PAPER Kernel-Based Regressors Equivalent to Stochastic Affine Estimators

Akira TANAKA^{†a)}, Masanari NAKAMURA^{†b)}, and Hideyuki IMAI^{†c)}, Members

SUMMARY The solution of the ordinary kernel ridge regression, based on the squared loss function and the squared norm-based regularizer, can be easily interpreted as a stochastic linear estimator by considering the autocorrelation prior for an unknown true function. As is well known, a stochastic affine estimator is one of the simplest extensions of the stochastic linear estimator. However, its corresponding kernel regression problem is not revealed so far. In this paper, we give a formulation of the kernel regression problem, whose solution is reduced to a stochastic affine estimator, and also give interpretations of the formulation.

key words: kernel regression, autocorrelation prior, linear estimators, affine estimators, optimization criterion

1. Introduction

The kernel ridge regression (KRR) [1]-[3] is still one of useful function estimators in the field of machine learning. The ordinary KRR is defined as the minimizer of the squared loss function for a given training data set and the squared norm-based regularizer defined in a certain reproducing kernel Hilbert space [4] to which an unknown true function is assumed to belong. This regularizer involves two hyperparameters, that specify a model. One is a reproducing kernel and the other is a regularization parameter. Selection of a model, that achieves good generalization performance, is one of crucial topics in this field. As shown in [5], assuming the existence of the autocorrelation function of an unknown true function, the optimal model of the KRR, in terms of the expected squared error, is specified by the autocorrelation function itself (as a kernel) and the variance of additive noise (as a regularization parameter), which immediately leads the equivalence of the KRR and a stochastic linear estimator [6]. Note that the KRR with this optimal model also agrees with the solution of the Gaussian process regression (GPR) [7]-[10] with the zero-mean assumption, which is formulated by an another approach, that is, the conditional expectation based on a given training data set. It is widely believed that the superiority of the GPR against the KRR is that the GPR can identify the variance of estimates. However, almost the same result was also obtained by the theory of reproducing kernel Hilbert spaces only without the assumptions of

Manuscript revised September 15, 2021.

Manuscript publicized October 5, 2021.

b) E-mail: masanari@ist.hokudai.ac.jp

DOI: 10.1587/transinf.2021EDP7156

probabilistic structures as shown in [11]. Therefore, the essential difference between the KRR and the GPR with the zero-mean assumption is their formulations only.

Needless to say, a stochastic affine estimator [6] is one of the simplest extensions of a stochastic linear estimator. In the framework of the GPR, the solution, corresponding to a stochastic affine estimator, is straightforwardly obtained by centering [10], that is, subtraction of mean, as the same with the general framework of a stochastic estimation [6]. On the other hand, the formulation of kernel regression problems, corresponding to a stochastic affine estimator, is not revealed so far. In this paper, we reveal the formulation of the kernel regression problem whose solution is reduced to that of a stochastic affine estimator. Since the formulation of kernelbased regression is quite different from those of stochastic estimators and the GPR, it is expected that our result can provide a novel aspect for further extensions of the kernel regression problems, while the GPR has limited room for extension due to its principle.

2. Overview of Kernel Ridge Regression

In this section, we give an overview of the KRR [1].

Firstly, we briefly review the theory of reproducing kernel Hilbert spaces [4], [12], [13].

Definition 1: [4] Let \mathbf{R}^d be a *d*-dimensional real vector space and let \mathcal{H} be a class of functions defined on a certain domain $\mathcal{D} \subset \mathbf{R}^d$, 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 the following two conditions hold.

1.
$$\forall \tilde{x} \in \mathcal{D}, \ K(\cdot, \tilde{x}) \in \mathcal{H}$$
 (1)

2.
$$\forall \tilde{x} \in \mathcal{D}, \forall f(\cdot) \in \mathcal{H}, f(\tilde{x}) = \langle f(\cdot), K(\cdot, \tilde{x}) \rangle_{\mathcal{H}},$$
 (2)

where $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ stands for the inner product of \mathcal{H} .

In the following contents, we use the term 'kernel' instead of 'reproducing kernel' for simplicity.

The Hilbert space \mathcal{H} that has a kernel is called a reproducing kernel Hilbert space (RKHS). Note that kernels are positive definite and symmetric [4]. If a kernel $K(\mathbf{x}, \tilde{\mathbf{x}})$ exists, it is unique [4]. Conversely, every positive definite function $K(\mathbf{x}, \tilde{\mathbf{x}})$ has the unique corresponding RKHS [4]. From Eqs. (1) and (2), it immediately follows that $\langle K(\cdot, \mathbf{x}), K(\cdot, \tilde{\mathbf{x}}) \rangle_{\mathcal{H}_{K}} = K(\mathbf{x}, \tilde{\mathbf{x}})$ for any $\mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{D}$, which plays a crucial role in kernel-based regression problems.

Let
$$T := \{(x_i, y_i) \mid i \in \{1, ..., n\}, x_i \in \mathcal{D}, y_i \in \mathbf{R}\}$$

Manuscript received July 19, 2021.

[†]The authors are with the Division of Computer Science and Information Technology, Faculty of Information Science and Technology, Hokkaido University, Sapporo-shi, 060–0814 Japan.

a) E-mail: takira@ist.hokudai.ac.jp

c) E-mail: imai@ist.hokudai.ac.jp

be a given training data set with *n* samples, where x_i and y_i denote an input vector and the corresponding output value, satisfying

$$y_i = f(\boldsymbol{x}_i) + \epsilon_i, \tag{3}$$

where $f(\cdot)$ denotes the unknown true function to be estimated and ϵ_i denotes an *i.i.d.* zero-mean additive noise whose variance is (unknown) $\sigma^2 > 0$. The ordinary KRR is formulated as the following problem [1]–[3].

Problem 1: Find a function $\hat{f}(\cdot)$ that minimizes

$$J_1(\hat{f}(\cdot)) := \sum_{i=1}^n (y_i - \hat{f}(\mathbf{x}_i))^2 + \gamma ||\hat{f}(\cdot)||_{\mathcal{H}_{\kappa}}^2,$$
(4)

where $\|\cdot\|_{\mathcal{H}_K}$ denotes the induced norm of the RKHS \mathcal{H}_K corresponding to an adopted kernel *K*, and γ denotes a positive regularization parameter.

