



Title	A Relationship Between Generalization Error and Training Samples in Kernel Regressors
Author(s)	Tanaka, Akira; Imai, Hideyuki; Kudo, Mineichi; Miyakoshi, Masaaki
Citation	2010 20th International Conference on Pattern Recognition (ICPR), 1421-1424 https://doi.org/10.1109/ICPR.2010.351
Issue Date	2010-08-23
Doc URL	http://hdl.handle.net/2115/46851
Rights	© 2010 IEEE. Reprinted, with permission, from Akira Tanaka, Hideyuki Imai, Mineichi Kudo, and Masaaki Miyakoshi, A Relationship Between Generalization Error and Training Samples in Kernel Regressors, 2010 International Conference on Pattern Recognition, Aug. 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	proceedings (author version)
File Information	PR_1421-1424.pdf



[Instructions for use](#)

A Relationship Between Generalization Error and Training Samples in Kernel Regressors

Akira Tanaka, Hideyuki Imai, Mineichi Kudo, and Masaaki Miyakoshi
Division of Computer Science
Hokkaido University
Sapporo, Japan
e-mail: {takira,imai,mine,miyakosi}@main.ist.hokudai.ac.jp

Abstract—A relationship between generalization error and training samples in kernel regressors is discussed in this paper. The generalization error can be decomposed into two components. One is a distance between an unknown true function and an adopted model space. The other is a distance between an estimated function and the orthogonal projection of the unknown true function onto the model space. In our previous work, we gave a framework to evaluate the first component. In this paper, we theoretically analyze the second one and show that a larger set of training samples usually causes a larger generalization error.

Keywords—kernel regressor; reproducing kernel Hilbert space; generalization error; sample points;

I. INTRODUCTION

Learning based on kernel machines[1], represented by the support vector machine[2] and the kernel ridge regression[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 the model selection, such as the leave-one-out cross-validation, are proposed, it is important to analyze the generalization error theoretically since it may be useful to improve the performance of the model selection methods. In kernel machines, a model space is specified by a linear space spanned by kernel functions corresponding to points in training data set. Theoretical analyses of the generalization error with respect to a kernel or its parameters is usually difficult since the metrics of the corresponding reproducing kernel Hilbert spaces may differ. Thus, we focus on theoretical analyses of the generalization error with respect to a training data set. There are two kinds of generalization error. One is a distance between an unknown true function and an adopted model space. The other is a distance between an estimated function and the orthogonal projection of the unknown true function onto the model space. In our previous work[3], we discussed the first generalization error and gave an upper bound of the absolute difference between the unknown true function and its orthogonal projection onto the model space at each point. According to the results,

it immediately follows that the first generalization error decreases when the number of training samples increases. In this paper, we theoretically analyze the second generalization error and show that a larger set of training samples usually causes a larger generalization error. Numerical examples are also given to confirm our theoretical results.

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[4].

Definition 1: [4] 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 $f \in \mathcal{H}$,

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

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 (*p.d.*) [4]:

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

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[4]. If a reproducing kernel $K(\mathbf{x}, \tilde{\mathbf{x}})$ exists, it is unique[4]. Conversely, every *p.d.* function $K(\mathbf{x}, \tilde{\mathbf{x}})$ has the unique corresponding RKHS [4].

Next, we introduce the Schatten product [5] that is a convenient tool to reveal the reproducing property of kernels.

Definition 2: [5] 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. PROBLEM FORMULATION OF LEARNING AND KERNEL REGRESSORS

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

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

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

In this paper, we assume that the unknown function f belongs to the RKHS \mathcal{H}_K corresponding to a certain kernel function K . If $f \in \mathcal{H}_K$, 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 transposition operator, 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

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

since $A_{K,X}$ is specified by the kernel K and the set of input vectors $X = \{\mathbf{x}_i \in \mathcal{D} \mid i \in \{1, \dots, \ell\}\}$. $A_{K,X}$ is a linear operator that maps an element of \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 and an output vector \mathbf{y} . Therefore, a machine learning problem can be interpreted as an inversion problem of the linear equation Eq.(10)[6].

