

Comparing Support Vector Machines with Gaussian Kernels to Radial Basis Function Classifiers

Bernhard Schölkopf, Kah-Kay Sung, Chris J. C. Burges, Federico Girosi,
Partha Niyogi, Tomaso Poggio, and Vladimir Vapnik

Abstract—The support vector (SV) machine is a novel type of learning machine, based on statistical learning theory, which contains polynomial classifiers, neural networks, and radial basis function (RBF) networks as special cases. In the RBF case, the SV algorithm automatically determines centers, weights, and threshold that minimize an upper bound on the expected test error.

The present study is devoted to an experimental comparison of these machines with a classical approach, where the centers are determined by k -means clustering, and the weights are computed using error backpropagation. We consider three machines, namely, a classical RBF machine, an SV machine with Gaussian kernel, and a hybrid system with the centers determined by the SV method and the weights trained by error backpropagation. Our results show that on the United States postal service database of handwritten digits, the SV machine achieves the highest recognition accuracy, followed by the hybrid system. The SV approach is thus not only theoretically well-founded but also superior in a practical application.

Index Terms—Clustering, pattern recognition, prototypes, radial basis function networks, support vector machines.

I. INTRODUCTION

CONSIDER Fig. 1. Suppose we want to construct a radial basis function classifier

$$f(\mathbf{x}) = \text{sgn} \left[\sum_{i=1}^{\ell} w_i \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{c_i} \right) + b \right] \quad (1)$$

(b and c_i being constants, the latter positive) separating balls from circles, i.e., taking different values on balls and circles. How do we choose the centers \mathbf{x}_i ? Two extreme cases are conceivable (Fig. 2).

Manuscript received June 20, 1997. Part of this work was done while B. Schölkopf was with AT&T Bell Laboratories and while K.-K. Sung and P. Niyogi were with the Massachusetts Institute of Technology. This work was supported in part by a grant from the National Science Foundation under Contract ASC-9217041, by the Studienstiftung des deutschen Volkes, and by ARPA under ONR Contract N00014-94-C-0186.

B. Schölkopf is with the Max-Planck-Institut für Biologische Kybernetik, Tübingen, Germany.

K.-K. Sung is with the Department of Information Systems and Computer Science, National University of Singapore, Singapore.

C. J. C. Burges is with Lucent Technologies, Bell Laboratories, Holmdel, NJ 07733 USA.

F. Girosi and T. Poggio are with the Center for Biological and Computational Learning, Massachusetts Institute of Technology, Cambridge, MA 02142 USA.

P. Niyogi is with Lucent Technologies, Bell Laboratories, Murray Hill, NJ 07974 USA.

V. Vapnik is with AT&T Research, Red Bank, NJ 07701 USA.

Publisher Item Identifier S 1053-587X(97)08062-8.

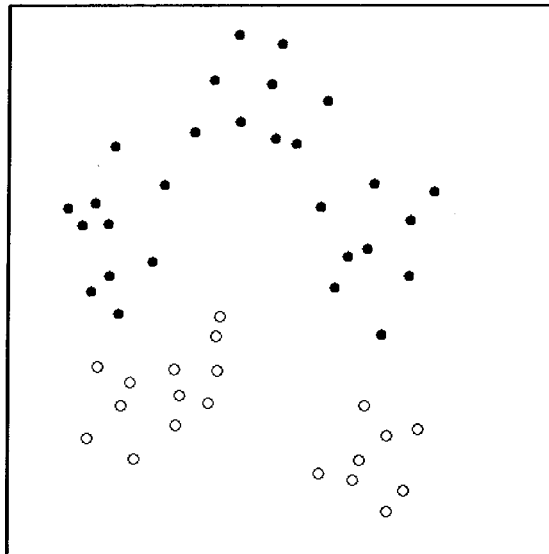


Fig. 1. Simple two-dimensional (2-D) classification problem: Find a decision function separating balls from circles. The box, as in all following figures, depicts the region $[-1, 1]^2$.

The first approach consists of choosing the centers for the two classes separately, irrespective of the classification task to be solved. The classical technique of finding the centers by some clustering technique (*before* tackling the classification problem) is such an approach. The weights w_i are then usually found by either error backpropagation [17] or the pseudo-inverse method (e.g., [15]).

An alternative approach consists of choosing, as centers, points that are *critical* for the classification task at hand. Recently, the *support vector algorithm* was developed [3], [7], [25], implementing the latter idea. It is a general algorithm, based on guaranteed risk bounds of statistical learning theory, which in particular allows the construction of RBF classifiers. This is done by simply choosing a suitable kernel function for the SV machine (see Section II-B). The SV training consists of a quadratic programming problem that can be solved efficiently and for which we are guaranteed to find a global extremum. The algorithm automatically computes the number and location of the above centers, the weights w_i , and the threshold b in the following way: By the use of a kernel function (in the present case, a Gaussian one), the patterns are mapped nonlinearly into a high-dimensional space. There, an optimal separating hyperplane is constructed, expressed in terms of those examples that are closest to the decision

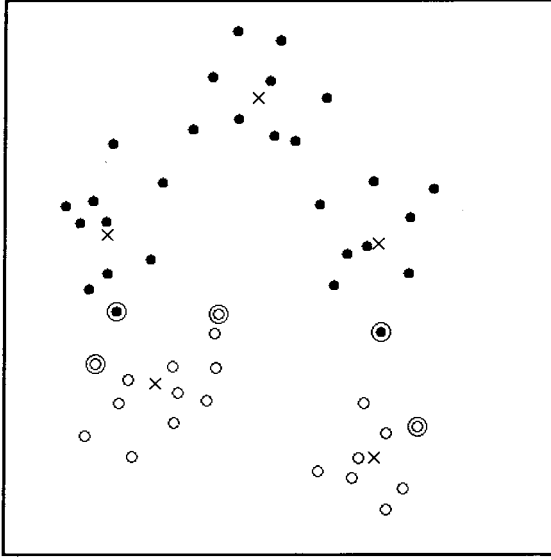


Fig. 2. RBF centers automatically found by the SV algorithm (indicated by extra circles), using $c_i = 1$ for all i [cf., (1)]. The number of SV centers accidentally coincides with the number of identifiable clusters (indicated by crosses found by k -means clustering with $k = 2$ and $k = 3$ for balls and circles, respectively), but the naive correspondence between clusters and centers is lost; indeed, three of the SV centers are circles, and only two of them are balls. Note that the SV centers are chosen with respect to the classification task to be solved.

boundary [24]. These are the *support vectors* that correspond to the centers in input space.

