

## New Support Vector Algorithms

**Bernhard Schölkopf\***

**Alex J. Smola**

*GMD FIRST, 12489 Berlin, Germany, and Department of Engineering, Australian National University, Canberra 0200, Australia*

**Robert C. Williamson**

*Department of Engineering, Australian National University, Canberra 0200, Australia*

**Peter L. Bartlett**

*RSISE, Australian National University, Canberra 0200, Australia*

We propose a new class of support vector algorithms for regression and classification. In these algorithms, a parameter  $\nu$  lets one effectively control the number of support vectors. While this can be useful in its own right, the parameterization has the additional benefit of enabling us to eliminate one of the other free parameters of the algorithm: the accuracy parameter  $\varepsilon$  in the regression case, and the regularization constant  $C$  in the classification case. We describe the algorithms, give some theoretical results concerning the meaning and the choice of  $\nu$ , and report experimental results.

### 1 Introduction ---

Support vector (SV) machines comprise a new class of learning algorithms, motivated by results of statistical learning theory (Vapnik, 1995). Originally developed for pattern recognition (Vapnik & Chervonenkis, 1974; Boser, Guyon, & Vapnik, 1992), they represent the decision boundary in terms of a typically small subset (Schölkopf, Burges, & Vapnik, 1995) of all training examples, called the support vectors. In order for this sparseness property to carry over to the case of SV Regression, Vapnik devised the so-called  $\varepsilon$ -insensitive loss function,

$$|y - f(\mathbf{x})|_\varepsilon = \max\{0, |y - f(\mathbf{x})| - \varepsilon\}, \quad (1.1)$$

which does not penalize errors below some  $\varepsilon > 0$ , chosen a priori. His algorithm, which we will henceforth call  $\varepsilon$ -SVR, seeks to estimate functions,

$$f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b, \quad \mathbf{w}, \mathbf{x} \in \mathbf{R}^N, b \in \mathbf{R}, \quad (1.2)$$

---

\* Present address: Microsoft Research, 1 Guildhall Street, Cambridge, U.K.

based on independent and identically distributed (i.i.d.) data,

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell) \in \mathbf{R}^N \times \mathbf{R}. \quad (1.3)$$

Here,  $\mathbf{R}^N$  is the space in which the input patterns live but most of the following also applies for inputs from a set  $\mathcal{X}$ . The goal of the learning process is to find a function  $f$  with a small risk (or test error),

$$R[f] = \int_{\mathcal{X}} l(f, \mathbf{x}, y) dP(\mathbf{x}, y), \quad (1.4)$$

where  $P$  is the probability measure, which is assumed to be responsible for the generation of the observations (see equation 1.3) and  $l$  is a loss function, for example,  $l(f, \mathbf{x}, y) = (f(\mathbf{x}) - y)^2$ , or many other choices (Smola & Schölkopf, 1998). The particular loss function for which we would like to minimize equation 1.4 depends on the specific regression estimation problem at hand. This does not necessarily have to coincide with the loss function used in our learning algorithm. First, there might be additional constraints that we would like our regression estimation to satisfy, for instance, that it have a sparse representation in terms of the training data. In the SV case, this is achieved through the insensitive zone in equation 1.1. Second, we cannot minimize equation 1.4 directly in the first place, since we do not know  $P$ . Instead, we are given the sample, equation 1.3, and we try to obtain a small risk by minimizing the regularized risk functional,

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \cdot R_{emp}^\varepsilon[f]. \quad (1.5)$$

Here,  $\|\mathbf{w}\|^2$  is a term that characterizes the model complexity,

$$R_{emp}^\varepsilon[f] := \frac{1}{\ell} \sum_{i=1}^{\ell} |y_i - f(\mathbf{x}_i)|_\varepsilon, \quad (1.6)$$

measures the  $\varepsilon$ -insensitive training error, and  $C$  is a constant determining the trade-off. In short, minimizing equation 1.5 captures the main insight of statistical learning theory, stating that in order to obtain a small risk, one needs to control both training error and model complexity—that is, explain the data with a simple model.

The minimization of equation 1.5 is equivalent to the following constrained optimization problem (see Figure 1):

$$\text{minimize } \tau(\mathbf{w}, \boldsymbol{\xi}^{(*)}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \cdot \frac{1}{\ell} \sum_{i=1}^{\ell} (\xi_i + \xi_i^*), \quad (1.7)$$

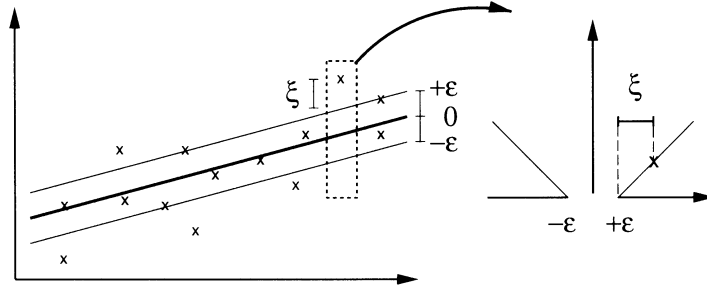


Figure 1: In SV regression, a desired accuracy  $\varepsilon$  is specified a priori. It is then attempted to fit a tube with radius  $\varepsilon$  to the data. The trade-off between model complexity and points lying outside the tube (with positive slack variables  $\xi$ ) is determined by minimizing the expression 1.5.

$$\text{subject to } ((\mathbf{w} \cdot \mathbf{x}_i) + b) - y_i \leq \varepsilon + \xi_i \quad (1.8)$$

$$y_i - ((\mathbf{w} \cdot \mathbf{x}_i) + b) \leq \varepsilon + \xi_i^* \quad (1.9)$$

$$\xi_i^{(*)} \geq 0. \quad (1.10)$$

Here and below, it is understood that  $i = 1, \dots, \ell$ , and that boldface Greek letters denote  $\ell$ -dimensional vectors of the corresponding variables;  $(*)$  is a shorthand implying both the variables with and without asterisks.

By using Lagrange multiplier techniques, one can show (Vapnik, 1995) that this leads to the following dual optimization problem. Maximize

$$\begin{aligned} W(\alpha, \alpha^*) = & -\varepsilon \sum_{i=1}^{\ell} (\alpha_i^* + \alpha_i) + \sum_{i=1}^{\ell} (\alpha_i^* - \alpha_i) y_i \\ & - \frac{1}{2} \sum_{i,j=1}^{\ell} (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j)(\mathbf{x}_i \cdot \mathbf{x}_j) \end{aligned} \quad (1.11)$$

$$\text{subject to } \sum_{i=1}^{\ell} (\alpha_i - \alpha_i^*) = 0 \quad (1.12)$$

$$\alpha_i^{(*)} \in \left[0, \frac{C}{\ell}\right]. \quad (1.13)$$

The resulting regression estimates are linear; however, the setting can be generalized to a nonlinear one by using the kernel method. As we will use precisely this method in the next section, we shall omit its exposition at this point.

To motivate the new algorithm that we shall propose, note that the parameter  $\varepsilon$  can be useful if the desired accuracy of the approximation can be specified beforehand. In some cases, however, we want the estimate to be as accurate as possible without having to commit ourselves to a specific level of accuracy a priori. In this work, we first describe a modification of the  $\varepsilon$ -SVR algorithm, called  $\nu$ -SVR, which automatically minimizes  $\varepsilon$ . Following this, we present two theoretical results on  $\nu$ -SVR concerning the connection to robust estimators (section 3) and the asymptotically optimal choice of the parameter  $\nu$  (section 4). Next, we extend the algorithm to handle parametric insensitivity models that allow taking into account prior knowledge about heteroscedasticity of the noise. As a bridge connecting this first theoretical part of the article to the second one, we then present a definition of a *margin* that both SV classification and SV regression algorithms maximize (section 6). In view of this close connection between both algorithms, it is not surprising that it is possible to formulate also a  $\nu$ -SV classification algorithm. This is done, including some theoretical analysis, in section 7. We conclude with experiments and a discussion.

## 2 $\nu$ -SV Regression

To estimate functions (see equation 1.2) from empirical data (see equation 1.3) we proceed as follows (Schölkopf, Bartlett, Smola, & Williamson, 1998). At each point  $\mathbf{x}_i$ , we allow an error of  $\varepsilon$ . Everything above  $\varepsilon$  is captured in slack variables  $\xi_i^{(*)}$ , which are penalized in the objective function via a regularization constant  $C$ , chosen a priori (Vapnik, 1995). The size of  $\varepsilon$  is traded off against model complexity and slack variables via a constant  $\nu \geq 0$ :

$$\text{minimize } \tau(\mathbf{w}, \boldsymbol{\xi}^{(*)}, \varepsilon) = \frac{1}{2} \|\mathbf{w}\|^2 + C \cdot \left( \nu \varepsilon + \frac{1}{\ell} \sum_{i=1}^{\ell} (\xi_i + \xi_i^*) \right) \quad (2.1)$$

$$\text{subject to } ((\mathbf{w} \cdot \mathbf{x}_i) + b) - y_i \leq \varepsilon + \xi_i \quad (2.2)$$

$$y_i - ((\mathbf{w} \cdot \mathbf{x}_i) + b) \leq \varepsilon + \xi_i^* \quad (2.3)$$

$$\xi_i^{(*)} \geq 0, \varepsilon \geq 0. \quad (2.4)$$

For the constraints, we introduce multipliers  $\alpha_i^{(*)}$ ,  $\eta_i^{(*)}$ ,  $\beta \geq 0$ , and obtain the Lagrangian

$$\begin{aligned} L(\mathbf{w}, b, \boldsymbol{\alpha}^{(*)}, \beta, \boldsymbol{\xi}^{(*)}, \varepsilon, \boldsymbol{\eta}^{(*)}) \\ = \frac{1}{2} \|\mathbf{w}\|^2 + C\nu\varepsilon + \frac{C}{\ell} \sum_{i=1}^{\ell} (\xi_i + \xi_i^*) - \beta\varepsilon - \sum_{i=1}^{\ell} (\eta_i \xi_i + \eta_i^* \xi_i^*) \end{aligned}$$

$$\begin{aligned}
& - \sum_{i=1}^{\ell} \alpha_i (\xi_i + y_i - (\mathbf{w} \cdot \mathbf{x}_i) - b + \varepsilon) \\
& - \sum_{i=1}^{\ell} \alpha_i^* (\xi_i^* + (\mathbf{w} \cdot \mathbf{x}_i) + b - y_i + \varepsilon).
\end{aligned} \tag{2.5}$$

To minimize the expression 2.1, we have to find the saddle point of  $L$ —that is, minimize over the primal variables  $\mathbf{w}$ ,  $\varepsilon$ ,  $b$ ,  $\xi_i^{(*)}$  and maximize over the dual variables  $\alpha_i^{(*)}$ ,  $\beta$ ,  $\eta_i^{(*)}$ . Setting the derivatives with respect to the primal variables equal to zero yields four equations:

$$\mathbf{w} = \sum_i (\alpha_i^* - \alpha_i) \mathbf{x}_i \tag{2.6}$$

$$C \cdot v - \sum_i (\alpha_i + \alpha_i^*) - \beta = 0 \tag{2.7}$$

$$\sum_{i=1}^{\ell} (\alpha_i - \alpha_i^*) = 0 \tag{2.8}$$

$$\frac{C}{\ell} - \alpha_i^{(*)} - \eta_i^{(*)} = 0. \tag{2.9}$$

In the SV expansion, equation 2.6, only those  $\alpha_i^{(*)}$  will be nonzero that correspond to a constraint, equations 2.2 or 2.3, which is precisely met; the corresponding patterns are called support vectors. This is due to the Karush-Kuhn-Tucker (KKT) conditions that apply to convex constrained optimization problems (Bertsekas, 1995). If we write the constraints as  $g(\mathbf{x}_i, y_i) \geq 0$ , with corresponding Lagrange multipliers  $\alpha_i$ , then the solution satisfies  $\alpha_i \cdot g(\mathbf{x}_i, y_i) = 0$  for all  $i$ .

Substituting the above four conditions into  $L$  leads to another optimization problem, called the Wolfe dual. Before stating it explicitly, we carry out one further modification. Following Boser et al. (1992), we substitute a kernel  $k$  for the dot product, corresponding to a dot product in some feature space related to input space via a nonlinear map  $\Phi$ ,

$$k(\mathbf{x}, \mathbf{y}) = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})). \tag{2.10}$$

By using  $k$ , we implicitly carry out all computations in the feature space that  $\Phi$  maps into, which can have a very high dimensionality. The feature space has the structure of a reproducing kernel Hilbert space (Wahba, 1999; Girosi, 1998; Schölkopf, 1997) and hence minimization of  $\|\mathbf{w}\|^2$  can be understood in the context of regularization operators (Smola, Schölkopf, & Müller, 1998).

The method is applicable whenever an algorithm can be cast in terms of dot products (Aizerman, Braverman, & Rozonoer, 1964; Boser et al., 1992; Schölkopf, Smola, & Müller, 1998). The choice of  $k$  is a research topic in its

own right that we shall not touch here (Williamson, Smola, & Schölkopf, 1998; Schölkopf, Shawe-Taylor, Smola, & Williamson, 1999); typical choices include gaussian kernels,  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/(2\sigma^2))$  and polynomial kernels,  $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^d$  ( $\sigma > 0, d \in \mathbb{N}$ ).