The minimum norm least squares solution for Eq.(10) is given by

$$\begin{aligned} \hat{f}(\cdot) &= A_{K,X}^+ \mathbf{y} = A_{K,X}^* G_{K,X}^+ \mathbf{y} \\ &= \sum_{k=1}^{\ell} \mathbf{y}' G_{K,X}^+ e_k^{(\ell)} K(\cdot, \mathbf{x}_k), \end{aligned} \quad (11)$$

where $G_{K,X}$ denotes the Gramian matrix of K with X and the superscript $+$ denotes the Moore-Penrose generalized inverse. Note that $A_{K,X}^+ A_{K,X}$ is the orthogonal projector onto $\mathcal{R}(A_{K,X}^*)$ (the range space of $A_{K,X}^*$, that is, the linear subspace spanned by $\{K(\cdot, \mathbf{x}_i) \mid i \in \{1, \dots, \ell\}\}$) and its closed form is given by

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

as shown in [7].

In practical problems, a solution by the kernel ridge regressor, which is a regularized version of Eq.(11) and is defined as

$$\hat{f}(\cdot) = \sum_{k=1}^{\ell} \mathbf{y}' (G_{K,X} + \mu I)^{-1} e_k^{(\ell)} K(\cdot, \mathbf{x}_k) \quad (13)$$

with $\mu > 0$ denoting a regularization parameter, is used instead of Eq.(11). However, theoretical analyses of a solution based on Eq.(11) can be an important basis of all other kernel machines including Eq.(13). Thus, we theoretically analyze the generalization error of the solution Eq.(11) in the following contents.

IV. GENERALIZATION ERROR OF A MODEL SPACE

In [3], we gave a framework to evaluate the generalization error of a model space, that is, $\mathcal{R}(A_{K,X}^*)$. In this section, we review results of [3] related to this paper.

Let $f \in \mathcal{H}_K$ be an unknown true function, then for any $\mathbf{x} \in \mathcal{D}$,

$$\begin{aligned} &|f(\mathbf{x}) - P_{K,X} f(\mathbf{x})| \\ &= |\langle f(\cdot), K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K} - \langle P_{K,X} f(\cdot), K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K}| \\ &= |\langle f(\cdot), K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K} - \langle f(\cdot), P_{K,X} K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K}| \\ &= |\langle f(\cdot), K(\cdot, \mathbf{x}) - P_{K,X} K(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_K}| \\ &\leq \|f\|_{\mathcal{H}_K} E_{K,X}^{1/2}(\mathbf{x}), \end{aligned}$$

holds, where

$$\begin{aligned} E_{K,X}(\mathbf{x}) &= \|K(\cdot, \mathbf{x}) - P_{K,X} K(\cdot, \mathbf{x})\|_{\mathcal{H}_K}^2 \\ &= \left(K(\mathbf{x}, \mathbf{x}) - \sum_{i,j=1}^{\ell} K(\mathbf{x}, \mathbf{x}_i) (G^+)_{i,j} K(\mathbf{x}, \mathbf{x}_j) \right), \end{aligned}$$

which implies that the absolute difference between f and $P_{K,X} f$ at a point $\mathbf{x} \in \mathcal{D}$ is proportional to $\|f\|_{\mathcal{H}_K}$ and $E_{K,X}^{1/2}(\mathbf{x})$ [3]. Let

$$J_{K,X}^{(1)} = \sup_{\mathbf{x} \in \mathcal{D}} (E_{K,X}(\mathbf{x})/K(\mathbf{x}, \mathbf{x}))^{1/2}, \quad (14)$$

then $0 \leq J_{K,X}^{(1)} \leq 1$ holds; and if $J_{K,X}^{(1)}$ is sufficiently close to zero, the model space $\mathcal{R}(A_{K,X}^*)$ has a sufficient ability to represent any $f \in \mathcal{H}_K$. Note that it is trivial that when $\tilde{X} \subset X$, $J_{K,X}^{(1)} \leq J_{K,\tilde{X}}^{(1)}$ holds. Thus, it is concluded that a larger set of training samples implies a smaller generalization error of an adopted model space with a fixed kernel.