The following theorem, called the nonparametric representer theorem, is a well known result on the representation of a solution of a certain class of kernel regression problems.

Theorem 1: [14] Let $c(\cdot) : (\mathcal{D} \times \mathbf{R}^2)^n \to R \cup \{\infty\}$ be an arbitrary cost function, and let $g(\cdot)$ be a strictly monotonically increasing function defined on $[0, \infty)$. Then, any minimizer $\hat{f}(\cdot) \in \mathcal{H}_K$ of the regularized risk functional

$$c((\mathbf{x}_1, y_1, \hat{f}(\mathbf{x}_1)), \dots, (\mathbf{x}_n, y_n, \hat{f}(\mathbf{x}_n))) + g(\|\hat{f}(\cdot)\|_{\mathcal{H}_K})$$
 (5)

can be represented by

$$\hat{f}(\cdot) := \sum_{i=1}^{n} \alpha_i K(\cdot, \mathbf{x}_i) = (\mathbf{g}_X^{(K)}(\cdot))' \boldsymbol{\alpha}, \tag{6}$$

where $g_X^{(K)}(\cdot) := [K(\cdot, \mathbf{x}_1), \dots, K(\cdot, \mathbf{x}_n)]' \in \mathbf{R}^n$ and $\alpha := [\alpha_1, \dots, \alpha_n]' \in \mathbf{R}^n$ with ' standing for the transposition operator, and $X := \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ being the set of training input vectors.

Since $J_1(\hat{f}(\cdot))$ in Problem 1 agrees to Eq. (5), the solution of Problem 1 can be represented by Eq. (6). It is well known [1]–[3] that the closed-form solution of the minimizer of Eq. (4) with the function model Eq. (6) is given by

$$\hat{f}^{(K,\gamma)}(\cdot) = (\boldsymbol{g}_{X}^{(K)}(\cdot))' (G_{XX}^{(K)} + \gamma I_{n})^{-1} \boldsymbol{y},$$
(7)

where $G_{XX}^{(K)} := (K(\mathbf{x}_i, \mathbf{x}_j)) \in \mathbf{R}^{n \times n}$ denotes the Gram matrix of the kernel *K* with *X*, I_n denotes the identity matrix of degree *n*, and $\mathbf{y} := [y_1, \dots, y_n]' \in \mathbf{R}^n$. Note that when $G_{XX}^{(K)}$ is non-singular, Eq. (7) is the

Note that when $G_{XX}^{(K)}$ is non-singular, Eq. (7) is the unique solution of Problem 1. On the other hand, when $G_{XX}^{(K)}$ is singular, Problem 1 may have many solutions including Eq. (7). Also note that the output vector corresponding to an arbitrary set of test input vectors $Z := \{z_1, \ldots, z_m\}, (z_i \in \mathcal{D})$ is given by

$$\hat{f}_{Z}^{(K,\gamma)} := [\hat{f}^{(K,\gamma)}(z_{1}), \dots, \hat{f}^{(K,\gamma)}(z_{m})]' = G_{ZX}^{(K)} (G_{XX}^{(K)} + \gamma I_{n})^{-1} \boldsymbol{y}$$
(8)

from Eq. (7), where $G_{ZX}^{(K)} := (K(\boldsymbol{z}_i, \boldsymbol{x}_j)) \in \mathbf{R}^{m \times n}$.

3. Kernel Ridge Regression with Autocorrelation Prior and Its relation to Stochastic Linear Estimator

In this section, we review the KRR with the autocorrelation prior introduced in our previous work [5], and discuss its relation to stochastic linear estimators.

We assume that the unknown true function $f(\cdot)$ is a realization of a random process whose autocorrelation function is defined as

$$R(\boldsymbol{x}, \tilde{\boldsymbol{x}}) := E_f[f(\boldsymbol{x})f(\tilde{\boldsymbol{x}})], \quad (\boldsymbol{x}, \tilde{\boldsymbol{x}} \in \mathcal{D}),$$
(9)

where E_f denotes the expectation operator over the random process. Let $Z := \{z_1, \ldots, z_m\}$ be an arbitrary finite subset of \mathcal{D} and let $f_Z := [f(z_1), \ldots, f(z_m)]'$. Since

$$R_{ZZ} := (R(z_i, z_j)) = (E[f(z_i)f(z_j)]) = E[f_Z f'_Z]$$
(10)

is trivially a non-negative definite (n.n.d.) matrix, it is concluded that the autocorrelation function $R(x, \tilde{x}), (x, \tilde{x} \in \mathcal{D})$ is also a kernel [15], [16]. Hereafter, we call $R(x, \tilde{x})$ the autocorrelation kernel. Since R_{ZZ} is identical to the Gram matrix of the autocorrelation kernel R with an arbitrary finite subset $Z \subset \mathcal{D}$, we use the symbol $G_{ZZ}^{(R)}$ instead of R_{ZZ} in the following contents.

We assume that $f(\mathbf{x}_i)$ and ϵ_j are uncorrelated. Let $\boldsymbol{\epsilon} := [\epsilon_1, \dots, \epsilon_n]'$ and $f_X := [f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)]'$, then the correlation matrix of the training output vector $\boldsymbol{y} = \boldsymbol{f}_X + \boldsymbol{\epsilon}$ is reduced to

$$E_{f,\boldsymbol{\epsilon}}[\boldsymbol{y}\boldsymbol{y}'] = E_{f,\boldsymbol{\epsilon}}[(\boldsymbol{f}_{X} + \boldsymbol{\epsilon})(\boldsymbol{f}_{X} + \boldsymbol{\epsilon})']$$

$$= E_{f}[\boldsymbol{f}_{X}\boldsymbol{f}'_{X}] + E_{\boldsymbol{\epsilon}}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}']$$

$$= G_{XX}^{(R)} + \sigma^{2}I_{n}.$$
(11)

We evaluate the generalization error of the KRR with a model (K, γ) for an arbitrary set of test input vectors $Z \subset \mathcal{D}$ by

$$L(K,\gamma;Z) := \|\hat{f}_{Z}^{(K,\gamma)} - f_{Z}\|^{2}.$$
(12)

The following theorem is one of main results in our previous work [5].

Theorem 2: [5] Let *K* be an arbitrary kernel defined on $\mathcal{D} \times \mathcal{D}$ and let γ be an arbitrary positive constant, then

$$E_{f,\boldsymbol{\epsilon}}L(K,\gamma;Z) \ge E_{f,\boldsymbol{\epsilon}}L(R,\sigma^{2};Z)$$
(13)

holds for any $Z \subset \mathcal{D}$.