The goal of the present study is to compare real-world results obtained with k -means clustering and classical RBF training to those obtained when the centers, weights, and threshold are automatically chosen by the SV algorithm. To this end, we decided to undertake a performance study by combining expertise on the SV algorithm (AT&T Bell Laboratories) and on the classical RBF networks (Massachusetts Institute of Technology). We report results obtained on a United States Postal Service database of handwritten digits.

We have organized the material as follows. In the next section, we describe the algorithms used to train the different types of RBF classifiers compared in this paper. Following that, we present an experimental comparison of the approaches. We conclude with a discussion of our findings.

II. CONSTRUCTING RADIAL BASIS FUNCTION CLASSIFIERS

We describe three RBF systems trained in different ways. In Section II-A, we discuss the first system trained along more classical lines. In Section II-B, we describe the SV algorithm, which constructs an RBF network whose parameters (centers, weights, threshold) are automatically optimized. In Section II-C, we finally use the SV algorithm merely to choose the centers of the RBF network and then optimize the weights separately.

A. Classical Spherical Gaussian RBF's

We begin by first describing the classical Gaussian RBF system. A N -dimensional spherical Gaussian RBF network

with K centers has the form

$$g(\mathbf{x}) = \sum_{i=1}^K w_i \mathcal{G}_i(\mathbf{x}) + b$$

$$= \sum_{i=1}^K w_i \frac{1}{(2\pi)^{N/2} \sigma_i^N} \exp\left(-\frac{\|\mathbf{x} - \mathbf{c}_i\|^2}{2\sigma_i^2}\right) + b$$

where \mathcal{G}_i is the i th Gaussian basis function with center \mathbf{c}_i and variance σ_i^2 . The weight coefficients w_i combine the Gaussian terms into a single output value, and b is a bias term. Building a Gaussian RBF network for a given learning task involves

- determining the total number of Gaussian basis functions to use for each output class and for the entire system;
- locating the Gaussian basis function centers;
- computing the cluster variance for each Gaussian basis function;
- solving for the weight coefficients and bias in the summation term.

One can implement a binary pattern classifier on input vectors \mathbf{x} as a Gaussian RBF network by defining an appropriate output threshold that separates the two pattern classes.

In this first system, we implement each individual digit recognizer as a spherical Gaussian RBF network trained with a classical RBF algorithm. Given a specified number of Gaussian basis functions for each digit class, the algorithm separately computes the Gaussian centers and variances for each of the ten digit classes to form the system's RBF kernels. The algorithm then solves for an optimal set of weight parameters between the RBF kernels and each output node to perform the desired digit recognition task. The training process constructs all ten digit recognizers in parallel so that one can reuse the same Gaussian basis functions among the ten digit recognizers. To avoid overfitting the available training data with an overly complex RBF classifier connected to every Gaussian kernel, we use a "bootstrap" like operation that selectively connects each recognizer's output node to only a "relevant" subset of all basis functions. The idea is similar to how we choose relevant near-miss clusters for each individual digit recognizer in the original system. The training procedure proceeds as follows (for further details, see [23]).

- The first training task is to determine an appropriate number k of Gaussian kernels for each digit class. This information is needed to initialize our clustering procedure for computing Gaussian RBF kernels. We opted for using the same numbers of Gaussian kernels as the ones automatically computed by the SV algorithm (see Table I).
- Our next task is to compute the Gaussian kernels for each digit class. We do this by separately performing classical k -means clustering (e.g., [11]) on each digit class in the training database. Each clustering operation returns a set of Gaussian centroids and their respective variances for the given digit class. Together, the Gaussian clusters from all ten digit classes form the system's RBF kernels.
- For each single-digit recognizer, we build an *initial* RBF network using only Gaussian kernels from its

TABLE I

NUMBERS OF CENTERS (SUPPORT VECTORS) AUTOMATICALLY EXTRACTED BY THE SV MACHINE. THE FIRST ROW GIVES THE TOTAL NUMBER FOR EACH BINARY CLASSIFIER, INCLUDING BOTH POSITIVE AND NEGATIVE EXAMPLES; IN THE SECOND ROW, WE ONLY COUNTED THE POSITIVE SUPPORT VECTORS. THE LATTER NUMBER WAS USED IN THE INITIALIZATION OF THE k -MEANS ALGORITHM, CF., SECTION II-A

digit class	0	1	2	3	4	5	6	7	8	9
# of Support Vectors	274	104	377	361	334	388	236	235	342	263
# of positive Support Vectors	172	77	217	179	211	231	147	133	194	166

target class, using on-line backpropagation of mean squared error to train the weights (the desired output is set to 1 or 0 for positive and negative examples, respectively). We then separately collect all the false positive mistakes each initial digit recognizer makes on the training database.

