
Advances in Large Margin Classifiers

Advances in Large Margin Classifiers

edited by
Alexander J. Smola
Peter Bartlett
Bernhard Schölkopf
Dale Schuurmans

The MIT Press
Cambridge, Massachusetts
London, England

©1999 Massachusetts Institute of Technology

All rights reserved. No part of this book may be reproduced in any form by any electronic or mechanical means (including photocopying, recording, or information storage and retrieval) without permission in writing from the publisher.

Printed and bound in the United States of America

Library of Congress Cataloging-in-Publication Data

Advances in large margin classifiers / edited by Alexander J. Smola, Peter Bartlett, Bernhard Schölkopf, Dale Schuurmans.

p. cm.

Includes bibliographical references and index.

ISBN 0-xxx-xxxxx-x (alk. paper)

1. Machine learning. 2. Algorithms. 3. Kernel functions

I. Smola, Alexander J. II. Bartlett, Peter. III. Schölkopf, Bernhard. IV. Schuurmans, Dale.

xxxx.x.xxx 1999

xxx.x'x-xxxx

99.xxxxx

CIP

Contents

Preface	vii
1 Introduction to Large Margin Classifiers <i>Alex J. Smola, Peter Bartlett, Bernhard Schölkopf, and Dale Schuurmans</i>	1
2 Adaptive Margin Support Vector Machines <i>Jason Weston</i>	29
References	44

Preface

Some good quote

who knows

some clever stuff ...
and some more visionary comments

Alexander J. Smola, Peter Bartlett, Bernhard Schölkopf, Dale Schuurmans

Berlin, Canberra, Waterloo, July 1999

1 Introduction to Large Margin Classifiers

The aim of this chapter is to provide a brief introduction to the basic concepts of large margin classifiers for readers unfamiliar with the topic. Moreover it is aimed at establishing a common basis in terms of notation and equations, upon which the subsequent chapters will build (and refer to) when dealing with more advanced issues.

1.1 A Simple Classification Problem

training data Assume that we are given a set of training data

$$X := \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathbb{R}^N \text{ where } m \in \mathbb{N} \quad (1.1)$$

labels together with corresponding labels

$$Y := \{y_1, \dots, y_m\} \subseteq \{-1, 1\}. \quad (1.2)$$

The goal is to find some decision function $g : \mathbb{R}^N \rightarrow \{-1, 1\}$ that accurately predicts the labels of unseen data points (\mathbf{x}, y) . That is, we seek a function g that minimizes the classification error, which is given by the probability that $g(\mathbf{x}) \neq y$. A common approach to representing decision functions is to use a real valued prediction function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ whose output is passed through a sign threshold to yield the final classification $g(\mathbf{x}) = \text{sgn}(f(\mathbf{x}))$. Let us start with a simple example: linear decision functions. In this case the unthresholded prediction is given by a simple linear function of the input vector \mathbf{x}

linear
decision
function

$$g(\mathbf{x}) := \text{sgn}(f(\mathbf{x})) \text{ where } f(\mathbf{x}) = (\mathbf{x} \cdot \mathbf{w}) + b \text{ for } \mathbf{w} \in \mathbb{R}^N \text{ and } b \in \mathbb{R}. \quad (1.3)$$

This gives a classification rule whose decision boundary $\{\mathbf{x} | f(\mathbf{x}) = 0\}$ is an $N - 1$ dimensional hyperplane separating the classes “+1” and “-1” from each other. Figure 1.1 depicts the situation. The problem of learning from data can be formulated as finding a set of parameters (\mathbf{w}, b) such that $\text{sgn}((\mathbf{w} \cdot \mathbf{x}_i) + b) = y_i$ for all $1 \leq i \leq m$. However, such a solution may not always exist, in particular if we are dealing with noisy data. For instance, consider Figure 1.1 with the triangle replaced by an open circle. This raises the question what to do in such a situation.

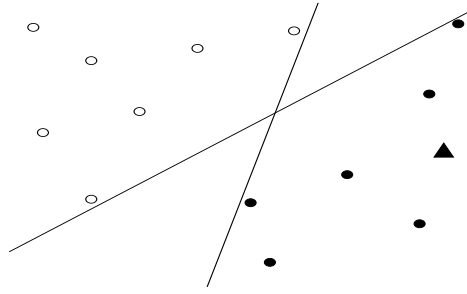


Figure 1.1 A linearly separable classification problem. Note that there may be several possible solutions as depicted by the two lines. The problem becomes non-separable if we replace the triangle by an open circle; in which case no solution (\mathbf{w}, b) exists.

1.1.1 Bayes Optimal Solution

Under the assumption that the data X, Y was generated from a probability distribution $p(\mathbf{x}, y)$ on $\mathbb{R}^N \times \{-1, 1\}$ and that p is known, it is straightforward to find a function that minimizes the probability of misclassification

$$R(g) := \int_{\mathbb{R}^N \times \{-1, 1\}} 1_{\{g(\mathbf{x}) \neq y\}} p(\mathbf{x}, y) d\mathbf{x} dy. \quad (1.4)$$

Bayes optimal
decision function

This function satisfies

$$g(\mathbf{x}) = \text{sgn}(p(\mathbf{x}, 1) - p(\mathbf{x}, -1)). \quad (1.5)$$

Consider a practical example.

Example 1.1 Two Gaussian Clusters

Assume that the two classes “+1” and “-1” are generated by two Gaussian clusters with the same covariance matrix Σ centered at $\boldsymbol{\mu}_+$ and $\boldsymbol{\mu}_-$ respectively

$$p(\mathbf{x}, y) = \frac{1}{2(2\sigma)^{N/2} |\Sigma|^{1/2}} \begin{cases} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_+)^{\top} \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_+)} & \text{if } y = +1 \\ e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_-)^{\top} \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_-)} & \text{if } y = -1. \end{cases} \quad (1.6)$$

Since the boundaries completely determine the decision function, we seek the set of points where $p(\mathbf{x}, +1) = p(\mathbf{x}, -1)$. In the case of (1.6) this is equivalent to seeking \mathbf{x} such that

$$(\mathbf{x} - \boldsymbol{\mu}_+)^{\top} \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_+) = (\mathbf{x} - \boldsymbol{\mu}_-)^{\top} \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_-). \quad (1.7)$$

By rearranging we find that this condition is equivalent to

$$\begin{aligned} \mathbf{x}^{\top} \Sigma^{-1} \mathbf{x} - 2\boldsymbol{\mu}_+^{\top} \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}_+^{\top} \Sigma^{-1} \boldsymbol{\mu}_+ - \mathbf{x}^{\top} \Sigma^{-1} \mathbf{x} + 2\boldsymbol{\mu}_-^{\top} \Sigma^{-1} \mathbf{x} - \boldsymbol{\mu}_-^{\top} \Sigma^{-1} \boldsymbol{\mu}_- &= 0 \\ 2(\boldsymbol{\mu}_+^{\top} \Sigma^{-1} - \boldsymbol{\mu}_-^{\top} \Sigma^{-1}) \mathbf{x} - (\boldsymbol{\mu}_+^{\top} \Sigma^{-1} \boldsymbol{\mu}_+ - \boldsymbol{\mu}_-^{\top} \Sigma^{-1} \boldsymbol{\mu}_-) &= 0 \end{aligned} \quad (1.8)$$

linear
discriminant

The latter form is equivalent to having a linear decision function determined by

$$f(\mathbf{x}) = ((\boldsymbol{\mu}_+ - \boldsymbol{\mu}_-)^T \Sigma^{-1}) \mathbf{x} - \frac{1}{2}(\boldsymbol{\mu}_+^T \Sigma^{-1} \boldsymbol{\mu}_+ - \boldsymbol{\mu}_-^T \Sigma^{-1} \boldsymbol{\mu}_-). \quad (1.9)$$

Hence in this simple example the Bayes optimal classification rule is linear.

Problems arise, however, if $p(\mathbf{x}, y)$ is not known (as generally happens in practice). In this case one has to obtain a good *estimate* of $g(\mathbf{x}) = \text{sgn}(f(\mathbf{x}))$ from the training data X, Y . A famous example of an algorithm for linear separation is the perceptron algorithm.

1.1.2 The Perceptron Algorithm

The *perceptron algorithm* is “incremental,” in the sense that small changes are made to the weight vector in response to each labelled example in turn. For any *learning rate* $\eta > 0$, the algorithm acts sequentially as shown in Table 1.1. Notice

```

argument: Training sample,  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathcal{X}$ ,  $Y = \{y_1, \dots, y_m\} \subset \{\pm 1\}$ 
             Learning rate,  $\eta$ 
returns: Weight vector  $\mathbf{w}$  and threshold  $b$ .
function Perceptron( $X, Y, \eta$ )
  initialize  $\mathbf{w}, b = 0$ 
  repeat
    for all  $i$  from  $i = 1, \dots, m$ 
      Compute  $g(\mathbf{x}_i) = \text{sgn}((\mathbf{w} \cdot \mathbf{x}_i) + b)$ 
      Update  $\mathbf{w}, b$  according to
           $\mathbf{w}' = \mathbf{w} + (\eta/2)(y_i - g(\mathbf{x}_i)) \mathbf{x}_i$ 
           $b' = b + (\eta/2)(y_i - g(\mathbf{x}_i))$ .
    endfor
  until for all  $1 \leq i \leq m$  we have  $g(\mathbf{x}_i) = y_i$ 
  return  $f : \mathbf{x} \mapsto (\mathbf{w} \cdot \mathbf{x}) + b$ 
end

```

Table 1.1 Basic Perceptron Algorithm.

perceptron
algorithm

that (\mathbf{w}, b) is only updated on a labelled example if the perceptron in state (\mathbf{w}, b) *misclassifies* the example. It is convenient to think of the algorithm as maintaining the hypothesis $g : \mathbf{x} \mapsto \text{sgn}((\mathbf{w} \cdot \mathbf{x}) + b)$, which is updated each time it misclassifies an example. The algorithm operates on a training sample by repeatedly cycling through the m examples, and when it has completed a cycle through the training data without updating its hypothesis, it returns that hypothesis.

