

Linear Dependency Between ϵ and the Input Noise in ϵ -Support Vector Regression

James T. Kwok and Ivor W. Tsang

Abstract—In using the ϵ -support vector regression (ϵ -SVR) algorithm, one has to decide a suitable value for the insensitivity parameter ϵ . Smola *et al.* considered its “optimal” choice by studying the statistical efficiency in a location parameter estimation problem. While they successfully predicted a linear scaling between the optimal ϵ and the noise in the data, their theoretically optimal value does not have a close match with its experimentally observed counterpart in the case of Gaussian noise. In this paper, we attempt to better explain their experimental results by studying the regression problem itself. Our resultant predicted choice of ϵ is much closer to the experimentally observed optimal value, while again demonstrating a linear trend with the input noise.

Index Terms—Support vector machines (SVMs), support vector regression.

I. INTRODUCTION

IN recent years, the use of support vector machines (SVMs) on various classification and regression problems have been increasingly popular. SVMs are motivated by results from statistical learning theory and, unlike other machine learning methods, their generalization performance does not depend on the dimensionality of the problem [3], [24]. In this paper, we focus on regression problems and consider the ϵ -support vector regression (ϵ -SVR) algorithm [4], [20] in particular. The ϵ -SVR has produced the best result on a timeseries prediction benchmark [11], as well as showing promising results in a number of different applications [7], [12], [22].

One issue about ϵ -SVR is how to set the insensitivity parameter ϵ . Data-resampling techniques such as cross-validation can be used [11], though they are usually very expensive in terms of computation and/or data. A more efficient approach is to use a variant of the SVR algorithm called ν -support vector regression (ν -SVR) [17]. By using another parameter ν to trade off ϵ with model complexity and training accuracy, ν -SVR allows the value of ϵ to be automatically determined. Moreover, it can be shown that ν represents an upper bound on the fraction of training errors and a lower bound on the fraction of support vectors. Thus, in situations where some prior knowledge on the value of ν is available, using ν -SVR may be more convenient than ϵ -SVR. Another approach is to consider the theoretically “optimal” choice of ϵ . Smola *et al.* [19] tackled this by studying the simpler location parameter estimation problem and derived the asymptotically optimal choice of ϵ by maximizing statistical

efficiency. They also showed that this optimal value scales linearly with the noise in the data, which is confirmed in the experiment. However, in the case of Gaussian noise, their predicted value of this optimal ϵ does not have a close match with their experimentally observed value.

In this paper, we attempt to better explain their experimental results. Instead of working on the location parameter estimation problem as in [19], our analysis will be based on the original ϵ -SVR formulation. The rest of this paper is organized as follows. Brief introduction to the ϵ -SVR is given in Section II. The analysis of the linear dependency between ϵ and the input noise level is given in Section III, while the last section gives some concluding remarks.

II. ϵ -SVR

In this section, we introduce some basic notations for ϵ -SVR. Interested readers are referred to [3], [20], [24] for more complete reviews.

Let the training set D be $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, with input $\mathbf{x}_i \in \mathbb{R}^n$ and output $y_i \in \mathbb{R}$. In ϵ -SVR, \mathbf{x} is first mapped to $\mathbf{z} = \psi(\mathbf{x})$ in a Hilbert space \mathcal{F} (with inner product $\langle \cdot, \cdot \rangle$) via a nonlinear map $\psi : \mathbb{R}^n \rightarrow \mathcal{F}$. This space \mathcal{F} is often called the *feature space* and its dimensionality is usually very high (sometimes infinite). Then, a linear function $f(\mathbf{x}) = \langle \mathbf{w}, \psi(\mathbf{x}) \rangle + b$ is constructed in \mathcal{F} such that it deviates least from the training data according to Vapnik’s ϵ -insensitive loss function

$$|y - f(\mathbf{x})|_\epsilon = \begin{cases} 0, & \text{if } |y - f(\mathbf{x})| \leq \epsilon \\ |y - f(\mathbf{x})| - \epsilon, & \text{otherwise} \end{cases}$$

while at the same time is as “flat” as possible (i.e., $\|\mathbf{w}\|$ is as small as possible). Mathematically, this means

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N (\xi_i + \xi_i^*) \\ & \text{subject to } \begin{cases} y_i - f(\mathbf{x}_i) \leq \epsilon + \xi_i \\ f(\mathbf{x}_i) - y_i \leq \epsilon + \xi_i^*, & i = 1, \dots, N \\ \xi_i, \xi_i^* \geq 0 \end{cases} \quad (1) \end{aligned}$$

where C is a user-defined constant. It is well known that (1) can be transformed to the following quadratic programming (QP) problem:

$$\begin{aligned} & \text{maximize } -\frac{1}{2} \sum_{i,j=1}^N (\gamma_i - \gamma_i^*) (\gamma_j - \gamma_j^*) \langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_j) \rangle \\ & \quad - \epsilon \sum_{i=1}^N (\gamma_i + \gamma_i^*) + \sum_{i=1}^N y_i (\gamma_i - \gamma_i^*) \quad (2) \end{aligned}$$

Manuscript received February 13, 2002; revised November 13, 2002.

The authors are with the Department of Computer Science, Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong (e-mail: jamesk@cs.ust.hk; ivor@cs.ust.hk).

Digital Object Identifier 10.1109/TNN.2003.810604

subject to

$$\sum_{i=1}^N (\gamma_i - \gamma_i^*) = 0 \quad \text{and} \quad \gamma_i, \gamma_i^* \in [0, C].$$

However, recall that the dimensionality of \mathcal{F} and thus also of $\psi(\mathbf{x}_i)$ and $\psi(\mathbf{x}_j)$, is usually very high. Hence, in order for this approach to be practical, a key characteristic of ϵ -SVR and kernel methods in general, is that one can obtain $\langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_j) \rangle$ in (2) without having to explicitly obtain $\psi(\mathbf{x}_i)$ and $\psi(\mathbf{x}_j)$ first. This is achieved by using a *kernel* function $K(\cdot, \cdot)$ such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_j) \rangle. \quad (3)$$

For example, the d -order polynomial kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j)^d$ corresponds to a map ψ into the space spanned by all products of exactly d order of \mathfrak{R}^n [24]. More generally, it can be shown that any function satisfying Mercer's theorem can be used as kernel and each will have an associated map ψ such that (3) holds [24]. Computationally, ϵ -SVR (and kernel methods in general) also has the important advantage that only quadratic programming¹ and not nonlinear optimization, is involved. Thus, the use of kernels provides an elegant nonlinear generalization of many existing linear algorithms [2], [3], [10], [16].

III. LINEAR DEPENDENCY BETWEEN ϵ AND THE INPUT NOISE PARAMETER

In order to derive a linear dependency between ϵ and the scale parameter of the input noise model, Smola *et al.* [19] considered the simpler problem of estimating the univariate location parameter w from a set of data points

$$y_i = w + \eta_i, \quad i = 1, \dots, N.$$

Here, η_i s are i.i.d. noise belonging to some distribution $\phi(\cdot)$. Using the Cramer-Rao information inequality for unbiased estimators, the maximum likelihood estimator of w with an "optimal" value of ϵ was then obtained by maximizing its statistical efficiency.