- 4) In the final training step, we augment each initial digit recognizer with additional Gaussian kernels from outside its target class to help reduce misclassification errors. We determine which Gaussian kernels are “relevant” for each recognizer as follows. For each *false positive* mistake the initial recognizer makes during the previous step, we look up the misclassified pattern’s actual digit class and include the nearest Gaussian kernel from its class in the “relevant” set. The *final* RBF network for each single-digit recognizer thus contains every Gaussian kernel from its target class, and several “relevant” kernels from the other nine digit classes, trained by error backpropagation. Because our final digit recognizers have fewer weight parameters than a naive system that fully connects all ten recognizers to every Gaussian kernel, we expect our system to generalize better on new data.

B. The Support Vector Machine

1) *Structural Risk Minimization*: For the case of two-class pattern recognition, the task of *learning from examples* can be formulated in the following way. Given a set of functions

$$\{f_\alpha: \alpha \in \Lambda\}, \quad f_\alpha: \mathbf{R}^N \rightarrow \{-1, +1\}$$

and a set of examples

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell), \quad \mathbf{x}_i \in \mathbf{R}^N, \quad y_i \in \{-1, +1\}$$

where each one is generated from an unknown probability distribution $P(\mathbf{x}, y)$, we want to find a function f_{α^*} that provides the smallest possible value for the average error committed on novel examples randomly drawn from P called the *risk*

$$R(\alpha) = \int \frac{1}{2} |f_\alpha(\mathbf{x}) - y| dP(\mathbf{x}, y).$$

The problem is that $R(\alpha)$ is unknown since $P(\mathbf{x}, y)$ is unknown. Therefore, an *induction principle* for risk minimization is necessary.

The straightforward approach to minimize the *empirical risk*

$$R_{\text{emp}}(\alpha) = \frac{1}{\ell} \sum_{i=1}^{\ell} \frac{1}{2} |f_\alpha(\mathbf{x}_i) - y_i|$$

turns out not to guarantee a small actual risk (i.e., a small error on the training set does not imply a small error on a

test set) if the number ℓ of training examples is limited. To make the most out of a limited amount of data, novel statistical techniques have been developed during the last 25 years. The *structural risk minimization* principle [24] is based on the fact that for the above learning problem, for any $\alpha \in \Lambda$ with a probability of at least $1 - \eta$, the bound

$$R(\alpha) \leq R_{\text{emp}}(\alpha) + \phi \left[\frac{h}{\ell} \right] \quad (2)$$

holds, where ϕ is defined as

$$\phi \left[\frac{h}{\ell} \right] = \sqrt{\frac{h \left(\log \frac{2\ell}{h} + 1 \right) - \log(\eta/4)}{\ell}}.$$

The parameter h is called the *Vapnik–Chervonenkis (VC) dimension* of a set of functions. It describes the capacity of a set of functions implementable by the learning machine. For binary classification, h is the maximal number of points that can be separated into two classes in all possible 2^h ways by using functions of the learning machine, i.e., for each possible separation, there exists a function that takes the value 1 on one class and -1 on the other class.

According to (2), given a fixed number ℓ of training examples, one can control the risk by controlling two quantities: $R_{\text{emp}}(\alpha)$ and $h(\{f_\alpha: \alpha \in \Lambda'\})$, with Λ' denoting some subset of the index set Λ . The empirical risk depends on the *function* chosen by the learning machine (i.e., on α), and it can be controlled by picking the right α . The VC dimension h depends on the *set of functions* $\{f_\alpha: \alpha \in \Lambda'\}$ that the learning machine can implement. To control h , one introduces a structure of nested subsets $S_n := \{f_\alpha: \alpha \in \Lambda_n\}$ of $\{f_\alpha: \alpha \in \Lambda\}$

$$S_1 \subset S_2 \subset \dots \subset S_n \subset \dots \quad (3)$$

whose VC dimensions, as a result, satisfy

$$h_1 \leq h_2 \leq \dots \leq h_n \leq \dots$$

For a given set of observations $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_\ell, y_\ell)$ the *structural risk minimization principle* chooses the function $f_{\alpha_\ell^*}$ in the subset $\{f_\alpha: \alpha \in \Lambda_n\}$ for which the guaranteed risk bound [the right-hand side of (2)] is minimal.

The remainder of this section follows [18] in briefly reviewing the SV algorithm. For details, the reader is referred to [25].

2) *A Structure on the Set of Hyperplanes*: Each particular choice of a structure (3) gives rise to a learning algorithm. The SV algorithm is based on a structure on the set of separating hyperplanes. To describe it, first note that given a dot product space \mathbf{Z} and a set of vectors $\mathbf{z}_1, \dots, \mathbf{z}_r \in \mathbf{Z}$, each hyperplane

$\{\mathbf{z} \in \mathbf{Z}: (\mathbf{w} \cdot \mathbf{z}) + b = 0\}$ corresponds to a canonical pair $(\mathbf{w}, b) \in \mathbf{Z} \times \mathbf{R}$ if we additionally require

$$\min_{i=1, \dots, r} |(\mathbf{w} \cdot \mathbf{z}_i) + b| = 1 \quad (4)$$

i.e., that the scaling of \mathbf{w} and b be such that the point closest to the hyperplane has a distance of $1/\|\mathbf{w}\|$. Let R be the radius of the smallest ball $B_R(\mathbf{a}) = \{\mathbf{z} \in \mathbf{Z}: \|\mathbf{z} - \mathbf{a}\| < R\}$ ($\mathbf{a} \in \mathbf{Z}$) containing the points $\mathbf{z}_1, \dots, \mathbf{z}_r$, and

$$f_{\mathbf{w}, b} = \text{sgn}[(\mathbf{w} \cdot \mathbf{z}) + b] \quad (5)$$

where the decision function is defined on these points. The possibility of introducing a structure on the set of hyperplanes is based on the result [25] that the set of canonical hyperplanes $\{f_{\mathbf{w}, b}: \|\mathbf{w}\| \leq A\}$ has a VC dimension h satisfying

$$h < R^2 A^2 + 1. \quad (6)$$

Note: Dropping the condition $\|\mathbf{w}\| \leq A$ leads to a set of functions whose VC dimension equals $N + 1$, where N is the dimensionality of \mathbf{Z} . Due to $\|\mathbf{w}\| \leq A$, we can get VC dimensions that are much smaller than N , enabling us to work in very high dimensional spaces.

