

On Predicting with Kernel Ridge Regression¹⁾

Changha Hwang²⁾

Abstract

Kernel machines are used widely in real-world regression tasks. Kernel ridge regressions(KRR) and support vector machines(SVM) are typical kernel machines. Here, we focus on two types of KRR. One is inductive KRR. The other is transductive KRR. In this paper, we study how differently they work in the interpolation and extrapolation areas. Furthermore, we study prediction interval estimation method for KRR. This turns out to be a reliable and practical measure of prediction interval and is essential in real-world tasks.

Keywords : Kernel ridge regression, prediction interval, transductive inference.

1. Introduction

In the case of noisy learning data, the use of traditional neural networks due to its learning method often leads to poor generalization and overfitting. Kernel machines such as support vector machine(SVM) and kernel ridge regression(KRR) were designed to overcome these problems. Foundations of SVM and KRR were established by Vapnik(1995) and Saunders et al. (1998), respectively. Kernel machines are used widely in real-world regression tasks. Here, we focus on two types of KRRs. One is inductive KRR. The other is transductive KRR. In this paper, we study how differently they work in terms of interpolation and extrapolation. Furthermore, we study prediction interval estimation method for

1) This research was supported by the Catholic University of Daegu research grants in 2003.

2) Associate Professor, Dept. of Statistical Information, Catholic University of Daegu.
E-mail: chhwang@cataegu.ac.kr

KRR. This turns out to be a reliable and practical measure of confidence bound and is essential in real-world tasks.

To review KRRs, we need to take lots of materials from Chapelle et al.(1999). Suppose there exists a function $y^* = f_0(\mathbf{x})$ from which we observe the measurements corrupted with noise

$$\{(\mathbf{x}_i, y_i), i = 1, \dots, l\}, y_i = y_i^* + \varepsilon_i. \quad (1)$$

Find an algorithm A that using both the given set of training data (1) and the given set of test data

$$\{\mathbf{x}_{l+1}, \dots, \mathbf{x}_{l+m}\} \quad (2)$$

selects from a set of functions $\{\mathbf{x} \mapsto f(\mathbf{x})\}$ a function

$$y = f(\mathbf{x}) = f_A(\mathbf{x} | \mathbf{x}_1, y_1, \dots, \mathbf{x}_l, y_l, \mathbf{x}_{l+1}, \dots, \mathbf{x}_{l+m}) \quad (3)$$

and minimizes at the points of interest the functional

$$R(A) = E\left(\sum_{i=l+1}^{l+m} (y_i^* - f_A(\mathbf{x} | \mathbf{x}_1, y_1, \dots, \mathbf{x}_l, y_l, \mathbf{x}_{l+1}, \dots, \mathbf{x}_{l+m}))^2\right) \quad (4)$$

where expectation is taken over \mathbf{x} and ε . For the training data we are given the vectors \mathbf{x} and the value y , for the test data we are only given \mathbf{x} .

Usually, the problem of estimating values of a function at points of interest is solved in two steps: first in a given set of functions $f(\mathbf{x}, \boldsymbol{\alpha})$, $\boldsymbol{\alpha} \in \Lambda$ one estimates the regression function which minimizes the functional

$$R(\boldsymbol{\alpha}) = \int ((y - f(\mathbf{x}, \boldsymbol{\alpha}))^2 dF(\mathbf{x}, y)) \quad (5)$$

(the **inductive step**) and then using the estimated function $y = f(\mathbf{x}, \boldsymbol{\alpha}_l)$ we calculate the values at points of interest

$$y_i^* = f(\mathbf{x}_i^*, \boldsymbol{\alpha}_l) \quad (6)$$

(the **deductive step**).

2. Kernel Ridge Regression and the Leave-One-Out procedure

For the discussion of the classical two-step (**inductive** plus **deductive**) KRR we consider the set of functions linear in their parameters

$$f(\mathbf{x}, \boldsymbol{\alpha}_l) = \sum_{i=1}^H \alpha_i \phi_i(\mathbf{x}). \quad (7)$$