In this paper, instead of working on the location parameter estimation problem, we study the regression problem of estimating the (possibly multivariate) weight parameter $\tilde{\mathbf{w}}$ given a set D of $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, with

$$y_i = \tilde{\mathbf{w}}^T \mathbf{x}_i + \eta_i, \quad i = 1, \dots, N. \quad (4)$$

Here, $\mathbf{x}_i \in \Omega$ follows distribution $p(\cdot)$ and η_i follows distribution $\phi(\cdot)$. The corresponding density function on y is denoted $p(y|\mathbf{x}) = \phi(y - \tilde{\mathbf{w}}^T \mathbf{x})$. Notice that the bias term has been dropped here for simplicity and this is equivalent to assuming that the \mathbf{x}_i s have zero mean. Moreover, using the notation in Section II, in general one can replace \mathbf{x}_i in (4) by $\mathbf{z}_i = \psi(\mathbf{x}_i)$ in feature space \mathcal{F} and thus recovers the original ϵ -SVR setting. Besides, while the work in [19] is based on maximum likelihood estimation, it is now well known that ϵ -SVR is related instead to

¹The SVM problem can also be formulated as a linear programming problem [9], [18] instead of a quadratic programming problem.

maximum *a posteriori* (MAP) estimation. As discussed in [14], [20], [21], the ϵ -insensitive loss function leads to the following probability density function on y

$$p(y_i|\mathbf{x}_i, \mathbf{w}, \beta, \epsilon) = \frac{\beta}{2(1 + \epsilon\beta)} \exp(-\beta|y_i - \mathbf{w}^T \mathbf{x}_i|_\epsilon). \quad (5)$$

Notice that [14], [20], [21] do not have the factor β in (5), but is introduced here to play the role of controlling the noise level² [6]. With the Gaussian prior³ on \mathbf{w}

$$p(\mathbf{w}|\alpha) = \sqrt{\frac{\alpha}{2\pi}} \exp\left(-\frac{\alpha}{2}\|\mathbf{w}\|^2\right)$$

and on applying the Bayes rule, $p(\mathbf{w}|D, \beta, \epsilon) \propto p(D|\mathbf{w}, \beta, \epsilon)p(\mathbf{w})$, we obtain

$$\begin{aligned} \log p(\mathbf{w}|D, \beta, \epsilon) = & -\frac{\alpha}{2}\|\mathbf{w}\|^2 - \beta \sum_{i=1}^N |y_i - \mathbf{w}^T \mathbf{x}_i|_\epsilon \\ & + N \log \frac{\beta}{2(1 + \epsilon\beta)} + \text{const.} \end{aligned} \quad (6)$$

On setting $C = \beta/\alpha$, the optimization problem in (1) can be interpreted as finding the MAP estimate of \mathbf{w} at given values of β and ϵ .

A. Estimating the Optimal Values of β and ϵ

In general, the MAP estimate $\hat{\mathbf{w}}$ in (6) depends on the particular training set and a closed-form solution is difficult to obtain. To simplify the analysis, we replace $1/N \sum_{i=1}^N |y_i - \mathbf{w}^T \mathbf{x}_i|_\epsilon$ in (6) by its expectation

$$\begin{aligned} E_{XY} (|y - \mathbf{w}^T \mathbf{x}|_\epsilon) &= \int_{\Omega} \int_{-\infty}^{\infty} |y - \mathbf{w}^T \mathbf{x}|_\epsilon p(y|\mathbf{x}) p(\mathbf{x}) dy d\mathbf{x} \\ &= \int_{\Omega} \left(\int_{-\infty}^{\mathbf{w}^T \mathbf{x} - \epsilon} (\mathbf{w}^T \mathbf{x} - \epsilon - y) p(y|\mathbf{x}) dy \right. \\ &\quad \left. + \int_{\mathbf{w}^T \mathbf{x} + \epsilon}^{\infty} (y - \mathbf{w}^T \mathbf{x} - \epsilon) p(y|\mathbf{x}) dy \right) \cdot p(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

Equation (6) thus becomes

$$\begin{aligned} M(\mathbf{w}, \beta, \epsilon) = & -\frac{\alpha}{2}\|\mathbf{w}\|^2 - \beta N E_{XY} (|y - \mathbf{w}^T \mathbf{x}|_\epsilon) \\ & + N \log \frac{\beta}{2(1 + \epsilon\beta)} + \text{const.} \end{aligned} \quad (7)$$

On setting its partial derivative with respect to \mathbf{w} to zero, it can be shown that $\hat{\mathbf{w}}$ has to satisfy

$$\begin{aligned} \alpha \hat{\mathbf{w}} + \beta N \int_{\Omega} \mathbf{x} \left(\int_{-\infty}^{\hat{\mathbf{w}}^T \mathbf{x} - \epsilon} p(y|\mathbf{x}) dy \right. \\ \left. - \int_{\hat{\mathbf{w}}^T \mathbf{x} + \epsilon}^{\infty} p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x} = 0. \end{aligned} \quad (8)$$

²To be more precise, we will show in Section III-B that the optimal value of β is inversely proportional to the noise level of the input noise model.

³In this paper, β and ϵ are regarded as hyperparameters while α is treated as a constant. Thus, dependence on α will not be explicitly mentioned in the sequel.

To find suitable values for the hyperparameters β and ϵ , we use the method of type II maximum likelihood prior⁴ (ML-II) in Bayesian statistics [1]. The ML-II solutions of β and ϵ are those that maximize $p(D|\beta, \epsilon)$, which can be obtained by integrating out \mathbf{w}

$$\begin{aligned} p(D|\beta, \epsilon) &= \int p(\mathbf{w}, D|\beta, \epsilon) d\mathbf{w} \\ &= \int p(D|\mathbf{w}, \beta, \epsilon) p(\mathbf{w}) d\mathbf{w} \\ &\simeq p(D|\hat{\mathbf{w}}, \beta, \epsilon) p(\hat{\mathbf{w}}) \Delta\mathbf{w} \end{aligned}$$

where $\Delta\mathbf{w}$ is the width of the integrand $p(D|\mathbf{w}, \beta, \epsilon)p(\mathbf{w})$ and measures the posterior uncertainty of \mathbf{w} . Assuming that the contributions due to $\Delta\mathbf{w}$ are comparable at different values of β and ϵ , maximizing $p(D|\beta, \epsilon)$ with respect to β and ϵ thus becomes maximizing $p(D|\hat{\mathbf{w}}, \beta, \epsilon)p(\hat{\mathbf{w}})$. This is the same as maximizing $\log p(\hat{\mathbf{w}}|D, \beta, \epsilon)$, or $M(\hat{\mathbf{w}}, \beta, \epsilon)$ approximately. Differentiating (7) with respect to ϵ and β and then using (8), we have

$$\begin{aligned} \frac{\partial M(\hat{\mathbf{w}}, \beta, \epsilon)}{\partial \epsilon} &= - \left(\alpha \hat{\mathbf{w}} + N\beta \int_{\Omega} \mathbf{x} \left(\int_{-\infty}^{\hat{\mathbf{w}}^T \mathbf{x} - \epsilon} p(y|\mathbf{x}) dy \right. \right. \\ &\quad \left. \left. - \int_{\hat{\mathbf{w}}^T \mathbf{x} + \epsilon}^{\infty} p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x} \right)^T \frac{\partial \hat{\mathbf{w}}}{\partial \epsilon} \\ &\quad - \frac{N\beta}{1 + \epsilon\beta} + N\beta \int_{\Omega} \left(\int_{-\infty}^{\hat{\mathbf{w}}^T \mathbf{x} - \epsilon} p(y|\mathbf{x}) dy \right. \\ &\quad \left. + \int_{\hat{\mathbf{w}}^T \mathbf{x} + \epsilon}^{\infty} p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x} \\ &= - \frac{N\beta}{1 + \epsilon\beta} + N\beta \int_{\Omega} \left(\int_{-\infty}^{\hat{\mathbf{w}}^T \mathbf{x} - \epsilon} p(y|\mathbf{x}) dy \right. \\ &\quad \left. + \int_{\hat{\mathbf{w}}^T \mathbf{x} + \epsilon}^{\infty} p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x} \quad (9) \end{aligned}$$

and

$$\begin{aligned} \frac{\partial M(\hat{\mathbf{w}}, \beta, \epsilon)}{\partial \beta} &= - \left(\alpha \hat{\mathbf{w}} + N\beta \int_{\Omega} \mathbf{x} \left(\int_{-\infty}^{\hat{\mathbf{w}}^T \mathbf{x} - \epsilon} p(y|\mathbf{x}) dy \right. \right. \\ &\quad \left. \left. - \int_{\hat{\mathbf{w}}^T \mathbf{x} + \epsilon}^{\infty} p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x} \right)^T \frac{\partial \hat{\mathbf{w}}}{\partial \beta} \\ &\quad + \frac{N}{\beta(1 + \epsilon\beta)} - NE_{XY}(|y - \hat{\mathbf{w}}^T \mathbf{x}|_{\epsilon}) \\ &= \frac{N}{\beta(1 + \epsilon\beta)} - NE_{XY}(|y - \hat{\mathbf{w}}^T \mathbf{x}|_{\epsilon}) \quad (10) \end{aligned}$$