The following result shows that if the training sample is consistent with some simple perceptron, then this algorithm converges after a finite number of iterations. In this theorem, \mathbf{w}^* and b^* define a decision boundary that correctly classifies all training points, and every training point is at least distance ρ from the decision boundary.

Theorem 1.1 Convergence of the Perceptron Algorithm

Suppose that there exists a $\rho > 0$, a weight vector \mathbf{w}^* satisfying $\|\mathbf{w}^*\| = 1$, and a threshold b^* such that

$$y_i((\mathbf{w}^* \cdot \mathbf{x}_i) + b^*) \geq \rho \text{ for all } 1 \leq i \leq m. \quad (1.10)$$

Then for all $\eta > 0$, the hypothesis maintained by the perceptron algorithm converges after no more than $(b^{*2} + 1)(R^2 + 1)/\rho^2$ updates, where $R = \max_i \|\mathbf{x}_i\|^2$. Clearly, the limiting hypothesis is consistent with the training data (X, Y) .

Proof Let (\mathbf{w}_j, b_j) be the state maintained immediately before the j th update occurring at, say, example (\mathbf{x}_i, y_i) . To measure the progress of the algorithm, we consider the evolution of the *angle* between (\mathbf{w}_j, b_j) and (\mathbf{w}^*, b^*) and note that the inner product $((\mathbf{w}_j, b_j) \cdot (\mathbf{w}^*, b^*))$ grows steadily with each update. To see this, note that (\mathbf{w}_j, b_j) is only updated when the corresponding hypothesis g_j misclassifies y_i , which implies that $y_i - g_j(\mathbf{x}_i) = 2y_i$. Therefore,

$$\begin{aligned} ((\mathbf{w}_{j+1}, b_{j+1}) \cdot (\mathbf{w}^*, b^*)) &= [((\mathbf{w}_j, b_j) + (\eta/2)(y_i - g_j(\mathbf{x}_i))(\mathbf{x}_i, 1)) \cdot (\mathbf{w}^*, b^*)] \\ &= ((\mathbf{w}_j, b_j) \cdot (\mathbf{w}^*, b^*)) + \eta y_i((\mathbf{x}_i, 1) \cdot (\mathbf{w}^*, b^*)) \\ &\geq ((\mathbf{w}_j, b_j) \cdot (\mathbf{w}^*, b^*)) + \eta\rho \\ &\geq j\eta\rho. \end{aligned}$$

On the other hand, the norm of (\mathbf{w}_j, b_j) cannot grow too fast, because on an update we have $y_i((\mathbf{w}_j \cdot \mathbf{x}_i) + b_j) < 0$, and therefore

$$\begin{aligned} \|(\mathbf{w}_{j+1}, b_{j+1})\|^2 &= \|(\mathbf{w}_j, b_j) + \eta y_i(\mathbf{x}_i, 1)\|^2 \\ &= \|(\mathbf{w}_j, b_j)\|^2 + 2\eta y_i((\mathbf{x}_i, 1) \cdot (\mathbf{w}_j, b_j)) + \eta^2 \|(\mathbf{x}_i, 1)\|^2 \\ &\leq \|(\mathbf{w}_j, b_j)\|^2 + \eta^2 \|(\mathbf{x}_i, 1)\|^2 \\ &\leq j\eta^2(R^2 + 1). \end{aligned}$$

Combining these two observations with the Cauchy-Schwarz inequality shows that

$$\begin{aligned} \sqrt{j\eta^2(R^2 + 1)} &\geq \|(\mathbf{w}_{j+1}, b_{j+1})\| \\ &\geq \frac{((\mathbf{w}_{j+1}, b_{j+1}) \cdot (\mathbf{w}^*, b^*))}{\sqrt{1 + b^{*2}}} \\ &\geq j\eta\rho, \end{aligned}$$

and thus $j \leq (1 + b^{*2})(R^2 + 1)/\rho^2$ as desired. ■

Since the perceptron algorithm makes an update at least once in every cycle through the training data, and each iteration involves $O(N)$ computation steps, this theorem implies that the perceptron algorithm has time complexity $O((R^2 + 1)mN/\rho^2)$.

1.1.3 Margins

The quantity ρ plays a crucial role in the previous theorem, since it determines how well the two classes can be separated and consequently how fast the perceptron

learning algorithm converges. This quantity ρ is what we shall henceforth call a *margin*.

Definition 1.1 Margin and Margin Errors

Denote by $f : \mathbb{R}^N \rightarrow \mathbb{R}$ a real valued hypothesis used for classification. Then

margin
$$\rho_f(\mathbf{x}, y) := yf(\mathbf{x}), \tag{1.11}$$

i.e. it is the margin by which the pattern \mathbf{x} is classified correctly (so that a negative value of $\rho_f(\mathbf{x}, y)$ corresponds to an incorrect classification). Moreover denote by

minimum margin
$$\rho_f := \min_{1 \leq i \leq m} \rho_f(\mathbf{x}_i, y_i) \tag{1.12}$$

the minimum margin over the whole sample. It is determined by the “worst” classification on the whole training set X, Y .

It appears to be desirable to have classifiers that achieve a large margin ρ_f since one might expect that an estimate that is “reliable” on the training set will also perform well on unseen examples. Moreover such an algorithm is more robust with respect to both patterns and parameters:

robustness in patterns

- Intuitively, for a pattern \mathbf{x} that is far from the decision boundary $\{\mathbf{x} | f(\mathbf{x}) = 0\}$ slight perturbations to \mathbf{x} will not change its classification $\text{sgn}(f(\mathbf{x}))$. To see this, note that if $f(\mathbf{x})$ is a continuous function in \mathbf{x} then small variations in \mathbf{x} will translate into small variations in $f(\mathbf{x})$. Therefore, if $y_i f(\mathbf{x}_i)$ is much larger than zero, $y_i f(\mathbf{x}_i \pm \varepsilon)$ will also be positive for small ε . (See, for example, Duda and Hart (1973).)

robustness in parameters

- Similarly, a slight perturbation to the function f will not affect any of the resulting classifications on the training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$. Assume that $f_{\mathbf{w}}(\mathbf{x})$ is continuous in its parameters \mathbf{w} . Then, again, if $y_i f_{\mathbf{w}}(\mathbf{x}_i)$ is much larger than zero, $y_i f_{\mathbf{w} \pm \varepsilon}(\mathbf{x}_i)$ will also be positive for small ε .

1.1.4 Maximum Margin Hyperplanes

As pointed out in the previous section, it is desirable to have an estimator with a large margin. This raises the question whether there exists an estimator with *maximum* margin, i.e. whether there exists some f^* with

$$f^* := \underset{f}{\operatorname{argmax}} \rho_f = \underset{f}{\operatorname{argmax}} \min_i y_i f(\mathbf{x}_i). \tag{1.13}$$

Without some constraint on the size of \mathbf{w} , this maximum does not exist. In Theorem 1.1, we constrained \mathbf{w}^* to have unit length. If we define $f : \mathbb{R}^N \rightarrow \mathbb{R}$ by

$$f(\mathbf{x}) = \frac{(\mathbf{w} \cdot \mathbf{x}) + b}{\|\mathbf{w}\|}, \tag{1.14}$$

optimal hyperplane then the maximum margin f is defined by the weight vector and threshold that satisfy

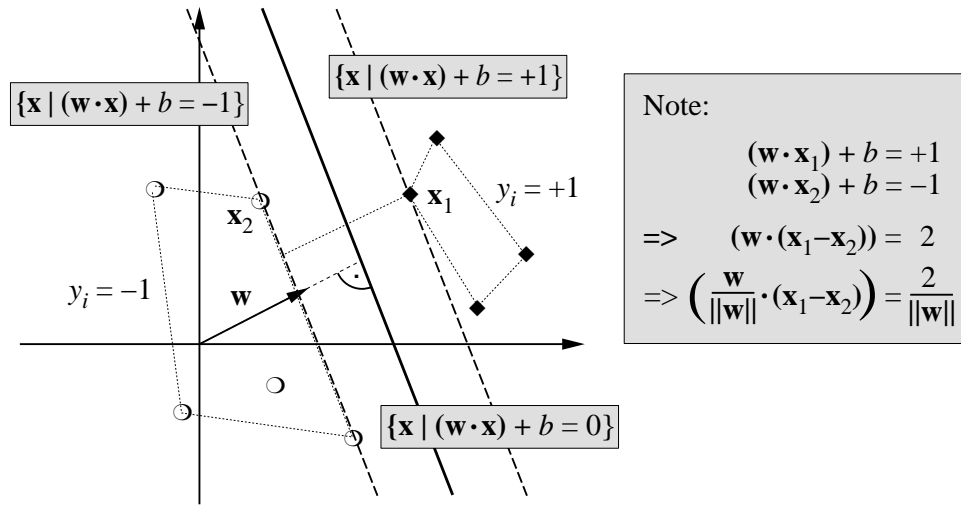


Figure 1.2 A binary classification toy problem: separate balls from diamonds. The *optimal hyperplane* is orthogonal to the shortest line connecting the convex hulls of the two classes (dotted), and intersects it half-way between the two classes. The problem being separable, there exists a weight vector \mathbf{w} and a threshold b such that $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) > 0$ ($i = 1, \dots, m$). Rescaling \mathbf{w} and b such that the point(s) closest to the hyperplane satisfy $|(\mathbf{w} \cdot \mathbf{x}_i) + b| = 1$, we obtain a *canonical form* (\mathbf{w}, b) of the hyperplane, satisfying $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1$. Note that in this case, the minimum Euclidean distance between the two classes (i.e. twice the margin), measured perpendicularly to the hyperplane, equals $2/\|\mathbf{w}\|$. This can be seen by considering two points $\mathbf{x}_1, \mathbf{x}_2$ on opposite sides of the margin, i.e. $(\mathbf{w} \cdot \mathbf{x}_1) + b = 1$, $(\mathbf{w} \cdot \mathbf{x}_2) + b = -1$, and projecting them onto the hyperplane normal vector $\mathbf{w}/\|\mathbf{w}\|$.