In the case of KRR we use a kernel function $K(\mathbf{x}_i, \mathbf{x})$ as basis function $\phi_i(\mathbf{x})$ and H equals to the sample size l . Notice that the first component of input vector \mathbf{x} is 1 and the rest of components are actual input variables. To

minimize the expected loss (5), we minimize the following empirical functional

$$R_{\text{emp}}(\mathbf{a}) = \frac{1}{l} \sum_{i=1}^l (y_i - f(\mathbf{x}_i, \mathbf{a}))^2 + \gamma \|\mathbf{a}\|^2 \quad (8)$$

where γ is a fixed positive constant, called the regularization parameter. The minimum is given by the vector of coefficients

$$\mathbf{a}_l = \mathbf{a}(\mathbf{x}_1, y_1, \dots, \mathbf{x}_l, y_l) = (K^t K + \gamma I_l)^{-1} K^t \mathbf{y} \quad (9)$$

where $\mathbf{y} = (y_1, \dots, y_l)^t$ and K is a matrix with elements,

$$K_{ij} = K(\mathbf{x}_i, \mathbf{x}_j), \quad i=1, \dots, l, \quad j=1, \dots, l. \quad (10)$$

Notice that the vector of coefficients in Saunders(1998) is given by

$$\mathbf{a}_l = 2\gamma(K + \gamma I_l)^{-1} \mathbf{y}. \quad (11)$$

These are a little different. Now, the problem is to choose the value γ which provides small expected loss for training on a sample $S_l = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)\}$. For this purpose, we would like to choose γ such that f_γ minimizing (8) also minimizes

$$R = \int (y^* - f_\gamma(\mathbf{x}^* | S_l))^2 dF(\mathbf{x}^*, y^*) dF(S_l) \quad (12)$$

Since $F(\mathbf{x}, y)$ is unknown, we cannot estimate this minimum directly. To solve this problem we use the leave-one-out(LOO) procedure, which is an unbiased estimator of (12). The LOO error of an algorithm on the training sample S_l is

$$T_{\text{loo}}(\gamma) = \frac{1}{l} \sum_{i=1}^l (y_i - f_\gamma(\mathbf{x}_i | S_l \setminus (\mathbf{x}_i, y_i)))^2 \quad (13)$$

The minimum over γ of (13) we consider as the minimum over γ of (12) since the expectation of (13) coincides with (12).

For KRR, we can derive a closed form expression for the LOO error. Denoting

$$A_\gamma^{-1} = (K^t K + \gamma I_l)^{-1} \quad (14)$$

the error incurred by the LOO procedure is

$$T_{\text{loo}}(\gamma) = \frac{1}{l} \sum_{i=1}^l \left(\frac{y_i - \mathbf{k}_i^t A_\gamma^{-1} K^t \mathbf{y}}{1 - \mathbf{k}_i^t A_\gamma^{-1} \mathbf{k}_i} \right)^2 \quad (15)$$

where $\mathbf{k}_i = (K(\mathbf{x}_1, \mathbf{x}_i), \dots, K(\mathbf{x}_l, \mathbf{x}_i))^t$. Let $\gamma = \gamma_0$ be the minimum of (15). Then the vector

$$\mathbf{y}_0 = \mathbf{K}^* (K^t K + \gamma_0 I_l)^{-1} K^t \mathbf{y} \quad (16)$$

where

$$\mathbf{K}^* = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_{l+1}) & \cdots & K(\mathbf{x}_l, \mathbf{x}_{l+1}) \\ \vdots & & \vdots \\ K(\mathbf{x}_1, \mathbf{x}_{l+m}) & \cdots & K(\mathbf{x}_l, \mathbf{x}_{l+m}) \end{pmatrix} \quad (17)$$

is the KRR estimate of the unknown values $(y_{l+1}^*, \dots, y_{l+m}^*)$.

3. Leave-One-Out Error for Transductive Inference

In transductive inference, our goal is to find an algorithm A which minimizes the functional (4) using both the training data (1) and the test data (2). We suggest the following method: predict $(y_{l+1}^*, \dots, y_{l+m}^*)$ by finding those values which minimize the LOO error of KRR training on the joint set

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l), (\mathbf{x}_{l+1}, y_{l+1}^*), \dots, (\mathbf{x}_{l+m}, y_{l+m}^*) \quad (18)$$

This is achieved in the following way. Suppose we treat the unknown values $\mathbf{y}_m^* = (y_{l+1}^*, \dots, y_{l+m}^*)$ as variables and for some fixed value of these variables we minimize the following empirical functional

$$R_{\text{emp}}(\boldsymbol{\alpha} | \mathbf{y}_m^*) = \frac{1}{l} \left(\sum_{i=1}^l (y_i - f(\mathbf{x}_i, \boldsymbol{\alpha}))^2 + \sum_{i=l+1}^{l+m} (y_i^* - f(\mathbf{x}_i, \boldsymbol{\alpha}))^2 \right) + \gamma \|\boldsymbol{\alpha}\|^2. \quad (19)$$

This functional differs only in the second term from the functional (8) and corresponds to performing KRR with the extra pairs $(\mathbf{x}_{l+1}, y_{l+1}^*), \dots, (\mathbf{x}_{l+m}, y_{l+m}^*)$.