Rewriting the constraints, noting that  $\beta, \eta_i^{(*)} \geq 0$  do not appear in the dual, we arrive at the  $\nu$ -SVR optimization problem: for  $\nu \geq 0, C > 0$ ,

$$\begin{aligned} \text{maximize } W(\boldsymbol{\alpha}^{(*)}) &= \sum_{i=1}^{\ell} (\alpha_i^* - \alpha_i) y_i \\ &\quad - \frac{1}{2} \sum_{i,j=1}^{\ell} (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) k(\mathbf{x}_i, \mathbf{x}_j) \end{aligned} \quad (2.11)$$

$$\text{subject to } \sum_{i=1}^{\ell} (\alpha_i - \alpha_i^*) = 0 \quad (2.12)$$

$$\alpha_i^{(*)} \in \left[0, \frac{C}{\ell}\right] \quad (2.13)$$

$$\sum_{i=1}^{\ell} (\alpha_i + \alpha_i^*) \leq C \cdot \nu. \quad (2.14)$$

The regression estimate then takes the form (cf. equations 1.2, 2.6, and 2.10),

$$f(\mathbf{x}) = \sum_{i=1}^{\ell} (\alpha_i^* - \alpha_i) k(\mathbf{x}_i, \mathbf{x}) + b, \quad (2.15)$$

where  $b$  (and  $\varepsilon$ ) can be computed by taking into account that equations 2.2 and 2.3 (substitution of  $\sum_j (\alpha_j^* - \alpha_j) k(\mathbf{x}_j, \mathbf{x})$  for  $(\mathbf{w} \cdot \mathbf{x})$  is understood; cf. equations 2.6 and 2.10) become equalities with  $\xi_i^{(*)} = 0$  for points with  $0 < \alpha_i^{(*)} < C/\ell$ , respectively, due to the KKT conditions.

Before we give theoretical results explaining the significance of the parameter  $\nu$ , the following observation concerning  $\varepsilon$  is helpful. If  $\nu > 1$ , then  $\varepsilon = 0$ , since it does not pay to increase  $\varepsilon$  (cf. equation 2.1). If  $\nu \leq 1$ , it can still happen that  $\varepsilon = 0$ —for example, if the data are noise free and can be perfectly interpolated with a low-capacity model. The case  $\varepsilon = 0$ , however, is not what we are interested in; it corresponds to plain  $L_1$ -loss regression.

We will use the term errors to refer to training points lying outside the tube<sup>1</sup> and the term fraction of errors or SVs to denote the relative numbers of

<sup>1</sup> For  $N > 1$ , the “tube” should actually be called a slab—the region between two parallel hyperplanes.

errors or SVs (i.e., divided by  $\ell$ ). In this proposition, we define the modulus of absolute continuity of a function  $f$  as the function  $\epsilon(\delta) = \sup \sum_i |f(b_i) - f(a_i)|$ , where the supremum is taken over all disjoint intervals  $(a_i, b_i)$  with  $a_i < b_i$  satisfying  $\sum_i (b_i - a_i) < \delta$ . Loosely speaking, the condition on the conditional density of  $y$  given  $\mathbf{x}$  asks that it be absolutely continuous “on average.”

**Proposition 1.** *Suppose  $v$ -SVR is applied to some data set, and the resulting  $\epsilon$  is nonzero. The following statements hold:*

- i.  $v$  is an upper bound on the fraction of errors.
- ii.  $v$  is a lower bound on the fraction of SVs.
- iii. *Suppose the data (see equation 1.3) were generated i.i.d. from a distribution  $P(\mathbf{x}, y) = P(\mathbf{x})P(y|\mathbf{x})$  with  $P(y|\mathbf{x})$  continuous and the expectation of the modulus of absolute continuity of its density satisfying  $\lim_{\delta \rightarrow 0} \mathbf{E}\epsilon(\delta) = 0$ . With probability 1, asymptotically,  $v$  equals both the fraction of SVs and the fraction of errors.*

**Proof.** **Ad (i).** The constraints, equations 2.13 and 2.14, imply that at most a fraction  $v$  of all examples can have  $\alpha_i^{(*)} = C/\ell$ . All examples with  $\xi_i^{(*)} > 0$  (i.e., those outside the tube) certainly satisfy  $\alpha_i^{(*)} = C/\ell$  (if not,  $\alpha_i^{(*)}$  could grow further to reduce  $\xi_i^{(*)}$ ).

**Ad (ii).** By the KKT conditions,  $\epsilon > 0$  implies  $\beta = 0$ . Hence, equation 2.14 becomes an equality (cf. equation 2.7).<sup>2</sup> Since SVs are those examples for which  $0 < \alpha_i^{(*)} \leq C/\ell$ , the result follows (using  $\alpha_i \cdot \alpha_i^* = 0$  for all  $i$ ; Vapnik, 1995).

**Ad (iii).** The strategy of proof is to show that asymptotically, the probability of a point is lying *on* the edge of the tube vanishes. The condition on  $P(y|\mathbf{x})$  means that

$$\sup_{f,t} \mathbf{E}P(|f(\mathbf{x}) + t - y| < \gamma | \mathbf{x}) < \delta(\gamma) \quad (2.16)$$

for some function  $\delta(\gamma)$  that approaches zero as  $\gamma \rightarrow 0$ . Since the class of SV regression estimates  $f$  has well-behaved covering numbers, we have (Anthony & Bartlett, 1999, chap. 21) that for all  $t$ ,

$$\Pr \left( \sup_f \left( \hat{P}_\ell(|f(\mathbf{x}) + t - y| < \gamma/2) < P(|f(\mathbf{x}) + t - y| < \gamma) \right) > \alpha \right) < c_1 c_2^{-\ell},$$

<sup>2</sup> In practice, one can alternatively work with equation 2.14 as an equality constraint.

where  $\hat{P}_\ell$  is the sample-based estimate of  $P$  (that is, the proportion of points that satisfy  $|f(\mathbf{x}) - y + t| < \gamma$ ), and  $c_1, c_2$  may depend on  $\gamma$  and  $\alpha$ . Discretizing the values of  $t$ , taking the union bound, and applying equation 2.16 shows that the supremum over  $f$  and  $t$  of  $\hat{P}_\ell(f(\mathbf{x}) - y + t = 0)$  converges to zero in probability. Thus, the fraction of points on the edge of the tube almost surely converges to 0. Hence the fraction of SVs equals that of errors. Combining statements i and ii then shows that both fractions converge almost surely to  $\nu$ .

Hence,  $0 \leq \nu \leq 1$  can be used to control the number of errors (note that for  $\nu \geq 1$ , equation 2.13 implies 2.14, since  $\alpha_i \cdot \alpha_i^* = 0$  for all  $i$  (Vapnik, 1995)). Moreover, since the constraint, equation 2.12, implies that equation 2.14 is equivalent to  $\sum_i \alpha_i^{(*)} \leq C\nu/2$ , we conclude that proposition 1 actually holds for the upper and the lower edges of the tube separately, with  $\nu/2$  each. (Note that by the same argument, the number of SVs at the two edges of the standard  $\varepsilon$ -SVR tube asymptotically agree.)

A more intuitive, albeit somewhat informal, explanation can be given in terms of the primal objective function (see equation 2.1). At the point of the solution, note that if  $\varepsilon > 0$ , we must have  $(\partial/\partial\varepsilon)\tau(\mathbf{w}, \varepsilon) = 0$ , that is,  $\nu + (\partial/\partial\varepsilon)R_{emp}^\varepsilon = 0$ , hence  $\nu = -(\partial/\partial\varepsilon)R_{emp}^\varepsilon$ . This is greater than or equal to the fraction of errors, since the points outside the tube certainly do contribute to a change in  $R_{emp}$  when  $\varepsilon$  is changed. Points at the *edge* of the tube possibly might also contribute. This is where the inequality comes from.

Note that this does not contradict our freedom to choose  $\nu > 1$ . In that case,  $\varepsilon = 0$ , since it does not pay to increase  $\varepsilon$  (cf. equation 2.1).

Let us briefly discuss how  $\nu$ -SVR relates to  $\varepsilon$ -SVR (see section 1). Both algorithms use the  $\varepsilon$ -insensitive loss function, but  $\nu$ -SVR automatically computes  $\varepsilon$ . In a Bayesian perspective, this automatic adaptation of the loss function could be interpreted as adapting the error model, controlled by the hyperparameter  $\nu$ . Comparing equation 1.11 (substitution of a kernel for the dot product is understood) and equation 2.11, we note that  $\varepsilon$ -SVR requires an additional term,  $-\varepsilon \sum_{i=1}^\ell (\alpha_i^* + \alpha_i)$ , which, for fixed  $\varepsilon > 0$ , encourages that some of the  $\alpha_i^{(*)}$  will turn out to be 0. Accordingly, the constraint (see equation 2.14), which appears in  $\nu$ -SVR, is not needed. The primal problems, equations 1.7 and 2.1, differ in the term  $\nu\varepsilon$ . If  $\nu = 0$ , then the optimization can grow  $\varepsilon$  arbitrarily large; hence zero empirical risk can be obtained even when all  $\alpha$ s are zero.

In the following sense,  $\nu$ -SVR includes  $\varepsilon$ -SVR. Note that in the general case, using kernels,  $\bar{\mathbf{w}}$  is a vector in feature space.

**Proposition 2.** *If  $\nu$ -SVR leads to the solution  $\bar{\varepsilon}, \bar{\mathbf{w}}, \bar{b}$ , then  $\varepsilon$ -SVR with  $\varepsilon$  set a priori to  $\bar{\varepsilon}$ , and the same value of  $C$ , has the solution  $\bar{\mathbf{w}}, \bar{b}$ .*



**Proof.** If we minimize equation 2.1, then fix  $\varepsilon$  and minimize only over the remaining variables. The solution does not change.

### 3 The Connection to Robust Estimators

---

Using the  $\varepsilon$ -insensitive loss function, only the patterns outside the  $\varepsilon$ -tube enter the empirical risk term, whereas the patterns closest to the actual regression have zero loss. This, however, does not mean that it is only the outliers that determine the regression. In fact, the contrary is the case.

**Proposition 3** (resistance of SV regression). *Using support vector regression with the  $\varepsilon$ -insensitive loss function (see equation 1.1), local movements of target values of points outside the tube do not influence the regression.*

**Proof.** Shifting  $y_i$  locally does not change the status of  $(\mathbf{x}_i, y_i)$  as being a point outside the tube. Then the dual solution  $\alpha^{(*)}$  is still feasible; it satisfies the constraints (the point still has  $\alpha_i^{(*)} = C/\ell$ ). Moreover, the primal solution, with  $\xi_i$  transformed according to the movement of  $y_i$ , is also feasible. Finally, the KKT conditions are still satisfied, as still  $\alpha_i^{(*)} = C/\ell$ . Thus (Bertsekas, 1995),  $\alpha^{(*)}$  is still the optimal solution.

The proof relies on the fact that everywhere outside the tube, the upper bound on the  $\alpha_i^{(*)}$  is the same. This is precisely the case if the loss function increases linearly outside the  $\varepsilon$ -tube (cf. Huber, 1981, for requirements on robust cost functions). Inside, we could use various functions, with a derivative smaller than the one of the linear part.

For the case of the  $\varepsilon$ -insensitive loss, proposition 3 implies that essentially, the regression is a generalization of an estimator for the mean of a random variable that does the following:

- Throws away the largest and smallest examples (a fraction  $\nu/2$  of either category; in section 2, it is shown that the sum constraint, equation 2.12, implies that proposition 1 can be applied separately for the two sides, using  $\nu/2$ ).
- Estimates the mean by taking the average of the two extremal ones of the remaining examples.

This resistance concerning outliers is close in spirit to robust estimators like the trimmed mean. In fact, we could get closer to the idea of the trimmed mean, which first throws away the largest and smallest points and then computes the mean of the remaining points, by using a quadratic loss inside the  $\varepsilon$ -tube. This would leave us with Huber's robust loss function.

Note, moreover, that the parameter  $\nu$  is related to the breakdown point of the corresponding robust estimator (Huber, 1981). Because it specifies the fraction of points that may be arbitrarily bad outliers,  $\nu$  is related to the fraction of some arbitrary distribution that may be added to a known noise model without the estimator failing.

Finally, we add that by a simple modification of the loss function (White, 1994)—weighting the slack variables  $\xi^{(*)}$  above and below the tube in the target function, equation 2.1, by  $2\lambda$  and  $2(1-\lambda)$ , with  $\lambda \in [0, 1]$ —respectively—one can estimate generalized *quantiles*. The argument proceeds as follows. Asymptotically, all patterns have multipliers at bound (cf. proposition 1). The  $\lambda$ , however, changes the upper bounds in the box constraints applying to the two different types of slack variables to  $2C\lambda/\ell$  and  $2C(1-\lambda)/\ell$ , respectively. The equality constraint, equation 2.8, then implies that  $(1-\lambda)$  and  $\lambda$  give the fractions of points (out of those which are outside the tube) that lie on the two sides of the tube, respectively.

#### 4 Asymptotically Optimal Choice of $\nu$

Using an analysis employing tools of information geometry (Murata, Yoshizawa, & Amari, 1994; Smola, Murata, Schölkopf, & Müller, 1998), we can derive the asymptotically optimal  $\nu$  for a given class of noise models in the sense of maximizing the statistical efficiency.<sup>3</sup>

**Remark.** Denote  $\mathfrak{p}$  a density with unit variance,<sup>4</sup> and  $\mathfrak{P}$  a family of noise models generated from  $\mathfrak{P}$  by  $\mathfrak{P} := \{p|p = \frac{1}{\sigma}\mathfrak{P}(\frac{y}{\sigma}), \sigma > 0\}$ . Moreover assume that the data were generated i.i.d. from a distribution  $p(x, y) = p(x)p(y - f(x))$  with  $p(y - f(x))$  continuous. Under the assumption that SV regression produces an estimate  $\hat{f}$  converging to the underlying functional dependency  $f$ , the asymptotically optimal  $\nu$ , for the estimation-of-location-parameter model of SV regression described in Smola, Murata, Schölkopf, & Müller (1998), is

$$\nu = 1 - \int_{-\varepsilon}^{\varepsilon} \mathfrak{P}(t) dt \text{ where} \quad \varepsilon := \operatorname{argmin}_{\tau} \frac{1}{(\mathfrak{P}(-\tau) + \mathfrak{P}(\tau))^2} \left( 1 - \int_{-\tau}^{\tau} \mathfrak{P}(t) dt \right) \quad (4.1)$$

<sup>3</sup> This section assumes familiarity with some concepts of information geometry. A more complete explanation of the model underlying the argument is given in Smola, Murata, Schölkopf, & Müller (1998) and can be downloaded from <http://svm.first.gmd.de>.