V. GENERALIZATION ERROR IN A MODEL SPACE

The minimum norm least squares solution for Eq.(10) given in Section III can be decomposed as

$$\hat{f}(\cdot) = A_{K,X}^+ \mathbf{y} = A_{K,X}^+ A_{K,X} f(\cdot) + A_{K,X}^+ \mathbf{n}. \quad (15)$$

The first term is the orthogonal projection of f onto $\mathcal{R}(A_{K,X}^*)$ and its generalization error was analyzed in the previous section. The second term is the generalization error in $\mathcal{R}(A_{K,X}^*)$, coming from the additive noise, whose closed form is given as

$$\begin{aligned} \hat{f}_n(\cdot) &= A_{K,X}^+ \mathbf{n} \\ &= \left(\sum_{k=1}^{\ell} K(\cdot, \mathbf{x}_i) \otimes \mathbf{e}_k^{(\ell)} \right) G_{K,X}^+ \mathbf{n} \\ &= \sum_{k=1}^{\ell} \mathbf{n}' G_{K,X}^+ \mathbf{e}_k^{(\ell)} K(\cdot, \mathbf{x}_k), \end{aligned} \quad (16)$$

and its squared norm is given as

$$\begin{aligned} J_{K,X}^{(2)} &= \|\hat{f}_n\|_{\mathcal{H}_K}^2 = \langle \hat{f}_n, \hat{f}_n \rangle_{\mathcal{H}_K} \\ &= \mathbf{n}' G_{K,X}^+ G_{K,X} G_{K,X}^+ \mathbf{n} = \mathbf{n}' G_{K,X}^+ \mathbf{n}. \end{aligned}$$

Lemma 1: Let

$$G = \begin{bmatrix} A & B \\ B' & C \end{bmatrix} \in \mathbf{R}^{(n+m) \times (n+m)} \quad (17)$$

be a *p.d.* matrix with $A \in \mathbf{R}^{n \times n}$, $C \in \mathbf{R}^{m \times m}$, and $B \in \mathbf{R}^{n \times m}$, then

$$Z = G^{-1} - \begin{bmatrix} A^{-1} & O_{n,m} \\ O_{m,n} & O_{m,m} \end{bmatrix} \quad (18)$$

is non-negative definite (*n.n.d.*), where $O_{m,n}$ denotes the zero matrix in $\mathbf{R}^{m \times n}$.

Proof: As shown in [8], G^{-1} can be represented as

$$G^{-1} = \begin{bmatrix} A^{-1} + FE^{-1}F' & -FE^{-1} \\ -E^{-1}F' & E^{-1} \end{bmatrix},$$

where $E = C - B'A^{-1}B$ and $F = A^{-1}B$. Thus, Eq.(18) is reduced to

$$Z = \begin{bmatrix} FE^{-1}F' & -FE^{-1} \\ -E^{-1}F' & E^{-1} \end{bmatrix}. \quad (19)$$

It is trivial that E^{-1} is *p.d.* since G , G^{-1} , and all their principal minors are *p.d.* Let $\mathbf{v}_1 \in \mathbf{R}^n$ and $\mathbf{v}_2 \in \mathbf{R}^m$ be

arbitrary vectors and let $\mathbf{v} = [\mathbf{v}_1' \ \mathbf{v}_2']'$, then

$$\begin{aligned} \mathbf{v}' Z \mathbf{v} &= [\mathbf{v}_1' \ \mathbf{v}_2'] \begin{bmatrix} FE^{-1}F' & -FE^{-1} \\ -E^{-1}F' & E^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \end{bmatrix} \\ &= \mathbf{v}_1' FE^{-1}F' \mathbf{v}_1 - \mathbf{v}_1' FE^{-1} \mathbf{v}_2 \\ &\quad - \mathbf{v}_2' E^{-1}F' \mathbf{v}_1 + \mathbf{v}_2' E^{-1} \mathbf{v}_2 \\ &= (F' \mathbf{v}_1 - \mathbf{v}_2)' E^{-1} (F' \mathbf{v}_1 - \mathbf{v}_2) \geq 0 \end{aligned}$$

holds, which concludes the proof. \blacksquare

Note that the non-singularity of G in this lemma is crucial. In fact, the singular matrix

$$G = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

gives a simple counter example since $G^+ - [1 \ 0]'([1 \ 0]G[1 \ 0]')^{-1}[1 \ 0]$ is not *n.n.d.*