Suppose that vector $\mathbf{y}^* = (y_{l+1}^*, \dots, y_{l+m}^*)^t$ is taken from some set $\mathbf{y}^* \in \mathcal{Y}$ such that the pairs $(\mathbf{x}_{l+1}, y_{l+1}^*), \dots, (\mathbf{x}_{l+m}, y_{l+m}^*)$ can be considered as a sample drawn from the distribution as the pairs $(\mathbf{x}_1, y_1^*), \dots, (\mathbf{x}_l, y_l^*)$. In this case the LOO error of minimizing (19) over the set (18) approximates the functional (4). Using the closed form (15) one obtains

$$T_{\text{loo}}(\gamma | y_{l+1}^*, \dots, y_{l+m}^*) = \frac{1}{l+m} \sum_{i=1}^{l+m} \left(\frac{\widehat{y}_i - \widehat{\mathbf{k}}_i^t \widehat{A}_\gamma^{-1} \widehat{K}^t \widehat{\mathbf{y}}}{1 - \widehat{\mathbf{k}}_i^t \widehat{A}_\gamma^{-1} \widehat{\mathbf{k}}_i} \right)^2 \quad (20)$$

where we denote $\widehat{\mathbf{x}} = (\mathbf{x}_1, \dots, \mathbf{x}_{l+m})^t$ and $\widehat{\mathbf{y}} = (y_1, \dots, y_l, y_{l+1}^*, \dots, y_{l+m}^*)^t$, and

$$\widehat{A}_\gamma^{-1} = (\widehat{K}^t \widehat{K} + \gamma I_{l+m})^{-1} \quad (21)$$

$$\widehat{K} = \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & \cdots & K(\mathbf{x}_1, \mathbf{x}_{l+m}) \\ \vdots & & \vdots \\ K(\mathbf{x}_{l+m}, \mathbf{x}_1) & \cdots & K(\mathbf{x}_{l+m}, \mathbf{x}_{l+m}) \end{pmatrix} \quad (22)$$

$$\widehat{\mathbf{k}}_i = (K(\mathbf{x}_i, \mathbf{x}_1) \cdots K(\mathbf{x}_i, \mathbf{x}_{l+m})) \quad (23)$$

Now let us rewrite the expression (20) in an equivalent form to separate the terms with $\widehat{\mathbf{y}}$ from the terms with \mathbf{x} . Introducing

$$C = I_{l+m} - \widehat{K} \widehat{A}_\gamma^{-1} \widehat{K}^t, \quad (24)$$

and the matrix M with elements

$$M_{ij} = \sum_{k=1}^{l+m} \frac{C_{ik} C_{kj}}{C_{kk}^2}. \quad (25)$$

We obtain the equivalent expression of (20)

$$T_{loo}(\gamma | y_{l+1}^*, \dots, y_{l+m}^*) = \frac{1}{l+m} (\widehat{\mathbf{y}}^t M \widehat{\mathbf{y}}). \quad (26)$$

In order for the \mathbf{y}^* which minimize the LOO procedure to be valid it is required that the pairs $(\mathbf{x}_{l+1}, y_{l+1}^*), \dots, (\mathbf{x}_{l+m}, y_{l+m}^*)$ are drawn from the same distribution as the pairs $(\mathbf{x}_1, y_1^*), \dots, (\mathbf{x}_l, y_l^*)$. To satisfy this constraint we choose vectors \mathbf{y}^* from the set

$$\mathcal{Y} = \{\mathbf{y}^*: \|\mathbf{y}^* - \mathbf{y}_0\|^2 \leq R\} \quad (27)$$

where the vector \mathbf{y}_0 is the solution obtained from classical ridge regression.