3) *The Support Vector Algorithm:* Now, suppose we want to find a decision function $f_{\mathbf{w}, b}$ with the property $f_{\mathbf{w}, b}(\mathbf{z}_i) = y_i$, $i = 1, \dots, \ell$. If this function exists, canonicity (4) implies

$$y_i[(\mathbf{w} \cdot \mathbf{z}_i) + b] \geq 1, \quad i = 1, \dots, \ell. \quad (7)$$

In practice, a separating hyperplane often does not exist. To allow for the possibility of examples violating (7), [7] introduces slack variables

$$\xi_i \geq 0, \quad i = 1, \dots, \ell \quad (8)$$

to get

$$y_i[(\mathbf{w} \cdot \mathbf{z}_i) + b] \geq 1 - \xi_i, \quad i = 1, \dots, \ell. \quad (9)$$

The SV approach to minimizing the guaranteed risk bound (2) consists of the following. Minimize

$$\tau(\mathbf{w}, \xi) = \frac{1}{2}(\mathbf{w} \cdot \mathbf{w}) + \gamma \sum_{i=1}^{\ell} \xi_i \quad (10)$$

subject to the constraints (8) and (9). Due to (6), minimizing the first term is related to minimizing the VC dimension of the considered class of learning machines, thereby minimizing the second term of the bound (2) [it also amounts to maximizing the separation margin, cf., the remark following (4)]. The term $\sum_{i=1}^{\ell} \xi_i$, on the other hand, is an upper bound on the number of misclassifications on the training set—this controls the empirical risk term in (2). For a suitable positive constant γ , this approach therefore constitutes a practical implementation of structural risk minimization on the given set of functions.

Introducing Lagrange multipliers α_i and using the Kuhn–Tucker theorem of optimization theory, the solution can be shown to have an expansion

$$\mathbf{w} = \sum_{i=1}^{\ell} y_i \alpha_i \mathbf{z}_i \quad (11)$$

with nonzero coefficients α_i only where the corresponding example (\mathbf{z}_i, y_i) precisely meets the constraint (9). These \mathbf{z}_i are called *support vectors*. All remaining examples of the training set are irrelevant. Their constraint (9) is satisfied automatically (with $\xi_i = 0$), and they do not appear in the expansion (11). The coefficients α_i are found by solving the following quadratic programming problem. Maximize

$$W(\alpha) = \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j (\mathbf{z}_i \cdot \mathbf{z}_j) \quad (12)$$

subject to

$$0 \leq \alpha_i \leq \gamma, \quad i = 1, \dots, \ell, \quad \text{and} \quad \sum_{i=1}^{\ell} \alpha_i y_i = 0. \quad (13)$$

By linearity of the dot product, the decision function (5) can thus be written as

$$f(\mathbf{z}) = \text{sgn} \left[\sum_{i=1}^{\ell} y_i \alpha_i \cdot (\mathbf{z} \cdot \mathbf{z}_i) + b \right]. \quad (14)$$

Thus far, we have described the case of linear decision surfaces [26]. To allow for much more general decision surfaces, one can first nonlinearly transform a set of input vectors $\mathbf{x}_1, \dots, \mathbf{x}_{\ell}$ into a high-dimensional feature space by a map $\Phi: \mathbf{x}_i \mapsto \Phi(\mathbf{x}_i)$ and then do a linear separation there. Maximizing (12) and evaluating (14) then requires the computation of dot products $(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i))$ in a high-dimensional space. We are interested in cases where these expensive calculations can be reduced significantly by using a suitable function K such that

$$(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i)) = K(\mathbf{x}, \mathbf{x}_i) \quad (15)$$

leading to decision functions of the form

$$f(\mathbf{x}) = \text{sgn} \left[\sum_{i=1}^{\ell} y_i \alpha_i \cdot K(\mathbf{x}, \mathbf{x}_i) + b \right]. \quad (16)$$

In practice, we need not worry about conceiving the map Φ . We will choose a K that is the continuous Kernel of a positive definite integral operator, and Mercer's theorem of functional analysis then tells us that K can be expanded in a uniformly convergent series in terms of Eigenfunctions ψ_j and positive Eigenvalues λ_j

$$K(\mathbf{x}, \mathbf{x}_i) = \sum_{j=1}^{\infty} \lambda_j \psi_j(\mathbf{x}) \psi_j(\mathbf{x}_i)$$

corresponding to a dot product $(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i))$ in the (possibly infinite-dimensional) range of

$$\Phi: \mathbf{x} \mapsto [\sqrt{\lambda_1} \psi_1(\mathbf{x}), \sqrt{\lambda_2} \psi_2(\mathbf{x}), \dots]$$

(see [3], [25]). Consequently, everything that has been said about the linear case also applies to nonlinear cases obtained by using a suitable kernel K instead of the Euclidean dot product (Fig. 3). We are now in a position to explain how the SV algorithm can construct RBF classifiers. We simply use

$$K(\mathbf{x}, \mathbf{x}_i) = \exp(-\|\mathbf{x} - \mathbf{x}_i\|^2/c) \quad (17)$$

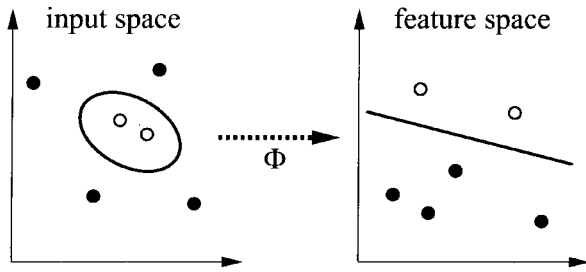


Fig. 3. By the use of a nonlinear kernel function (15), it is possible to compute a separating hyperplane with maximum margin in a feature space without explicitly mapping into that space. The feature space is related to input space via a nonlinear map Φ , causing the decision surface to be nonlinear in input space (e.g., of a Gaussian RBF type).