respectively. Setting (9) and (10) to zero, we obtain

$$\begin{aligned} \frac{1}{1 + \epsilon\beta} &= \int_{\Omega} \left(\int_{-\infty}^{\hat{\mathbf{w}}^T \mathbf{x} - \epsilon} p(y|\mathbf{x}) dy \right. \\ &\quad \left. + \int_{\hat{\mathbf{w}}^T \mathbf{x} + \epsilon}^{\infty} p(y|\mathbf{x}) dy \right) p(\mathbf{x}) d\mathbf{x} \\ &= \int_{\Omega} \int_{y \notin [\hat{\mathbf{w}}^T \mathbf{x} - \epsilon, \hat{\mathbf{w}}^T \mathbf{x} + \epsilon]} p(\mathbf{x}, y) dy d\mathbf{x} \quad (11) \end{aligned}$$

⁴Notice that MacKay's evidence framework [8], which has been popularly used in the neural networks community, is computationally equivalent to the ML-II method.

$$\frac{1}{\beta(1 + \epsilon\beta)} = E_{XY}(|y - \hat{\mathbf{w}}^T \mathbf{x}|_{\epsilon}). \quad (12)$$

These can then be used to solve for β and ϵ .

B. Applications to Some Common Noise Models

Recall that the ϵ -insensitive loss function implicitly corresponds to the noise model in (5). Of course, in cases where the underlying noise model of a particular data set is known, the loss function should be chosen such that it has a close match with this known noise model, while at the same time ensuring that the resultant optimization problem can still be solved efficiently. It is thus possible that the ϵ -insensitive loss function may not be best. For example, when the noise model is known to be Gaussian, the corresponding loss function is the squared loss function and the resultant model becomes a regularization network [5], [13].⁵ However, an important advantage of ϵ -SVR using the ϵ -insensitive loss function, just like SVMs using the hinge loss in classification problems, is that sparseness of the dual variables can be ensured. On the contrary, sparseness will be lost if the squared loss function is used instead. Thus, even for Gaussian noise, the ϵ -insensitive loss function is sometimes still desirable and the resultant performance is often very competitive. In this section, by solving (11) and (12), we will show that there is always a linear relationship between the optimal value of ϵ and the noise level, even when the true noise model is different from (5). As in [19], three commonly used noise models, namely the Gaussian, Laplacian, and uniform models, will be studied.

1) *Gaussian Noise:* For the Gaussian noise model

$$\phi(\eta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\eta^2}{2\sigma^2}\right)$$

where σ is the standard deviation. Define $\delta(\mathbf{x}) = (\hat{\mathbf{w}} - \hat{\mathbf{w}})^T \mathbf{x}$, $b_1(\mathbf{x}) = 1/\sqrt{2}(\epsilon/\sigma - \delta(\mathbf{x})/\sigma)$ and $b_2(\mathbf{x}) = 1/\sqrt{2}(\epsilon/\sigma + \delta(\mathbf{x})/\sigma)$. It can be shown that (11) and (12) reduce to

$$\frac{1}{1 + \epsilon\beta} = \int_{\Omega} \frac{\text{erfc}(b_1) + \text{erfc}(b_2)}{2} p(\mathbf{x}) d\mathbf{x} \quad (13)$$

$$\begin{aligned} \frac{1}{\beta(1 + \epsilon\beta)} &= \int_{\Omega} \left(-\frac{\sigma b_1}{\sqrt{2}} \text{erfc}(b_1) - \frac{\sigma b_2}{\sqrt{2}} \text{erfc}(b_2) \right. \\ &\quad \left. + \frac{\sigma}{\sqrt{2\pi}} \exp(-b_1^2) + \frac{\sigma}{\sqrt{2\pi}} \exp(-b_2^2) \right) \\ &\quad \cdot p(\mathbf{x}) d\mathbf{x} \quad (14) \end{aligned}$$

where $\text{erfc}(x) = 2/\sqrt{\pi} \int_x^{\infty} \exp(-t^2) dt$ is the complementary error function [15]. Substituting (13) into (14) and after performing the integration,⁶ it can be shown that ϵ always appears as the factor ϵ/σ in the solution and, thus, there is always a linear relationship between the optimal values of ϵ and σ .

⁵This is sometimes also called the least-squares SVM [23] in the SVM literature.

⁶Here, the integration is performed by using the symbolic math toolbox of Matlab.

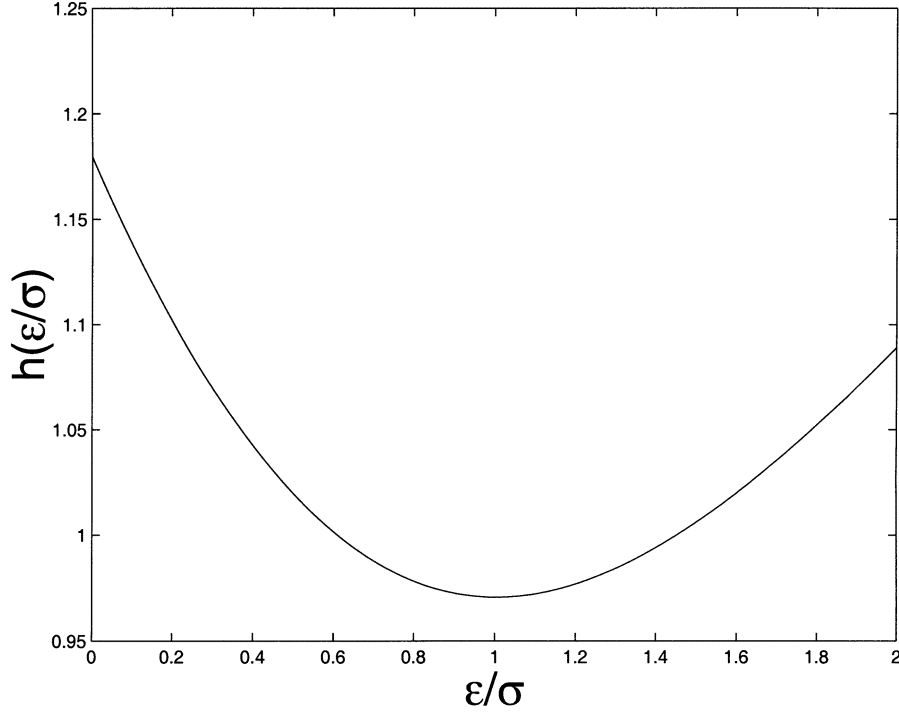


Fig. 1. Objective function $h(\epsilon/\sigma)$ to be minimized in the case of Gaussian noise.