According to Theorem 2, it is concluded that the model (R, σ^2) is optimal for the KRR in terms of the expected generalization error, which agrees to the optimal model (R, σ^2) in the GPR, specified by the conditional expectation based on the training data set.

Next, let us consider the linear estimation model

$$\hat{\boldsymbol{f}}_{\boldsymbol{Z}}^{(L)} = \boldsymbol{B}\boldsymbol{y},\tag{14}$$

where $B \in \mathbf{R}^{m \times n}$, and also consider the following problem for this estimation model.

Problem 2: Find the matrix $B \in \mathbb{R}^{m \times n}$ that minimizes

$$J_2(B) := E_{f,\boldsymbol{\epsilon}} ||B\boldsymbol{y} - \boldsymbol{f}_Z||^2.$$
(15)

From Eq. (11) and

.

$$E_f[\boldsymbol{f}_Z \boldsymbol{f}_Z'] = G_{ZZ}^{(R)},$$

$$E_{f,\boldsymbol{\epsilon}}[\boldsymbol{y}\boldsymbol{f}_Z'] = E_f[\boldsymbol{f}_X \boldsymbol{f}_Z'] = G_{XZ}^{(R)}$$

where $G_{ZZ}^{(R)} := (R(z_i, z_j)) \in \mathbf{R}^{m \times m}, \ G_{XZ}^{(R)} := (R(x_i, z_j)) \in \mathbf{R}^{n \times m},$

$$J_2(B) := E_{f,\epsilon} ||By - f_Z||^2$$

= tr[B(G_{XX}^{(R)} + \sigma^2 I_n)B' - 2BG_{XZ}^{(R)} + G_{ZZ}^{(R)}]

is obtained and its first order differential is reduced to

$$dJ_2(B) = 2 \operatorname{tr}[dB(G_{XX}^{(R)} + \sigma^2 I_n)B' - dBG_{XZ}^{(R)}]$$

= 2 tr[dB((G_{XX}^{(R)} + \sigma^2 I_n)B' - G_{XZ}^{(R)})].

Since $J_2(B)$ is a non-negative quadratic function with respect to B, the stationary point of $J_2(B)$ is its minimizer. Thus, the solution of Problem 2 is immediately obtained by

$$\hat{B} = G_{ZX}^{(R)} (G_{XX}^{(R)} + \sigma^2 I_n)^{-1}$$
(16)

and the estimate of f_Z by this optimal matrix \hat{B} is written by

$$\hat{f}_{Z}^{(L)} = G_{ZX}^{(R)} (G_{XX}^{(R)} + \sigma^2 I_n)^{-1} \boldsymbol{y},$$
(17)

where $G_{ZX}^{(R)} := (G_{XZ}^{(R)})' = (R(z_i, x_j)) \in \mathbf{R}^{m \times n}$, which is identical to Eq. (8) with the optimal model (R, σ^2) . Accordingly, it is concluded that the KRR with the optimal model (R, σ^2) is identical to the stochastic linear estimator defined by Problem 2. Note that we do not assume $E_f[f_X] = \mathbf{0}_n$ and $E_f[f_{\mathcal{I}}] = \mathbf{0}_m$, where $\mathbf{0}_n \in \mathbf{R}^n$ denotes the *n*-dimensional zero vector, which implies that the Gram matrices are not always reduced to covariance matrices. When $E_f[f_x] = \mathbf{0}_n$ and $E_f[f_{z}] = \mathbf{0}_m$ hold, the above discussion, concerned with a stochastic linear estimator, is the same with that given in [6].

Stochastic Affine Estimator and Its Properties 4.

In this section, we discuss a stochastic affine estimator, which is one of the simplest extensions of a stochastic linear estimator. Note that almost all of results of this section is given in [6]. However, we adopt a different way to obtain the same results in order for theoretical analyses concerned with kernel regression problems given in the next section.

Let us consider the following affine estimation model:

$$\hat{f}_{Z}^{(A)} = B \begin{bmatrix} \mathbf{y} \\ 1 \end{bmatrix}, \tag{18}$$

where $B \in \mathbf{R}^{m \times (n+1)}$. As the same with the linear case, we

consider the following problem.

Problem 3: Find the matrix $B \in \mathbb{R}^{m \times (n+1)}$ that minimizes

$$J_{3}(B) := E_{f,\boldsymbol{\epsilon}} \left\| B \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix} - \boldsymbol{f}_{Z} \right\|^{2}.$$
(19)

Let $\mu_X := E_f[f_X], \mu_Z := E_f[f_Z], \text{ and } M := G_{XX}^{(R)} + \sigma^2 I_n,$ then

$$E_{f,\boldsymbol{\epsilon}}\left[\left[\begin{array}{c}\boldsymbol{y}\\1\end{array}\right]\left[\begin{array}{c}\boldsymbol{y}'&1\end{array}\right]\right] = \left[\begin{array}{c}\boldsymbol{M} & \boldsymbol{\mu}_{X}\\\boldsymbol{\mu}_{X}'&1\end{array}\right] =: M_{A}$$

and

$$E_{f,\boldsymbol{\epsilon}} \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix} \boldsymbol{f}'_{Z} = \begin{bmatrix} G_{XZ}^{(R)} \\ \boldsymbol{\mu}'_{Z} \end{bmatrix} =: Q$$

hold. Therefore

$$J_{3}(B) := E_{f,\epsilon} \left\| B \begin{bmatrix} \mathbf{y} \\ 1 \end{bmatrix} - \mathbf{f}_{Z} \right\|^{2}$$
$$= \operatorname{tr} \left[BM_{A}B' - 2BQ + G_{ZZ}^{(R)} \right]$$

is obtained and its first order differential is reduced to

$$dJ_3(B) = 2 \operatorname{tr} \left[dBM_AB' - dBQ \right]$$

= 2 tr [dB (M_AB' - Q))].

Theorem 3: The matrix M_A is non-singular.