(see [1]). Other possible choices of K include $K(\mathbf{x}, \mathbf{x}_i) = (\mathbf{x} \cdot \mathbf{x}_i)^d$, yielding polynomial classifiers of degree d , and $K(\mathbf{x}, \mathbf{x}_i) = \tanh[\kappa \cdot (\mathbf{x} \cdot \mathbf{x}_i) + \Theta]$ for constructing neural networks.

To find the decision function (16), we maximize [cf., (12)]

$$W(\alpha) = \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) \quad (18)$$

subject to the constraints (13). Since K is required to satisfy Mercer's conditions, it corresponds to a dot product in another space (15); thus, $[y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)]_{ij}$ is a positive matrix, providing us with a problem that can be solved efficiently. To compute the threshold b , one takes into account that due to (9), for support vectors \mathbf{x}_j for which $\xi_j = 0$, we have $\sum_{i=1}^{\ell} y_i \alpha_i \cdot K(\mathbf{x}_j, \mathbf{x}_i) + b = y_j$.

4) *Further Remarks on SV Research:* The SV algorithm has been empirically shown to exhibit good generalization ability [7]. This applies for different types of SV machines obtained by using different kernel functions [cf., (17)]; moreover, different machines have been found to use largely the same SV's, e.g., most of the centers of an SV machine with Gaussian kernel coincide with the first-layer weights of neural network SV classifiers [18].

The accuracy can be further improved by incorporating invariances of the problem at hand, as with the virtual SV method. There, one generates artificial examples by transforming SV's and retrains on these in a second training run with a computational complexity comparable to the first one [19]. In addition, the decision rule (16), which requires the evaluation of kernel functions for the test example and all SV's, can be sped up with the reduced set technique [4] by approximating (16) with an expansion using fewer terms. These methods have led to substantial improvements for polynomial SV machines ([5] cut the error rate on a character recognition task by a third and obtained a speed-up by a factor of 50), and they are directly applicable also to Gaussian SV machines.

Besides pattern recognition, modified versions of the SV algorithm have been used for regression estimation and time series prediction [12], [13], [25], [27]. Moreover, the kernel method for computing dot products in high-dimensional fea-

TABLE II BINARY CLASSIFICATION: NUMBERS OF TEST ERRORS (OUT OF 2007 TEST PATTERNS) FOR THE SYSTEMS DESCRIBED IN SECTIONS II-A-C										
digit class	0	1	2	3	4	5	6	7	8	9
classical RBF	20	16	43	38	46	31	15	18	37	26
RBF with SV centers	9	12	27	24	32	24	19	16	26	16
full SV machine	16	8	25	19	29	23	14	12	25	16

ture spaces has been applied in other domains as, for instance, nonlinear principal component analysis [20].¹

C. A Hybrid System: RBF Network with SV Centers

The previous section discusses how one can train RBF like networks using the SV algorithm [cf., (16), (17)]. This involves the choice of an appropriate kernel function K and solving the quadratic optimization problem (18). The SV algorithm thus automatically determines the centers (which are the support vectors), the weights (given by $y_i \alpha_i$), and the threshold b for the RBF machine.

To assess the relative influence of the automatic SV center choice and the SV weight optimization, respectively, we built another RBF system constructed with centers that are simply the support vectors arising from the SV optimization and with the weights trained separately.

D. Addenda

1) *Computational Complexity:* By construction, the resulting classifiers after training will have the same architecture and comparable sizes. Thus, the three machines are comparable in classification speed and memory requirements.

Differences are, however, noticeable in training. Regarding training time, the SV machine was faster than the RBF system by about an order of magnitude due to the formulation as a quadratic programming problem that can be solved efficiently. The optimization, however, requires us to work with potentially large matrices. In the implementation that we used, the training data is processed in chunks, and matrix sizes were of the order 500×500 . For problems with very large numbers of SV's, a modified training algorithm has recently been proposed by [14].

2) *Different Ways of Training an RBF Classifier:* Due to (9) and (10), the SV algorithm puts emphasis on correctly separating the training data. In this respect, it is different from the classical RBF approach of training in the least-squares metric, which is more concerned with the general problem of estimating posterior probabilities than with directly solving a classification task at hand. There exist, however, studies investigating the question of how to select RBF centers or exemplars to minimize the number of misclassifications; see, for instance, [2], [6], [8], and [16]. A classical RBF system could also be made more discriminant by using moving centers (e.g., [15]) or a different cost function as the classification figure of merit [9]. In fact, it can be shown that Gaussian RBF regularization networks are equivalent to SV machines

¹ More information on SV methods can be obtained via <http://www.mpik-tueb.mpg.de/people/personal/bs/svm.html>.