<sup>4</sup>  $\mathfrak{p}$  is a prototype generating the class of densities  $\mathfrak{P}$ . Normalization assumptions are made for ease of notation.

To see this, note that under the assumptions stated above, the probability of a deviation larger than  $\varepsilon$ ,  $Pr\{|y - \hat{f}(x)| > \varepsilon\}$ , will converge to

$$\begin{aligned} Pr\{|y - f(x)| > \varepsilon\} &= \int_{\mathcal{X} \times \{\mathbb{R} \setminus [-\varepsilon, \varepsilon]\}} p(x)p(\xi) dx d\xi \\ &= 1 - \int_{-\varepsilon}^{\varepsilon} p(\xi) d\xi. \end{aligned} \quad (4.2)$$

This is also the fraction of samples that will (asymptotically) become SVs (proposition 1, iii). Therefore an algorithm generating a fraction  $\nu = 1 - \int_{-\varepsilon}^{\varepsilon} p(\xi) d\xi$  SVs will correspond to an algorithm with a tube of size  $\varepsilon$ . The consequence is that given a noise model  $p(\xi)$ , one can compute the optimal  $\varepsilon$  for it, and then, by using equation 4.2, compute the corresponding optimal value  $\nu$ .

To this end, one exploits the linear scaling behavior between the standard deviation  $\sigma$  of a distribution  $p$  and the optimal  $\varepsilon$ . This result, established in Smola, Murata, Schölkopf, & Müller (1998) and Smola (1998), cannot be proved here; instead, we shall merely try to give a flavor of the argument. The basic idea is to consider the estimation of a location parameter using the  $\varepsilon$ -insensitive loss, with the goal of maximizing the statistical efficiency. Using the Cramér-Rao bound and a result of Murata et al. (1994), the efficiency is found to be

$$e\left(\frac{\varepsilon}{\sigma}\right) = \frac{Q^2}{GI}. \quad (4.3)$$

Here,  $I$  is the Fisher information, while  $Q$  and  $G$  are information geometrical quantities computed from the loss function and the noise model.

This means that one only has to consider distributions of unit variance, say,  $\mathfrak{P}$ , to compute an optimal value of  $\nu$  that holds for the whole class of distributions  $\mathfrak{P}$ . Using equation 4.3, one arrives at

$$\frac{1}{e(\varepsilon)} \propto \frac{G}{Q^2} = \frac{1}{(\mathfrak{P}(-\varepsilon) + \mathfrak{P}(\varepsilon))^2} \left(1 - \int_{-\varepsilon}^{\varepsilon} \mathfrak{P}(t) dt\right). \quad (4.4)$$

The minimum of equation 4.4 yields the optimal choice of  $\varepsilon$ , which allows computation of the corresponding  $\nu$  and thus leads to equation 4.1.

Consider now an example: arbitrary polynomial noise models ( $\propto e^{-|\xi|^p}$ ) with unit variance can be written as

$$\mathfrak{P}(\xi) = \frac{1}{2} \sqrt{\frac{\Gamma(3/p)}{\Gamma(1/p)}} \frac{p}{\Gamma(1/p)} \exp\left(-\left(\sqrt{\frac{\Gamma(3/p)}{\Gamma(1/p)}} |\xi|\right)^p\right) \quad (4.5)$$

where  $\Gamma$  denotes the gamma function. Table 1 shows the optimal value of  $\nu$  for different polynomial degrees. Observe that the more “lighter-tailed”

Table 1: Optimal  $\nu$  for Various Degrees of Polynomial Additive Noise.

Polynomial degree $p$	1	2	3	4	5	6	7	8
Optimal $\nu$	1.00	0.54	0.29	0.19	0.14	0.11	0.09	0.07

the distribution becomes, the smaller  $\nu$  are optimal—that is, the tube width increases. This is reasonable as only for very long tails of the distribution (data with many outliers) it appears reasonable to use an early cutoff on the influence of the data (by basically giving all data equal influence via  $\alpha_i = C/\ell$ ). The extreme case of Laplacian noise ( $\nu = 1$ ) leads to a tube width of 0, that is, to  $L_1$  regression.

We conclude this section with three caveats: first, we have only made an asymptotic statement; second, for nonzero  $\varepsilon$ , the SV regression need not necessarily converge to the target  $f$ : measured using  $|\cdot|_\varepsilon$ , many other functions are just as good as  $f$  itself; third, the proportionality between  $\varepsilon$  and  $\sigma$  has only been established in the estimation-of-location-parameter context, which is not quite SV regression.

## 5 Parametric Insensitivity Models

We now return to the algorithm described in section 2. We generalized  $\varepsilon$ -SVR by estimating the width of the tube rather than taking it as given a priori. What we have so far retained is the assumption that the  $\varepsilon$ -insensitive zone has a tube (or slab) shape. We now go one step further and use parametric models of arbitrary shape. This can be useful in situations where the noise is heteroscedastic, that is, where it depends on  $\mathbf{x}$ .

Let  $\{\zeta_q^{(*)}\}$  (here and below,  $q = 1, \dots, p$  is understood) be a set of  $2p$  positive functions on the input space  $\mathcal{X}$ . Consider the following quadratic program: for given  $\nu_1^{(*)}, \dots, \nu_p^{(*)} \geq 0$ , minimize

$$\begin{aligned} \tau(\mathbf{w}, \boldsymbol{\xi}^{(*)}, \boldsymbol{\varepsilon}^{(*)}) = & \|\mathbf{w}\|^2/2 \\ & + C \cdot \left( \sum_{q=1}^p (\nu_q \varepsilon_q + \nu_q^* \varepsilon_q^*) + \frac{1}{\ell} \sum_{i=1}^{\ell} (\xi_i + \xi_i^*) \right) \end{aligned} \quad (5.1)$$

$$\text{subject to} \quad ((\mathbf{w} \cdot \mathbf{x}_i) + b) - y_i \leq \sum_q \varepsilon_q \zeta_q(\mathbf{x}_i) + \xi_i \quad (5.2)$$

$$y_i - ((\mathbf{w} \cdot \mathbf{x}_i) + b) \leq \sum_q \varepsilon_q^* \zeta_q^*(\mathbf{x}_i) + \xi_i^* \quad (5.3)$$

$$\xi_i^{(*)} \geq 0, \varepsilon_q^{(*)} \geq 0. \quad (5.4)$$

A calculation analogous to that in section 2 shows that the Wolfe dual consists of maximizing the expression 2.11 subject to the constraints 2.12 and 2.13, and, instead of 2.14, the modified constraints, still linear in  $\alpha^{(*)}$ ,

$$\sum_{i=1}^{\ell} \alpha_i^{(*)} \zeta_q^{(*)}(\mathbf{x}_i) \leq C \cdot v_q^{(*)}. \quad (5.5)$$

In the experiments in section 8, we use a simplified version of this optimization problem, where we drop the term  $v_q^* \varepsilon_q^*$  from the objective function, equation 5.1, and use  $\varepsilon_q$  and  $\zeta_q$  in equation 5.3. By this, we render the problem symmetric with respect to the two edges of the tube. In addition, we use  $p = 1$ . This leads to the same Wolfe dual, except for the last constraint, which becomes (cf. equation 2.14),

$$\sum_{i=1}^{\ell} (\alpha_i + \alpha_i^*) \zeta(\mathbf{x}_i) \leq C \cdot v. \quad (5.6)$$

Note that the optimization problem of section 2 can be recovered by using the constant function  $\zeta \equiv 1$ .<sup>5</sup>

The advantage of this setting is that since the same  $v$  is used for both sides of the tube, the computation of  $\varepsilon$ ,  $b$  is straightforward: for instance, by solving a linear system, using two conditions as those described following equation 2.15. Otherwise, general statements are harder to make; the linear system can have a zero determinant, depending on whether the functions  $\zeta_p^{(*)}$ , evaluated on the  $\mathbf{x}_i$  with  $0 < \alpha_i^{(*)} < C/\ell$ , are linearly dependent. The latter occurs, for instance, if we use constant functions  $\zeta^{(*)} \equiv 1$ . In this case, it is pointless to use two different values  $v, v^*$ , for the constraint (see equation 2.12) then implies that *both* sums  $\sum_{i=1}^{\ell} \alpha_i^{(*)}$  will be bounded by  $C \cdot \min\{v, v^*\}$ . We conclude this section by giving, without proof, a generalization of proposition 1 to the optimization problem with constraint (see equation 5.6):

**Proposition 4.** *Suppose we run the above algorithm on a data set with the result that  $\varepsilon > 0$ . Then*

- i.  $\frac{v\ell}{\sum_i \zeta(\mathbf{x}_i)}$  is an upper bound on the fraction of errors.
- ii.  $\frac{v\ell}{\sum_i \zeta(\mathbf{x}_i)}$  is an upper bound on the fraction of SVs.

---

<sup>5</sup> Observe the similarity to semiparametric SV models (Smola, Frieß, & Schölkopf, 1999) where a modification of the expansion of  $f$  led to similar additional constraints. The important difference in the present setting is that the Lagrange multipliers  $\alpha_i$  and  $\alpha_i^*$  are treated equally and not with different signs, as in semiparametric modeling.

- iii. Suppose the data in equation 1.3 were generated i.i.d. from a distribution  $P(\mathbf{x}, y) = P(\mathbf{x})P(y|\mathbf{x})$  with  $P(y|\mathbf{x})$  continuous and the expectation of its modulus of continuity satisfying  $\lim_{\delta \rightarrow 0} \mathbf{E}\epsilon(\delta) = 0$ . With probability 1, asymptotically, the fractions of SVs and errors equal  $v \cdot (\int \xi(\mathbf{x}) d\tilde{P}(\mathbf{x}))^{-1}$ , where  $\tilde{P}$  is the asymptotic distribution of SVs over  $\mathbf{x}$ .

## 6 Margins in Regression and Classification

The SV algorithm was first proposed for the case of pattern recognition (Boser et al., 1992), and then generalized to regression (Vapnik, 1995). Conceptually, however, one can take the view that the latter case is actually the simpler one, providing a posterior justification as to why we started this article with the regression case. To explain this, we will introduce a suitable definition of a *margin* that is maximized in both cases.

At first glance, the two variants of the algorithm seem conceptually different. In the case of pattern recognition, a margin of separation between two pattern classes is maximized, and the SVs are those examples that lie closest to this margin. In the simplest case, where the training error is fixed to 0, this is done by minimizing  $\|\mathbf{w}\|^2$  subject to  $y_i \cdot ((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1$  (note that in pattern recognition, the targets  $y_i$  are in  $\{\pm 1\}$ ).

In regression estimation, on the other hand, a tube of radius  $\varepsilon$  is fitted to the data, in the space of the target values, with the property that it corresponds to the flattest function in feature space. Here, the SVs lie at the edge of the tube. The parameter  $\varepsilon$  does not occur in the pattern recognition case.

We will show how these seemingly different problems are identical (cf. also Vapnik, 1995; Pontil, Rifkin, & Evgeniou, 1999), how this naturally leads to the concept of *canonical* hyperplanes (Vapnik, 1995), and how it suggests different generalizations to the estimation of vector-valued functions.

**Definition 1** ( $\varepsilon$ -margin). Let  $(E, \|\cdot\|_E)$ ,  $(F, \|\cdot\|_F)$  be normed spaces, and  $\mathcal{X} \subset E$ . We define the  $\varepsilon$ -margin of a function  $f: \mathcal{X} \rightarrow F$  as

$$m_\varepsilon(f) := \inf\{\|\mathbf{x} - \mathbf{y}\|_E: \mathbf{x}, \mathbf{y} \in \mathcal{X}, \|f(\mathbf{x}) - f(\mathbf{y})\|_F \geq 2\varepsilon\}. \quad (6.1)$$

$m_\varepsilon(f)$  can be zero, even for continuous functions, an example being  $f(x) = 1/x$  on  $\mathcal{X} = \mathbf{R}^+$ . There,  $m_\varepsilon(f) = 0$  for all  $\varepsilon > 0$ .

Note that the  $\varepsilon$ -margin is related (albeit not identical) to the traditional modulus of continuity of a function: given  $\delta > 0$ , the latter measures the largest difference in function values that can be obtained using points within a distance  $\delta$  in  $E$ .

The following observations characterize the functions for which the margin is strictly positive.

**Lemma 1** (uniformly continuous functions). With the above notations,  $m_\varepsilon(f)$  is positive for all  $\varepsilon > 0$  if and only if  $f$  is uniformly continuous.

**Proof.** By definition of  $m_\varepsilon$ , we have

$$(\|f(\mathbf{x}) - f(\mathbf{y})\|_F \geq 2\varepsilon \implies \|\mathbf{x} - \mathbf{y}\|_E \geq m_\varepsilon(f)) \quad (6.2)$$

$$\iff (\|\mathbf{x} - \mathbf{y}\|_E < m_\varepsilon(f) \implies \|f(\mathbf{x}) - f(\mathbf{y})\|_F < 2\varepsilon), \quad (6.3)$$

that is, if  $m_\varepsilon(f) > 0$ , then  $f$  is uniformly continuous. Similarly, if  $f$  is uniformly continuous, then for each  $\varepsilon > 0$ , we can find a  $\delta > 0$  such that  $\|f(\mathbf{x}) - f(\mathbf{y})\|_F \geq 2\varepsilon$  implies  $\|\mathbf{x} - \mathbf{y}\|_E \geq \delta$ . Since the latter holds uniformly, we can take the infimum to get  $m_\varepsilon(f) \geq \delta > 0$ .

We next specialize to a particular set of uniformly continuous functions.