If we assume that

$$\begin{aligned}
 E_X (\delta^2(\mathbf{x})) &= \int_{\Omega} \delta^2(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\
 &= E_X \left((\hat{\mathbf{w}}^T \mathbf{x} - \tilde{\mathbf{w}}^T \mathbf{x})^2 \right) \\
 &\simeq E_X \left((y - \tilde{\mathbf{w}}^T \mathbf{x})^2 \right) \\
 &= \sigma^2
 \end{aligned} \tag{15}$$

and also that the variation of $\hat{\mathbf{w}}$ with ϵ is small,⁷ then maximizing M in (7) is effectively the same as minimizing $\beta E_{XY}(|y - \tilde{\mathbf{w}}^T \mathbf{x}|_{\epsilon}) - \log(\beta/2(1 + \epsilon\beta))$, which is the same as minimizing⁸

$$\begin{aligned}
 h\left(\frac{\epsilon}{\sigma}\right) &= \sqrt{2\pi} \exp\left(\frac{\epsilon^2}{2\sigma^2}\right) \left(\left(\frac{\epsilon}{\sigma}\right)^2 + 3\right)^{-1} \\
 &\cdot \left(-\frac{\epsilon}{\sigma} \operatorname{erfc}\left(\frac{\epsilon}{\sqrt{2}\sigma}\right) + \frac{3}{\sqrt{2\pi}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right)\right) \\
 &- \log\left(\sqrt{2\pi} \exp\left(\frac{\epsilon^2}{2\sigma^2}\right) \left(\left(\frac{\epsilon}{\sigma}\right)^2 + 3\right)^{-1}\right) \\
 &- \log\left(\operatorname{erfc}\left(\frac{\epsilon}{\sqrt{2}\sigma}\right) + \frac{\epsilon}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right)\right).
 \end{aligned} \tag{16}$$

⁷Notice that from (8), we have

$$\begin{aligned}
 \frac{\partial \hat{\mathbf{w}}}{\partial \epsilon} &= \left(\frac{\alpha}{\beta N} \mathbf{I} + \int_{\Omega} \mathbf{x} \mathbf{x}^T (p(\hat{\mathbf{w}}^T \mathbf{x} - \epsilon|\mathbf{x}) \right. \\
 &\quad \left. + p(\hat{\mathbf{w}}^T \mathbf{x} + \epsilon|\mathbf{x})) p(\mathbf{x}) d\mathbf{x} \right)^{-1} \\
 &\cdot \int_{\Omega} \mathbf{x} (p(\hat{\mathbf{w}}^T \mathbf{x} - \epsilon|\mathbf{x}) \\
 &\quad - p(\hat{\mathbf{w}}^T \mathbf{x} + \epsilon|\mathbf{x})) p(\mathbf{x}) d\mathbf{x}
 \end{aligned}$$

where \mathbf{I} is the identity matrix. $\partial \hat{\mathbf{w}} / \partial \epsilon$ thus involves the difference of two very similar terms and will be small.

⁸More detailed derivations are in the Appendix.

A plot of $h(\epsilon/\sigma)$ is shown in Fig. 1. Its minimum can be obtained numerically as

$$\epsilon = 1.0043\sigma. \tag{17}$$

By substituting (15), (17) into (29), we can also obtain the optimal value of β as $1/\beta = 0.9658\sigma$. This confirms the role of β in controlling the noise variance, as mentioned in Section III.

In the following, we repeat an experiment in [19] to verify this ratio of ϵ/σ . The target function to be learned is $f(x) = 0.9\operatorname{sinc}(10x/\pi)$. The training set consists of 200 x s drawn independently from a uniform distribution on $[-1, 1]$ and the corresponding outputs have Gaussian noise $N(0, \sigma^2)$ added. Testing error is obtained by directly integrating the squared difference between the true target function and the estimated function over the range $[-1, 1]$. As in [19], the model selection problem for the regularization parameter C in (1) is side-stepped by always choosing the value of C that yields the smallest testing error. The whole experiment is repeated 40 times and the spline kernel with an infinite number of knots is used.

Fig. 2 shows the linear relationship of $\epsilon = 0.9846\sigma$ obtained from the experiment. This is close to the value of $\epsilon = 0.9\sigma$ obtained in the experiment in [19] and is also in good agreement with our predicted ratio of $\epsilon = 1.0043\sigma$. In comparison, the theoretical results in [19] suggest the ‘‘optimal’’ choice of $\epsilon = 0.6166\sigma$.

2) *Laplacian Noise*: Next, we consider the Laplacian noise model

$$\phi(\eta) = \frac{1}{2\sigma} \exp\left(-\frac{|\eta|}{\sigma}\right).$$

Here, we only consider the simpler case when \mathbf{x} is one-dimensional, with uniform density over the range $[-L, L]$ (i.e., $X \sim$

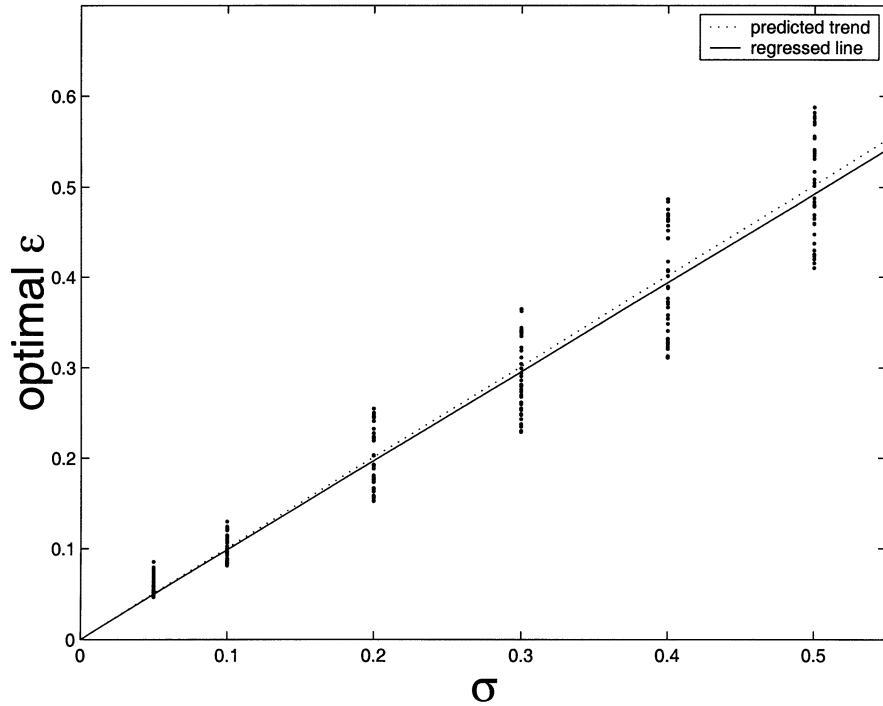


Fig. 2. Relationship between the optimal value of ϵ and the Gaussian noise level σ .