$$\mathbf{w}^*, b^* = \underset{\mathbf{w}, b}{\operatorname{argmax}} \min_{i=1}^m \frac{y_i((\mathbf{w} \cdot \mathbf{x}_i) + b)}{\|\mathbf{w}\|} \quad (1.15)$$

$$= \underset{\mathbf{w}, b}{\operatorname{argmax}} \min_{i=1}^m y_i \operatorname{sgn}((\mathbf{w} \cdot \mathbf{x}_i) + b) \left\| \frac{(\mathbf{w} \cdot \mathbf{x}_i)}{\|\mathbf{w}\|^2} \mathbf{w} + \frac{b}{\|\mathbf{w}\|^2} \mathbf{w} \right\| \quad (1.16)$$

Euclidean
Margin

The formulation (1.16) has a simple geometric interpretation: $-\mathbf{b}\mathbf{w}/\|\mathbf{w}\|^2$ is the vector in direction \mathbf{w} that ends right on the decision hyperplane (since $(\mathbf{w} \cdot (-\mathbf{b}\mathbf{w}/\|\mathbf{w}\|^2)) = -b$), and for a vector \mathbf{x}_i , $(\mathbf{w} \cdot \mathbf{x}_i)\mathbf{w}/\|\mathbf{w}\|^2$ is the projection of \mathbf{x}_i onto \mathbf{w} . Therefore, we are interested in maximizing the length of the vector differences $(\mathbf{w} \cdot \mathbf{x}_i)\mathbf{w}/\|\mathbf{w}\|^2 - (-\mathbf{b}\mathbf{w}/\|\mathbf{w}\|^2)$ appropriately signed by $y_i g(\mathbf{x}_i)$.

The maxi-min problem (1.15) can be easily transformed into an equivalent constrained optimization task by conjecturing a lower bound on the margin, ρ , and maximizing ρ subject to the constraint that it really is a lower bound:

$$\mathbf{w}^*, b^*, \rho^*$$

optimization
problems

$$= \operatorname{argmax}_{\mathbf{w}, b, \rho} \rho \quad \text{subject to} \quad \frac{y_i((\mathbf{w} \cdot \mathbf{x}_i) + b)}{\|\mathbf{w}\|} \geq \rho \text{ for } 1 \leq i \leq m \quad (1.17)$$

$$= \operatorname{argmax}_{\mathbf{w}, b, \rho} \rho \quad \text{subject to} \quad \|\mathbf{w}\| = 1 \text{ and } y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq \rho \text{ for } 1 \leq i \leq m \quad (1.18)$$

$$= \operatorname{argmin}_{\mathbf{w}, b} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1 \text{ for } 1 \leq i \leq m \quad (1.19)$$

quadratic
program

This last formulation is in the form of a quadratic programming problem, which can be easily handled using standard numerical routines (Luenberger, 1973; Bertsekas, 1995).

Notice that (1.18) is in a particularly intuitive form. This formulation states that we are seeking a weight vector \mathbf{w} that obtains large dot products $y_i(\mathbf{w} \cdot \mathbf{x}_i)$, but constrain the weight vector to lie on the unit sphere to prevent obtaining such large dot products “for free” by scaling up \mathbf{w} . Interesting variants of problem (1.18) are obtained by choosing different norms to constrain the length of the weight vector. For example, constraining \mathbf{w} to lie on the unit ℓ_1 sphere instead of the unit ℓ_2 sphere gives the problem of determining

$$\mathbf{w}^*, b^*, \rho^* \\ = \operatorname{argmax}_{\mathbf{w}, b, \rho} \rho \quad \text{subject to} \quad \|\mathbf{w}\|_1 = 1 \text{ and } y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq \rho \text{ for } 1 \leq i \leq m \quad (1.20)$$

ℓ_∞ margin

which can easily be shown to be in the form of a linear programming problem. Mangasarian (1997) shows that this is equivalent to finding the weight vector and threshold that maximize the minimum ℓ_∞ distance between the training patterns and the decision hyperplane, in a direct analogue to the original Euclidean formulation (1.15).

Similarly, the constraint that \mathbf{w} lie on the unit ℓ_∞ sphere yields the problem

$$\mathbf{w}^*, b^*, \rho^* \\ = \operatorname{argmax}_{\mathbf{w}, b, \rho} \rho \quad \text{subject to} \quad \|\mathbf{w}\|_\infty = 1 \text{ and } y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq \rho \text{ for } 1 \leq i \leq m \quad (1.21)$$

ℓ_1 margin

which is also a linear programming problem, but now equivalent to finding the weight vector and threshold that maximize the minimum ℓ_1 distance between the training patterns and the decision hyperplane. In general, constraining \mathbf{w} to lie on the unit ℓ_p sphere yields a convex programming problem

$$\mathbf{w}^*, b^*, \rho^* \\ = \operatorname{argmax}_{\mathbf{w}, b, \rho} \rho \quad \text{subject to} \quad \|\mathbf{w}\|_p = 1 \text{ and } y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq \rho \text{ for } 1 \leq i \leq m \quad (1.22)$$

ℓ_q margin

which is equivalent to finding the weight vector and threshold that maximize the minimum ℓ_q distance between the training patterns and the decision hyperplane, where ℓ_p and ℓ_q are conjugate norms, i.e. such that $\frac{1}{p} + \frac{1}{q} = 1$ (Mangasarian, 1997).

In solving any of these constrained optimization problems, there is a notion of *critical constraints*; i.e. those inequality constraints that are satisfied as equalities by the optimal solution. In our setting, constraints correspond to training examples (\mathbf{x}_i, y_i) , $1 \leq i \leq m$, and the *critical* constraints are given by those training

Support Vectors examples that lie right on the margin a distance ρ from the optimal hyperplane (cf. Figure 1.2). These critical training patterns are called *Support Vectors*.

Notice that all the remaining examples of the training set are irrelevant: for non-critical examples the corresponding constraint $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1$ in (1.19) does not play a role in the optimization, and therefore these points could be removed from the training set without affecting the results. This nicely captures our intuition of the problem: the hyperplane (cf. Figure 1.2) is completely determined by the patterns closest to it, the solution should not depend on the other examples.

soft margin hyperplane

In practice, a separating hyperplane may not exist, e.g. if a high noise level causes a large overlap of the classes. The previous maximum margin algorithms perform poorly in this case because the maximum achievable minimum margin is negative, and this means the critical constraints are the mislabelled patterns that are furthest from the decision hyperplane. That is, the solution hyperplane is determined entirely by misclassified examples! To overcome the sensitivity to noisy training patterns, a standard approach is to allow for the possibility of examples violating the constraint in (1.19) by introducing *slack variables* (Cortes and Vapnik, 1995; Vapnik, 1995)

slack variables

$$\xi_i \geq 0, \text{ for all } i = 1, \dots, m, \quad (1.23)$$

along with relaxed constraints

$$y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1 - \xi_i, \text{ for all } i = 1, \dots, m. \quad (1.24)$$

A classifier which generalizes well is then found by controlling both the size of \mathbf{w} and the number of training errors, minimizing the objective function

$$\tau(\mathbf{w}, \boldsymbol{\xi}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \quad (1.25)$$

subject to the constraints (1.23) and (1.24), for some value of the constant $C > 0$.

In the following section, we shall see why the size of \mathbf{w} is a good measure of the complexity of the classifier.

1.2 Theory

In order to provide a theoretical analysis of the learning problem we have to introduce a few definitions and assumptions about the process generating the data.

1.2.1 Basic Assumptions

independently
identically
distributed

We assume that the training data X, Y is drawn independently and identically distributed (iid) according to some probability measure $p(\mathbf{x}, y)$. This means that all examples (\mathbf{x}_i, y_i) are drawn from $p(\mathbf{x}, y)$ regardless of the other examples or the index i .

This assumption is stronger than it may appear at first glance. For instance, time series data fails to satisfy the condition, since the observations are typically dependent, and their statistics might depend on the index i .

In (1.4), we defined the functional $R(g)$ of a decision function g as the probability of misclassification. We can generalize this definition to apply to prediction functions f as well as thresholded decision functions g . This yields what we call the risk functional.