TABLE III
TEN-CLASS DIGIT RECOGNITION ERROR RATES FOR THREE RBF CLASSIFIERS CONSTRUCTED WITH DIFFERENT ALGORITHMS. THE FIRST SYSTEM IS A CLASSICAL ONE, CHOOSING ITS CENTERS BY k -MEANS CLUSTERING. IN THE SECOND SYSTEM, THE SUPPORT VECTORS WERE USED AS CENTERS, AND IN THE THIRD ONE, THE ENTIRE NETWORK WAS TRAINED USING THE SV ALGORITHM

USPS Database	Classification Error Rate		
	classical RBF	RBF with SV centers	full SV machine
Training (7291 patterns)	1.7%	0.0%	0.0%
Test (2007 patterns)	6.7%	4.9%	4.2%

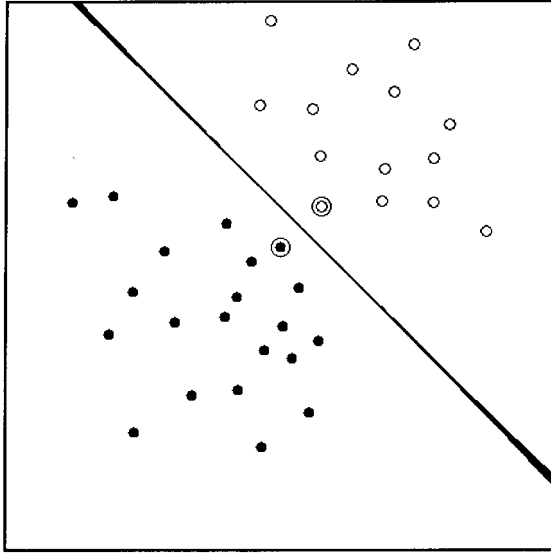


Fig. 4. Simple two-class classification problem as solved by the SV algorithm [$c_i = 1$ for all i ; cf., (1)]. Note that the RBF centers (indicated by extra circles) are closest to the decision boundary.

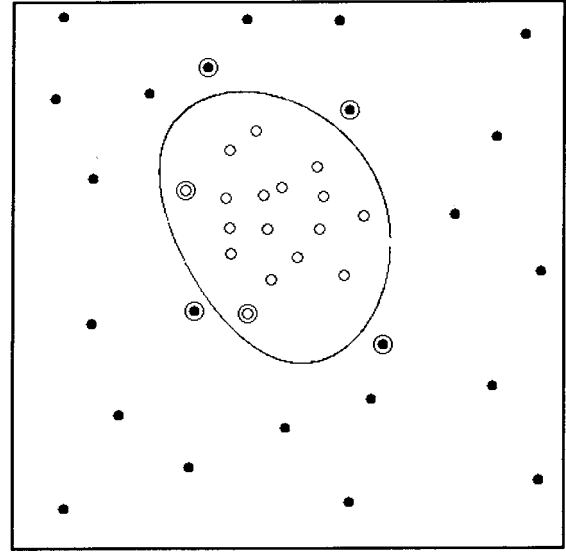


Fig. 5. Two-class classification problem solved by the SV algorithm [$c_i = 1$ for all i ; cf., (1)].

if the regularization operator and the cost function are chosen appropriately [22].

It is important to stress that the SV machine does not minimize the empirical risk (misclassification error on the training set) alone. Instead, it minimizes the sum of an upper bound on the empirical risk and a penalty term that depends on the complexity of the classifier used.

III. EXPERIMENTAL RESULTS

A. Toy Examples

What are the support vectors? They are elements of the data set that are “important” in separating the two classes from each other. Support vectors with zero slack variables (8) lie on the boundary of the decision surface, as they precisely satisfy the inequality (9) in the high-dimensional space. Figs. 4 and 5 illustrate that for the used Gaussian kernel, this is also the case in input space. This raises an interesting question from the point of view of interpreting the structure of trained RBF networks. The traditional view of RBF networks has been one where the centers were regarded as “templates” or stereotypical patterns. It is this point of view that leads to the clustering heuristic for training RBF networks. In contrast, the SV machine posits an alternate point of view, with the centers being those examples that are critical for a given classification task.

B. United States Postal Service Database

We used the USPS database of 9298 handwritten digits (7291 for training, 2007 for testing) collected from mail envelopes in Buffalo, NY (cf., [10]). Each digit is a 16×16 image represented as a 256-dimensional vector with entries between -1 and 1 . Preprocessing consisted of smoothing with a Gaussian kernel of width $\sigma = 0.75$. The SV machine results reported in the following were obtained with $\gamma = 10$ and $c = 0.3 \cdot N$ [cf., (10), (17)], where $N = 256$ is the dimensionality of input space.² In all experiments, we used the SV algorithm with standard quadratic programming techniques (conjugate gradient descent).

1) *Two-Class Classification*: Table I shows the numbers of support vectors, i.e., RBF centers, extracted by the SV algorithm. Table II gives the results of binary classifiers separating single digits from the rest for the systems described in Sections II-A–C.

2) *Ten-Class Classification*: For each test pattern, the arbitration procedure in all three systems simply returns the digit class whose recognizer gives the strongest response.³ Table III

²The SV machine is rather insensitive to different choices of c . For all values in $0.1, 0.2, \dots, 1.0$, the performance is about the same (in the area of 4–4.5%).

³In the SV case, we constructed ten two-class classifiers, each trained to separate a given digit from the other nine. Ten-class classification was done according to the maximal output before applying the sgn function, among the two-class classifiers. These outputs can be used for reject decisions.

shows the ten-class digit recognition error rates for our original system and the two RBF-based systems.

The fully automatic SV machine exhibits the highest test accuracy of the three systems.⁴ Using the SV algorithm to choose the centers for the RBF network is also better than the baseline procedure of choosing the centers by a clustering heuristic. It can be seen that in contrast to the k -means cluster centers, the centers chosen by the SV algorithm allow zero training error rates.