Proof The covariance matrix Σ_X of f_X is represented by

$$\Sigma_{X} := E_{f}[(f_{X} - \mu_{X})(f_{X} - \mu_{X})']$$

= $E_{f}[f_{X}f'_{X}] - \mu_{X}\mu'_{X}$
= $G_{XX}^{(R)} - \mu_{X}\mu'_{X}.$ (20)

Therefore.

$$M = \Sigma_X + \sigma^2 I_n + \mu_X \mu'_X$$

and

$$M_A = T_1 + T_2$$

are obtained with

$$T_1 := \begin{bmatrix} \Sigma_X + \sigma^2 I_n & \mathbf{0}_n \\ \mathbf{0}'_n & \mathbf{0} \end{bmatrix}, \quad T_2 := \begin{bmatrix} \boldsymbol{\mu}_X \\ 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\mu}'_X & 1 \end{bmatrix}.$$

Let $[\mathbf{u}' \ v]' \in \mathbf{R}^{n+1}$ be an arbitrary vector in the null space of M_A , written as $\mathcal{N}(M_A)$, then

$$\begin{bmatrix} \boldsymbol{u}' & \boldsymbol{v} \end{bmatrix} T_1 \begin{bmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{bmatrix} = 0, \tag{21}$$

$$\begin{bmatrix} \boldsymbol{u}' & \boldsymbol{v} \end{bmatrix} T_2 \begin{bmatrix} \boldsymbol{u} \\ \boldsymbol{v} \end{bmatrix} = 0 \tag{22}$$

must hold since T_1 and T_2 are *n.n.d.* Therefore, Eq. (21) immediately yields $\boldsymbol{u} = \boldsymbol{0}_n$ since $\Sigma_X + \sigma^2 I_n$ is non-singular, and then Eq. (22) yields v = 0, which implies $\mathcal{N}(M_A) =$ $\{\mathbf{0}_{n+1}\}.$ Since $J_3(B)$ is a non-negative quadratic function with respect to *B*, the stationary point of $J_3(B)$ is the minimizer of $J_3(B)$. Thus, the solution of Problem 3 is immediately obtained by

$$\hat{B} = Q' M_A^{-1}. \tag{23}$$

By substituting Eq. (23) into Eq. (18), the estimate of f_Z by the solution of Problem 3 is reduced to

$$\hat{\boldsymbol{f}}_{Z}^{(A)} = \boldsymbol{Q}' \boldsymbol{M}_{A}^{-1} \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix}$$
$$= \begin{bmatrix} \boldsymbol{G}_{ZX}^{(R)} & \boldsymbol{\mu}_{Z} \end{bmatrix}$$
$$\times \begin{bmatrix} \boldsymbol{G}_{XX}^{(R)} + \sigma^{2} \boldsymbol{I}_{n} & \boldsymbol{\mu}_{X} \\ \boldsymbol{\mu}_{X}' & 1 \end{bmatrix}^{-1} \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix}.$$
(24)

5. Kernel Regressors Equivalent to Stochastic Affine Estimators

In this section, we discuss the kernel regression problem, whose solution is reduced to the stochastic affine estimator $\hat{f}_Z^{(A)}$ in Eq. (24).

The function with the model (K, γ) that gives Eq. (24) by substituting the points in *Z* can be represented by

$$\hat{f}^{(K,\gamma)}(\cdot) = \begin{bmatrix} (\boldsymbol{g}_{X}^{(K)}(\cdot))' & \mu(\cdot) \end{bmatrix} \times \begin{bmatrix} G_{XX}^{(K)} + \gamma I_{n} & \boldsymbol{\mu}_{X} \\ \boldsymbol{\mu}_{X}^{(K)} & 1 \end{bmatrix}^{-1} \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix}, \quad (25)$$

where $\mu(\cdot)$ is a certain function, that yields $\mu_X := [\mu(\mathbf{x}_1) \dots, \mu(\mathbf{x}_n)]'$, corresponding to the mean vector in the stochastic affine estimator. Thus, we consider the following function model:

$$\hat{f}(\cdot) := (\boldsymbol{g}_{\boldsymbol{X}}^{(K)}(\cdot))'\boldsymbol{\alpha} + \boldsymbol{\beta}\boldsymbol{\mu}(\cdot), \tag{26}$$

where $\alpha \in \mathbf{R}^n$ and $\beta \in \mathbf{R}$.

Let $K_{\mu}(\mathbf{x}, \tilde{\mathbf{x}}) := \mu(\mathbf{x})\mu(\tilde{\mathbf{x}}), (\mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{D})$, then $K_{\mu}(\mathbf{x}, \tilde{\mathbf{x}})$ is the kernel whose corresponding RKHS is consisting of $a\mu(\cdot), (a \in \mathbf{R})$ as shown in [13], which trivially implies $\mu(\cdot) \in \mathcal{H}_{K_{u}}$. Note that since

$$\mu(\mathbf{x}) = \langle \mu(\cdot), K_{\mu}(\cdot, \mathbf{x}) \rangle_{\mathcal{H}_{K_{\mu}}} \\ = \langle \mu(\cdot), \mu(\cdot)\mu(\mathbf{x}) \rangle_{\mathcal{H}_{K_{\mu}}} \\ = \mu(\mathbf{x}) \langle \mu(\cdot), \mu(\cdot) \rangle_{\mathcal{H}_{K_{\mu}}} \\ = \mu(\mathbf{x}) ||\mu(\cdot)||_{\mathcal{H}_{K_{\mu}}}^{2}$$

holds for any $\mathbf{x} \in \mathcal{D}$, we have $\|\mu(\cdot)\|_{\mathcal{H}_{K_u}}^2 = 1$.

Theorem 4: [4] If K_i is the reproducing kernel of the class F_i with the norm $\|\cdot\|_i$, then $K = K_1 + K_2$ is the reproducing kernel of the class F of all functions $f(\cdot) = f_1(\cdot) + f_2(\cdot)$ with $f_i(\cdot) \in F_i$, and with the norm defined by

$$\|f(\cdot)\|^{2} = \min\left[\|f_{1}(\cdot)\|_{1}^{2} + \|f_{2}(\cdot)\|_{2}^{2}\right],$$
(27)

the minimum taken for all the decompositions $f(\cdot) = f_1(\cdot) + f_2(\cdot)$ with $f_i(\cdot) \in F_i$.