Definition 1.2 Risk Functional

Expected Risk

Denote by $c(\mathbf{x}, y, f(\mathbf{x})) : \mathbb{R}^N \times \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$ a cost function and by $p(\mathbf{x}, y)$ a probability measure as described above. Then the risk functional for a function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ is defined as

$$R(f) := \int_{\mathbb{R}^N \times \mathbb{R}} c(\mathbf{x}, y, f(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x}dy. \quad (1.26)$$

Moreover the *empirical* risk functional for an m -sample X, Y is given by

Empirical Risk

$$R_{\text{emp}}(f) := \frac{1}{m} \sum_{i=1}^m c(\mathbf{x}_i, y_i, f(\mathbf{x}_i)). \quad (1.27)$$

For thresholded decision functions $g : \mathbb{R}^N \rightarrow \{-1, 1\}$ we often use 0–1 classification error as the cost function $c(\mathbf{x}, y, g(\mathbf{x})) = 1_{\{g(\mathbf{x}) \neq y\}}$. In this case we obtain the risk functional defined in (1.4) (the probability of misclassification),

$$R(g) := \Pr\{g(\mathbf{x}) \neq y\}. \quad (1.28)$$

In this case, the empirical risk functional is

$$R_{\text{emp}}(g) := \frac{1}{m} \sum_{i=1}^m 1_{\{g(\mathbf{x}_i) \neq y_i\}}, \quad (1.29)$$

which is just the training error.

margin error

Finally we need a quantity called the *margin error* which is given by the proportion of training points that have margin less than ρ , i.e.

$$R_{\rho}(f) := \frac{1}{m} \sum_{i=1}^m 1_{\{y_i f(\mathbf{x}_i) < \rho\}}. \quad (1.30)$$

This empirical estimate of risk counts a point as an error if it is either incorrectly classified or correctly classified by with margin less than ρ .

While one wants to minimize the risk $R(g)$ this is hardly ever possible since $p(\mathbf{x}, y)$ is unknown. Hence one may only resort to minimizing $R_{\text{emp}}(g)$ which is based on the training data. This, however, is not an effective method by itself—just consider an estimator that memorizes all the training data X, Y and generates random outputs for any other data. This clearly would have an empirical risk $R_{\text{emp}}(g) = 0$ but would obtain a true risk $R(g) = 0.5$ (assuming the finite training sample has measure 0). The solution is to take the complexity of the estimate g into account as well, which will be discussed in the following sections.

1.2.2 Error Bounds for Thresholded Decision Functions

VC dimension

The central result of this analysis is to relate the number of training examples, the training set error, and the complexity of the hypothesis space to the generalization error. For thresholded decision functions, an appropriate measure for the complexity of the hypothesis space is the Vapnik-Chervonenkis (VC) dimension.

Definition 1.3 VC dimension

The VC dimension h of a space of $\{-1, 1\}$ -valued functions, G , is the size of the largest subset of domain points that can be labelled arbitrarily by choosing functions only from G (Vapnik and Chervonenkis, 1971).

The VC dimension can be used to prove high probability bounds on the error of a hypothesis chosen from a class of decision functions G —this is the famous result of Vapnik and Chervonenkis (1971). The bounds have since been improved slightly by Talagrand (1994)—see also (Alexander, 1984).

Theorem 1.2 VC Upper Bound

Let G be a class of decision functions mapping \mathbb{R}^N to $\{-1, 1\}$ that has VC dimension h . For any probability distribution $p(\mathbf{x}, y)$ on $\mathbb{R}^N \times \{-1, 1\}$, with probability at least $1 - \delta$ over m random examples \mathbf{x} , for any hypothesis g in G the risk functional with 0–1 loss is bounded by

$$R(g) \leq R_{\text{emp}}(g) + \sqrt{\frac{c}{m} \left(h + \ln \left(\frac{1}{\delta} \right) \right)} \quad (1.31)$$

where c is a universal constant. Furthermore, if $g^* \in G$ minimizes $R_{\text{emp}}(\cdot)$, then with probability $1 - \delta$

$$R(g^*) \leq \inf_{g \in G} R(g) + \sqrt{\frac{c}{m} \left(h + \ln \left(\frac{1}{\delta} \right) \right)} \quad (1.32)$$

(A short proof of this result is given by Long (1998), but with worse constants than Talagrand's.)

These upper bounds are asymptotically close to the best possible, since there is also a lower bound with the same form:

Theorem 1.3 VC Lower Bound

Let G be a hypothesis space with finite VC dimension $h \geq 1$. Then for any learning algorithm there exist distributions such that with probability at least δ over m random examples, the error of its hypothesis g satisfies

$$R(g) \geq \inf_{g' \in G} R(g') + \sqrt{\frac{c}{m} \left(h + \ln \left(\frac{1}{\delta} \right) \right)} \quad (1.33)$$

where c is a universal constant.

(Results of this form have been given by Devroye and Lugosi (1995); Simon (1996); Anthony and Bartlett (1999), using ideas from Ehrenfeucht et al. (1989).)

Theorems 1.2 and 1.3 give a fairly complete characterization of the generalization error that can be achieved by choosing decision functions from a class G . However, this characterization suffers from two drawbacks.

- The first drawback is that the VC dimension must actually be determined (or at least bounded) for the class of interest—and this is often not easy to do. (However, bounds on the VC dimension h have been computed for many natural decision function classes, including parametric classes involving standard arithmetic and boolean operations. See Anthony and Bartlett (1999) for a review of these results.)
- The second (more serious) drawback is that the analysis ignores the *structure* of the mapping from training samples to hypotheses, and concentrates solely on the *range* of the learner’s possible outputs. Ignoring the details of the learning map can omit many of the factors that are *crucial* for determining the success of the learning algorithm in real situations.

For example, consider learning algorithms that operate by first computing a real valued prediction function f from some class F and then thresholding this hypothesis to obtain the final decision function $g(\mathbf{x}) = \text{sgn}(f(\mathbf{x}))$. Here, the VC dimension is a particularly weak method for measuring the representational capacity of the resulting function class $G = \text{sgn}(F)$.

One reason is that the VC dimension of G is not sensitive to the *scale* of F at the accuracy level of interest. That is, it does not pay attention to whether the complexity of the hypothesis class is at a scale that is relevant for the outcome of the predictions.

The first step towards a more refined analysis that takes scale into account is given by Vapnik (1979). Consider a set $X_0 \subset \mathbb{R}^N$ of input points with norm bounded by $R > 0$ (that is, $\|\mathbf{x}_i\| \leq R$ for $\mathbf{x} \in X_0$), and the set F of bounded linear functions defined on X_0 ,

$$F = \{\mathbf{x} \mapsto (\mathbf{w} \cdot \mathbf{x}) \mid \|\mathbf{w}\| \leq 1, \mathbf{x} \in X_0\} \quad (1.34)$$

satisfying $|f(\mathbf{x})| \geq \rho$ for all patterns \mathbf{x} in X_0 . Then if we consider the set G of linear decision functions obtained by thresholding functions in F , Vapnik (1979) shows

$$\text{VCdim}(G) \leq \min\{R^2/\rho^2, N\} + 1. \quad (1.35)$$

Note that this can be much smaller than the VC dimension of $\text{sgn}(F)$ obtained without taking ρ into account, which is $N + 1$ in this case. Therefore, one could hope to obtain significant benefits by using scale sensitive bounds which give much tighter results for large margin classifiers. Unfortunately, the bound (1.35) does not yet suffice for our purposes, because note that it requires that *all* points (including the test points) satisfy the margin condition, and therefore theorem 1.2 does not apply in this case. Rigorously obtaining these scale sensitive improvements is the topic we now address. In the following section, we consider scale-sensitive versions of the VC dimension, and obtain upper and lower bounds on risk in terms of these dimensions.

1.2.3 Margin Dependent Error Bounds for Real Valued Predictors

Definition 1.4 Fat Shattering Dimension

Let F be a set of real valued functions. We say that a set of points $S \subset \mathcal{X}$, which we will index as a vector $\mathbf{x} \in \mathcal{X}^{|S|}$, is ρ -shattered by F if there is a vector of real numbers $\mathbf{b} \in \mathbb{R}^{|S|}$ such that for any choice of signs $\mathbf{y} \in \{-1, 1\}^{|S|}$ there is a function f in F that satisfies

$$y_i(f(x_i) - b_i) \geq \rho \text{ for } 1 \leq i \leq |S|. \quad (1.36)$$

(That is, $f(x_i) \geq b_i + \rho$ if $y_i = 1$, and $f(x_i) \leq b_i - \rho$ if $y_i = -1$, for all x_i in S . Notice how similar this is to the notion of a minimum margin defined by (1.12).)

fat shattering

The *fat shattering dimension* $\text{fat}_F(\rho)$ of the set F is a function from the positive real numbers to the integers which maps a value ρ to the size of the largest ρ -shattered set, if this is finite, or infinity otherwise.

We may think of the fat-shattering dimension of a set of real-valued functions as the VC dimension obtained by thresholding but requiring that outputs are ρ above the threshold for positive classification and ρ below for negative.

The fat-shattering dimension is closely related to a more basic quantity, the covering number of a class of functions.