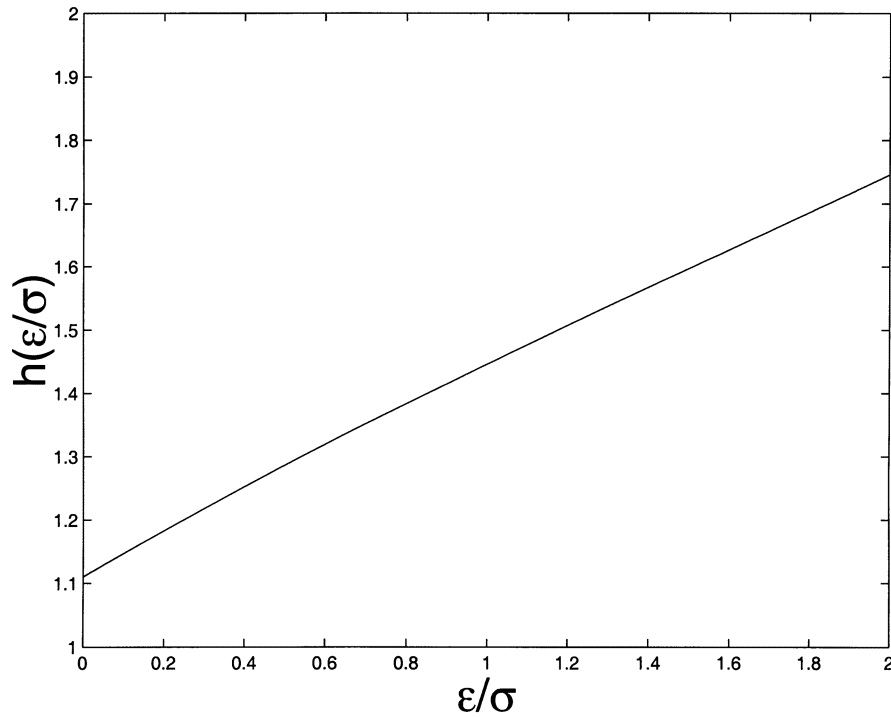


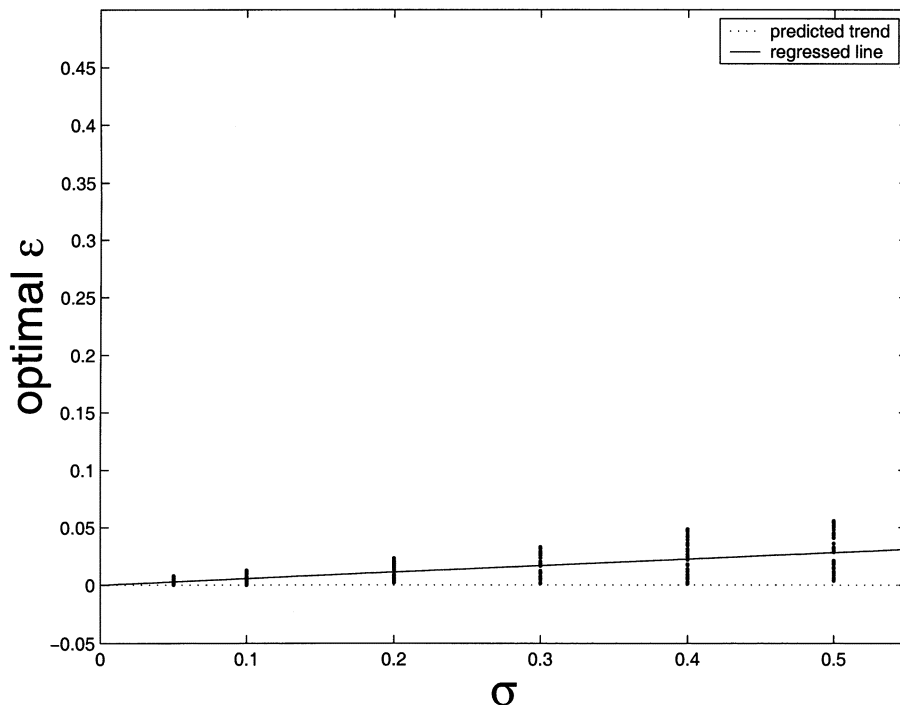
Fig. 3. Objective function $h(\epsilon/\sigma)$ to be minimized in the case of Laplacian noise.

$U([-L, L])$). After tedious computation, it can be shown that (11) and (12) reduce to

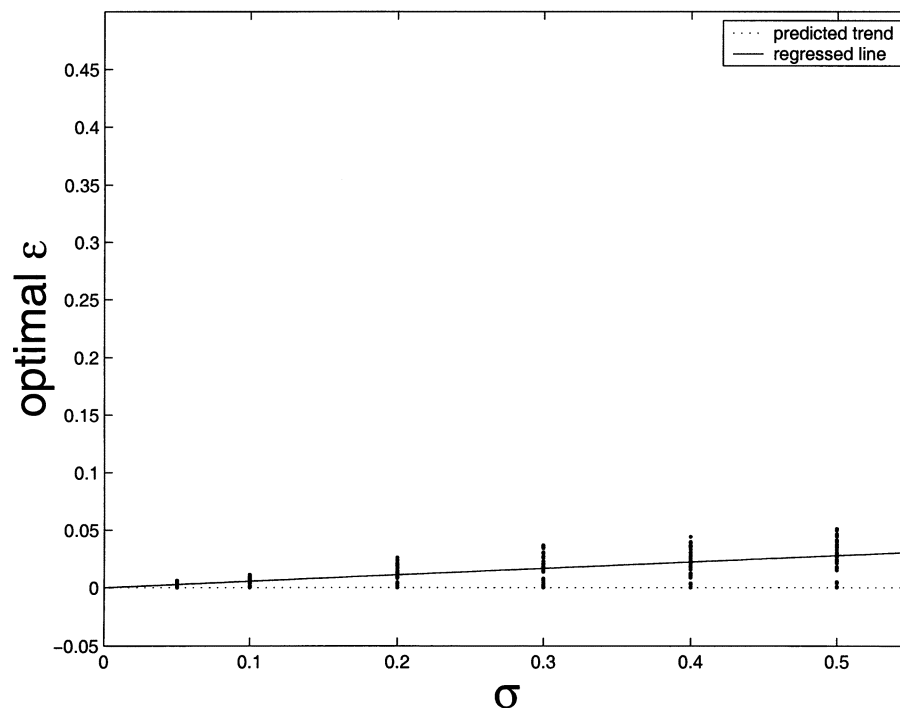
$$\frac{1}{1 + \epsilon\beta} = 1 - \frac{\epsilon}{L(\tilde{w} - \hat{w})} + \frac{\sigma}{2L(\tilde{w} - \hat{w})} \cdot \exp\left(-\frac{L(\tilde{w} - \hat{w})}{\sigma}\right) \cdot \left(\exp\left(\frac{\epsilon}{\sigma}\right) - \exp\left(-\frac{\epsilon}{\sigma}\right)\right) \quad (18)$$

and

$$\frac{1}{\beta} = \frac{\sigma^2}{L(\tilde{w} - \hat{w})} + \frac{\sigma}{2L(\tilde{w} - \hat{w})} \cdot \exp\left(-\frac{L(\tilde{w} - \hat{w})}{\sigma}\right) \cdot \left((\epsilon - \sigma) \exp\left(\frac{\epsilon}{\sigma}\right) - (\epsilon + \sigma) \exp\left(-\frac{\epsilon}{\sigma}\right)\right) + \frac{(\tilde{w} - \hat{w})L}{2} - \frac{\epsilon^2}{2L(\tilde{w} - \hat{w})}. \quad (19)$$



(a)



(b)

Fig. 4. Relationship between the optimal value of ϵ and the Laplacian noise level σ . (a) One-dimensional x . (b) Two-dimensional x .

Substituting (19) back into (18), we notice again that ϵ always appears in the factor ϵ/σ and thus there is always a linear relationship between the optimal value of ϵ and σ under the Laplacian noise. Recall that the variance for the uniform distribution $U(a, b)$ is $(b - a)^2/12$ and that for the exponential distribution⁹

⁹The density function of the exponential distribution $\theta \sim \text{Expon}(\beta)$ is $p(\theta) = \beta \exp(-\beta\theta)$ (where $\beta > 0, \theta \geq 0$).

$\text{Expon}(\beta)$ is $1/\beta^2$. As in Section III-B1, we consider $\delta(x) = (\hat{w} - \tilde{w})x$ and assume that $E_X(\delta^2(x)) \simeq E_X((y - \hat{w}x)^2) = 2\sigma^2$. With $X \sim U([-L, L])$, we also obtain $\text{var}(X) = L^2/3$. Consequently

$$E_X(\delta^2(x)) \simeq \text{var}((\tilde{w} - \hat{w})X) \simeq \frac{L^2(\tilde{w} - \hat{w})^2}{3} \quad (20)$$

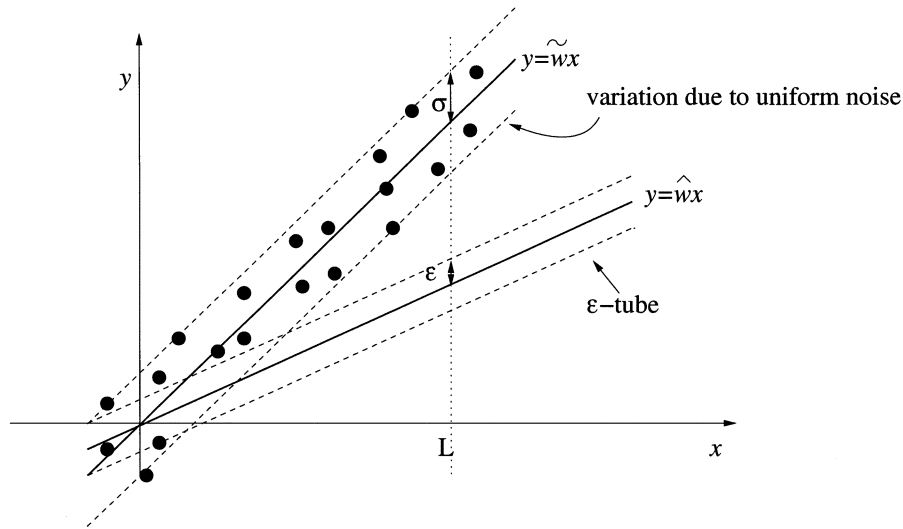


Fig. 5. Scenario when $L > \sigma + \epsilon/\tilde{w} - \hat{w}$ in the case of uniform noise.

and hence

$$\sigma = \frac{L(\tilde{w} - \hat{w})}{\sqrt{6}}. \quad (21)$$

Analogous to Section III-B1, by plugging in (18), (19), and (21), it can be shown that the problem reduces to minimizing (Fig. 3) as shown in the equation at the bottom of the page. Again, the solution can be obtained numerically as $\epsilon = 0$, which was also obtained in [19]. Notice that this is intuitively correct as the density function in (5) degenerates to the Laplacian density when $\epsilon = 0$. Moreover, we can also obtain the optimal value of β as $1/\beta = 1.5977\sigma$, again showing that β is inversely proportional to the noise level of the Laplacian noise model.