The considered recognition task is known to be rather hard—the human error rate is 2.5%, which is almost matched by a memory-based Tangent-distance classifier (2.6%; see [21]). Other results on this data base include a Euclidean distance nearest neighbor classifier (5.9%; see [21]), a two-layer perceptron (5.9%), and a convolutional neural network (5.0%; see [10]). By incorporating translational and rotational invariance using the Virtual SV technique ([19]; cf., our remark at the end of Section II-B3), we were able to improve the performance of the considered Gaussian kernel SV machine (same values of γ and c) from 4.2 to 3.2% error.

IV. SUMMARY AND DISCUSSION

The SV algorithm provides a principled way of choosing the number and the locations of RBF centers. Our experiments on a real-world pattern recognition problem have shown that *in contrast to a corresponding number of centers chosen by k -means, the centers chosen by the SV algorithm allowed a training error of zero*, even if the weights were trained by classical RBF methods. The interpretation of this finding is that the SV centers are specifically chosen for the classification task at hand, whereas k -means does not care about picking those centers that will make a problem separable.⁵

In addition, *the SV centers yielded lower test error rates than k -means*. It is interesting to note that using SV centers, while sticking to the classical procedure for training the weights, improved training and test error rates by approximately the same margin (2%). In view of the guaranteed risk bound (2), this can be understood in the following way. The improvement in test error (risk) was solely due to the lower value of the training error (empirical risk); the confidence term [the second term on the right-hand side of (2)], depending on the VC dimension and, thus, on the norm of the weight vector, did not change, as we stuck to the classical weight training procedure. However, when we also trained the weights with the SV algorithm, we minimized the norm of the weight vector [see (6) and (10)] in feature space and, thus, the confidence term, while still keeping the training error zero. Thus, consistent with (2), *the support vector machine achieved the highest test accuracy of the three systems*.

⁴ An analysis of the errors showed that about 85% of the errors committed by the SV machine were also made by the other systems. This makes the differences in error rates very reliable.

⁵ As pointed out in Section II-D2, the classical system could be made more discriminant as well by using different cost functions, e.g., the classification figure of merit [9], or by using moving centers (e.g., [15]).

ACKNOWLEDGMENT

The authors would like to thank K.-R. Müller, P. Simard, and A. Smola for useful discussions. B. Schölkopf thanks the Massachusetts Institute of Technology for its hospitality during a three-week visit in March 1995, where this work was started.

REFERENCES

- [1] M. Aizerman, E. Braverman, and L. Rozonoer, "Theoretical foundations of the potential function method in pattern recognition learning," *Automat. Remote Contr.*, vol. 25, pp. 821–837, 1964.
- [2] A. Barron, "Predicted squared error: A criterion for automatic model selection," in *Self-Organizing Methods in Modeling*, S. Farlow, Ed. New York: Marcel Dekker, 1984.
- [3] B. E. Boser, I. M. Guyon, and V. Vapnik, "A training algorithm for optimal margin classifiers," in *Fifth Annu. Workshop Comput. Learning Theory*, Pittsburgh, PA, 1992, pp. 144–152.
- [4] C. J. C. Burges, "Simplified support vector decision rules," in *Proc. 13th Int. Conf. Machine Learning*, L. Saitta, Ed. San Mateo, CA: Morgan Kaufmann, 1996.
- [5] C. J. C. Burges and B. Schölkopf, "Improving the accuracy and speed of support vector machines," in *Advances in Neural Information Processing Systems 9*, M. Mozer, M. Jordan, and T. Petsche, Eds. Cambridge, MA: MIT Press, 1997.
- [6] E. I. Chang and R. L. Lippmann, "A boundary hunting radial basis function classifier which allocates centers constructively," in *Advances in Neural Information Processing Systems 5*, S. J. Hanson, J. D. Cowan, and C. L. Giles, Eds. San Mateo, CA: Morgan Kaufmann, 1993.
- [7] C. Cortes and V. Vapnik, "Support vector networks," *Machine Learning*, vol. 20, pp. 273–297, 1995.
- [8] R. O. Duda and P. E. Hart, *Pattern Classification and Scene Analysis*. New York: Wiley, 1973.
- [9] J. B. Hampshire and A. Waibel, "A novel objective function for improved phoneme recognition using time-delay neural networks," *IEEE Trans. Neural Networks*, vol. 1, pp. 216–228, 1990.
- [10] Y. Le Cun et al., "Backpropagation applied to handwritten zip code recognition," *Neural Comput.*, vol. 1, pp. 541–551, 1989.
- [11] S. P. Lloyd, "Least squares quantization in PCM," *IEEE Trans. Inform. Theory*, vol. IT-28, pp. 129–137, 1982.
- [12] S. Mukherjee, E. Osuna, and F. Girosi, "Nonlinear prediction of chaotic time series using a support vector machine," in *Proc. NNSP*, 1997.
- [13] K.-R. Müller, A. Smola, G. Rätsch, B. Schölkopf, J. Kohlmorgen, and V. Vapnik, "Predicting time series with support vector machines," in *Proc. ICANN*, 1997.
- [14] E. Osuna, R. Freund, and F. Girosi, "Improved training algorithm for support vector machines," in *Proc. NNSP*, 1997.
- [15] T. Poggio and F. Girosi, "Networks for approximation and learning," *Proc. IEEE*, vol. 78, pp. 1481–1497, 1990.
- [16] D. Reilly, L. N. Cooper, and C. Elbaum, "A neural model for category learning," *Biol. Cybern.*, vol. 45, pp. 35–41, 1982.
- [17] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, pp. 533–536, 1986.
- [18] B. Schölkopf, C. Burges, and V. Vapnik, "Extracting support data for a given task," in *Proc. First Int. Conf. Knowledge Discovery Data Mining*, U. M. Fayyad and R. Uthurusamy, Eds. Menlo Park, CA: AAAI, 1995.
- [19] ———, "Incorporating invariances in support vector learning machines," in *Proc. Artificial Neural Networks-ICANN*, C. von der Malsburg, W. von Seelen, J. C. Vorbrüggen, and B. Sendhoff, Eds. Berlin, Germany: Springer, 1996, vol. 1112, pp. 47–52.
- [20] B. Schölkopf, A. Smola, and K.-R. Müller, "Nonlinear component analysis as a kernel eigenvalue problem," Tech. Rep. 44, Max-Planck-Institut für biologische Kybernetik, Tübingen, Germany, 1996.
- [21] P. Simard, Y. Le Cun, and J. Denker, "Efficient pattern recognition using a new transformation distance," in *Advances in Neural Information Processing Systems 5*, J. E. Moody, S. J. Hanson, and R. P. Lippmann, Eds. San Mateo, CA: Morgan Kaufmann, 1993.
- [22] A. Smola and B. Schölkopf, "On a kernel-based method for pattern recognition, regression, approximation and operator inversion," Tech. Rep. 1064, GMD, Berlin, 1997.
- [23] K. Sung, "Learning and example selection for object and pattern detection," Ph.D. thesis, Mass. Inst. Technol., Cambridge, 1995.
- [24] V. Vapnik, *Estimation of Dependences Based on Empirical Data*. Moscow: Nauka, 1979, in Russian; English translation: New York: Springer-Verlag, 1982.