Definition 1.5 Covering Numbers of a Set

covering number

Denote by (S, d) a pseudometric space, $B_r(\mathbf{x})$ the closed ball in S centred at \mathbf{x} with radius r , T a subset of S , and ε some positive constant. Then the covering number $\mathcal{N}(\varepsilon, T)$ is defined as the minimum cardinality (that is, number of elements) of a set of points $T' \subset S$ such that

$$T \subseteq \bigcup_{\mathbf{x}_i \in T'} B_\varepsilon(\mathbf{x}_i), \quad (1.37)$$

i.e. such that the maximum difference of any element in T and the closest element in T' is less than or equal to ε .

Covering a class of functions F with an ε -cover means that one is able to approximately represent F (which may be of infinite cardinality) by a finite set. For learning, it turns out that it suffices to approximate the restrictions of functions in a class F to finite samples. For a subset X of some domain \mathcal{X} , define the pseudometric $\ell_{\infty, X}$ by

$$\ell_{\infty, X}(f, f') = \max_{\mathbf{x} \in X} |f(\mathbf{x}) - f'(\mathbf{x})| \quad (1.38)$$

where f and f' are real-valued functions defined on \mathcal{X} . Let $\mathcal{N}(\varepsilon, F, m)$ denote the maximum, over all $X \subset \mathcal{X}$ of size $|X| = m$, of the covering number $\mathcal{N}(\varepsilon, F)$ with respect to $\ell_{\infty, X}$. The following theorem shows that the fat-shattering dimension is intimately related to these covering numbers. (The upper bound is due to Alon et al. (1997), and the lower bound to Bartlett et al. (1997).)

Theorem 1.4 Bounds on \mathcal{N} in terms of fat_F

Let F be a set of real functions from a domain \mathcal{X} to the bounded interval $[0, B]$. Let $\varepsilon > 0$ and let $m \geq \text{fat}_F(\varepsilon/4)$. Then

$$\frac{\log_2 e}{8} \text{fat}_F(16\varepsilon) \leq \log_2 \mathcal{N}(\varepsilon, F, m) \leq 3 \text{fat}_F\left(\frac{\varepsilon}{4}\right) \log_2^2 \left(\frac{4eBm}{\varepsilon} \right). \quad (1.39)$$

Unfortunately, directly bounding \mathcal{N} can be quite difficult in general. Useful tools from functional analysis (which deal with the functional inverse of \mathcal{N} wrt. ε , the so called entropy number) for obtaining these bounds have been developed for classes of functions F defined by linear mappings from Hilbert spaces (Carl and Stephani, 1990), and linear functions over kernel expansions (Williamson et al., 1998b).

The following result shows that we can use covering numbers to obtain upper bounds on risk in terms of margin error (Shawe-Taylor et al., 1998; Bartlett, 1998).

Theorem 1.5 Bounds on $R(f)$ in terms of \mathcal{N} and ρ

Suppose that F is a set of real-valued functions defined on \mathcal{X} , $\varepsilon \in (0, 1)$ and $\rho > 0$. Fix a probability distribution on $\mathcal{X} \times \{-1, 1\}$ and a sample size m . Then the probability that some f in F has $R_\rho(f) = 0$ but $R(f) \geq \varepsilon$ is no more than

$$2 \mathcal{N}\left(\frac{\rho}{2}, F, 2m\right) 2^{-\varepsilon m/2}. \quad (1.40)$$

Furthermore,

$$\Pr(\text{“some } f \text{ in } F \text{ has } R(f) \geq R_\rho(f) + \varepsilon\text{”}) \leq 2 \mathcal{N}\left(\frac{\rho}{2}, F, 2m\right) e^{-\varepsilon^2 m/8}. \quad (1.41)$$

In fact, it is possible to obtain a similar result that depends only on the behaviour of functions in F near the threshold (see (Anthony and Bartlett, 1999) for details).

anatomy of a
uniform conver-
gence bound

Let us have a close look at the bound (1.41) on the probability of excessive error. The factor $e^{-\varepsilon^2 m/8}$ in (1.41) stems from a bound of Hoeffding (1963) on the probability of a large deviation of a sum of random variables from its mean. The factor $\mathcal{N}\left(\frac{\rho}{2}, F, 2m\right)$ stems from the fact that the continuous class of functions F was approximated (to accuracy $\rho/2$) by a finite number of functions. The $2m$ is due to the use of a symmetrization argument which is needed to make the overall argument work. Theorem 1.4 shows that this term is bounded by an exponential function of the fat-shattering dimension at scale $\rho/8$.

Interestingly, a similar result holds in regression. (For a review of these uniform convergence results, see (Anthony and Bartlett, 1999)).

Theorem 1.6 Bounds on $R(f)$ for Regression

Suppose that F is a set of functions defined on a domain \mathcal{X} and mapping into the real interval $[0, 1]$. Let p be any probability distribution on $\mathcal{X} \times [0, 1]$, ε any real number between 0 and 1, and $m \in \mathbb{N}$. Then for the quadratic cost function $c(\mathbf{x}, y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$ we have

$$\Pr\left(\sup_{f \in F} |R(f) - R_{\text{emp}}(f)| \geq \varepsilon\right) \leq 4 \mathcal{N}\left(\frac{\varepsilon}{16}, F, 2m\right) e^{-\varepsilon^2 m/32}. \quad (1.42)$$

Comparing with (1.41), notice that the scale of the covering number depends on the desired accuracy ε , whereas in (1.41) it depends on the scale ρ at which the margins are examined.

1.2.4 Error Bounds for Linear Decision Functions

The following result, due to Bartlett and Shawe-Taylor (1999), gives a bound on the fat-shattering dimension of large margin linear classifiers. It has a similar form to the bound (1.35) on the VC dimension of linear functions restricted to certain sets. It improves on a straightforward corollary of that result, and on a result of Gurvits (1997).

Theorem 1.7 Fat Shattering Dimension for Linear Classifiers

Suppose that B_R is the ℓ_2 ball of radius R in \mathbb{R}^n , centered at the origin, and consider the set

$$F := \{f_{\mathbf{w}} \mid f_{\mathbf{w}}(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) \text{ with } \|\mathbf{w}\| \leq 1, \mathbf{x} \in B_R\}. \quad (1.43)$$

Then

$$\text{fat}_F(\rho) \leq \left(\frac{R}{\rho}\right)^2. \quad (1.44)$$

Using this result together with Theorems 1.4 and 1.5 gives the following theorem.

Theorem 1.8 Error Bounds for Linear Classifiers

Define the class F of real-valued functions on the ball of radius R as in (1.43). There is a constant c such that, for all probability distributions, with probability at least $1 - \delta$ over m independently generated training examples, every $\rho > 0$ and every function $f \in F$ with margin at least ρ on all training examples (i.e. $R_\rho(f) = 0$) satisfies

$$R(f) \leq \frac{c}{m} \left(\frac{R^2}{\rho^2} \log^2 \left(\frac{m}{\rho} \right) + \log \left(\frac{1}{\delta} \right) \right). \quad (1.45)$$

Furthermore, with probability at least $1 - \delta$, for all $\rho > 0$, every function f in F has error

$$R(f) \leq R_\rho(f) + \sqrt{\frac{c}{m} \left(\frac{R^2}{\rho^2} \log^2 \left(\frac{m}{\rho} \right) + \log \left(\frac{1}{\delta} \right) \right)}. \quad (1.46)$$

For estimators using a linear programming approach as in (Mangasarian, 1968) one may state the following result, which then, via Theorem 1.4 can be transformed into a generalization bound as well.

Theorem 1.9 Capacity Bounds for Linear Classifiers

There is a constant c such that for the class

$$F_R = \{\mathbf{x} \mapsto \mathbf{w}^T \mathbf{x} \mid \|\mathbf{x}\|_\infty \leq 1, \|\mathbf{w}\|_1 \leq R\} \quad (1.47)$$

we have

$$\text{fat}_{F_R}(\varepsilon) \leq c \left(\frac{R}{\varepsilon} \right)^2 \ln(2N + 2). \quad (1.48)$$

Finally, we can obtain bounds for convex combinations of arbitrary hypotheses from a class G of $\{-1, 1\}$ -valued functions,

$$\text{co}(G) = \left\{ \sum_i \alpha_i g_i \mid \alpha_i > 0, \sum_i \alpha_i = 1, g_i \in G \right\}. \quad (1.49)$$

See (Schapire et al., 1998). These bounds are useful in analysing boosting algorithms; see Section 1.4.

Theorem 1.10 Bounds for Convex Combinations of Hypotheses

Let $p(\mathbf{x}, y)$ be a distribution over $\mathcal{X} \times \{-1, 1\}$, and let X be a sample of m examples chosen iid according to p . Suppose the base-hypothesis space G has VC dimension h , and let $\delta > 0$. Then with probability at least $1 - \delta$ over the random choice of the training set X, Y , every convex combination of functions $f \in \text{co}(G)$ satisfies the following bound for all $\rho > 0$.

$$R(f) \leq R_\rho(f) + \sqrt{\frac{c}{m} \left(\frac{h \log^2(m/h)}{\rho^2} + \log \left(\frac{1}{\delta} \right) \right)} \quad (1.50)$$

1.3 Support Vector Machines

1.3.1 Optimization Problem

To construct the *Optimal Hyperplane* (cf. Figure 1.2), one solves the following optimization problem:

$$\text{minimize} \quad \tau(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 \quad (1.51)$$

$$\text{subject to} \quad y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \geq 1, \text{ for all } i = 1, \dots, m. \quad (1.52)$$

Lagrangian

This constrained optimization problem is dealt with by introducing Lagrange multipliers $\alpha_i \geq 0$ and a Lagrangian

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i (y_i((\mathbf{x}_i \cdot \mathbf{w}) + b) - 1). \quad (1.53)$$

The Lagrangian L has to be minimized with respect to the *primal variables* \mathbf{w} and b and maximized with respect to the *dual variables* α_i (i.e. a saddle point has to be found). Let us try to get some intuition for this. If a constraint (1.52) is violated, then $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) - 1 < 0$, in which case L can be increased by increasing the corresponding α_i . At the same time, \mathbf{w} and b will have to change such that L decreases. To prevent $-\alpha_i (y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) - 1)$ from becoming arbitrarily large, the change in \mathbf{w} and b will ensure that, provided the problem is separable, the

constraint will eventually be satisfied.