To verify the ratio of ϵ/σ , we repeat the experiment in Section III-B1 but with the Laplacian noise $L(0, \sigma)$ instead. Moreover, as our analysis applies only when x is one dimensional, we also investigate experimentally the optimal value of ϵ in the two-dimensional case. The ratios obtained for the one- and two-dimensional cases are $\epsilon = 0.0565\sigma$ and $\epsilon = 0.0559\sigma$, respectively (Fig. 4), which are close to our prediction of $\epsilon = 0$.

3) *Uniform Noise*: Finally, we consider the uniform noise model

$$\phi(\eta) = \frac{1}{2\sigma} I_{[-\sigma, \sigma]}(\eta) \quad (22)$$

where $I_A(\eta) = \begin{cases} 1, & \text{if } \eta \in A \\ 0, & \text{otherwise.} \end{cases}$ is the indicator function. Again, we only consider the simpler case when \mathbf{x} is one-dimen-

sional, with uniform density over the range $[-L, L]$. Moreover, only the case when $(\sigma - \epsilon)/(\tilde{w} - \hat{w}) \leq L \leq (\sigma + \epsilon)/(\tilde{w} - \hat{w})$ will be discussed here. Derivation for the case when $L < (\sigma - \epsilon)/(\tilde{w} - \hat{w})$ is similar. Whereas, for $L > \sigma + \epsilon/\tilde{w} - \hat{w}$ (i.e., $\tilde{w}L - \sigma > \hat{w}L + \epsilon$), the corresponding SVR solution will be very poor (Fig. 5) and thus usually not the case of interest.

It can be shown that (11) and (12) reduce to

$$\frac{1}{1 + \epsilon\beta} = \frac{((\sigma - \epsilon) + (\tilde{w} - \hat{w})L)^2}{4\sigma(\tilde{w} - \hat{w})L}$$

and

$$\frac{1}{\beta(1 + \epsilon\beta)} = \frac{((\tilde{w} - \hat{w})L + (\sigma - \epsilon))^3}{12\sigma(\tilde{w} - \hat{w})L}. \quad (23)$$

Combining and simplifying, we have $1/\beta = ((\tilde{w} - \hat{w})L + (\sigma - \epsilon))/3$. Substituting back into (23), we obtain

$$\frac{L}{\sigma}(\tilde{w} - \hat{w}) \left(2 - \frac{\epsilon}{\sigma} - (\tilde{w} - \hat{w}) \frac{L}{\sigma} \right) = 1 - 2 \left(\frac{\epsilon}{\sigma} \right)^2 + \frac{\epsilon}{\sigma}. \quad (24)$$

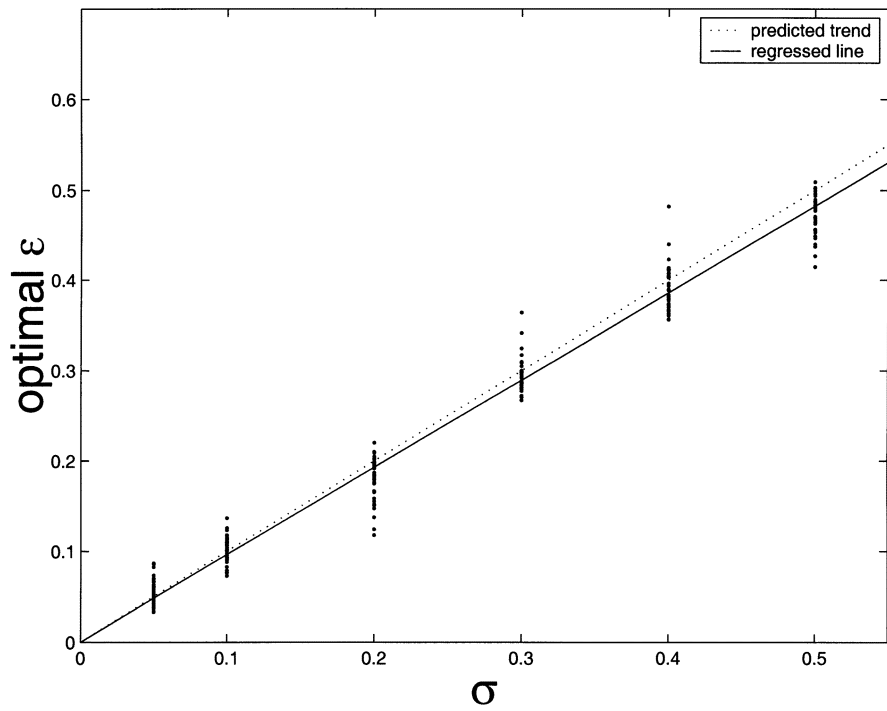
Again, we notice that ϵ always appears in the factor ϵ/σ and, thus, there is also a linear relationship between the optimal value of ϵ and σ in the case of uniform noise.

As in Section III-B1, consider $\delta(x) = (\hat{w} - \tilde{w})x$ and assume that $E_X(\delta^2(x)) \simeq E_X((y - \hat{w}x)^2) = \sigma^2/3$. Using (20) again, we obtain $\sigma = L(\tilde{w} - \hat{w})$. On substituting back into (24) and proceed as in Section III-B1, we obtain, after tedious computation, $\epsilon = \sigma$. This is intuitively correct since the density function

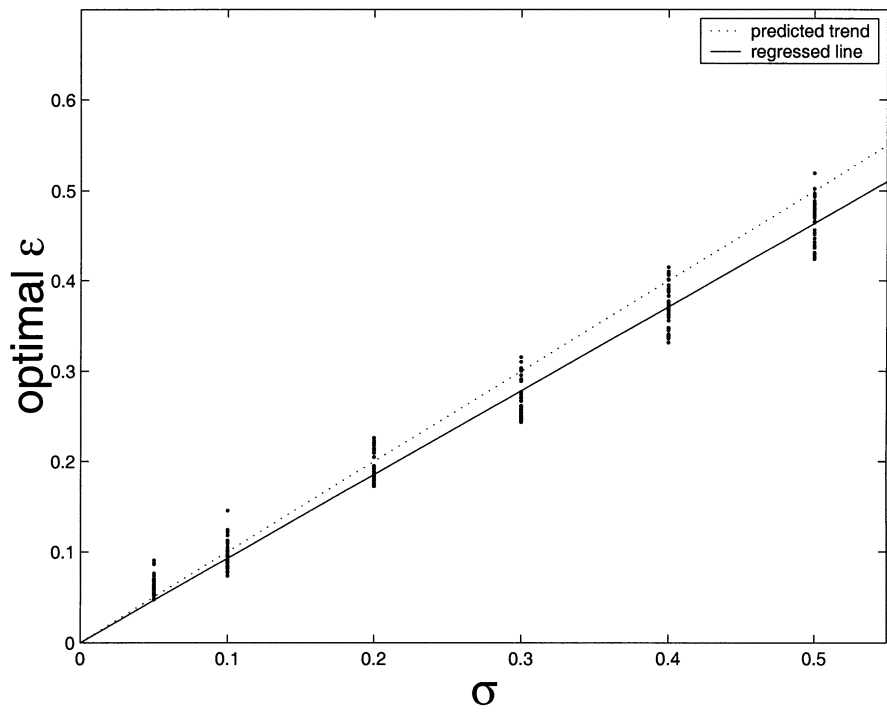
$$h\left(\frac{\epsilon}{\sigma}\right) = \frac{-\frac{\epsilon}{\sigma} + \frac{1}{\sqrt{6}} \left(\left(\frac{\epsilon}{\sigma}\right)^2 + 1 \right) - \frac{\exp(-\sqrt{6})}{2\sqrt{6}} \left(\exp\left(\frac{\epsilon}{\sigma}\right) + \exp\left(-\frac{\epsilon}{\sigma}\right) \right) + \frac{\sqrt{6}}{2} - \frac{1}{2\sqrt{6}} \left(\frac{\epsilon}{\sigma}\right)^2}{\frac{1}{\sqrt{6}} + \frac{\exp(-\sqrt{6})}{2\sqrt{6}} \left(\left(\frac{\epsilon}{\sigma} - 1\right) \exp\left(\frac{\epsilon}{\sigma}\right) - \left(\frac{\epsilon}{\sigma} + 1\right) \exp\left(-\frac{\epsilon}{\sigma}\right) \right) + \frac{\sqrt{6}}{2} - \frac{1}{2\sqrt{6}} \left(\frac{\epsilon}{\sigma}\right)^2}$$