Hereafter, we assume that the kernel, used in the following contents, is of the form:

$$K(\boldsymbol{x}, \tilde{\boldsymbol{x}}) := K_c(\boldsymbol{x}, \tilde{\boldsymbol{x}}) + K_{\mu}(\boldsymbol{x}, \tilde{\boldsymbol{x}}), \qquad (28)$$

where K_c is a certain kernel, whose Gram matrix with *X* corresponds to a covariance matrix Σ_X in the stochastic affine estimator as in Eq. (20). We call K_c the covariance part of *K* and also call K_{μ} the mean part of *K*. It should be noted that $\mu(\cdot) \in \mathcal{H}_K$ holds due to Theorem 4 and the fact that $\mu(\cdot) \in \mathcal{H}_{K_{\mu}}$. We assume $\mu(\cdot) \notin \mathcal{H}_{K_c}^{\dagger}$, which implies $\mathcal{H}_{K_c} \cap \mathcal{H}_{K_{\mu}} = \{0\}$, and then $||\mu(\cdot)||^2_{\mathcal{H}_K} = 1$ is obtained, since the decomposition in Theorem 4 is unique. We also assume that $\mu(\cdot) \notin \mathcal{L}_X^{(K)} := \operatorname{span}\{K(\cdot, \mathbf{x}_i) \mid i \in \{1, \ldots, n\}\}$, since when $\mu(\cdot) \in \mathcal{L}_X^{(K)}$ holds, Eq. (26) is reduced to Eq. (6). Note that the Gram matrix of *K* with *X* can be represented as

$$G_{XX}^{(K)} = G_{XX}^{(K_c)} + G_{XX}^{(K_{\mu})} = G_{XX}^{(K_c)} + \mu_X \mu'_X.$$

Here, we define some subsets of \mathcal{H}_K . Let $\mathcal{L}_{\mu} := \operatorname{span}\{\mu(\cdot)\}$ and let $\mathcal{L}_{\mu}^{\perp}$ be its orthogonal complement, then, we define the linear manifold

$$\mathcal{M}_{\mu} := \{ \mu^{\perp}(\cdot) + \mu(\cdot) \mid \mu^{\perp}(\cdot) \in \mathcal{L}_{\mu}^{\perp} \} \subset \mathcal{H}_{K}.$$

Under these preparations, we define the following problem.

Problem 4: Find the function $\hat{f}(\cdot) \in \mathcal{M}_{\mu}$ that minimizes

$$J_4(\hat{f}(\cdot)) := \sum_{i=1}^n (y_i - \hat{f}(\boldsymbol{x}_i))^2 + \gamma \|\hat{f}(\cdot)\|_{\mathcal{H}_K}^2$$
(29)

with a positive constant γ .

It should be noted that the constraint

$$\hat{c}(\cdot) \in \mathcal{M}_{\mu} \tag{30}$$

can be represented as

$$\mu(\cdot) = P_{\mu}\hat{f}(\cdot),\tag{31}$$

where P_{μ} denotes the orthogonal projector onto \mathcal{L}_{μ} , since $\hat{f}(\cdot) \in \mathcal{M}_{\mu}$ can be represented as

$$\hat{f}(\cdot) = f^{\perp}(\cdot) + \mu(\cdot)$$

with $f^{\perp}(\cdot) \in \mathcal{L}^{\perp}_{\mu}$ and

$$P_{\mu}\hat{f}(\cdot) = P_{\mu}f^{\perp}(\cdot) + P_{\mu}\mu(\cdot) = 0 + P_{\mu}\mu(\cdot) = \mu(\cdot)$$

is immediately followed.

The first issue to be resolved is whether the solution of Problem 4 can be represented by Eq. (26) or not. In [14], the representer theorem concerned with the function model Eq. (26), which is called the semiparametric representer theorem, is given as follows.

Theorem 5: [14] In addition to the assumptions of Theorem 1, we assume that a set of n_P real-valued function

[†]This is the only arbitrary assumption that the framework of the stochastic affine estimators has no counterpart.

 $\{\psi_p(\cdot)\}, (p \in \{1, \ldots, n_P\})$ defined on \mathcal{D} with the property that the matrix $(\psi_p(\mathbf{x}_i)) \in \mathbf{R}^{n_P \times n}$ has rank n_P is given. Then, any $\hat{f}(\cdot) := f(\cdot) + h(\cdot)$ with $f(\cdot) \in \mathcal{H}_K$ and $h(\cdot) \in \text{span}\{\psi_1(\cdot), \ldots, \psi_{n_P}(\cdot)\}$, minimizing the regularized risk functional

$$c((\mathbf{x}_1, y_1, \hat{f}(\mathbf{x}_1)), \dots, (\mathbf{x}_n, y_n, \hat{f}(\mathbf{x}_n))) + g(||f(\cdot)||_{\mathcal{H}_K})$$
 (32)

can be represented by

$$\hat{f}(\cdot) := \sum_{i=1}^{n} \alpha_i K(\cdot, \mathbf{x}_i) + \sum_{p=1}^{n_p} \beta_p \psi_p(\cdot)$$
(33)

with unique coefficients $\beta_p \in \mathbf{R}$ for all $p \in \{1, \ldots, n_P\}$.

The function model Eq. (26) agrees to Eq. (33) with $n_P = 1$. However, the optimization criterion and the constraint in Problem 4 does not agree to Eq. (32), which implies that Theorem 5 can not be applied to problem 4. In order to positively resolve this issue, we give the following theorem.

Theorem 6: Any function $\hat{f}(\cdot) \in \mathcal{M}_{\mu}$ that minimizes Eq. (29) with a positive γ can be represented by Eq. (26).