KKT
conditions

Similarly, one can understand that for all constraints which are not precisely met as equalities, i.e. for which $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) - 1 > 0$, the corresponding α_i must be 0: this is the value of α_i that maximizes L . The latter is the statement of the Karush-Kuhn-Tucker complementarity conditions of optimization theory (Karush, 1939; Kuhn and Tucker, 1951; Bertsekas, 1995).

The condition that at the saddle point, the derivatives of L with respect to the primal variables must vanish,

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \boldsymbol{\alpha}) = 0 \text{ and } \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \boldsymbol{\alpha}) = 0, \quad (1.54)$$

leads to

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad (1.55)$$

and

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i. \quad (1.56)$$

support vector
expansion

The solution vector thus has an expansion in terms of a subset of the training patterns, namely those patterns whose Lagrange multiplier α_i is non-zero. By the Karush-Kuhn-Tucker complementarity conditions these training patterns are the ones for which

$$\alpha_i (y_i((\mathbf{x}_i \cdot \mathbf{w}) + b) - 1) = 0, \quad i = 1, \dots, m, \quad (1.57)$$

and therefore they correspond precisely to the *Support Vectors* (i.e. critical constraints) discussed in Section 1.1.4. Thus we have the satisfying result that the Support Vectors are the only training patterns that determine the optimal decision hyperplane; all other training patterns are irrelevant and do not appear in the expansion (1.56).

dual
optimization
problem

By substituting (1.55) and (1.56) into L , one eliminates the primal variables and arrives at the Wolfe dual of the optimization problem (e.g. Bertsekas, 1995): find multipliers α_i which

$$\text{maximize} \quad W(\boldsymbol{\alpha}) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) \quad (1.58)$$

$$\text{subject to} \quad \alpha_i \geq 0 \text{ for all } i = 1, \dots, m, \text{ and } \sum_{i=1}^m \alpha_i y_i = 0. \quad (1.59)$$

The hyperplane decision function can thus be written as

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^m y_i \alpha_i (\mathbf{x} \cdot \mathbf{x}_i) + b \right) \quad (1.60)$$

where b is computed using (1.57).

The structure of the optimization problem closely resembles those that typically arise in Lagrange's formulation of mechanics (e.g. Goldstein, 1986). In that case also, it is often only a subset of the constraints that are active. For instance, if we keep a ball in a box, then it will typically roll into one of the corners. The constraints corresponding to the walls which are not touched by the ball are irrelevant, the walls could just as well be removed.

Seen in this light, it is not too surprising that it is possible to give a mechanical interpretation of optimal margin hyperplanes (Burges and Schölkopf, 1997): If we assume that each support vector \mathbf{x}_i exerts a perpendicular force of size α_i and sign y_i on a solid plane sheet lying along the hyperplane, then the solution satisfies the requirements of mechanical stability. The constraint (1.55) states that the forces on the sheet sum to zero; and (1.56) implies that the torques also sum to zero, via $\sum_i \mathbf{x}_i \times y_i \alpha_i \mathbf{w} / \|\mathbf{w}\| = \mathbf{w} \times \mathbf{w} / \|\mathbf{w}\| = 0$.

1.3.2 Feature Spaces and Kernels

feature space

To construct *Support Vector Machines*, the optimal hyperplane algorithm is augmented by a method for computing dot products in feature spaces that are *nonlinearly* related to input space (Aizerman et al., 1964; Boser et al., 1992). The basic idea is to map the data into some other dot product space (called the *feature space*) \mathcal{F} via a nonlinear map

$$\Phi : \mathbb{R}^N \rightarrow \mathcal{F}, \quad (1.61)$$

and then in the space \mathcal{F} perform the linear algorithm described above.

For instance, suppose we are given patterns $\mathbf{x} \in \mathbb{R}^N$ where most information is contained in the d -th order products (monomials) of entries x_j of \mathbf{x} , i.e. $x_{j_1} x_{j_2} \cdots x_{j_d}$, where $j_1, \dots, j_d \in \{1, \dots, N\}$. In that case, we might prefer to extract these monomial features first, and work in the feature space \mathcal{F} of all products of d entries.

This approach, however, fails for realistically sized problems: for N -dimensional input patterns, there exist $(N + d - 1)! / (d!(N - 1)!)$ different monomials. Already 16×16 pixel input images (e.g. in character recognition) and a monomial degree $d = 5$ yield a dimensionality of 10^{10} .

This problem can be overcome by noticing that both the construction of the optimal hyperplane in \mathcal{F} (cf. (1.58)) and the evaluation of the corresponding decision function (1.60) only require the evaluation of dot products $(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}'))$, and never require the mapped patterns $\Phi(\mathbf{x})$ in explicit form. This is crucial, since in some cases, the dot products can be evaluated by a simple kernel (Aizerman et al., 1964; Boser et al., 1992).

Mercer kernel

$$k(\mathbf{x}, \mathbf{x}') = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')). \quad (1.62)$$

polynomial
kernel

For instance, the polynomial kernel

$$k(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}')^d \quad (1.63)$$

can be shown to correspond to a map Φ into the space spanned by all products of exactly d dimensions of \mathbb{R}^N (Poggio (1975); Boser et al. (1992)). For a proof, see Schölkopf (1997). For $d = 2$ and $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^2$, for example, we have (Vapnik, 1995)

$$(\mathbf{x} \cdot \mathbf{x}')^2 = (x_1^2, x_2^2, \sqrt{2} x_1 x_2)(y_1^2, y_2^2, \sqrt{2} y_1 y_2)^\top = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')), \quad (1.64)$$

defining $\Phi(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2} x_1 x_2)$.

By using $k(\mathbf{x}, \mathbf{x}') = ((\mathbf{x} \cdot \mathbf{x}') + c)^d$ with $c > 0$, we can take into account all product of order up to d (i.e. including those of order smaller than d).

More generally, the following theorem of functional analysis shows that kernels k of positive integral operators give rise to maps Φ such that (1.62) holds (Mercer, 1909; Aizerman et al., 1964; Boser et al., 1992; Dunford and Schwartz, 1963):

Theorem 1.11 Mercer

If k is a continuous symmetric kernel of a positive integral operator T , i.e.

$$(Tf)(\mathbf{x}') = \int_{\mathcal{X}} k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) d\mathbf{x} \quad (1.65)$$

with

$$\int_{\mathcal{X} \times \mathcal{X}} k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \geq 0 \quad (1.66)$$

for all $f \in L_2(\mathcal{X})$ (\mathcal{X} being a compact subset of \mathbb{R}^N), it can be expanded in a uniformly convergent series (on $\mathcal{X} \times \mathcal{X}$) in terms of T 's eigenfunctions ψ_j and positive eigenvalues λ_j ,

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^{N_{\mathcal{F}}} \lambda_j \psi_j(\mathbf{x}) \psi_j(\mathbf{x}'), \quad (1.67)$$

where $N_{\mathcal{F}} \leq \infty$ is the number of positive eigenvalues.

An equivalent way to characterize Mercer kernels is that they give rise to positive matrices $K_{ij} := k(\mathbf{x}_i, \mathbf{x}_j)$ for all $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ (Saitoh, 1988).

From (1.67), it is straightforward to construct a map Φ into a potentially infinite-dimensional l_2 space which satisfies (1.62). For instance, we may use

$$\Phi(\mathbf{x}) = (\sqrt{\lambda_1} \psi_1(\mathbf{x}), \sqrt{\lambda_2} \psi_2(\mathbf{x}), \dots). \quad (1.68)$$

Rather than thinking of the feature space as an l_2 space, we can alternatively represent it as the Hilbert space \mathcal{H}_k containing all linear combinations of the functions $f(\cdot) = k(\mathbf{x}_i, \cdot)$ ($\mathbf{x}_i \in \mathcal{X}$). To ensure that the map $\Phi : \mathcal{X} \rightarrow \mathcal{H}_k$, which in this case is defined as

$$\Phi(\mathbf{x}) = k(\mathbf{x}, \cdot), \quad (1.69)$$

satisfies (1.62), we need to endow \mathcal{H}_k with a suitable dot product $\langle \cdot, \cdot \rangle$. In view of the definition of Φ , this dot product needs to satisfy

$$\langle k(\mathbf{x}, \cdot), k(\mathbf{x}', \cdot) \rangle = k(\mathbf{x}, \mathbf{x}'), \quad (1.70)$$

positive
integral
operator

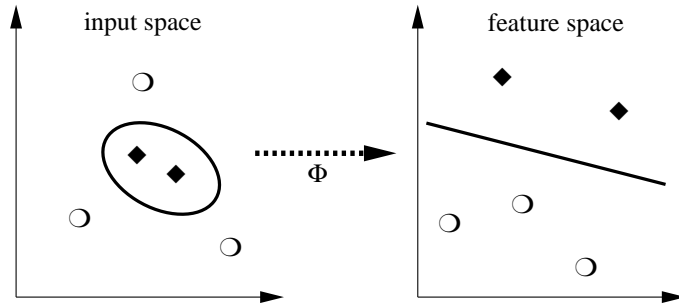


Figure 1.3 The idea of SV machines: map the training data nonlinearly into a higher-dimensional feature space via Φ , and construct a separating hyperplane with maximum margin there. This yields a nonlinear decision boundary in input space. By the use of a kernel function (1.62), it is possible to compute the separating hyperplane without explicitly carrying out the map into the feature space.