**Lemma 2** (Lipschitz-continuous functions). *If there exists some  $L > 0$  such that for all  $\mathbf{x}, \mathbf{y} \in \mathcal{X}$ ,  $\|f(\mathbf{x}) - f(\mathbf{y})\|_F \leq L \cdot \|\mathbf{x} - \mathbf{y}\|_E$ , then  $m_\varepsilon \geq \frac{2\varepsilon}{L}$ .*

**Proof.** Take the infimum over  $\|\mathbf{x} - \mathbf{y}\|_E \geq \frac{\|f(\mathbf{x}) - f(\mathbf{y})\|_F}{L} \geq \frac{2\varepsilon}{L}$ .

**Example 1** (SV regression estimation). *Suppose that  $E$  is endowed with a dot product  $(\cdot, \cdot)$  (generating the norm  $\|\cdot\|_E$ ). For linear functions (see equation 1.2) the margin takes the form  $m_\varepsilon(f) = \frac{2\varepsilon}{\|\mathbf{w}\|}$ . To see this, note that since  $|f(\mathbf{x}) - f(\mathbf{y})| = |(\mathbf{w} \cdot (\mathbf{x} - \mathbf{y}))|$ , the distance  $\|\mathbf{x} - \mathbf{y}\|$  will be smallest given  $|(\mathbf{w} \cdot (\mathbf{x} - \mathbf{y}))| = 2\varepsilon$ , when  $\mathbf{x} - \mathbf{y}$  is parallel to  $\mathbf{w}$  (due to Cauchy-Schwartz), i.e. if  $\mathbf{x} - \mathbf{y} = \pm 2\varepsilon \mathbf{w} / \|\mathbf{w}\|^2$ . In that case,  $\|\mathbf{x} - \mathbf{y}\| = 2\varepsilon / \|\mathbf{w}\|$ . For fixed  $\varepsilon > 0$ , maximizing the margin hence amounts to minimizing  $\|\mathbf{w}\|$ , as done in SV regression: in the simplest form (cf. equation 1.7 without slack variables  $\xi_i$ ) the training on data (equation 1.3) consists of minimizing  $\|\mathbf{w}\|^2$  subject to*

$$|f(\mathbf{x}_i) - y_i| \leq \varepsilon. \quad (6.4)$$

**Example 2** (SV pattern recognition; see Figure 2). *We specialize the setting of example 1 to the case where  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}$ . Then  $m_1(f) = \frac{2}{\|\mathbf{w}\|}$  is equal to the margin defined for Vapnik's canonical hyperplane (Vapnik, 1995). The latter is a way in which, given the data set  $\mathcal{X}$ , an oriented hyperplane in  $E$  can be uniquely expressed by a linear function (see equation 1.2) requiring that*

$$\min\{|f(\mathbf{x})| : \mathbf{x} \in \mathcal{X}\} = 1. \quad (6.5)$$

*Vapnik gives a bound on the VC-dimension of canonical hyperplanes in terms of  $\|\mathbf{w}\|$ . An optimal margin SV machine for pattern recognition can be constructed from data,*

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell) \in \mathcal{X} \times \{\pm 1\} \quad (6.6)$$

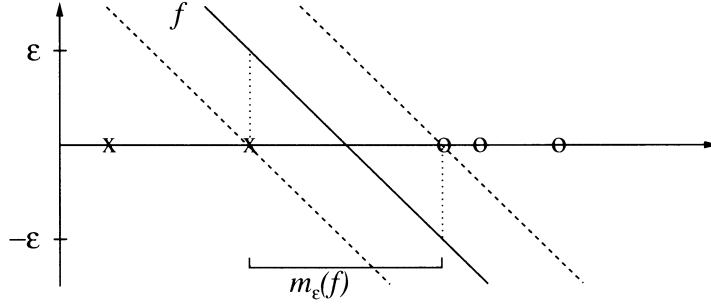


Figure 2: 1D toy problem. Separate x from o. The SV classification algorithm constructs a linear function  $f(x) = w \cdot x + b$  satisfying equation 6.5 ( $\epsilon = 1$ ). To maximize the margin  $m_\epsilon(f)$ , one has to minimize  $|w|$ .

as follows (Boser et al., 1992):

$$\text{minimize } \|\mathbf{w}\|^2 \text{ subject to } y_i \cdot f(\mathbf{x}_i) \geq 1. \quad (6.7)$$

The decision function used for classification takes the form

$$f^*(\mathbf{x}) = \text{sgn}((\mathbf{w} \cdot \mathbf{x}) + b). \quad (6.8)$$

The parameter  $\epsilon$  is superfluous in pattern recognition, as the resulting decision function,

$$f^*(\mathbf{x}) = \text{sgn}((\mathbf{w} \cdot \mathbf{x}) + b), \quad (6.9)$$

will not change if we minimize  $\|\mathbf{w}\|^2$  subject to

$$y_i \cdot f(\mathbf{x}_i) \geq \epsilon. \quad (6.10)$$

Finally, to understand why the constraint (see equation 6.7) looks different from equation 6.4 (e.g., one is multiplicative, the other one additive), note that in regression, the points  $(\mathbf{x}_i, y_i)$  are required to lie within a tube of radius  $\epsilon$ , whereas in pattern recognition, they are required to lie outside the tube (see Figure 2), and on the correct side. For the points on the tube, we have  $1 = y_i \cdot f(\mathbf{x}_i) = 1 - |f(\mathbf{x}_i) - y_i|$ .

So far, we have interpreted known algorithms only in terms of maximizing  $m_\epsilon$ . Next, we consider whether we can use the latter as a guide for constructing more general algorithms.

**Example 3** (SV regression for vector-valued functions). Assume  $E = \mathbf{R}^N$ . For linear functions  $f(\mathbf{x}) = W\mathbf{x} + \mathbf{b}$ , with  $W$  being an  $N \times N$  matrix, and  $\mathbf{b} \in \mathbf{R}^N$ ,



we have, as a consequence of lemma 1,

$$m_\varepsilon(f) \geq \frac{2\varepsilon}{\|W\|}, \quad (6.11)$$

where  $\|W\|$  is any matrix norm of  $W$  that is compatible (Horn & Johnson, 1985) with  $\|\cdot\|_E$ . If the matrix norm is the one induced by  $\|\cdot\|_E$ , that is, there exists a unit vector  $\mathbf{z} \in E$  such that  $\|W\mathbf{z}\|_E = \|W\|$ , then equality holds in 6.11. To see the latter, we use the same argument as in example 1, setting  $\mathbf{x} - \mathbf{y} = 2\varepsilon\mathbf{z}/\|W\|$ .

For the Hilbert-Schmidt norm  $\|W\|_2 = \sqrt{\sum_{i,j=1}^N W_{ij}^2}$ , which is compatible with the vector norm  $\|\cdot\|_2$ , the problem of minimizing  $\|W\|$  subject to separate constraints for each output dimension separates into  $N$  regression problems.

In Smola, Williamson, Mika, & Schölkopf (1999), it is shown that one can specify invariance requirements, which imply that the regularizers act on the output dimensions separately and identically (i.e., in a scalar fashion). In particular, it turns out that under the assumption of quadratic homogeneity and permutation symmetry, the Hilbert-Schmidt norm is the only admissible one.

## 7 $\nu$ -SV Classification

We saw that  $\nu$ -SVR differs from  $\varepsilon$ -SVR in that it uses the parameters  $\nu$  and  $C$  instead of  $\varepsilon$  and  $C$ . In many cases, this is a useful reparameterization of the original algorithm, and thus it is worthwhile to ask whether a similar change could be incorporated in the original SV classification algorithm (for brevity, we call it C-SVC). There, the primal optimization problem is to minimize (Cortes & Vapnik, 1995)

$$\tau(\mathbf{w}, \boldsymbol{\xi}) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{C}{\ell} \sum_i \xi_i \quad (7.1)$$

subject to

$$y_i \cdot ((\mathbf{x}_i \cdot \mathbf{w}) + b) \geq 1 - \xi_i, \quad \xi_i \geq 0. \quad (7.2)$$

The goal of the learning process is to estimate a function  $f^*$  (see equation 6.9) such that the probability of misclassification on an independent test set, the risk  $R[f^*]$ , is small.<sup>6</sup>

Here, the only parameter that we can dispose of is the regularization constant  $C$ . To substitute it by a parameter similar to the  $\nu$  used in the regression case, we proceed as follows. As a primal problem for  $\nu$ -SVC, we

<sup>6</sup> Implicitly we make use of the  $\{0, 1\}$  loss function; hence the risk equals the probability of misclassification.

consider the minimization of

$$\tau(\mathbf{w}, \boldsymbol{\xi}, \rho) = \frac{1}{2} \|\mathbf{w}\|^2 - \nu\rho + \frac{1}{\ell} \sum_i \xi_i \quad (7.3)$$

subject to (cf. equation 6.10)

$$y_i \cdot ((\mathbf{x}_i \cdot \mathbf{w}) + b) \geq \rho - \xi_i, \quad (7.4)$$

$$\xi_i \geq 0, \quad \rho \geq 0. \quad (7.5)$$

For reasons we shall explain, no constant  $C$  appears in this formulation. To understand the role of  $\rho$ , note that for  $\boldsymbol{\xi} = 0$ , the constraint (see 7.4) simply states that the two classes are separated by the *margin*  $2\rho/\|\mathbf{w}\|$  (cf. example 2).

To derive the dual, we consider the Lagrangian

$$\begin{aligned} L(\mathbf{w}, \boldsymbol{\xi}, b, \rho, \boldsymbol{\alpha}, \boldsymbol{\beta}, \delta) &= \frac{1}{2} \|\mathbf{w}\|^2 - \nu\rho + \frac{1}{\ell} \sum_i \xi_i \\ &\quad - \sum_i (\alpha_i (y_i ((\mathbf{x}_i \cdot \mathbf{w}) + b) - \rho + \xi_i) + \beta_i \xi_i) \\ &\quad - \delta\rho, \end{aligned} \quad (7.6)$$

using multipliers  $\alpha_i, \beta_i, \delta \geq 0$ . This function has to be minimized with respect to the primal variables  $\mathbf{w}, \boldsymbol{\xi}, b, \rho$  and maximized with respect to the dual variables  $\boldsymbol{\alpha}, \boldsymbol{\beta}, \delta$ . To eliminate the former, we compute the corresponding partial derivatives and set them to 0, obtaining the following conditions:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i \quad (7.7)$$

$$\alpha_i + \beta_i = 1/\ell, \quad 0 = \sum_i \alpha_i y_i, \quad \sum_i \alpha_i - \delta = \nu. \quad (7.8)$$

In the SV expansion (see equation 7.7), only those  $\alpha_i$  can be nonzero that correspond to a constraint (see 7.4) that is precisely met (KKT conditions; cf. Vapnik, 1995).

Substituting equations 7.7 and 7.8 into  $L$ , using  $\alpha_i, \beta_i, \delta \geq 0$ , and incorporating kernels for dot products leaves us with the following quadratic optimization problem: maximize

$$W(\boldsymbol{\alpha}) = -\frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \quad (7.9)$$

subject to

$$0 \leq \alpha_i \leq 1/\ell \quad (7.10)$$

$$0 = \sum_i \alpha_i y_i \quad (7.11)$$

$$\sum_i \alpha_i \geq \nu. \quad (7.12)$$

The resulting decision function can be shown to take the form

$$f^*(\mathbf{x}) = \text{sgn} \left( \sum_i \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i) + b \right). \quad (7.13)$$

Compared to the original dual (Boser et al., 1992; Vapnik, 1995), there are two differences. First, there is an additional constraint, 7.12, similar to the regression case, 2.14. Second, the linear term  $\sum_i \alpha_i$  of Boser et al. (1992) no longer appears in the objective function 7.9. This has an interesting consequence: 7.9 is now quadratically homogeneous in  $\alpha$ . It is straightforward to verify that one obtains exactly the same objective function if one starts with the primal function  $\tau(\mathbf{w}, \xi, \rho) = \|\mathbf{w}\|^2/2 + C \cdot (-\nu\rho + (1/\ell) \sum_i \xi_i)$  (i.e., if one does use C), the only difference being that the constraints, 7.10 and 7.12 would have an extra factor C on the right-hand side. In that case, due to the homogeneity, the solution of the dual would be scaled by C; however, it is straightforward to see that the corresponding decision function will not change. Hence we may set  $C = 1$ .

To compute  $b$  and  $\rho$ , we consider two sets  $S_{\pm}$ , of identical size  $s > 0$ , containing SVs  $\mathbf{x}_i$  with  $0 < \alpha_i < 1$  and  $y_i = \pm 1$ , respectively. Then, due to the KKT conditions, 7.4 becomes an equality with  $\xi_i = 0$ . Hence, in terms of kernels,

$$b = -\frac{1}{2s} \sum_{\mathbf{x} \in S_+ \cup S_-} \sum_j \alpha_j y_j k(\mathbf{x}, \mathbf{x}_j), \quad (7.14)$$

$$\rho = \frac{1}{2s} \left( \sum_{\mathbf{x} \in S_+} \sum_j \alpha_j y_j k(\mathbf{x}, \mathbf{x}_j) - \sum_{\mathbf{x} \in S_-} \sum_j \alpha_j y_j k(\mathbf{x}, \mathbf{x}_j) \right). \quad (7.15)$$

As in the regression case, the  $\nu$  parameter has a more natural interpretation than the one we removed,  $C$ . To formulate it, let us first define the term *margin error*. By this, we denote points with  $\xi_i > 0$ —that is, that are either errors or lie within the margin. Formally, the fraction of margin errors is

$$R_{emp}^\rho[f] := \frac{1}{\ell} |\{i: y_i \cdot f(\mathbf{x}_i) < \rho\}|. \quad (7.16)$$

Here,  $f$  is used to denote the argument of the  $\text{sgn}$  in the decision function, equation 7.13, that is,  $f^* = \text{sgn} \circ f$ .