$$+ \log \left(\frac{1}{\sqrt{6}} + \frac{\exp(-\sqrt{6})}{2\sqrt{6}} \left(\left(\frac{\epsilon}{\sigma} - 1\right) \exp\left(\frac{\epsilon}{\sigma}\right) - \left(\frac{\epsilon}{\sigma} + 1\right) \exp\left(-\frac{\epsilon}{\sigma}\right) \right) + \frac{\sqrt{6}}{2} - \frac{1}{2\sqrt{6}} \left(\frac{\epsilon}{\sigma}\right)^2 \right)$$

$$- \log \left(1 - \frac{1}{\sqrt{6}} \frac{\epsilon}{\sigma} + \frac{\exp(-\sqrt{6})}{2\sqrt{6}} \left(\exp\left(\frac{\epsilon}{\sigma}\right) - \exp\left(-\frac{\epsilon}{\sigma}\right) \right) \right).$$



(a)



(b)

Fig. 6. Relationship between the optimal value of ϵ and the uniform noise level σ . (a) One-dimensional x . (b) Two-dimensional x .

in (5) becomes effectively the same as (22) when $\epsilon = \sigma$, as noise with magnitude larger than the size of the ϵ -tube then becomes impossible. Moreover, as a side-product, we can also obtain the optimal value of β as $1/\beta = \sigma/3$, showing that β is again inversely proportional to σ .

Fig. 6 shows the linear relationships obtained from the experiment, with $\epsilon = 0.9646\sigma$ and $\epsilon = 0.9275\sigma$ for the one- and

two-dimensional cases, respectively. These are again in close match with our prediction of $\epsilon = \sigma$.

IV. CONCLUSION

In this paper, we study the ϵ -SVR problem and derive the optimal choice of ϵ at a given value of the input noise parameter.

While the results in [19] considered only the problem of location parameter estimation using maximum likelihood, our analysis is based on the original ϵ -SVR formulation and corresponds to MAP estimation. Consequently, our predicted ratio of ϵ/σ in the case of Gaussian noise is much closer to the experimentally observed optimal value. Besides, we accord with [19] in that ϵ scales linearly with the scale parameter under the Gaussian, Laplacian, and uniform noise models.

In order to apply these linear dependency results in practical applications, one has to first arrive at an estimate of the noise level σ . One way to obtain this is by using Bayesian methods (e.g., [6]). In the future, we will investigate the integration of these two and also its applications in some real-world problems. Besides, the work here is also useful in designing simulation experiments. Typically, a researcher/practitioner may be experimenting a new technique on ϵ -SVR, while not directly addressing the issue of finding the optimal value of ϵ . The results here can then be used to ascertain that a suitable value/range of ϵ has been chosen in the simulation. Finally, our analyses under the Laplacian and uniform noise models are restricted to the case when the input \mathbf{x} is one dimensional and with uniform density over a certain range. Extensions to the multivariate case and to other noise models will also be investigated in the future.

APPENDIX

Recall the following notations introduced in Section III-B1:

$$\begin{aligned}\delta(\mathbf{x}) &= (\hat{\mathbf{w}} - \tilde{\mathbf{w}})^T \mathbf{x} \\ b_1(\mathbf{x}) &= \frac{1}{\sqrt{2}} \left(\frac{\epsilon}{\sigma} - \frac{\delta(\mathbf{x})}{\sigma} \right), \\ b_2(\mathbf{x}) &= \frac{1}{\sqrt{2}} \left(\frac{\epsilon}{\sigma} + \frac{\delta(\mathbf{x})}{\sigma} \right).\end{aligned}$$

In the following, we will use the second-order Taylor series expansions for $\text{erfc}(x)$ and $\exp(-x^2)$

$$\begin{aligned}\text{erfc}(x+h) &= \text{erfc}(x) - \frac{2}{\sqrt{\pi}} e^{-x^2} h + \frac{2}{\sqrt{\pi}} x e^{-x^2} h^2 \\ &\quad + O(h^3)\end{aligned}\quad (25)$$

$$\begin{aligned}\exp(-(x+h)^2) &= \exp(-x^2) (1 - 2xh + (2x^2 - 1)h^2) \\ &\quad + O(h^3).\end{aligned}\quad (26)$$

Using (25), (13) then becomes

$$\begin{aligned}\frac{1}{1+\epsilon\beta} &= \int_{\Omega} \frac{\text{erfc}(b_1) + \text{erfc}(b_2)}{2} p(\mathbf{x}) d\mathbf{x} \\ &\simeq \int_{\Omega} \left(\text{erfc}\left(\frac{\epsilon}{\sqrt{2}\sigma}\right) + \frac{2}{\sqrt{\pi}} \frac{\epsilon}{\sqrt{2}\sigma} \right. \\ &\quad \cdot \left. \left(\exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \right) \frac{\delta^2(\mathbf{x})}{2\sigma^2} \right) p(\mathbf{x}) d\mathbf{x} \\ &= \text{erfc}\left(\frac{\epsilon}{\sqrt{2}\sigma}\right) + \frac{\epsilon}{\sqrt{2\pi}\sigma^3} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \\ &\quad \cdot E_X(\delta^2(\mathbf{x})).\end{aligned}\quad (27)$$

Similarly, using (26), (14) becomes

$$\begin{aligned}\frac{1}{\beta(1+\epsilon\beta)} &= \int_{\Omega} \left(-\frac{\sigma b_1}{\sqrt{2}} \text{erfc}(b_1) - \frac{\sigma b_2}{\sqrt{2}} \text{erfc}(b_2) \right. \\ &\quad \left. + \frac{\sigma}{\sqrt{2\pi}} \exp(-b_1^2) + \frac{\sigma}{\sqrt{2\pi}} \exp(-b_2^2) \right) p(\mathbf{x}) d\mathbf{x} \\ &\simeq \int_{\Omega} \left[-\epsilon \left(\frac{\text{erfc}(b_1) + \text{erfc}(b_2)}{2} \right) \right. \\ &\quad \left. + \frac{\delta(\mathbf{x})}{2} (\text{erfc}(b_1) - \text{erfc}(b_2)) \right. \\ &\quad \left. + \frac{\sigma}{\sqrt{2\pi}} \cdot 2 \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \right. \\ &\quad \left. \cdot \left(1 + \left(\frac{\epsilon^2}{\sigma^2} - 1 \right) \frac{\delta^2(\mathbf{x})}{2\sigma^2} \right) \right] p(\mathbf{x}) d\mathbf{x}.\end{aligned}$$