Proof Let $v(\cdot) \in \mathcal{H}_K$ be an arbitrary function orthogonal to $\mathcal{L}_{\mu} + \mathcal{L}_X^{(K)}$. Without loss of generality, any function $\hat{f}(\cdot) \in \mathcal{H}_K$ can be represented by

$$\hat{f}(\cdot) = (\boldsymbol{g}_X^{(K)}(\cdot))'\boldsymbol{\alpha} + \beta\boldsymbol{\mu}(\cdot) + \boldsymbol{v}(\cdot)$$

Since $v(\cdot) \in \left(\mathcal{L}_{\mu} + \mathcal{L}_{X}^{(K)}\right)^{\perp}$, we have $\langle \mu(\cdot), v(\cdot) \rangle_{\mathcal{H}_{K}} = 0$ and $\langle v(\cdot), K(\cdot, \mathbf{x}_{i}) \rangle_{\mathcal{H}_{K}} = 0$ for any $i \in \{1, \ldots, n\}$. The former leads

$$P_{\mu}\hat{f}(\cdot) = P_{\mu}((\boldsymbol{g}_{X}^{(K)}(\cdot))'\alpha + \beta\mu(\cdot)),$$

which implies that the constraint Eq. (30) is independent of $v(\cdot)$. Similarly, the latter leads

$$\begin{split} \hat{f}(\boldsymbol{x}_i) &= \langle \hat{f}(\cdot), K(\cdot, \boldsymbol{x}_i) \rangle_{\mathcal{H}_K} \\ &= \langle (\boldsymbol{g}_X^{(K)}(\cdot))' \boldsymbol{\alpha} + \beta \boldsymbol{\mu}(\cdot) + \boldsymbol{v}(\cdot), K(\cdot, \boldsymbol{x}_i) \rangle_{\mathcal{H}_K} \\ &= \langle (\boldsymbol{g}_X^{(K)}(\cdot))' \boldsymbol{\alpha} + \beta \boldsymbol{\mu}(\cdot), K(\cdot, \boldsymbol{x}_i) \rangle_{\mathcal{H}_K}, \end{split}$$

which implies that the first term of Eq. (29) is also independent of $v(\cdot)$. The second term of Eq. (29) satisfies

$$\begin{split} \gamma \|\hat{f}(\cdot)\|_{\mathcal{H}_{K}}^{2} &= \gamma \langle \hat{f}(\cdot), \hat{f}(\cdot) \rangle_{\mathcal{H}_{K}} \\ &= \gamma \langle (\boldsymbol{g}_{X}^{(K)}(\cdot))' \alpha + \beta \mu(\cdot) + v(\cdot), \\ & (\boldsymbol{g}_{X}^{(K)}(\cdot))' \alpha + \beta \mu(\cdot) + v(\cdot) \rangle_{\mathcal{H}_{K}} \\ &= \gamma \|(\boldsymbol{g}_{X}^{(K)}(\cdot))' \alpha + \beta \mu(\cdot)\|_{\mathcal{H}_{K}}^{2} + \gamma \|v(\cdot)\|_{\mathcal{H}_{\mu}}^{2} \\ &\geq \gamma \|(\boldsymbol{g}_{X}^{(K)}(\cdot))' \alpha + \beta \mu(\cdot)\|_{\mathcal{H}_{K}}^{2} \end{split}$$

since $\gamma > 0$, and equality is attained if and only if $v(\cdot) = 0$. Accordingly, in order for $\hat{f}(\cdot)$ to be a minimizer of Eq. (29), $v(\cdot) = 0$ must hold, which concludes the proof.

According to Theorem 6, it is enough to consider the function model Eq. (26) as a candidate of the solution of Problem 4.

Since $\hat{f}(\cdot)$ in Problem 4 is specified by α and β only, we use the notation $J_4(\alpha,\beta)$ as the criterion of Problem 4 instead of $J_4(\hat{f}(\cdot))$. Let $S_i(\alpha,\beta)$, $(i \in \{1,2\})$ be the *i*-th term of Eq. (29). It is easy to show that

$$S_1(\boldsymbol{\alpha},\boldsymbol{\beta}) = \left\| \boldsymbol{y} - \begin{bmatrix} G_{XX}^{(K)} & \boldsymbol{\mu}_X \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2.$$

Also, we have

$$S_{2}(\alpha,\beta) = \gamma \|\hat{f}(\cdot)\|_{\mathcal{H}_{K}}^{2}$$

$$= \gamma \|(\boldsymbol{g}_{X}^{(K)}(\cdot))'\alpha + \beta \mu(\cdot)\|_{\mathcal{H}_{K}}^{2}$$

$$= \gamma (\alpha' G_{XX}^{(K)}\alpha + 2\beta \alpha' \boldsymbol{\mu}_{X} + \beta^{2})$$

$$= \gamma \left[\begin{array}{cc} \alpha' & \beta \end{array} \right] \left[\begin{array}{cc} G_{XX}^{(K)} & \boldsymbol{\mu}_{X} \\ \boldsymbol{\mu}_{X}' & 1 \end{array} \right] \left[\begin{array}{cc} \alpha \\ \beta \end{array} \right]$$

for $\hat{f}(\cdot)$ specified by the function model Eq. (26).

As mentioned before, the linear manifold constraint Eq. (30) is identical to Eq. (31), which can be also represented as

$$\mu(\cdot) = P_{\mu}f(\cdot)$$

$$= \langle \hat{f}(\cdot), \mu(\cdot) \rangle_{\mathcal{H}_{K}}\mu(\cdot)$$

$$= \langle (\boldsymbol{g}_{X}^{(K)}(\cdot))'\boldsymbol{\alpha} + \beta\mu(\cdot), \mu(\cdot) \rangle_{\mathcal{H}_{K}}\mu(\cdot)$$

$$= (\boldsymbol{\mu}_{X}'\boldsymbol{\alpha} + \beta)\,\mu(\cdot). \tag{34}$$

Since $\|\mu(\cdot)\|_{\mathcal{H}_{\nu}}^2 = 1 \neq 0$ holds, Eq. (34) is identical to

$$1 = \mu'_X \alpha + \beta, \tag{35}$$

which can be also represented as

$$1 = \left[\begin{array}{c} \mu_X' & 1 \end{array} \right] \left[\begin{array}{c} \alpha \\ \beta \end{array} \right] \tag{36}$$

and

$$\begin{bmatrix} \boldsymbol{\mu}_X \\ 1 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_X \\ 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\mu}'_X & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix}, \quad (37)$$

since $[\mu'_X 1]'$ is an $(n + 1) \times 1$ full column rank matrix.

The following theorem is the main result of this paper.

Theorem 7: Eq. (25) is a solution of Problem 4.