We are now in a position to modify proposition 1 for the case of pattern recognition:

**Proposition 5.** *Suppose  $k$  is a real analytic kernel function, and we run  $v$ -SVC with  $k$  on some data with the result that  $\rho > 0$ . Then*

- i.  *$v$  is an upper bound on the fraction of margin errors.*
- ii.  *$v$  is a lower bound on the fraction of SVs.*
- iii. *Suppose the data (see equation 6.6) were generated i.i.d. from a distribution  $P(\mathbf{x}, y) = P(\mathbf{x})P(y|\mathbf{x})$  such that neither  $P(\mathbf{x}, y = 1)$  nor  $P(\mathbf{x}, y = -1)$  contains any discrete component. Suppose, moreover, that the kernel is analytic and non-constant. With probability 1, asymptotically,  $v$  equals both the fraction of SVs and the fraction of errors.*

**Proof.** **Ad (i).** By the KKT conditions,  $\rho > 0$  implies  $\delta = 0$ . Hence, inequality 7.12 becomes an equality (cf. equations 7.8). Thus, at most a fraction  $v$  of all examples can have  $\alpha_i = 1/\ell$ . All examples with  $\xi_i > 0$  do satisfy  $\alpha_i = 1/\ell$  (if not,  $\alpha_i$  could grow further to reduce  $\xi_i$ ).

**Ad (ii).** SVs can contribute at most  $1/\ell$  to the left-hand side of 7.12; hence there must be at least  $v\ell$  of them.

**Ad (iii).** It follows from the condition on  $P(\mathbf{x}, y)$  that apart from some set of measure zero (arising from possible singular components), the two class distributions are absolutely continuous and can be written as integrals over distribution functions. Because the kernel is analytic and non-constant, it cannot be constant in any open set; otherwise it would be constant everywhere. Therefore, functions  $f$  constituting the argument of the  $\text{sgn}$  in the SV decision function (see equation 7.13) essentially functions in the class of SV regression functions) transform the distribution over  $\mathbf{x}$  into distributions such that for all  $f$ , and all  $t \in \mathbf{R}$ ,  $\lim_{\gamma \rightarrow 0} P(|f(\mathbf{x}) + t| < \gamma) = 0$ . At the same time, we know that the class of these functions has well-behaved covering numbers; hence we get uniform convergence: for all  $\gamma > 0$ ,  $\sup_f |P(|f(\mathbf{x}) + t| < \gamma) - \hat{P}_\ell(|f(\mathbf{x}) + t| < \gamma)|$  converges to zero in probability, where  $\hat{P}_\ell$  is the sample-based estimate of  $P$  (that is, the proportion of points that satisfy  $|f(\mathbf{x}) + t| < \gamma$ ). But then for all  $\alpha > 0$ ,  $\lim_{\gamma \rightarrow 0} \lim_{\ell \rightarrow \infty} P(\sup_f \hat{P}_\ell(|f(\mathbf{x}) + t| < \gamma) > \alpha) = 0$ . Hence,  $\sup_f \hat{P}_\ell(|f(\mathbf{x}) + t| = 0)$  converges to zero in probability. Using  $t = \pm\rho$  thus shows that almost surely the fraction of points exactly on the margin tends to zero; hence the fraction of SVs equals that of margin errors.

Combining (i) and (ii) shows that both fractions converge almost surely to  $\nu$ .

Moreover, since equation 7.11 means that the sums over the coefficients of positive and negative SVs respectively are equal, we conclude that proposition 5 actually holds for both classes separately, with  $\nu/2$ . (Note that by the same argument, the number of SVs at the two sides of the margin asymptotically agree.)

A connection to standard SV classification, and a somewhat surprising interpretation of the regularization parameter  $C$ , is described by the following result:

**Proposition 6.** *If  $\nu$ -SV classification leads to  $\rho > 0$ , then  $C$ -SV classification, with  $C$  set a priori to  $1/\rho$ , leads to the same decision function.*

**Proof.** If one minimizes the function 7.3 and then fixes  $\rho$  to minimize only over the remaining variables, nothing will change. Hence the obtained solution  $\mathbf{w}_0, b_0, \xi_0$  minimizes the function 7.1 for  $C = 1$ , subject to the constraint 7.4. To recover the constraint 7.2, we rescale to the set of variables  $\mathbf{w}' = \mathbf{w}/\rho, b' = b/\rho, \xi' = \xi/\rho$ . This leaves us, up to a constant scaling factor  $\rho^2$ , with the objective function 7.1 using  $C = 1/\rho$ .

As in the case of regression estimation (see proposition 3), linearity of the target function in the slack variables  $\xi^{(*)}$  leads to “outlier” resistance of the estimator in pattern recognition. The exact statement, however, differs from the one in regression in two respects. First, the perturbation of the point is carried out in feature space. What it precisely corresponds to in input space therefore depends on the specific kernel chosen. Second, instead of referring to points outside the  $\varepsilon$ -tube, it refers to margin error points—points that are misclassified or fall into the margin. Below, we use the shorthand  $\mathbf{z}_i$  for  $\Phi(\mathbf{x}_i)$ .

**Proposition 7** (resistance of SV classification). *Suppose  $\mathbf{w}$  can be expressed in terms of the SVs that are not at bound, that is,*

$$\mathbf{w} = \sum_i \gamma_i \mathbf{z}_i, \quad (7.17)$$

*with  $\gamma_i \neq 0$  only if  $\alpha_i \in (0, 1/\ell)$  (where the  $\alpha_i$  are the coefficients of the dual solution). Then local movements of any margin error  $\mathbf{z}_m$  parallel to  $\mathbf{w}$  do not change the hyperplane.*

**Proof.** Since the slack variable of  $\mathbf{z}_m$  satisfies  $\xi_m > 0$ , the KKT conditions (e.g., Bertsekas, 1995) imply  $\alpha_m = 1/\ell$ . If  $\delta$  is sufficiently small, then transforming the point into  $\mathbf{z}'_m := \mathbf{z}_m + \delta \cdot \mathbf{w}$  results in a slack that is still nonzero,

that is,  $\xi'_m > 0$ ; hence we have  $\alpha'_m = 1/\ell = \alpha_m$ . Updating the  $\xi_m$  and keeping all other primal variables unchanged, we obtain a modified set of primal variables that is still feasible.

We next show how to obtain a corresponding set of feasible dual variables. To keep  $\mathbf{w}$  unchanged, we need to satisfy

$$\sum_i \alpha_i y_i \mathbf{z}_i = \sum_{i \neq m} \alpha'_i y_i \mathbf{z}_i + \alpha_m y_m \mathbf{z}'_m.$$

Substituting  $\mathbf{z}'_m = \mathbf{z}_m + \delta \cdot \mathbf{w}$  and equation 7.17, we note that a sufficient condition for this to hold is that for all  $i \neq m$ ,

$$\alpha'_i = \alpha_i - \delta \gamma_i y_i \alpha_m y_m.$$

Since by assumption  $\gamma_i$  is nonzero only if  $\alpha_i \in (0, 1/\ell)$ ,  $\alpha'_i$  will be in  $(0, 1/\ell)$  if  $\alpha_i$  is, provided  $\delta$  is sufficiently small, and it will equal  $1/\ell$  if  $\alpha_i$  does. In both cases, we end up with a feasible solution  $\alpha'$ , and the KKT conditions are still satisfied. Thus (Bertsekas, 1995),  $(\mathbf{w}, b)$  are still the hyperplane parameters of the solution.

Note that the assumption (7.17) is not as restrictive as it may seem. Although the SV expansion of the solution,  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{z}_i$ , often contains many multipliers  $\alpha_i$  that are at bound, it is nevertheless conceivable that, especially when discarding the requirement that the coefficients be bounded, we can obtain an expansion (see equation 7.17) in terms of a subset of the original vectors. For instance, if we have a 2D problem that we solve directly in input space, with  $k(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})$ , then it already suffices to have two linearly independent SVs that are not at bound in order to express  $\mathbf{w}$ . This holds for any overlap of the two classes—even if there are many SVs at the upper bound.

For the selection of  $C$ , several methods have been proposed that could probably be adapted for  $\nu$  (Schölkopf, 1997; Shawe-Taylor & Cristianini, 1999). In practice, most researchers have so far used cross validation. Clearly, this could be done also for  $\nu$ -SVC. Nevertheless, we shall propose a method that takes into account specific properties of  $\nu$ -SVC.

The parameter  $\nu$  lets us control the number of margin errors, the crucial quantity in a class of bounds on the generalization error of classifiers using covering numbers to measure the classifier capacity. We can use this connection to give a generalization error bound for  $\nu$ -SVC in terms of  $\nu$ . There are a number of complications in doing this the best possible way, and so here we will indicate the simplest one. It is based on the following result:

**Proposition 8** (Bartlett, 1998). *Suppose  $\rho > 0$ ,  $0 < \delta < \frac{1}{2}$ ,  $P$  is a probability distribution on  $\mathcal{X} \times \{-1, 1\}$  from which the training set, equation 6.6, is drawn. Then with probability at least  $1 - \delta$  for every  $f$  in some function class  $\mathcal{F}$ , the*

probability of error of the classification function  $f^* = \text{sgn} \circ f$  on an independent test set is bounded according to

$$R[f^*] \leq R_{\text{emp}}^\rho[f] + \sqrt{\frac{2}{\ell} (\ln \mathcal{N}(\mathcal{F}, l_\infty^\ell, \rho/2) + \ln(2/\delta))}, \quad (7.18)$$

where  $\mathcal{N}(\mathcal{F}, l_\infty^\ell, \rho) = \sup_{\mathbf{x}_1, \dots, \mathbf{x}_\ell} \mathcal{N}(\mathcal{F}|_{\mathbf{x}}, l_\infty, \rho)$ ,  $\mathcal{F}|_{\mathbf{x}} = \{(f(\mathbf{x}_1), \dots, f(\mathbf{x}_\ell)) : f \in \mathcal{F}\}$ ,  $\mathcal{N}(\mathcal{F}_{\mathbf{x}}, l_\infty, \rho)$  is the  $\rho$ -covering number of  $\mathcal{F}_{\mathbf{x}}$  with respect to  $l_\infty$ , the usual  $l_\infty$  metric on a set of vectors.

To obtain the generalization bound for  $\nu$ -SVC, we simply substitute the bound  $R_{\text{emp}}^\rho[f] \leq \nu$  (proposition 5, i) and some estimate of the covering numbers in terms of the margin. The best available bounds are stated in terms of the functional inverse of  $\mathcal{N}$ , hence the slightly complicated expressions in the following.

**Proposition 9** (Williamson et al., 1998). Denote  $B_R$  the ball of radius  $R$  around the origin in some Hilbert space  $F$ . Then the covering number  $\mathcal{N}$  of the class of functions

$$\mathcal{F} = \{\mathbf{x} \mapsto (\mathbf{w} \cdot \mathbf{x}) : \|\mathbf{w}\| \leq 1, \mathbf{x} \in B_R\} \quad (7.19)$$

at scale  $\rho$  satisfies

$$\log_2 \mathcal{N}(\mathcal{F}, l_\infty^\ell, \rho) \leq \inf \left\{ n \left\lceil \frac{c^2 R^2}{\rho^2} \frac{1}{n} \log_2 \left( 1 + \frac{\ell}{n} \right) \geq 1 \right\rceil - 1, \right. \quad (7.20)$$

where  $c < 103$  is a constant.

This is a consequence of a fundamental theorem due to Maurey. For  $\ell \geq 2$  one thus obtains

$$\log_2 \mathcal{N}(\mathcal{F}, l_\infty^\ell, \rho) \leq \frac{c^2 R^2}{\rho^2} \log_2 \ell - 1. \quad (7.21)$$

To apply these results to  $\nu$ -SVC, we rescale  $\mathbf{w}$  to length 1, thus obtaining a margin  $\rho/\|\mathbf{w}\|$  (cf. equation 7.4). Moreover, we have to combine propositions 8 and 9. Using  $\rho/2$  instead of  $\rho$  in the latter yields the following result.

**Proposition 10.** Suppose  $\nu$ -SVC is used with a kernel of the form  $k(\mathbf{x}, \mathbf{y}) = k(\|\mathbf{x} - \mathbf{y}\|)$  with  $k(0) = 1$ . Then all the data points  $\Phi(\mathbf{x}_i)$  in feature space live in a ball of radius 1 centered at the origin. Consequently with probability at least  $1 - \delta$  over the training set (see equation 6.6), the  $\nu$ -SVC decision function  $f^* = \text{sgn} \circ f$ ,

with  $f(\mathbf{x}) = \sum_i \alpha_i y_i k(\mathbf{x}, \mathbf{x}_i)$  (cf. equation 7.13), has a probability of test error bounded according to

$$\begin{aligned} R[f^*] &\leq R_{\text{emp}}^\rho[f] + \sqrt{\frac{2}{\ell} \left( \frac{4c^2 \|\mathbf{w}\|^2}{\rho^2} \log_2(2\ell) - 1 + \ln(2/\delta) \right)} \\ &\leq \nu + \sqrt{\frac{2}{\ell} \left( \frac{4c^2 \|\mathbf{w}\|^2}{\rho^2} \log_2(2\ell) - 1 + \ln(2/\delta) \right)}. \end{aligned}$$

Notice that in general,  $\|\mathbf{w}\|$  is a vector in feature space.

Note that the set of functions in the proposition differs from support vector decision functions (see equation 7.13) in that it comes without the  $+b$  term. This leads to a minor modification (for details, see Williamson et al., 1998).

Better bounds can be obtained by estimating the radius or even optimizing the choice of the center of the ball (cf. the procedure described by Schölkopf et al. 1995; Burges, 1998). However, in order to get a theorem of the above form in that case, a more complex argument is necessary (see Shawe-Taylor, Bartlett, Williamson, & Anthony, 1998, sec. VI for an indication).

We conclude this section by noting that a straightforward extension of the  $\nu$ -SVC algorithm is to include parametric models  $\zeta_k(\mathbf{x})$  for the margin, and thus to use  $\sum_q \rho_q \zeta_q(\mathbf{x}_i)$  instead of  $\rho$  in the constraint (see equation 7.4)—in complete analogy to the regression case discussed in section 5.

