which agrees with the theoretical results obtained in the previous section.

Figure 3 shows the squared difference of $f(\cdot)$ and $P_{K, X} f(\cdot)$ with X_S and X_L .

According to this result, there exist points $x \in \mathcal{D}$, where the squared difference of $f(x)$ and $P_{K, X_L} f(x)$ is larger than that of $f(x)$ and $P_{K, X_S} f(x)$. This is a counter example for Eq.(22). Thus, our evaluation based on upper bound of the squared difference seems to be crucial.

VI. CONCLUSIONS

In this paper, we theoretically analyzed the generalization ability of a model space in kernel machines with respect to a set of training samples; and showed that a larger set of training samples achieves a smaller upper bound of the squared

difference of the unknown true function and the orthogonal projection of it onto the model space compared with a smaller set of training samples. We also gave numerical examples to confirm our theoretical results. Applying our result to practical model selection is one of future works that should be resolved.

ACKNOWLEDGMENTS

This work was partially supported by Grant-in-Aid No.21700001 for Young Scientists (B) from the Ministry of Education, Culture, Sports and Technology of Japan.

REFERENCES

- [1] K. Muller, S. Mika, G. Ratsch, K. Tsuda, and B. Scholkopf, "An introduction to kernel-based learning algorithms," *IEEE Transactions on Neural Networks*, vol. 12, pp. 181–201, 2001.
- [2] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines and other kernel-based learning methods*. Cambridge: Cambridge University Press, 2000.
- [3] M. Sugiyama and H. Ogawa, "Subspace Information Criterion for Model Selection," *Neural Computation*, vol. 13, no. 8, pp. 1863–1889, 2001.
- [4] M. Sugiyama, M. Kawanabe, and K. Muller, "Trading variance reduction with unbiasedness: The regularized subspace information criterion for robust model selection in kernel regression," *Neural Computation*, vol. 16, no. 5, pp. 1077–1104, 2004.
- [5] A. Tanaka, H. Imai, M. Kudo, and M. Miyakoshi, "Optimal kernel in a class of kernels with an invariant metric," in *Proc. S&SSPR2008*, 2008, pp. 530–539.
- [6] N. Aronszajn, "Theory of Reproducing Kernels," *Transactions of the American Mathematical Society*, vol. 68, no. 3, pp. 337–404, 1950.
- [7] J. Mercer, "Functions of Positive and Negative Type and Their Connection with The Theory of Integral Equations," *Transactions of the London Philosophical Society*, vol. A, no. 209, pp. 415–446, 1909.
- [8] R. Schatten, *Norm Ideals of Completely Continuous Operators*. Springer-Verlag, Berlin, 1960.
- [9] H. Ogawa, "What Can we See behind Sampling Theorems?" *IEICE Transactions on Fundamentals*, vol. E92-A, no. 3, pp. 688–707, 2009.
- [10] C. R. Rao and S. K. Mitra, *Generalized Inverse of Matrices and its Applications*. John Wiley & Sons, 1971.
- [11] A. Tanaka, H. Imai, and M. Miyakoshi, "Kernel-induced sampling theorem," *IEEE Transactions on Signal Processing (in printing)*.