Proof Let

$$\mathscr{L}(\alpha,\beta,\lambda) = \sum_{i=1}^{2} S_i(\alpha,\beta) - 2\lambda(\mu'_X\alpha + \beta - 1)$$
(38)

be the Lagrangian function induced from Problem 4 with a Lagrange multiplier λ . Since

$$\begin{aligned} \mathscr{L}(\boldsymbol{\alpha},\boldsymbol{\beta},\boldsymbol{\lambda}) \\ &= \|\boldsymbol{y}\|^2 + \begin{bmatrix} \boldsymbol{\alpha}' & \boldsymbol{\beta} \end{bmatrix} \begin{bmatrix} \boldsymbol{G}_{XX}^{(K)} \\ \boldsymbol{\mu}'_X \end{bmatrix} \begin{bmatrix} \boldsymbol{G}_{XX}^{(K)} & \boldsymbol{\mu}_X \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} \\ &-2\begin{bmatrix} \boldsymbol{\alpha}' & \boldsymbol{\beta} \end{bmatrix} \begin{bmatrix} \boldsymbol{G}_{XX}^{(K)} \\ \boldsymbol{\mu}'_X \end{bmatrix} \boldsymbol{y} \end{aligned}$$

$$+\gamma \begin{bmatrix} \alpha' & \beta \end{bmatrix} \begin{bmatrix} G_{XX}^{(K)} & \mu_X \\ \mu'_X & 1 \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}$$
$$-2\lambda \left(\begin{bmatrix} \alpha' & \beta \end{bmatrix} \begin{bmatrix} \mu_X \\ 1 \end{bmatrix} - 1 \right)$$

holds, its first order differential with respect to $[\alpha' \beta]'$ is reduced to

$$d\mathscr{L}(\alpha,\beta,\lambda) = 2d\left[\begin{array}{cc} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] y \\ +2\gamma d\left[\begin{array}{c} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] y \\ +2\gamma d\left[\begin{array}{c} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} \alpha \\ \beta\end{array}\right] \\ -2\lambda d\left[\begin{array}{c} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} \alpha \\ \beta\end{array}\right] \\ = 2d\left[\begin{array}{c} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} \alpha \\ \beta\end{array}\right] \\ +2\gamma d\left[\begin{array}{c} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} \alpha \\ \beta\end{array}\right] \\ -2d\left[\begin{array}{c} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} \alpha \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} \alpha \\ \beta\end{array}\right] \\ -2d\left[\begin{array}{c} \alpha' & \beta\end{array}\right] \left[\begin{array}{c} G_{XX}^{(K)} \\ \mu'_{X} \end{array}\right] \left[\begin{array}{c} y \\ \mu'_{X} \end{array}\right] .$$
(39)

Since $J_4(\alpha,\beta)$ is a non-negative quadratic function with respect to $[\alpha'\beta]'$, Eq. (38) is also a lower-bounded quadratic function with respect to $[\alpha'\beta]'$ under the constraint Eq. (30). Thus, stationary points of Eq. (38) is its minimizer. Therefore,

$$\begin{pmatrix} \begin{bmatrix} G_{XX}^{(K)} \\ \mu_X' \end{bmatrix} \begin{bmatrix} G_{XX}^{(K)} & \mu_X \end{bmatrix} + \gamma \begin{bmatrix} G_{XX}^{(K)} & \mu_X \\ \mu_X' & 1 \end{bmatrix} \begin{pmatrix} \alpha \\ \beta \end{bmatrix} - \begin{bmatrix} G_{XX}^{(K)} & \mu_X \\ \mu_X' & 1 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \lambda \end{bmatrix} = \mathbf{0}_{n+1}$$
(40)

must hold, under the constraint Eq. (30). Note that Eq. (30) can be represented by

$$\begin{bmatrix} \boldsymbol{\mu}_{X} \\ 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\mu}'_{X} & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} - \begin{bmatrix} \boldsymbol{\mu}_{X} \\ 1 \end{bmatrix} = \boldsymbol{0}_{n+1}$$
(41)

from Eq. (37). Since Eqs. (40) and (41) must hold simultaneously, the sum of their left sides must be zero, which is represented by

$$\begin{pmatrix} \begin{bmatrix} G_{XX}^{(K)} & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_X^{\prime} & 1 \end{bmatrix}^2 + \gamma \begin{bmatrix} G_{XX}^{(K)} & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_X^{\prime} & 1 \end{bmatrix} \end{pmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} - \begin{bmatrix} G_{XX}^{(K)} & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_X^{\prime} & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{y} \\ \lambda + 1 \end{bmatrix} = \boldsymbol{0}_{n+1},$$
 (42)

that can be also represented as

$$\begin{bmatrix} G_{XX}^{(K)} & \boldsymbol{\mu}_{X} \\ \boldsymbol{\mu}_{X}^{\prime} & 1 \end{bmatrix} \begin{bmatrix} G_{XX}^{(K)} + \gamma I_{n} & \boldsymbol{\mu}_{X} \\ \boldsymbol{\mu}_{X}^{\prime} & 1 + \gamma \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix}$$
$$-\begin{bmatrix} G_{XX}^{(K)} & \boldsymbol{\mu}_{X} \\ \boldsymbol{\mu}_{X}^{\prime} & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{y} \\ \boldsymbol{\lambda} + 1 \end{bmatrix} = \boldsymbol{0}_{n+1}.$$
(43)

Therefore, the solution of the linear equation

$$\begin{bmatrix} G_{XX}^{(K)} + \gamma I_n & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_X' & 1 + \gamma \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} - \begin{bmatrix} \boldsymbol{y} \\ \lambda + 1 \end{bmatrix} = \boldsymbol{0}_{n+1} \quad (44)$$

is a stationary point of Eq. (38). Since the last row of Eq. (44) is reduced to

$$\begin{aligned} \mu'_X \alpha + (1+\gamma)\beta - (\lambda+1) \\ &= 1+\gamma\beta - (\lambda+1) = \gamma\beta - \lambda = 0 \end{aligned}$$

due to the constraint Eq. (35), $\lambda = \gamma\beta$ is immediately followed. Accordingly, the second term of the left side of Eq. (44) can be represented by

$$\begin{bmatrix} \boldsymbol{y} \\ \gamma\beta+1 \end{bmatrix} = \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix} + \begin{bmatrix} O_n & \boldsymbol{0}_n \\ \boldsymbol{0}'_n & \gamma \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix}, \quad (45)$$

where O_n stands for the $n \times n$ zero matrix, and substituting it to Eq. (44) yields

$$\begin{bmatrix} G_{XX}^{(K)} + \gamma I_n & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_X' & 1 \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} - \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix} = \boldsymbol{0}_{n+1}.$$
(46)

It is easy to show that the matrix

$$M_B := \begin{bmatrix} G_{XX}^{(K)} + \gamma I_n & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_X' & 1 \end{bmatrix}$$

is non-singular by Theorem 3. Thus,

$$\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} G_{XX}^{(K)} + \gamma I_n & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}'_X & 1 \end{bmatrix}^{-1} \begin{bmatrix} \boldsymbol{y} \\ 1 \end{bmatrix}$$
(47)

surely exists and it is a minimizer of Eq. (38) under the constraint Eq. (30) and then Eq. (25) is a solution of Problem 4.