## 8 Experiments

**8.1 Regression Estimation.** In the experiments, we used the optimizer LOQO.<sup>7</sup> This has the serendipitous advantage that the primal variables  $b$  and  $\varepsilon$  can be recovered as the dual variables of the Wolfe dual (see equation 2.11) (i.e., the double dual variables) fed into the optimizer.

**8.1.1 Toy Examples.** The first task was to estimate a noisy sinc function, given  $\ell$  examples  $(x_i, y_i)$ , with  $x_i$  drawn uniformly from  $[-3, 3]$ , and  $y_i = \sin(\pi x_i)/(\pi x_i) + v_i$ , where the  $v_i$  were drawn from a gaussian with zero mean and variance  $\sigma^2$ . Unless stated otherwise, we used the radial basis function (RBF) kernel  $k(x, x') = \exp(-|x - x'|^2)$ ,  $\ell = 50$ ,  $C = 100$ ,  $\nu = 0.2$ , and  $\sigma = 0.2$ . Whenever standard deviation error bars are given, the results were obtained from 100 trials. Finally, the risk (or test error) of a regression estimate  $f$  was computed with respect to the sinc function without noise, as  $\frac{1}{6} \int_{-3}^3 |f(x) - \sin(\pi x)/(\pi x)| dx$ . Results are given in Table 2 and Figures 3 through 9.

<sup>7</sup> Available online at <http://www.princeton.edu/~rvdb/>.



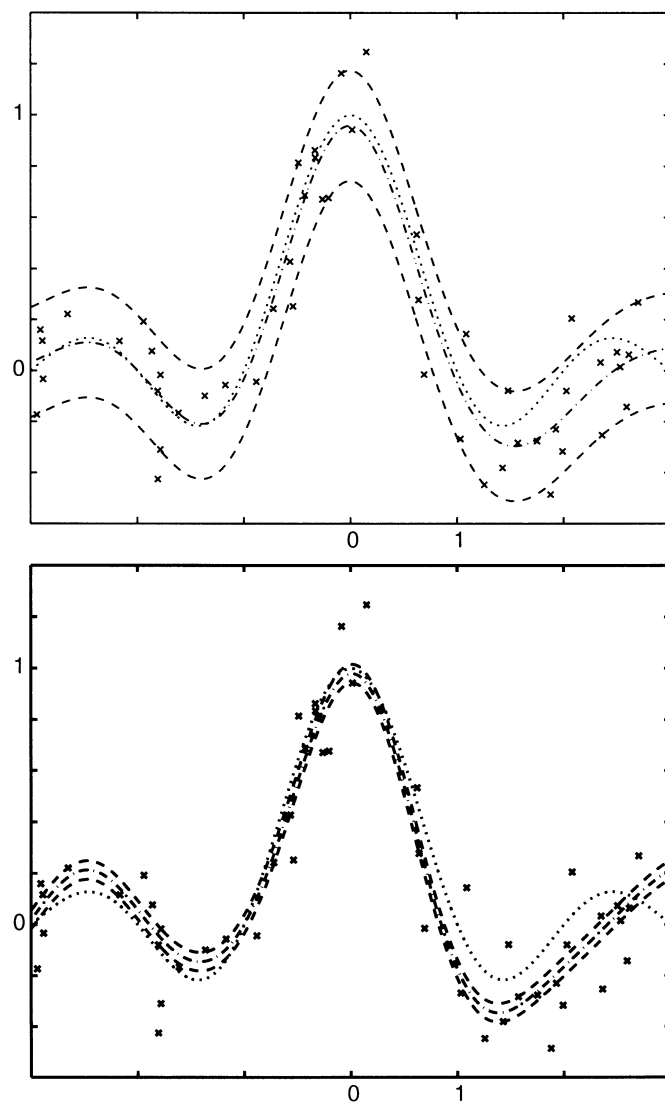


Figure 3:  $\nu$ -SV regression with  $\nu = 0.2$  (top) and  $\nu = 0.8$  (bottom). The larger  $\nu$  allows more points to lie outside the tube (see section 2). The algorithm automatically adjusts  $\varepsilon$  to 0.22 (top) and 0.04 (bottom). Shown are the sinc function (dotted), the regression  $f$ , and the tube  $f \pm \varepsilon$ .

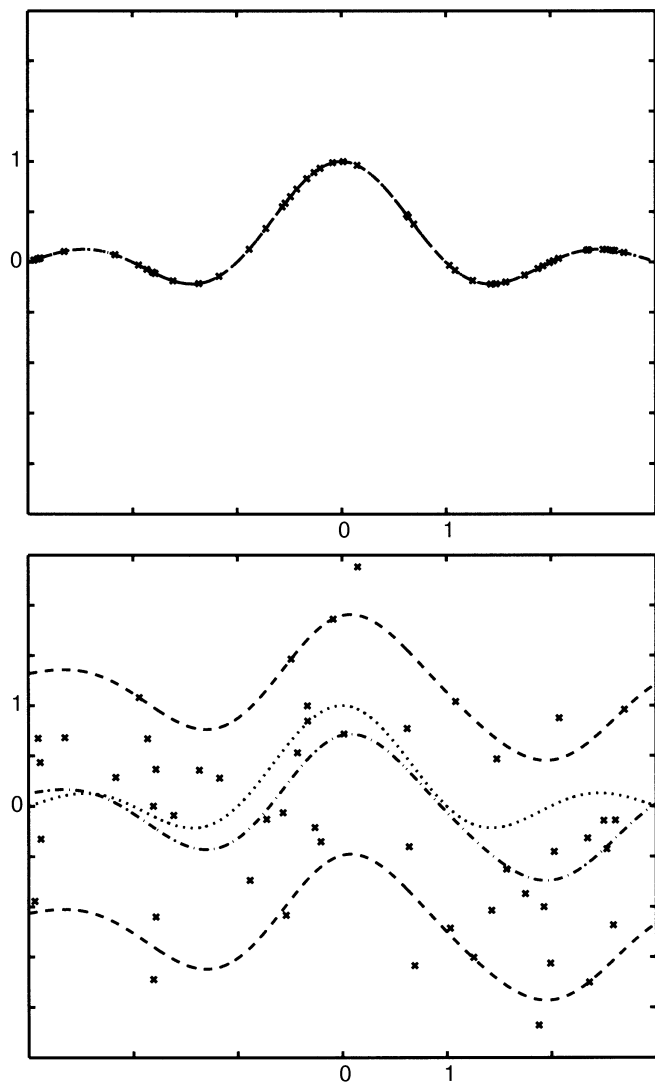


Figure 4:  $\nu$ -SV regression on data with noise  $\sigma = 0$  (top) and  $\sigma = 1$ ; (bottom). In both cases,  $\nu = 0.2$ . The tube width automatically adjusts to the noise (top:  $\varepsilon = 0$ ; bottom:  $\varepsilon = 1.19$ ).

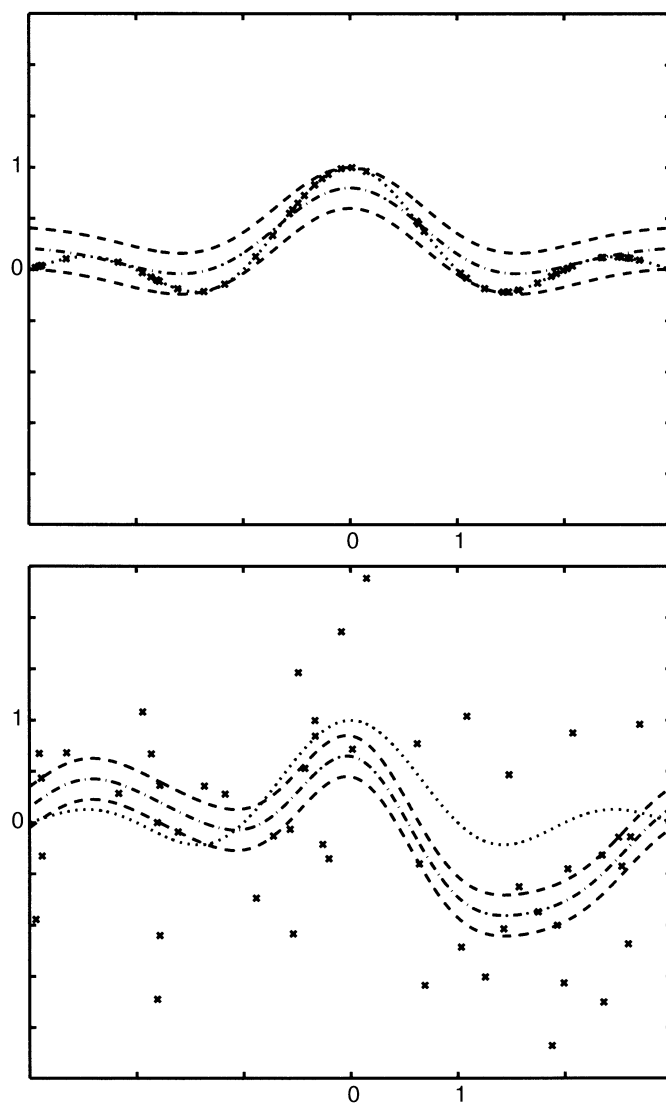


Figure 5:  $\varepsilon$ -SV regression (Vapnik, 1995) on data with noise  $\sigma = 0$  (top) and  $\sigma = 1$  (bottom). In both cases,  $\varepsilon = 0.2$ . This choice, which has to be specified a priori, is ideal for neither case. In the upper figure, the regression estimate is biased; in the lower figure,  $\varepsilon$  does not match the external noise (Smola, Murata, Schölkopf, & Müller, 1998).

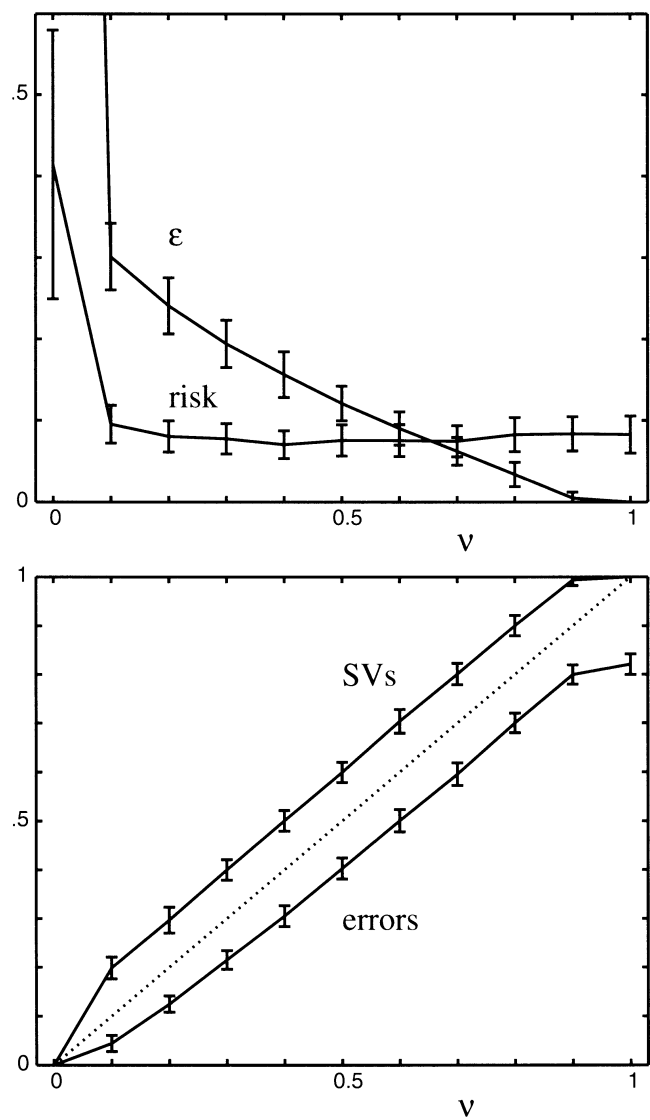


Figure 6:  $\nu$ -SVR for different values of the error constant  $\nu$ . Notice how  $\varepsilon$  decreases when more errors are allowed (large  $\nu$ ), and that over a large range of  $\nu$ , the test error (risk) is insensitive toward changes in  $\nu$ .

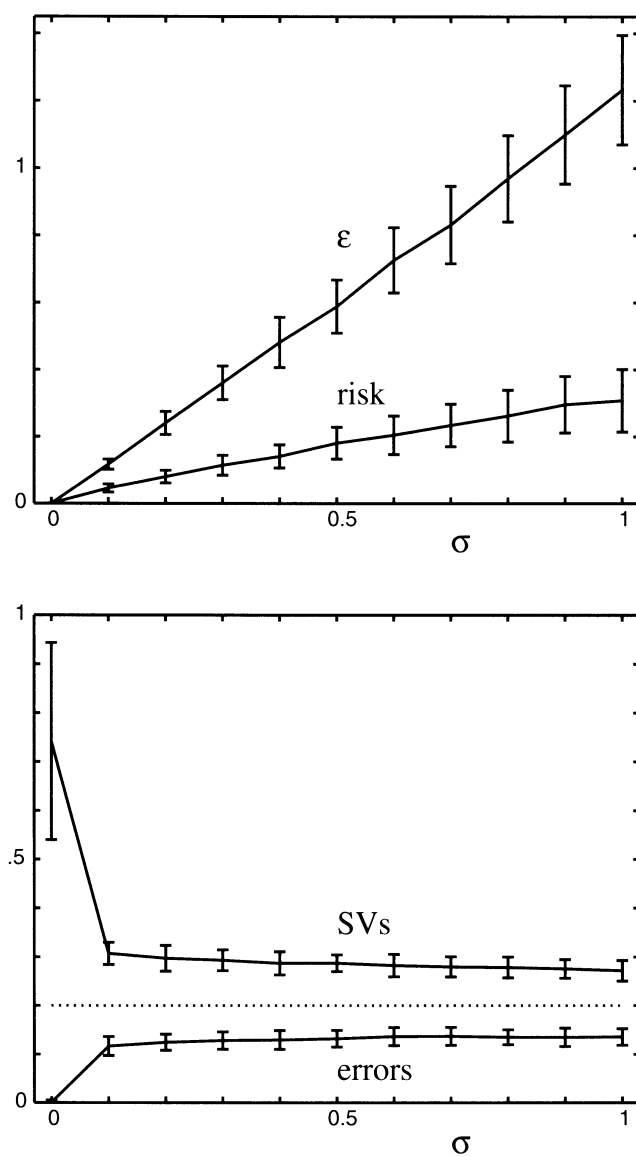


Figure 7:  $\nu$ -SVR for different values of the noise  $\sigma$ . The tube radius  $\varepsilon$  increases linearly with  $\sigma$  (largely due to the fact that both  $\varepsilon$  and the  $\xi_i^{(*)}$  enter the cost function linearly). Due to the automatic adaptation of  $\varepsilon$ , the number of SVs and points outside the tube (errors) are, except for the noise-free case  $\sigma = 0$ , largely independent of  $\sigma$ .