reproducing kernel

which amounts to saying that k is a *reproducing kernel* for \mathcal{H}_k . For a Mercer kernel (1.67), such a dot product does exist. Since k is symmetric, the ψ_i ($i = 1, \dots, N_{\mathcal{F}}$) can be chosen to be orthogonal with respect to the dot product in $L_2(C)$, i.e. $(\psi_j, \psi_n)_{L_2(C)} = \delta_{jn}$, using the Kronecker δ_{jn} . From this, we can construct $\langle \cdot, \cdot \rangle$ such that

$$\langle \sqrt{\lambda_j} \psi_j, \sqrt{\lambda_n} \psi_n \rangle = \delta_{jn}. \tag{1.71}$$

Substituting (1.67) into (1.70) then proves the desired equality (for further details, see Aronszajn (1950); Wahba (1973); Girosi (1998); Schölkopf (1997)).

sigmoid kernel

Besides (1.63), SV practitioners use sigmoid kernels

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\kappa(\mathbf{x} \cdot \mathbf{x}') + \Theta) \tag{1.72}$$

Gaussian RBF kernel

for suitable values of gain κ and threshold Θ , and radial basis function kernels, as for instance (Aizerman et al., 1964; Boser et al., 1992; Schölkopf et al., 1997)

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / (2\sigma^2)), \tag{1.73}$$

with $\sigma > 0$. Note that when using Gaussian kernels, for instance, the feature space \mathcal{H}_k thus contains all superpositions of Gaussians on \mathcal{X} (plus limit points), whereas by definition of Φ (1.69), only single bumps $k(\mathbf{x}, \cdot)$ do have pre-images under Φ .

The main lesson from the study of kernel functions, is that the use of kernels can turn any algorithm that only depends on dot products into a nonlinear algorithm which is linear in feature space. In the time since this was explicitly pointed out (Schölkopf et al., 1998c) a number of such algorithms have been proposed: until then the applications of the kernel trick were a proof of the convergence of rbf network training by (Aizerman et al., 1964) and the nonlinear variant of the SV algorithm by Boser et al. (1992) (see Figure 1.3). To construct SV machines, one computes an optimal hyperplane in feature space. To this end, we substitute $\Phi(\mathbf{x}_i)$ for each training example \mathbf{x}_i . The weight vector (cf. (1.56)) then becomes an expansion in

feature space. Note that \mathbf{w} will typically no more correspond to the image of just a single vector from input space (cf. Schölkopf et al. (1998a) for a formula to compute the pre-image if it exists), in other words, \mathbf{w} may not be directly accessible any more. However, since all patterns only occur in dot products, one can substitute Mercer kernels k for the dot products (Boser et al., 1992; Guyon et al., 1993), leading to decision functions of the more general form (cf. (1.60))

decision
function

$$g(\mathbf{x}) = \operatorname{sgn} \left(\sum_{i=1}^m y_i \alpha_i (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}_i)) + b \right) = \operatorname{sgn} \left(\sum_{i=1}^m y_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b \right) \quad (1.74)$$

and the following quadratic program (cf. (1.58)):

$$\text{maximize} \quad W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \quad (1.75)$$

$$\text{subject to} \quad \alpha_i \geq 0, \quad i = 1, \dots, m, \quad \text{and} \quad \sum_{i=1}^m \alpha_i y_i = 0. \quad (1.76)$$

soft margin
and kernels

Recall that, as discussed in Section 1.1.4 a separating hyperplane may not always exist, even in the expanded feature space \mathcal{F} . To cope with this difficulty, slack variables were introduced to yield the *soft margin* optimal hyperplane problem (1.25). Incorporating kernels, and rewriting (1.25) in terms of Lagrange multipliers, this again leads to the problem of maximizing (1.75), but now subject to the constraints

$$0 \leq \alpha_i \leq C, \quad i = 1, \dots, m, \quad \text{and} \quad \sum_{i=1}^m \alpha_i y_i = 0. \quad (1.77)$$

The only difference from the separable case (1.76) is the upper bound C on the Lagrange multipliers α_i . This way, the influence of the individual patterns (which could always be outliers) gets limited. As above, the solution takes the form (1.74). The threshold b can be computed by exploiting the fact that for all SVs \mathbf{x}_i with $\alpha_i < C$, the slack variable ξ_i is zero (this again follows from the Karush-Kuhn-Tucker complementarity conditions), and hence

$$\sum_{j=1}^m y_j \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + b = y_i. \quad (1.78)$$

If one uses an optimizer that works with the double dual (e.g. Vanderbei, 1997), one can also recover the value of the primal variable b directly from the corresponding double dual variable.

Finally, the algorithm can be modified such that it does not require the regularization constant C . Instead, one specifies an upper bound $0 \leq \nu \leq 1$ on the fraction of points allowed to lie in the margin (asymptotically, the number of SVs) (Schölkopf et al., 1998d). This leaves us with a homogeneous target function made up by the quadratic part of (1.75), and an additional lower bound constraint on the sum over all Lagrange multipliers.

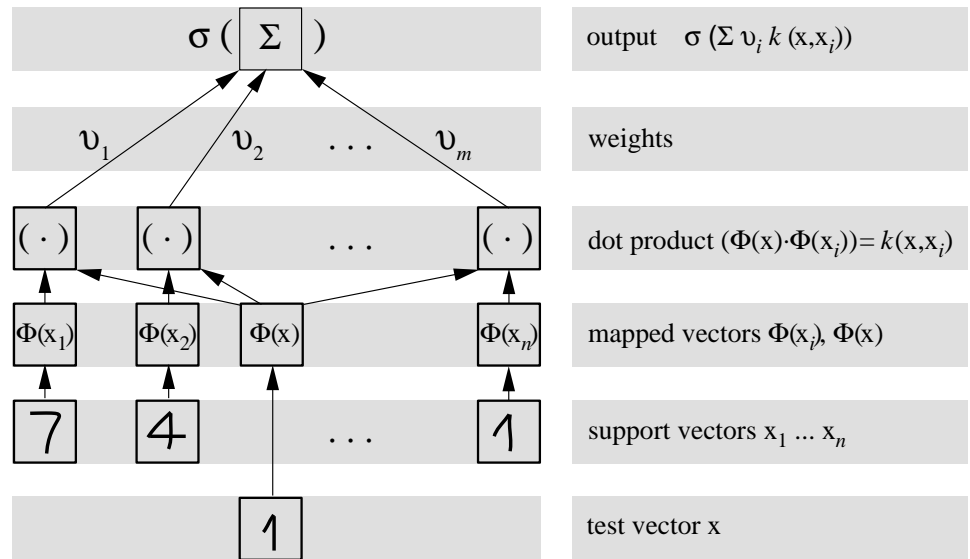


Figure 1.4 Architecture of SV machines. The input x and the Support Vectors x_i are nonlinearly mapped (by Φ) into a feature space \mathcal{F} , where dot products are computed. By the use of the kernel k , these two layers are in practice computed in one single step. The results are linearly combined by weights v_i , found by solving a quadratic program (in pattern recognition, $v_i = y_i \alpha_i$; in regression estimation, $v_i = \alpha_i^* - \alpha_i$). The linear combination is fed into the function σ (in pattern recognition, $\sigma(x) = \text{sgn}(x + b)$; in regression estimation, $\sigma(x) = x + b$).

1.3.3 Smoothness and Regularization

For kernel-based function expansions, one can show (Smola and Schölkopf, 1998b) that given a regularization operator P mapping the functions of the learning machine into some dot product space, the problem of minimizing the regularized risk

regularized risk
$$R_{\text{reg}}(f) := R_{\text{emp}}(f) + \frac{\lambda}{2} \|Pf\|^2 \tag{1.79}$$

(with a regularization parameter $\lambda \geq 0$) can be written as a constrained optimization problem. For particular choices of the loss function, it further reduces to a SV type quadratic programming problem. The latter thus is not specific to SV machines, but is common to a much wider class of approaches. What gets lost in the general case, however, is the fact that the solution can usually be expressed in terms of a small number of SVs (cf. also Girosi (1998), who establishes a connection between SV machines and basis pursuit denoising (Chen et al., 1995)). This specific feature of SV machines is due to the fact that the type of regularization and the

class of functions that the estimate is chosen from are intimately related (Girosi et al., 1993; Smola and Schölkopf, 1998a; Smola et al., 1998): the SV algorithm is equivalent to minimizing the regularized risk $R_{\text{reg}}(f)$ on the set of functions

$$f(\mathbf{x}) = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x}) + b, \quad (1.80)$$

provided that k and P are interrelated by

$$k(\mathbf{x}_i, \mathbf{x}_j) = ((Pk)(\mathbf{x}_i, \cdot) \cdot (Pk)(\mathbf{x}_j, \cdot)). \quad (1.81)$$

To this end, k is chosen as a Green's function of P^*P , for in that case, the right hand side of (1.81) equals $(k(\mathbf{x}_i, \cdot) \cdot (P^*Pk)(\mathbf{x}_j, \cdot)) = (k(\mathbf{x}_i, \cdot) \cdot \delta_{\mathbf{x}_j}(\cdot)) = k(\mathbf{x}_i, \mathbf{x}_j)$. For instance, an RBF kernel corresponds to regularization with a functional containing a specific differential operator.