Theorem 1: Let $X = \{\mathbf{x}_i \mid i \in \{1, \dots, \ell\}\}$ be a set of training input vectors and let $\tilde{X} \subset X$ be a subset of X with $|\tilde{X}| = m < \ell$. Let $\mathbf{n} \in \mathbf{R}^\ell$ be the noise vector for X and let $\tilde{\mathbf{n}} \in \mathbf{R}^m$ be the noise vector for \tilde{X} defined by $\tilde{\mathbf{n}} = P_1' \mathbf{n}$, where $P_1' \in \mathbf{R}^{m \times \ell}$ denotes the full row-rank matrix, such that all components are zero except for one component being unity in each row, that extracts components corresponding to \tilde{X} from those of X . Then, if $G_{K,X}$ is non-singular,

$$J_{K,X}^{(2)} - J_{K,\tilde{X}}^{(2)} \geq 0 \quad (20)$$

holds.

Proof: Let $P_2' \in \mathbf{R}^{(\ell-m) \times \ell}$ be a full row-rank matrix, such that all components are zero except for one components being unity in each row, satisfying $P_1' P_2 = O_{m, \ell-m}$. It is obvious that $[P_1' \ P_2']$ is a permutation matrix, which implies that it is also an orthogonal matrix. Then,

$$\begin{aligned} J_{K,X}^{(2)} - J_{K,\tilde{X}}^{(2)} &= \mathbf{n}' G_{K,X}^{-1} \mathbf{n} - \tilde{\mathbf{n}}' G_{K,\tilde{X}}^{-1} \tilde{\mathbf{n}} \\ &= \mathbf{n}' [P_1' \ P_2'] [P_1' \ P_2']' G_{K,X}^{-1} [P_1' \ P_2'] [P_1' \ P_2']' \mathbf{n} \\ &\quad - \mathbf{n}' P_1' (P_1' G_{K,X} P_1)^{-1} P_1' \mathbf{n} \\ &= \mathbf{n}' [P_1' \ P_2'] ([P_1' \ P_2']' G_{K,X} [P_1' \ P_2'])^{-1} [P_1' \ P_2']' \mathbf{n} \\ &\quad - \mathbf{n}' P_1' (P_1' G_{K,X} P_1)^{-1} P_1' \mathbf{n} \\ &= \mathbf{n}' [P_1' \ P_2'] \begin{bmatrix} P_1' G_{K,X} P_1 & P_1' G_{K,X} P_2 \\ P_2' G_{K,X} P_1 & P_2' G_{K,X} P_2 \end{bmatrix}^{-1} \\ &\quad \times [P_1' \ P_2']' \mathbf{n} \\ &\quad - \mathbf{n}' [P_1' \ P_2'] \begin{bmatrix} (P_1' G_{K,X} P_1)^{-1} & O_{m, \ell-m} \\ O_{\ell-m, m} & O_{\ell-m, \ell-m} \end{bmatrix} \\ &\quad \times [P_1' \ P_2']' \mathbf{n} \geq 0 \end{aligned}$$

holds from Lemma 1. \blacksquare

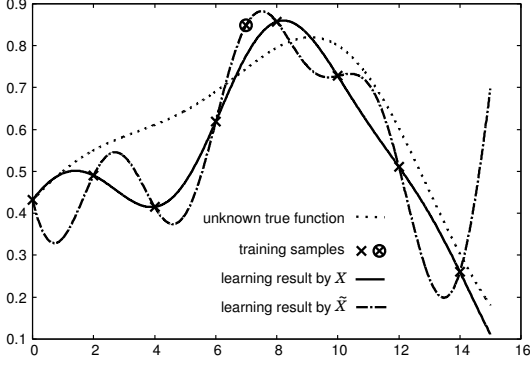


Figure 1. The unknown true function, the training samples, and the learning results by \tilde{X} and X .