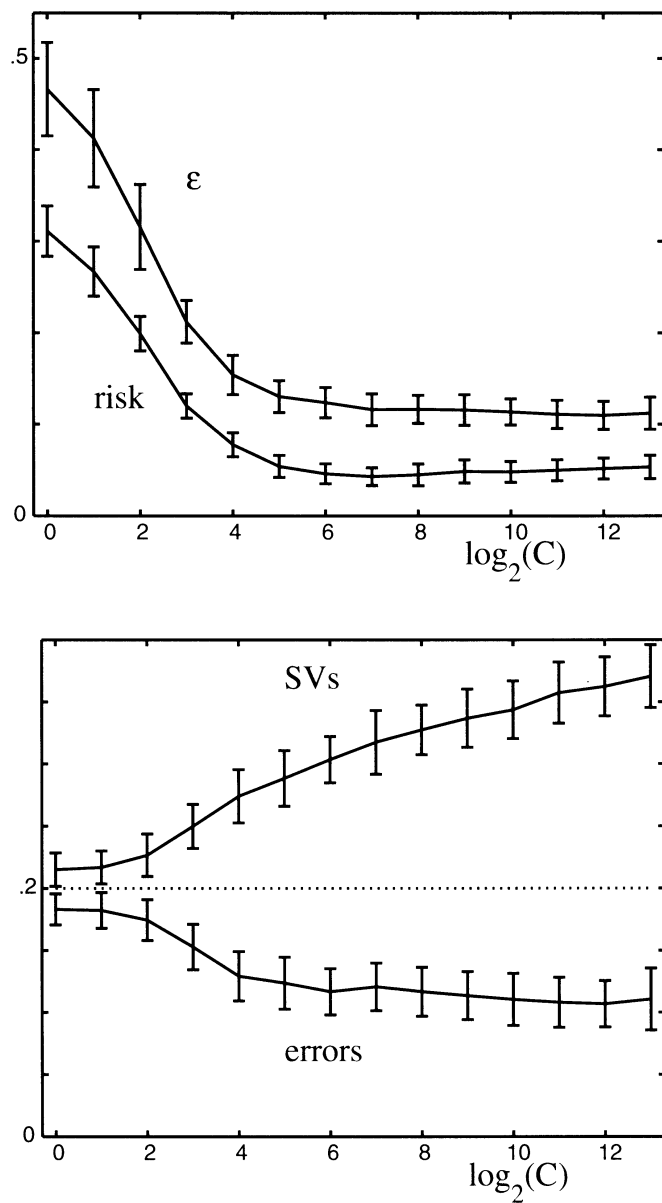


Figure 8:  $\nu$ -SVR for different values of the constant  $C$ . (Top)  $\epsilon$  decreases when the regularization is decreased (large  $C$ ). Only very little, if any, overfitting occurs. (Bottom)  $\nu$  upper bounds the fraction of errors, and lower bounds the fraction of SVs (cf. proposition 1). The bound gets looser as  $C$  increases; this corresponds to a smaller number of examples  $\ell$  relative to  $C$  (cf. Table 2).

Table 2: Asymptotic Behavior of the Fraction of Errors and SVs.

$\ell$	10	50	100	200	500	1000	1500	2000
$\varepsilon$	0.27	0.22	0.23	0.25	0.26	0.26	0.26	0.26
Fraction of errors	0.00	0.10	0.14	0.18	0.19	0.20	0.20	0.20
Fraction of SVs	0.40	0.28	0.24	0.23	0.21	0.21	0.20	0.20

Notes: The  $\varepsilon$  found by  $\nu$ -SV regression is largely independent of the sample size  $\ell$ . The fraction of SVs and the fraction of errors approach  $\nu = 0.2$  from above and below, respectively, as the number of training examples  $\ell$  increases (cf. proposition 1).

Figure 10 gives an illustration of how one can make use of parametric insensitivity models as proposed in section 5. Using the proper model, the estimate gets much better. In the parametric case, we used  $\nu = 0.1$  and  $\zeta(x) = \sin^2((2\pi/3)x)$ , which, due to  $\int \zeta(x) dP(x) = 1/2$ , corresponds to our standard choice  $\nu = 0.2$  in  $\nu$ -SVR (cf. proposition 4). Although this relies on the assumption that the SVs are uniformly distributed, the experimental findings are consistent with the asymptotics predicted theoretically: for  $\ell = 200$ , we got 0.24 and 0.19 for the fraction of SVs and errors, respectively.

**8.1.2 Boston Housing Benchmark.** Empirical studies using  $\varepsilon$ -SVR have reported excellent performance on the widely used Boston housing regression benchmark set (Stitson et al., 1999). Due to proposition 2, the only difference between  $\nu$ -SVR and standard  $\varepsilon$ -SVR lies in the fact that different parameters,  $\varepsilon$  versus  $\nu$ , have to be specified a priori. Accordingly, the goal of the following experiment was not to show that  $\nu$ -SVR is better than  $\varepsilon$ -SVR, but that  $\nu$  is a useful parameter to select. Consequently, we are interested only in  $\nu$  and  $\varepsilon$ , and hence kept the remaining parameters fixed. We adjusted  $C$  and the width  $2\sigma^2$  in  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2/(2\sigma^2))$  as in Schölkopf et al. (1997). We used  $2\sigma^2 = 0.3 \cdot N$ , where  $N = 13$  is the input dimensionality, and  $C/\ell = 10 \cdot 50$  (i.e., the original value of 10 was corrected since in the present case, the maximal  $y$ -value is 50 rather than 1). We performed 100 runs, where each time the overall set of 506 examples was randomly split into a training set of  $\ell = 481$  examples and a test set of 25 examples (cf. Stitson et al., 1999). Table 3 shows that over a wide range of  $\nu$  (note that only  $0 \leq \nu \leq 1$  makes sense), we obtained performances that are close to the best performances that can be achieved by selecting  $\varepsilon$  a priori by looking at the test set. Finally, although we did not use validation techniques to select the optimal values for  $C$  and  $2\sigma^2$ , the performances are state of the art (Stitson et al., 1999, report an MSE of 7.6 for  $\varepsilon$ -SVR using ANOVA kernels, and 11.7 for Bagging regression trees). Table 3, moreover, shows that in this real-world application,  $\nu$  can be used to control the fraction of SVs/errors.

Table 3: Results for the Boston Housing Benchmark (top:  $\nu$ -SVR; bottom:  $\varepsilon$ -SVR).

$\nu$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
automatic $\varepsilon$	2.6	1.7	1.2	0.8	0.6	0.3	0.0	0.0	0.0	0.0
MSE	9.4	8.7	9.3	9.5	10.0	10.6	11.3	11.3	11.3	11.3
STD	6.4	6.8	7.6	7.9	8.4	9.0	9.6	9.5	9.5	9.5
Errors	0.0	0.1	0.2	0.2	0.3	0.4	0.5	0.5	0.5	0.5
SVs	0.3	0.4	0.6	0.7	0.8	0.9	1.0	1.0	1.0	1.0

$\varepsilon$	0	1	2	3	4	5	6	7	8	9	10
MSE	11.3	9.5	8.8	9.7	11.2	13.1	15.6	18.2	22.1	27.0	34.3
STD	9.5	7.7	6.8	6.2	6.3	6.0	6.1	6.2	6.6	7.3	8.4
Errors	0.5	0.2	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
SVs	1.0	0.6	0.4	0.3	0.2	0.1	0.1	0.1	0.1	0.1	0.1

Note: MSE: Mean squared errors; STD: standard deviations thereof (100 trials); errors: fraction of training points outside the tube; SVs: fraction of training points that are SVs.

**8.2 Classification.** As in the regression case, the difference between  $C$ -SVC and  $\nu$ -SVC lies in the fact that we have to select a different parameter a priori. If we are able to do this well, we obtain identical performances. In other words,  $\nu$ -SVC could be used to reproduce the excellent results obtained on various data sets using  $C$ -SVC (for an overview; see Schölkopf, Burges, & Smola, 1999). This would certainly be a worthwhile project; however, we restrict ourselves here to showing some toy examples illustrating the influence of  $\nu$  (see Figure 11). The corresponding fractions of SVs and margin errors are listed in Table 4.

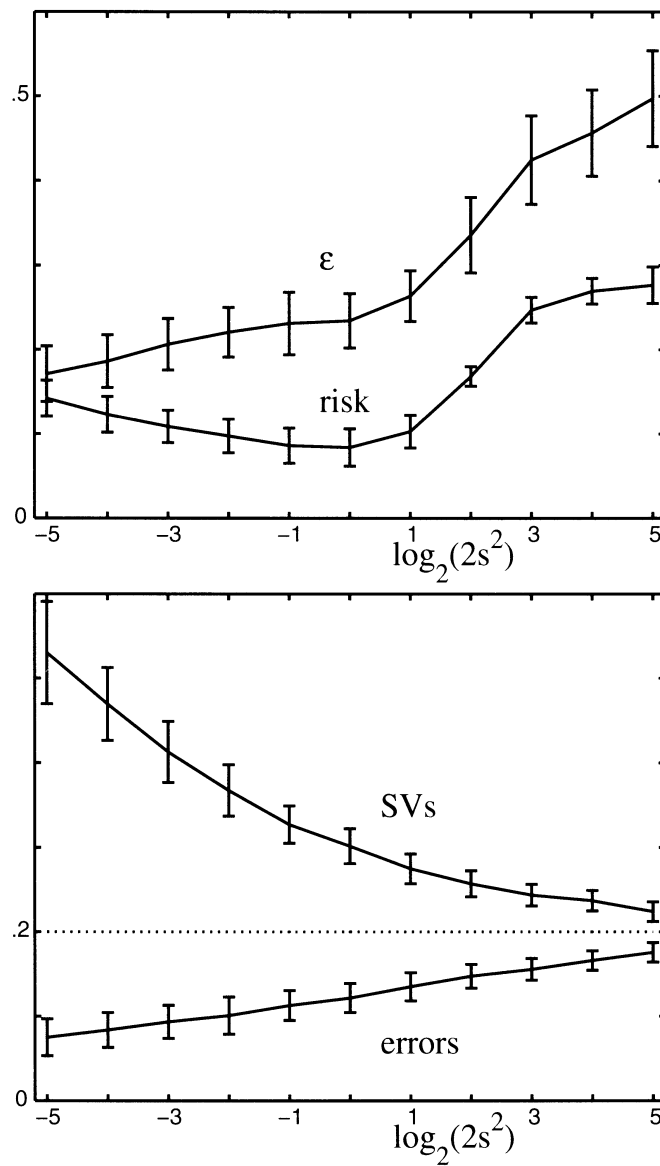
## 9 Discussion

We have presented a new class of SV algorithms, which are parameterized by a quantity  $\nu$  that lets one control the number of SVs and errors. We described  $\nu$ -SVR, a new regression algorithm that has been shown to be rather

Figure 9: Facing page.  $\nu$ -SVR for different values of the gaussian kernel width  $2s^2$ , using  $k(x, x') = \exp(-|x - x'|^2/(2s^2))$ . Using a kernel that is too wide results in underfitting; moreover, since the tube becomes too rigid as  $2s^2$  gets larger than 1, the  $\varepsilon$  needed to accomodate a fraction  $(1 - \nu)$  of the points, increases significantly. In the bottom figure, it can again be seen that the speed of the uniform convergence responsible for the asymptotic statement given in proposition 1 depends on the capacity of the underlying model. Increasing the kernel width leads to smaller covering numbers (Williamson et al., 1998) and therefore faster convergence.



useful in practice. We gave theoretical results concerning the meaning and the choice of the parameter  $\nu$ . Moreover, we have applied the idea underlying  $\nu$ -SV regression to develop a  $\nu$ -SV classification algorithm. Just like its regression counterpart, the algorithm is interesting from both a prac-



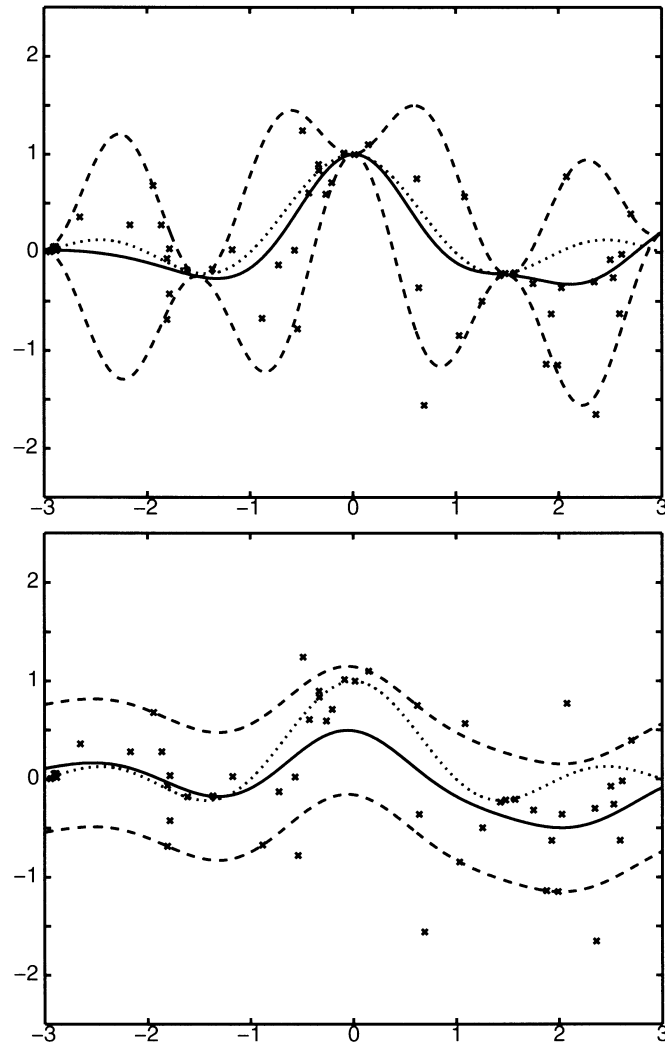


Figure 10: Toy example, using prior knowledge about an  $x$ -dependence of the noise. Additive noise ( $\sigma = 1$ ) was multiplied by the function  $\sin^2((2\pi/3)x)$ . (Top) The same function was used as  $\zeta$  as a parametric insensitivity tube (section 5). (Bottom)  $\nu$ -SVR with standard tube.