Using (13) and (25), this simplifies to (28), shown at the bottom of the page, from which we obtain

$$\beta = \sqrt{\frac{\pi}{2}} \exp\left(\frac{\epsilon^2}{2\sigma^2}\right) \frac{2\sigma}{2\sigma^2 + \left(\frac{\epsilon^2}{\sigma^2} + 1\right) E_X(\delta^2(\mathbf{x}))}. \quad (29)$$

Substituting (27) into (28) and simplifying, we have

$$\begin{aligned}E_{XY}(|y - \tilde{\mathbf{w}}^T \mathbf{x}|_{\epsilon}) &= -\epsilon \text{erfc}\left(\frac{\epsilon}{\sqrt{2}\sigma}\right) \\ &\quad + \sqrt{\frac{2}{\pi}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \left(\sigma + \frac{E_X(\delta^2(\mathbf{x}))}{2\sigma} \right).\end{aligned}\quad (30)$$

$$\begin{aligned}E_{XY}(|y - \tilde{\mathbf{w}}^T \mathbf{x}|_{\epsilon}) &= \frac{1}{\beta(1+\epsilon\beta)} \\ &= -\frac{\epsilon}{1+\epsilon\beta} + \int_{\Omega} \left[\delta(\mathbf{x}) \cdot \frac{2}{\sqrt{\pi}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \cdot \frac{\delta(\mathbf{x})}{\sqrt{2}\sigma} + \sqrt{\frac{2}{\pi}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \left(\sigma + \left(\frac{\epsilon^2}{\sigma^2} - 1 \right) \frac{\delta^2(\mathbf{x})}{2\sigma} \right) \right] p(\mathbf{x}) d\mathbf{x} \\ &= -\frac{\epsilon}{1+\epsilon\beta} + \sqrt{\frac{2}{\pi}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \left(\sigma + \int_{\Omega} \left(\frac{\epsilon^2}{\sigma^2} + 1 \right) \frac{\delta^2(\mathbf{x})}{2\sigma} p(\mathbf{x}) d\mathbf{x} \right) \\ &= -\frac{\epsilon}{1+\epsilon\beta} + \sqrt{\frac{2}{\pi}} \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \left(\sigma + \frac{1}{2\sigma} \left(\frac{\epsilon^2}{\sigma^2} + 1 \right) E_X(\delta^2(\mathbf{x})) \right)\end{aligned}\quad (28)$$

Now, maximizing M in (7) is thus the same as minimizing $\beta E_{XY}(|y - \mathbf{w}^T \mathbf{x}|_\epsilon) - \log \beta / 2(1 + \epsilon\beta)$. Substituting in (15), (27), (29), and (30), this becomes minimizing $h(\epsilon/\sigma)$ in (16). \square

REFERENCES

- [1] J. O. Berger, *Statistical Decision Theory and Bayesian Analysis*, 2nd ed, ser. Springer Series in Statistics. New York: Springer-Verlag, 1985.
- [2] V. Cherkassky and F. Mulier, *Learning From Data: Concepts, Theory and Methods*. New York: Wiley, 1998.
- [3] N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines*. New York: Cambridge Univ. Press, 2000.
- [4] H. D. Drucker, C. J. C. Burges, L. Kaufman, A. Smola, and V. Vapnik, "Support vector regression machines," in *Advances in Neural Information Processing Systems 9*, M. C. Mozer, M. I. Jordan, and T. Petsche, Eds. San Mateo, CA: Morgan Kaufmann, 1997, pp. 155–161.
- [5] T. Evgenious, M. Pontil, and T. Poggio, "Regularization networks and support vector machines," in *Advances in Computational Mathematics*, 1999.
- [6] M. H. Law and J. T. Kwok, "Bayesian support vector regression," in *Proc. 8th Int. Workshop Artificial Intelligence Statistics*, Key West, FL, 2001, pp. 239–244.
- [7] Y. Li, S. Gong, and H. Liddell, "Support vector regression and classification based multi-view face detection and recognition," in *Proc. 4th IEEE Int. Conf. Face Gesture Recognition*, Grenoble, France, 2000, pp. 300–305.
- [8] D. J. C. MacKay, "Bayesian interpolation," *Neural Comput.*, vol. 4, no. 3, pp. 415–447, May 1992.
- [9] O. L. Mangasarian and D. R. Musicant, "Robust linear and support vector regression," *IEEE Trans. Pattern Anal. Machine Learning*, vol. 22, Sept. 2000.
- [10] K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf, "An introduction to kernel-based learning algorithms," *IEEE Trans. Neural Networks*, vol. 12, pp. 181–202, Mar. 2001.
- [11] K. R. Müller, A. J. Smola, G. Rätsch, B. Schölkopf, J. Kohlmorgen, and V. Vapnik, "Predicting time series with support vector machines," in *Proc. Int. Conf. Artificial Neural Networks*, 1997, pp. 999–1004.
- [12] C. Papageorgiou, F. Girosi, and T. Poggio, "Sparse correlation kernel reconstruction," in *Proc. Int. Conf. Acoust., Speech, Signal Processing*, Phoenix, AZ, 1999, pp. 1633–1636.
- [13] T. Poggio and F. Girosi, *A Theory of Networks for Approximation and Learning*. Cambridge, MA: MIT Press, 1989.
- [14] M. Pontil, S. Mukherjee, and F. Girosi, *On the Noise Model of Support Vector Machine Regression*. Cambridge, MA: MIT Press, 1998.
- [15] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, *Numerical Recipes in C*, 2nd ed. New York: Cambridge Univ. Press, 1992.
- [16] B. Schölkopf, C. Burges, and A. e. Smola, *Advances in Kernel Methods—Support Vector Learning*. Cambridge, MA: MIT Press, 1998.

- [17] B. Schölkopf, A. Smola, and R. C. Williamson, "New support vector algorithms," *Neural Comput.*, vol. 12, no. 4, pp. 1207–1245, 2000.
- [18] A. Smola, B. Schölkopf, and G. Rätsch, "Linear programs for automatic accuracy control in regression," in *Proc. Int. Conf. Artificial Neural Networks*, 1999, pp. 575–580.
- [19] A. J. Smola, N. Murata, B. Schölkopf, and K.-R. Müller, "Asymptotically optimal choice of ϵ -loss for support vector machines," in *Proc. Int. Conf. Artificial Neural Networks*, 1998, pp. 105–110.
- [20] A. J. Smola and B. Schölkopf, "A Tutorial on Support Vector Regression," Royal Holloway College, London, U.K., NeuroCOLT2 Tech. Rep. NC2-TR-1998-030, 1998.
- [21] A. J. Smola, B. Schölkopf, and K.-R. Müller, "General cost functions for support vector regression," in *Proc. Australian Congr. Neural Networks*, 1998, pp. 78–83.
- [22] M. Stitson, A. Gammerman, V. N. Vapnik, V. Vovk, C. Watkins, and J. Weston, "Support vector regression with ANOVA decomposition kernels," in *Advances in Kernel Methods—Support Vector Learning*, B. Schölkopf, C. Burges, and A. Smola, Eds. Cambridge, MA: MIT Press, 1999.
- [23] J. A. K. Suykens, "Least squares support vector machines for classification and nonlinear modeling," in *Proc. 7th Int. Workshop Parallel Applications Statistics Economics*, 2000, pp. 29–48.
- [24] V. Vapnik, *Statistical Learning Theory*. New York: Wiley, 1998.



James T. Kwok received the Ph.D. degree in computer science from the Hong Kong University of Science and Technology in 1996. He then joined the Department of Computer Science, Hong Kong Baptist University as an Assistant Professor. He returned to the Hong Kong University of Science and Technology in 2000 and is now an Assistant Professor in the Department of Computer Science. His research interests include kernel methods, artificial neural networks, pattern recognition and machine learning.



Ivor W. Tsang received his B.Eng. degree in computer science from the Hong Kong University of Science and Technology (HKUST) in 2001. He is currently a master student in HKUST. He was the Honor Outstanding student in 2001. His research interests include machine learning and kernel methods.