To minimize (26) under constraint (27) we use the functional

$$T_{loo}(\gamma | y_{l+1}^*, \dots, y_{l+m}^*) = \widehat{\mathbf{y}}^t M \widehat{\mathbf{y}} + \gamma^* \|\mathbf{y}^* - \mathbf{y}_0\|^2 \quad (28)$$

where γ^* is a constant depending on R .

Now, all that remains is to find the minimum of (28) in \mathbf{y}^* . Note that the matrix M is obtained using only the vectors $\widehat{\mathbf{x}}$. Therefore, to find the minimum of this functional we rewrite (28) as

$$T_{loo}(\gamma) = \mathbf{y}^t M_0 \mathbf{y} + 2\mathbf{y}^{*t} M_1^t \mathbf{y} + \mathbf{y}^{*t} M_2 \mathbf{y}^* + \gamma^* \|\mathbf{y}^* - \mathbf{y}_0\|^2 \quad (29)$$

where

$$M = \begin{pmatrix} M_0 & M_1 \\ M_1^t & M_2 \end{pmatrix} \quad (30)$$

and M_0 is a $l \times l$ matrix, M_1 is a $l \times m$ matrix and M_2 is a $m \times m$ matrix.

Taking the derivative of (29) in \mathbf{y}^* we obtain the condition for the solution

$$2M_1^t \mathbf{y} + 2M_2 \mathbf{y}^* - 2\gamma^* \mathbf{y}_0 + 2\gamma^* \mathbf{y}^* = \mathbf{0} \quad (31)$$

which gives the predictions

$$\mathbf{y}^* = (\gamma^* I_m + M_2)^{-1} (-M_1^t \mathbf{y} + \gamma^* \mathbf{y}_0) \quad (32)$$

In this transductive KRR we have two parameters to control: γ and γ^* . The choice of γ can be found using the LOO estimator (15) for KRR. This leaves γ^* as the only free parameter.

4. Prediction Interval Estimation for KRR

Standard methods for computing prediction intervals in nonlinear regression can be effectively applied to neural networks when the number of training points is large. However, De Veaux et al.(1998) presented an approach to estimating prediction intervals which uses weight decay to fit the network and show that this method is effective on a wide range of problems. Since KRR uses weight decay, we can estimate prediction intervals for KRR in the same way as De Veaux et al.(1998) did. By the way, Seok et al.(2002) considered to estimating prediction intervals for standard SVM in a different way.

It is commonly assumed that the KRR satisfies the nonlinear regression model

$$y = f(\mathbf{x}, \mathbf{a}^*) + \varepsilon \quad (33)$$

where \mathbf{x} represent the inputs, y the outputs, \mathbf{a}^* the true values of the set of parameters, and ε is the error associated with the function f in modeling the system. Let $\widehat{\mathbf{a}}$ be the least squares estimate of \mathbf{a}^* obtained by minimizing the error function

$$S(\mathbf{a}) = \sum_{i=1}^l (y_i - f(\mathbf{x}_i, \mathbf{a}))^2 \quad (34)$$

for a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_l, y_l)$. The predicted output of the input \mathbf{x}_0 is

$$\widehat{y}_0 = f(\mathbf{x}_0, \widehat{\mathbf{a}}). \quad (35)$$

Assume that ε is independently and normally distributed with zero means. The $100(1-\alpha)\%$ confidence interval for the predicted value \widehat{y}_0 is $\widehat{y}_0 \pm c$, where c is

$$c = t_{l-p^*}^{\frac{\alpha}{2}} s (1 + \mathbf{f}_0^t (F^t F)^{-1} \mathbf{f}_0)^{\frac{1}{2}}. \quad (36)$$

Here, $t_{l-p^*}^{\frac{\alpha}{2}}$ is the inverse of the Student t cumulative distribution function with $l-p^*$ degrees of freedom, evaluated at $\alpha/2$, p^* is the effective number of parameters, and $s^2 = S(\mathbf{a})/(l-p^*)$. In effect, the vector \mathbf{f}_0 is given by

$$\begin{aligned} \mathbf{f}_0 &= \left(\frac{\partial f(\mathbf{x}_0, \mathbf{a}^*)}{\partial a_1^*} \quad \frac{\partial f(\mathbf{x}_0, \mathbf{a}^*)}{\partial a_2^*} \quad \dots \quad \frac{\partial f(\mathbf{x}_0, \mathbf{a}^*)}{\partial a_l^*} \right)^t \\ &= (K(\mathbf{x}_1, \mathbf{x}_0) \quad \dots \quad K(\mathbf{x}_l, \mathbf{x}_0)) \end{aligned} \quad (37)$$