According to Theorem 1, given a realization of an additive noise vector \mathbf{n} for X , the norm of the noise component in a learning result is larger than that based on $\tilde{\mathbf{n}}$ for \tilde{X} . It trivially holds that

$$E\mathbf{n}J_{K,X}^{(2)} - E\tilde{\mathbf{n}}J_{K,\tilde{X}}^{(2)} \geq 0 \quad (21)$$

when components of \mathbf{n} and $\tilde{\mathbf{n}}$ are *i.i.d.* with the variance σ^2 since $E\mathbf{n}J_{K,X}^{(2)} = E\mathbf{n}'G_{K,X}^+\mathbf{n} = \sigma^2\text{tr}(G_{K,X}^{-1})$.

VI. NUMERICAL EXAMPLE

In this section, we verify the behavior of the norm of the noise component in the solution Eq.(11) with an artificial data. We adopt the Gaussian kernel given by

$$K(x, y) = \exp\left(-\frac{(x-y)^2}{16}\right) \quad (22)$$

as a kernel.

Figure 1 shows the unknown true function in the corresponding RKHS, training data set with $\sigma = 0.1$ where $\tilde{X} = \{0, 2, 4, 6, 8, 10, 12, 14\}$ denoted by 'x' and $X = \tilde{X} \cup X_e$, where $X_e = \{7\}$ denoted by 'o', and the learning results based on \tilde{X} and X . $J_{K,\tilde{X}}^{(1)} = 0.139$ and $J_{K,X}^{(1)} = 0.116$, which implies that the model space $\mathcal{R}(A_{K,X}^*)$ is slightly better than $\mathcal{R}(A_{K,\tilde{X}}^*)$. On the other hand, $J_{K,\tilde{X}}^{(2)} = 0.408$ and $J_{K,X}^{(2)} = 58.227$, which agrees with the theoretical analyses given in the previous section.

Figure 2 shows the normalized histogram of

$$d(x_r) = \log\left(1 + J_{K,X}^{(2)} - J_{K,\tilde{X}}^{(2)}\right) \quad (23)$$

over 1,000 trials with $X_e = \{x_r\}$ in which x_r is randomly selected from $[0, 15]$. According to Fig. 2, it is confirmed that $J_{K,X}^{(2)}$ is larger than $J_{K,\tilde{X}}^{(2)}$ in all cases.

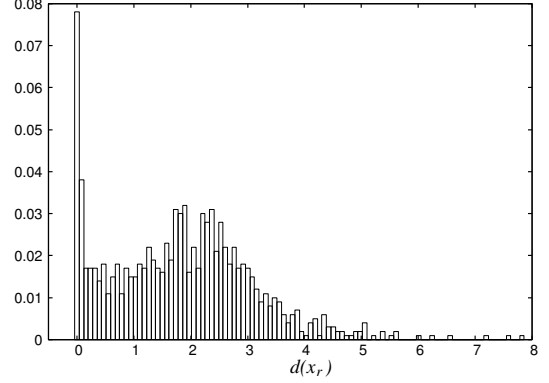


Figure 2. Normalized histogram of $d(x_r)$.

VII. CONCLUSION

In this paper, we investigated the relationship between the generalization error and training samples in kernel regressors and showed that a larger set of training samples causes a larger generalization error corresponding to noise components. Extending our result for other practical kernel machines is one of future works that should be resolved.

ACKNOWLEDGMENT

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] A. Tanaka, H. Imai, and M. Miyakoshi, "Kernel-induced sampling theorem," *IEEE Transactions on Signal Processing (in printing)*.
- [4] N. Aronszajn, "Theory of Reproducing Kernels," *Transactions of the American Mathematical Society*, vol. 68, no. 3, pp. 337–404, 1950.
- [5] R. Schatten, *Norm Ideals of Completely Continuous Operators*. Springer-Verlag, Berlin, 1960.
- [6] H. Ogawa, "Neural Networks and Generalization Ability," *IEICE Technical Report*, vol. NC95-8, pp. 57–64, 1995 (in Japanese).
- [7] 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.
- [8] H. Yanai and K. Takeuchi, *Projection Matrices, Generalized Inverses, and Singular Value Decomposition (in Japanese)*. Tokyo: University of Tokyo Press, 1983.