According to Theorem 7, it is concluded that the model (R, σ^2) is also optimal for Problem 4 with a given function $\mu(\cdot)$ in terms of the expected squared error under the autocorrelation prior, as the same with the ordinary KRR. Note that when the matrix

$$M_C := \left[\begin{array}{cc} G_{XX}^{(K)} & \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_X' & 1 \end{array} \right]$$

is non-singular, Eq. (47) gives the unique solution of Problem 4. On the other hand, when M_C is singular, problem 4 may have many solutions including the learning result based on Eq. (47) as the same with the ordinary KRR.

Here, we give some remarks on the optimization criterion of Problem 4. As the same with the ordinary KRR, $S_1(\alpha,\beta)$ in $J_4(\alpha,\beta)$ represents the fidelity of the estimated function to the given training data set. However, the minimizer of $S_1(\alpha,\beta)$ has remaining degree of freedom caused by the singularity of the coefficient matrix in the normal equation obtained from the first order differential of $S_1(\alpha,\beta)$. The solution Eq. (47) is obtained by vanishing the degree of freedom by the linear manifold constraint Eq. (30), instead of the minimum norm constraint, which is usually adopted in solving an underdetermined linear system. It seems that $S_2(\alpha,\beta)$ in $J_4(\alpha,\beta)$ plays the same role of the regularization term in Problem 1. However, it should be noted that the effect of the regularization parameter is restricted to the Gram matrix concerned with the term $(g_X^{(K)}(\cdot))'\alpha$ in the function model Eq. (26) due to the linear manifold constraint. Therefore, it is concluded that the linear manifold constraint $\hat{f}(\cdot) \in \mathcal{M}_{\mu}$ is crucial in obtaining the kernel-based regressor Eq. (25)

It is expected that replacing $S_2(\alpha,\beta)$ in $J_4(\alpha,\beta)$ to another alternatives yields various kernel regressors represented as affine estimators together with the linear manifold constraint, and our formulation may be useful to construct a kernel-based regression problem with multiple given functions instead of a single $\mu(\cdot)$, while these extensions can not be obtained in the framework of the GPR due to its principle.

6. Conclusion

In this paper, we discussed the kernel-based regression problems by considering the relation to stochastic estimators, and obtained the formulation, which is equivalent to stochastic affine estimators, and gave the interpretation of it. Moreover, we gave a novel representer theorem concerned with the formulation. It is expected that our result give a novel aspect for further extensions of kernel-based regression problems.

Acknowledgments

This work was partially supported by JSPS KAKENHI Grant Number JP20H04238.

References

- J. Shawe-Taylor and N. Cristianini, Kernel Methods for Pattern Analysis, Cambridge University Press, Cambridge, 2004.
- [2] N. Cristianini and J. Shawe-Taylor, An Introduction to Support Vector Machines and Other Kernel-Based Learning Methods, Cambridge University Press, Cambridge, 2000.
- [3] K.-R. Muller, S. Mika, G. Ratsch, K. Tsuda, and B. Scholkopf, "An introduction to kernel-based learning algorithms," IEEE Trans. Neural Netw., vol.12, no.2, pp.181–201, 2001.
- [4] N. Aronszajn, "Theory of reproducing kernels," Trans. American Mathematical Society, vol.68, no.3, pp.337–404, 1950.
- [5] A. Tanaka and H. Imai, "Kernel ridge regression with autocorrelation prior: Optimal model and cross-validation," 2020 IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP2020), pp.3867–3871, 2020.
- [6] T. Kailath, A.H. Sayed, and B. Hassibi, Linear Estimation, Prentice-Hall, Upper Saddle River, NJ, 2000.
- [7] C.E. Rasmussen, "Gaussian processes in machine learning," Summer School on Machine Learning, Lecture Notes in Computer Science, vol.3176, pp.63–71, Springer 2004.
- [8] C.E. Rasmussen and C.K.I. Williams, Gaussian processes for machine learning, MIT Press, Cambridge, MA, 2006.
- [9] K.P. Murphy, Machine learning: A probabilistic perspective, MIT Press, 2012.
- [10] M. Kanagawa, P. Hennig, D. Sejdinovic, and B.K. Sriperumbudur, "Gaussian process and kernel methods: A review on connections and equivalences," arXiv:1807.02582 [stat.ML], 2018.

- [11] A. Tanaka, H. Imai, and M. Miyakoshi, "Kernel-induced sampling theorem," IEEE Trans. Signal Process., vol.58, no.7, pp.3569–3577, 2010.
- [12] J. Mercer, "Functions of positive and negative type, and their connection with the theory of integral equations," Philosophical Trans. Royal Society A, vol.209, no.441-458, pp.415–446, 1909.
- [13] S. Saitoh, Integral Transforms, Reproducing Kernels and Their Applications, Addison Wesley Longman Ltd, UK, 1997.
- [14] B. Schölkopf, R. Herbrich, and A. Smola, "A generalized representer theorem," Computational Learning Theory, Lecture Notes in Computer Science, vol.2111, pp.416–426, Springer, 2001.
- [15] R. Kong and B. Zhang, "Autocorrelation kernel functions for support vector machines," The 3rd International Conference on Natural Computation, 2007.
- [16] Y. Saito and A. Tanaka, "Optimal kernel in kernel regression problems with autocorrelation prior," 2017 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC), pp.969–972, IEEE, 2017.



Akira Tanaka received the D.E. from Hokkaido University in 2000. He joined the Faculty of Information Science and Technology, Hokkaido University. His research interests include image processing, acoustic signal processing, and machine learning theory.



Masanari Nakamura received the D.S. from Hokkaido University in 2018. He joined the Faculty of Information Science and Technology, Hokkaido University. His research interests include acoustic signal processing, positioning, and target tracking.



Hideyuki Imai received the D.E. from Hokkaido University in 1999. He joined the Graduate School of Information Science and Technology, Hokkaido University. His research interests include statistical inference.