and the Jacobian matrix F is given by

$$F = \begin{pmatrix} \frac{\partial f(\mathbf{x}_1, \hat{\mathbf{a}})}{\partial \hat{a}_1} & \frac{\partial f(\mathbf{x}_1, \hat{\mathbf{a}})}{\partial \hat{a}_2} & \dots & \frac{\partial f(\mathbf{x}_1, \hat{\mathbf{a}})}{\partial \hat{a}_l} \\ \frac{\partial f(\mathbf{x}_2, \hat{\mathbf{a}})}{\partial \hat{a}_1} & \frac{\partial f(\mathbf{x}_2, \hat{\mathbf{a}})}{\partial \hat{a}_2} & \dots & \frac{\partial f(\mathbf{x}_2, \hat{\mathbf{a}})}{\partial \hat{a}_l} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f(\mathbf{x}_l, \hat{\mathbf{a}})}{\partial \hat{a}_1} & \frac{\partial f(\mathbf{x}_l, \hat{\mathbf{a}})}{\partial \hat{a}_2} & \dots & \frac{\partial f(\mathbf{x}_l, \hat{\mathbf{a}})}{\partial \hat{a}_l} \end{pmatrix} = K. \quad (38)$$

Notice that

$$p^* = \sum_{i=1}^l \frac{\lambda_i}{\lambda_i + \gamma}, \quad (39)$$

where λ_i is the eigenvalues of F . See for details Vapnik(1998). Replacing \mathbf{a}^* in (37) by $\hat{\mathbf{a}}$, we can estimate the prediction interval straightforwardly.

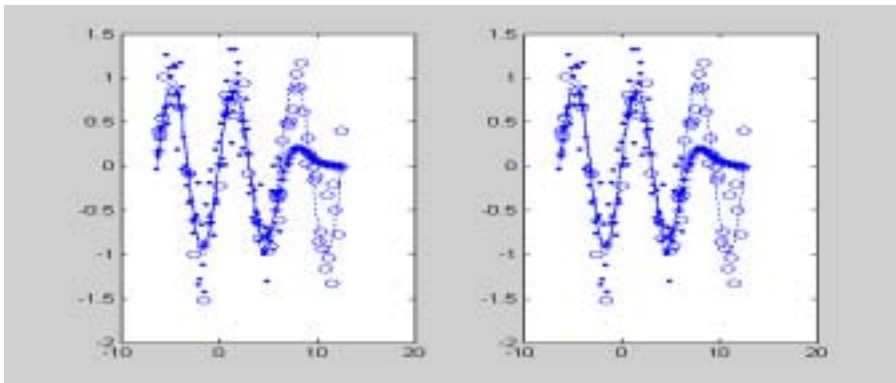
De Veaux et al.(1998) showed that the above method for the prediction interval works well when training data set is large. However, when the training data set is small and the network is trained to convergence, the matrix $F^t F$ can be nearly singular. In this case, the estimated prediction intervals are unreliable. They suggested using the weight decay method, i.e minimizing the error function

$$S(\mathbf{a}) + \gamma \|\mathbf{a}\|^2 \quad (40)$$

instead of $S(\mathbf{a})$. Following De Veaux et al.(1998), for KRR we get finally the prediction interval given by

$$c = t_{n-p}^{\frac{\alpha}{2}} s (1 + \mathbf{f}_0^t (F^t F + \gamma I)^{-1} F^t F (F^t F + \gamma I)^{-1} \mathbf{f}_0)^{\frac{1}{2}}. \quad (41)$$

5. Experiments



(a) The Result of Inductive KRR (b) The Result of Transductive KRR
Figure 1. A Comparison of Transductive KRR to Inductive KRR

We first examine the shape of the estimated regression functions to compare the one-step transductive KRR with the classical two-step inductive KRR in terms of interpolation and extrapolation. Artificial data is generated by a simple function $y = \sin x$ which is corrupted by Gaussian noise with variance 0.1. This function is well used in the papers, for example Gao et al.(2001) related to confidence bound. The training data points of size 100 are evenly distributed between -2π and 2π . The test data points of size 20 are unevenly distributed between -2π and 2π . The test data points of size 30 are evenly distributed between 2π and 4π . When $x > 2\pi$, two KRRs extrapolate the training data. Typical results are shown in Figure 1. Two KRRs are implemented with an Gaussian kernel function with $\sigma = 1.8$, $\gamma = 0.08$ for inductive KRR and $\gamma^* = 1.2$ for transductive KRR. The LOO estimator was used to determine these values.