Table 4: Fractions of Errors and SVs, Along with the Margins of Class Separation, for the Toy Example Depicted in Figure 11.

$\nu$	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8
Fraction of errors	0.00	0.07	0.25	0.32	0.39	0.50	0.61	0.71
Fraction of SVs	0.29	0.36	0.43	0.46	0.57	0.68	0.79	0.86
Margin $2\rho/\ \mathbf{w}\ $	0.009	0.035	0.229	0.312	0.727	0.837	0.922	1.092

Note:  $\nu$  upper bounds the fraction of errors and lower bounds the fraction of SVs, and that increasing  $\nu$ , i.e. allowing more errors, increases the margin.

tical and a theoretical point of view. Controlling the number of SVs has consequences for (1) run-time complexity, since the evaluation time of the estimated function scales linearly with the number of SVs (Burges, 1998); (2) training time, e.g., when using a chunking algorithm (Vapnik, 1979) whose complexity increases with the number of SVs; (3) possible data compression applications— $\nu$  characterizes the compression ratio: it suffices to train the algorithm only on the SVs, leading to the same solution (Schölkopf et al., 1995); and (4) generalization error bounds: the algorithm directly optimizes a quantity using which one can give generalization bounds. These, in turn, could be used to perform structural risk minimization over  $\nu$ . Moreover, asymptotically,  $\nu$  directly controls the number of support vectors, and the latter can be used to give a leave-one-out generalization bound (Vapnik, 1995).

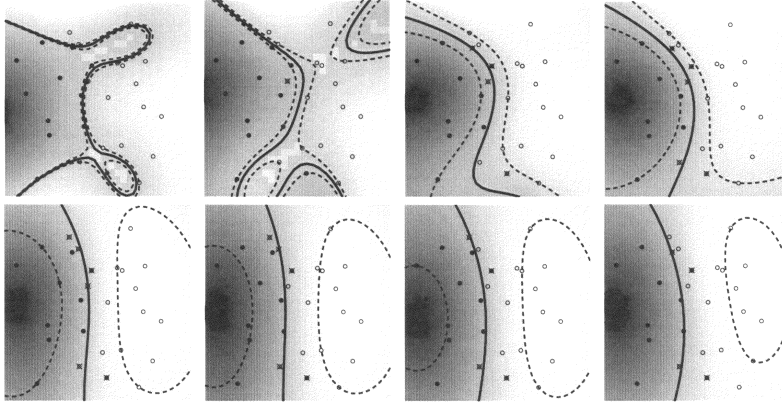


Figure 11: Toy problem (task: separate circles from disks) solved using  $\nu$ -SV classification, using parameter values ranging from  $\nu = 0.1$  (top left) to  $\nu = 0.8$  (bottom right). The larger we select  $\nu$ , the more points are allowed to lie inside the margin (depicted by dotted lines). As a kernel, we used the gaussian  $k(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2)$ .

In both the regression and the pattern recognition case, the introduction of  $\nu$  has enabled us to dispose of another parameter. In the regression case, this was the accuracy parameter  $\varepsilon$ ; in pattern recognition, it was the regularization constant  $C$ . Whether we could have as well abolished  $C$  in the regression case is an open problem.

Note that the algorithms are not fundamentally different from previous SV algorithms; in fact, we showed that for certain parameter settings, the results coincide. Nevertheless, we believe there are practical applications where it is more convenient to specify a fraction of points that is allowed to become errors, rather than quantities that are either hard to adjust a priori (such as the accuracy  $\varepsilon$ ) or do not have an intuitive interpretation (such as  $C$ ). On the other hand, desirable properties of previous SV algorithms, including the formulation as a definite quadratic program, and the sparse SV representation of the solution, are retained. We are optimistic that in many applications, the new algorithms will prove to be quite robust. Among these should be the reduced set algorithm of Osuna and Girosi (1999), which approximates the SV pattern recognition decision surface by  $\varepsilon$ -SVR. Here,  $\nu$ -SVR should give a direct handle on the desired speed-up.

Future work includes the experimental test of the asymptotic predictions of section 4 and an experimental evaluation of  $\nu$ -SV classification on real-world problems. Moreover, the formulation of efficient chunking algorithms for the  $\nu$ -SV case should be studied (cf. Platt, 1999). Finally, the additional freedom to use parametric error models has not been exploited yet. We expect that this new capability of the algorithms could be very useful in situations where the noise is heteroscedastic, such as in many problems of financial data analysis, and general time-series analysis applications (Müller et al., 1999; Mattera & Haykin, 1999). If a priori knowledge about the noise is available, it can be incorporated into an error model  $\zeta$ ; if not, we can try to estimate the model directly from the data, for example, by using a variance estimator (e.g., Seifert, Gasser, & Wolf, 1993) or quantile estimator (section 3).

## Acknowledgments

---

This work was supported in part by grants of the Australian Research Council and the DFG (Ja 379/7-1 and Ja 379/9-1). Thanks to S. Ben-David, A. Elisseeff, T. Jaakkola, K. Müller, J. Platt, R. von Sachs, and V. Vapnik for discussions and to L. Almeida for pointing us to White's work. Jason Weston has independently performed experiments using a sum inequality constraint on the Lagrange multipliers, but declined an offer of coauthorship.

## References

---

- Aizerman, M., Braverman, E., & Rozonoer, L. (1964). Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, 25, 821–837.

- Anthony, M., & Bartlett, P. L. (1999). *Neural network learning: Theoretical foundations*. Cambridge: Cambridge University Press.
- Bartlett, P. L. (1998). The sample complexity of pattern classification with neural networks: The size of the weights is more important than the size of the network. *IEEE Transactions on Information Theory*, 44(2), 525–536.
- Bertsekas, D. P. (1995). *Nonlinear programming*. Belmont, MA: Athena Scientific.
- Boser, B. E., Guyon, I. M., & Vapnik, V. N. (1992). A training algorithm for optimal margin classifiers. In D. Haussler (Ed.), *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory* (pp. 144–152). Pittsburgh, PA: ACM Press.
- Burges, C. J. C. (1998). A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, 2(2), 1–47.
- Cortes, C., & Vapnik, V. (1995). Support vector networks. *Machine Learning*, 20, 273–297.
- Girosi, F. (1998). An equivalence between sparse approximation and support vector machines. *Neural Computation*, 10(6), 1455–1480.
- Horn, R. A., & Johnson, C. R. (1985). *Matrix analysis*. Cambridge: Cambridge University Press.
- Huber, P. J. (1981). *Robust statistics*. New York: Wiley.
- Mattera, D., & Haykin, S. (1999). Support vector machines for dynamic reconstruction of a chaotic system. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods—Support vector learning* (pp. 211–241). Cambridge, MA: MIT Press.
- Müller, K.-R., Smola, A., Rätsch, G., Schölkopf, B., Kohlmorgen, J., and Vapnik, V. (1999). Predicting time series with support vector machines. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods—Support vector learning* (pp. 243–253). Cambridge, MA: MIT Press.
- Murata, N., Yoshizawa, S., & Amari, S. (1994). Network information criterion—determining the number of hidden units for artificial neural network models. *IEEE Transactions on Neural Networks*, 5, 865–872.
- Osuna, E., & Girosi, F. (1999). Reducing run-time complexity in support vector machines. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods—Support vector learning* (pp. 271–283). Cambridge, MA: MIT Press.
- Platt, J. (1999). Fast training of SVMs using sequential minimal optimization. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods—Support vector learning* (pp. 185–208). Cambridge, MA: MIT Press.
- Pontil, M., Rifkin, R., & Evgeniou, T. (1999). From regression to classification in support vector machines. In M. Verleysen (Ed.), *Proceedings ESANN* (pp. 225–230). Brussels: D Facto.
- Schölkopf, B. (1997). *Support vector learning*. Munich: R. Oldenbourg Verlag.
- Schölkopf, B., Bartlett, P. L., Smola, A., & Williamson, R. C. (1998). Support vector regression with automatic accuracy control. In L. Niklasson, M. Bodén, & T. Ziemke (Eds.), *Proceedings of the 8th International Conference on Artificial Neural Networks* (pp. 111–116). Berlin: Springer-Verlag.

- Schölkopf, B., Burges, C. J. C., & Smola, A. J. (1999). *Advances in kernel methods—Support vector learning*. Cambridge, MA: MIT Press.
- Schölkopf, B., Burges, C., & Vapnik, V. (1995). Extracting support data for a given task. In U. M. Fayyad & R. Uthurusamy (Eds.), *Proceedings, First International Conference on Knowledge Discovery and Data Mining*. Menlo Park, CA: AAAI Press.
- Schölkopf, B., Shawe-Taylor, J., Smola, A. J., & Williamson, R. C. (1999). Kernel-dependent support vector error bounds. In *Ninth International Conference on Artificial Neural Networks* (pp. 103–108). London: IEE.
- Schölkopf, B., Smola, A., & Müller, K.-R. (1998). Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10, 1299–1319.
- Schölkopf, B., Sung, K., Burges, C., Girosi, F., Niyogi, P., Poggio, T., & Vapnik, V. (1997). Comparing support vector machines with gaussian kernels to radial basis function classifiers. *IEEE Trans. Sign. Processing*, 45, 2758–2765.
- Seifert, B., Gasser, T., & Wolf, A. (1993). Nonparametric estimation of residual variance revisited. *Biometrika*, 80, 373–383.
- Shawe-Taylor, J., Bartlett, P. L., Williamson, R. C., & Anthony, M. (1998). Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 44(5), 1926–1940.
- Shawe-Taylor, J., & Cristianini, N. (1999). Margin distribution bounds on generalization. In *Computational Learning Theory: 4th European Conference* (pp. 263–273). New York: Springer.
- Smola, A. J. (1998). *Learning with kernels*. Doctoral dissertation, Technische Universität Berlin. Also: *GMD Research Series No. 25*, Birlinghoven, Germany.
- Smola, A., Frieß, T., & Schölkopf, B. (1999). Semiparametric support vector and linear programming machines. In M. S. Kearns, S. A. Solla, & D. A. Cohn (Eds.), *Advances in neural information processing systems*, 11 (pp. 585–591). Cambridge, MA: MIT Press.
- Smola, A., Murata, N., Schölkopf, B., & Müller, K.-R. (1998). Asymptotically optimal choice of  $\varepsilon$ -loss for support vector machines. In L. Niklasson, M. Bodén, & T. Ziemke (Eds.), *Proceedings of the 8th International Conference on Artificial Neural Networks* (pp. 105–110). Berlin: Springer-Verlag.
- Smola, A., & Schölkopf, B. (1998). On a kernel-based method for pattern recognition, regression, approximation and operator inversion. *Algorithmica*, 22, 211–231.
- Smola, A., Schölkopf, B., & Müller, K.-R. (1998). The connection between regularization operators and support vector kernels. *Neural Networks*, 11, 637–649.
- Smola, A., Williamson, R. C., Mika, S., & Schölkopf, B. (1999). Regularized principal manifolds. In *Computational Learning Theory: 4th European Conference* (pp. 214–229). Berlin: Springer-Verlag.
- Stitson, M., Gammerman, A., Vapnik, V., Vovk, V., Watkins, C., & Weston, J. (1999). Support vector regression with ANOVA decomposition kernels. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods—Support vector learning* (pp. 285–291). Cambridge, MA: MIT Press.

- Vapnik, V. (1979). *Estimation of dependences based on empirical data* [in Russian]. Nauka: Moscow. (English translation: Springer-Verlag, New York, 1982).
- Vapnik, V. (1995). *The nature of statistical learning theory*. New York: Springer-Verlag.
- Vapnik, V., & Chervonenkis, A. (1974). *Theory of pattern recognition* [in Russian]. Nauka: Moscow. (German Translation: W. Wapnik & A. Tscherwonienkis, *Theorie der Zeichenerkennung*, Akademie-Verlag, Berlin, 1979).
- Wahba, G. (1999). Support vector machines, reproducing kernel Hilbert spaces and the randomized GACV. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods—Support vector learning* (pp. 69–88). Cambridge, MA: MIT Press.
- White, H. (1994). Parametric statistical estimation with artificial neural networks: A condensed discussion. In V. Cherkassky, J. H. Friedman, & H. Wechsler (Eds.), *From statistics to neural networks*. Berlin: Springer.
- Williamson, R. C., Smola, A. J., & Schölkopf, B. (1998). *Generalization performance of regularization networks and support vector machines via entropy numbers of compact operators* (Tech. Rep. 19 Neurocolt Series). London: Royal Holloway College. Available online at <http://www.neurocolt.com>.

---

Received December 2, 1998; accepted May 14, 1999.