- [25] ———, *The Nature of Statistical Learning Theory*. New York: Springer-Verlag, 1995.
- [26] V. Vapnik and A. Chervonenkis, *Theory of Pattern Recognition*. Moscow, Russia: Nauka, 1974, in Russian. German translation: W. Vapnik and A. Tschervonenkis, *Theorie der Zeichenerkennung*. Berlin, Germany: Akademie-Verlag, 1979.
- [27] V. Vapnik, S. Golowich, and A. Smola, "Support vector method for function approximation, regression estimation, and signal processing," in *Advances in Neural Information Processing Systems 9*, M. Mozer, M. Jordan, and T. Petsche, Eds. Cambridge, MA: MIT Press, 1997.



Bernhard Schölkopf received the M.Sc. degree in mathematics and the Lionel Cooper Memorial Prize from the University of London, U.K., in 1992. In 1994, he received the Diplom in physics from the Eberhard-Karls-Universität, Tübingen, Germany.

He is a doctoral student with H. Bülthoff (Max-Planck-Institut für Biologische Kybernetik) and V. Vapnik (AT&T Research). His scientific interests include machine learning and perception.

Kah-Kay Sung received the Ph.D. degree in electrical engineering and computer science from the Massachusetts Institute of Technology, Cambridge, in 1995.

He is currently a Lecturer at the Department of Information Systems and Computer Science, National University of Singapore. His research interests include computer vision and machine learning.



Chris J. C. Burges received the B.A. degree with first class honors in physics from Oxford University, Oxford, U.K., and the Ph.D. degree in particle physics from Brandeis University, Waltham, MA.

After a two-year postdoctoral position at the Massachusetts Institute of Technology, Cambridge, in theoretical particle physics, he joined AT&T Bell Laboratories, (now Lucent Technologies), Holmdel, NJ, and developed the routing algorithm AT&T now uses to ensure physical diversity of its signaling links. He then started working on applying neural

networks to handwriting recognition and developed several methods now used by banks in automated check readers. For the last two years, he has been working on developing support vector algorithms with V. Vapnik and others, concentrating on the pattern recognition problem, methods to speed up support vector machines, and the theory of kernel-based methods.

Federico Girosi received the Ph.D. degree in theoretical physics from the University of Genoa, Genoa, Italy.

He is a Principal Research Scientist in the Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology (MIT), Cambridge, and is Associate Director of the Center for Biological and Computational Learning at MIT. He has worked in the fields of neural networks, approximation theory, and computer vision. His current research topics include the study of neural network architectures, the complexity of the approximation problem, and support vector machines and their applications to computer vision.

Partha Niyogi received the B.Tech. degree from the Indian Institute of Technology, New Delhi, and the Ph.D. degree in electrical engineering and computer science from the Massachusetts Institute of Technology (MIT), Cambridge.

His research interests are in pattern recognition and learning theory and its application to problems in speech and language processing. He is currently with Bell Laboratories, Murray Hill, NJ.

Tomaso Poggio received the Ph.D. degree in theoretical physics from the University of Genoa, Genoa, Italy, in 1970.

He is the Uncas and Helen Whitaker Professor of the Department of Brain and Cognitive Sciences at the Massachusetts Institute of Technology (MIT), Cambridge. He is currently doing research in computational learning and vision at the MIT Center for Biological and Computational Learning, of which he is Co-Director. He is also a member of the Artificial Intelligence Laboratory. He is the author of several papers in areas ranging from psychophysics and biophysics to information processing in man and machine, artificial intelligence, and machine vision and learning.

Dr. Poggio has received several awards, is on the editorial boards of a number of interdisciplinary journals, is a fellow of the American Association for Artificial Intelligence, as well as the American Academy of Arts and Sciences, and is an Honorary Associate of the Neuroscience Research Program at Rockefeller University.



Vladimir Vapnik was born in Russia and received the Ph.D. degree in theoretical cybernetics from the Institute of Control Sciences, Academy of Science of the USSR, Moscow, in 1964.

Since 1991, he has been working for AT&T Bell Laboratories (since 1996, AT&T Labs Research), Red Bank, NJ. His research interests include statistical learning theory, theoretical and applied statistics, theory and methods of solving stochastic ill-posed problems, and methods of multidimensional function approximation. His main results in the last three

years are related to the development of the support vector method. He is author of many publications, including seven monographs on various problems of statistical learning theory.