In Figure 1, the dotted line indicates the true regression function, and the solid line indicates the estimated regression function based on 100 training data points between -2π and 2π . The points marked with dot and circle are the training and test data points, respectively. The points marked with asterisk are the estimated values of test data points. As seen from Figure 1, we notice that for training data points and test data points in the interpolation area two KRRs work very equally well. However, for test data points in the extrapolation area two KRRs do not work well. As a whole, two KRRs work in the exactly same way. Therefore, at least for this particular data set we can not argue that for test data set transductive KRR works better than inductive KRR does, even though transductive KRR is developed to improve estimation proficiency for test data. Actually, a series of experiments has been conducted, but only results for sine function are reported here.

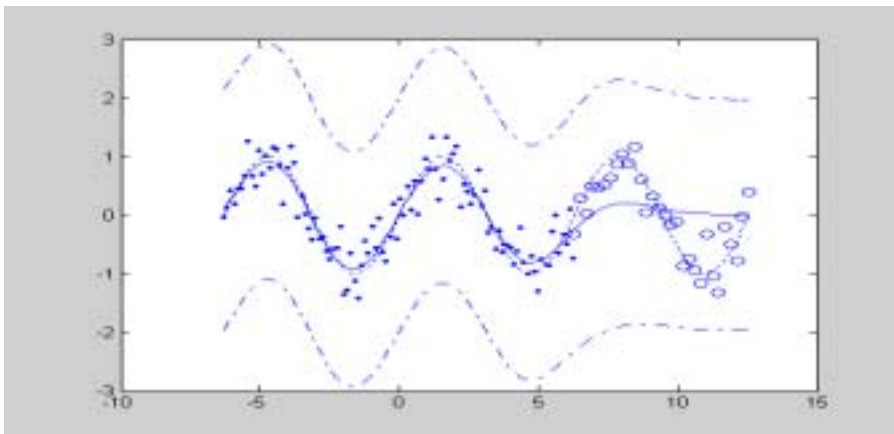


Figure 2 illustrates the prediction intervals for KRR. The points marked with dot and circle are the training and test data points, respectively. The test data

points are only in the extrapolations area. As seen from this figure, we notice that prediction intervals are somewhat wide and are wider in the extrapolation area. We see that prediction intervals are smooth. It is because KRR uses regularization method. This phenomenon happens to the neural networks using weight decay.

6. Conclusions

In this paper we perform transductive inference in the problem of estimating values of functions at the points of interest. We demonstrate that estimating the unknown values via a one-step transductive KRR is not necessarily more accurate than the traditional two-step (inductive plus deductive) KRR.

Like De Veaux et al.(1998), we also study an approach to estimating prediction intervals for KRR, which uses weight decay to fit the network and show that this method is especially simple and effective for kernel machines such as SVM and KRR.

References

1. Chapelle, O., Vapnik, V. and Weston, J. (1999). Transductive Inference for Estimating Values of Functions, *Advances in Neural Information Processing Systems*, 12.
2. DeVeaux, R., Schumi, J., Schweinsberg, J., Shellington, D. and Ungar, L.H. (1998). Prediction Intervals for Neural Networks via Nonlinear Regression, *Technometrics*, 40, 4, 273-282.
3. Gao, J. B., Gunn, S. R., Harris, C. J. and Brown, B. (2001). A Probabilistic Framework for SVM Regression and Error Bar Estimation, *ISIS Technical Report*, U. of Southampton.
4. Saunders, C., Gammerman, A. and Vo $\check{\text{z}}\text{a}$, V. (1998). Ridge Regression Learning Algorithm in Dual Variables, *Proceedings of the 15th International Conference on Machine Learning*, 515-521.
5. Seok, K., Hwang, C. and Cho, D. (2002). Prediction Intervals for Support Vector Machine Regression, *Communications in Statistics: Theory and Methods*, 31, 10, 1887-1898.
6. Vapnik, V. (1995). *The Nature of Statistical Learning Theory*, Springer, New York.
7. Vapnik, V. (1998). *Statistical Learning Theory*, John Wiley & Sons, New York.

[received date : Dec. 2002, accepted date : Feb. 2003]