In SV machines, the kernel thus plays a dual role: firstly, it determines the class of functions (1.80) that the solution is taken from; secondly, via (1.81), the kernel determines the type of regularization that is used. The next question, naturally, is what type of regularization (i.e. kernel) we should use in order to get the best generalization performance. Using bounds on covering numbers of Hilbert spaces (Carl and Stephani, 1990), one can show (Williamson et al., 1998b,a; Schölkopf et al., 1999) that the eigenspectrum of the matrix $k(x_i, x_j)$ is closely connected to the latter and also to the eigenspectrum of the kernel k .

regularization
networks

For arbitrary expansions of f into basis functions, say f_i , the considerations about smoothness of the estimate still hold, provided $\|Pf\|$ is a norm in the space spanned by the basis functions f_i (otherwise one could find functions $f \in \text{span}\{f_i\}$ with $\|Pf\| = 0$, however $f \neq 0$). In this case the existing bounds for kernel expansions can be readily applied to regularization networks as well (cf. e.g. (Williamson et al., 1998b; Smola, 1998) for details). However, one can show (Kimeldorf and Wahba, 1971; Cox and O'Sullivan, 1990), that such an expansion may not fully minimize the regularized risk functional (1.79). This is one of the reasons why often only kernel expansions are considered.

Gaussian
processes

Finally it is worth while pointing out the connection between Gaussian Processes and Support Vector machines. The similarity is most obvious in regression, where the Support Vector solution is the maximum a posteriori estimate of the corresponding Bayesian inference scheme (Williams, 1998). In particular, the kernel k of Support Vector machines plays the role of a covariance function such that the prior probability of a function $f = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x})$ is given by

$$P(f) \propto \exp\left(-\frac{1}{2}\|Pf\|^2\right) = \exp\left(-\frac{1}{2}\sum_{i,j} \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)\right). \quad (1.82)$$

Bayesian methods, however, require averaging over the posterior distribution $P(f|X, Y)$ in order to obtain the final estimate and to derive error bounds. In classification the situation is even more complicated, since we have Bernoulli distributed random variables for the labels of the classifier. See (Williams, 1998) for

more details on this subject.

1.3.4 A Bound on the Leave-One-Out Estimate

Besides the bounds directly involving large margins, which are useful for stating uniform convergence results, one may also try to estimate $R(f)$ by using leave-one-out estimates. Denote by f_i the estimate obtained from $X \setminus \{\mathbf{x}_i\}, Y \setminus \{y_i\}$. Then

$$R_{\text{out}}(f) := \frac{1}{m} \sum_{i=1}^m c(\mathbf{x}_i, y_i, f_i(\mathbf{x}_i)) \quad (1.83)$$

One can show (cf. e.g. (Vapnik, 1979)) that the latter is an unbiased estimator of $R(f)$. Unfortunately, $R_{\text{out}}(f)$ is hard to compute and thus rarely used. In the case of Support Vector classification, however, an upper bound on $R_{\text{out}}(f)$ is not too difficult to obtain. Vapnik (1995) showed that the fraction of Support Vectors is an upper bound on $R_{\text{out}}(f)$. Jaakkola and Haussler (1999) have generalized this result as follows

$$\begin{aligned} R_{\text{out}}(f) &\leq \frac{1}{m} \sum_{i=1}^m 1_{\{y_i \sum_{j \neq i} \alpha_j y_j k(\mathbf{x}_j, \mathbf{x}_i) + y_i b > 0\}} \\ &= \frac{1}{m} \sum_{i=1}^m 1_{\{y_i f(\mathbf{x}_i) - \alpha_i k(\mathbf{x}_i, \mathbf{x}_i) > 0\}}. \end{aligned} \quad (1.84)$$

The latter can be obtained easily without explicitly solving the optimization problem again for the reduced samples. In particular, for kernels with $k(\mathbf{x}, \mathbf{x}) = 1$ like many RBF kernels the condition reduces to testing whether $y_i f(\mathbf{x}_i) - \alpha_i > 0$. The remaining problem is that $R_{\text{out}}(f)$ itself is a random variable and thus it does not immediately give a *bound* on $R(f)$. See also chapters 2 and ?? for further details on how to exploit these bounds in practical cases.

1.4 Boosting

Freund and Schapire (1995) proposed the AdaBoost algorithm for combining classifiers produced by other learning algorithms. AdaBoost has been very successful in practical applications (see Section 1.5). It turns out that it is also a large margin technique.

Table 1.2 gives the pseudocode for the algorithm. It returns a convex combination of classifiers from a class G , by using a learning algorithm L that takes as input a training sample X, Y and a distribution D on X (not to be confused with the true distribution p), and returns a classifier from G . The algorithm L aims to minimize training error on X, Y , weighted according to D . That is, it aims to minimize

$$\sum_{i=1}^m D_i 1_{\{h(\mathbf{x}_i) \neq y_i\}}. \quad (1.85)$$

argument: Training sample, $X = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subset \mathcal{X}$, $Y = \{y_1, \dots, y_m\} \subset \{\pm 1\}$
Number of iterations, T

returns: Convex combination of functions from G , $f = \sum_{t=1}^T \alpha_t g_t$.

function AdaBoost(X, Y, T)
for all i **from** $i = 1, \dots, m$
 $D_1(i) := 1/m$
endfor
for all t **from** $\{1, \dots, T\}$
 $g_t := L(X, Y, D_t)$
 $\varepsilon_t := \sum_{i=1}^m D_t(i) 1_{g_t(x_i) \neq y_i}$
 $\alpha_t := \frac{1}{2} \ln \left(\frac{1 - \varepsilon_t}{\varepsilon_t} \right)$
 $Z_t := 2\sqrt{\varepsilon_t(1 - \varepsilon_t)}$
for all i **from** $i = 1, \dots, m$
 $D_{t+1}(i) := \begin{cases} D_t(i)e^{-\alpha_t}/Z_t & \text{if } y_i = g_t(x_i) \\ D_t(i)e^{\alpha_t}/Z_t & \text{otherwise,} \end{cases}$
endfor
endfor
return $f = \frac{\sum_{t=1}^T \alpha_t g_t}{\sum_{i=1}^T \alpha_t}$.
end

Table 1.2 Pseudocode for the Adaboost algorithm. (L is a learning algorithm that chooses a classifier from G to minimize weighted training error.)

AdaBoost iteratively combines the classifiers returned by L . The idea behind AdaBoost is to start with a uniform weighting over the training sample, and progressively adjust the weights to emphasize the examples that have been frequently misclassified by the classifiers returned by L . These classifiers are combined with convex coefficients that depend on their respective weighted errors. The following theorem shows that Adaboost produces a large margin classifier, provided L is successful at finding classifiers with small weighted training error. See (Schapire et al., 1998). Recall (1.30) that the margin error of a function f with respect to ρ on a sample X, Y is $R_\rho(f) = \frac{1}{m} \sum_{i=1}^m 1_{\{y_i f(\mathbf{x}_i) < \rho\}}$.

Theorem 1.12 Margin Error of AdaBoost

If, at iteration t , L returns a function with weighted training error $\varepsilon_t < 1/2$, then AdaBoost returns a function f that satisfies

$$R_\rho(f) \leq 2^T \prod_{t=1}^T \sqrt{\varepsilon_t^{1-\rho} (1 - \varepsilon_t)^{1+\rho}}. \quad (1.86)$$

In particular, if $\varepsilon_t \leq 1/2 - 2\rho$, then

$$R_\rho(f) < (1 - \rho^2)^{T/2}, \quad (1.87)$$

and this is less than ε for $T \geq (2/\rho^2) \ln(1/\varepsilon)$.

1.5 Empirical Results, Implementations, and Further Developments

Large margin classifiers are not only promising from the theoretical point of view. They also have proven to be competitive or superior to other learning algorithms in practical applications. In the following we will give references to such situations.

1.5.1 Boosting

Experimental results show that boosting is able to improve the performance of classifiers significantly. Extensive studies on the UC Irvine dataset, carried out by Freund and Schapire (1996) and Quinlan (1996) with tree classifiers show the performance of such methods. However, also other learning algorithms can benefit from boosting. Schwenk and Bengio (1998) achieve record performance on an OCR task on the UC Irvine database, using neural networks as the base classifiers. See Rätsch (1998) and chapter ?? for further results on the performance of improved versions of boosted classifiers.

1.5.2 Support Vector Machines

SV Machines perform particularly well in feature rich highdimensional problems. Schölkopf et al. (1995); Schölkopf et al. (1996, 1998b); Burges and Schölkopf (1997); Schölkopf (1997) achieve state of the art, or even record performance in several Optical Character Recognition (OCR) tasks such as the digit databases of the United Postal Service (USPS) and the National Institute of Standards and Technology (NIST). The latter can be obtained at

<http://www.research.att.com/~yann/ocr/mnist/>

Similar results have been obtained for face recognition by Oren et al. (1997); Osuna et al. (1997b) and object recognition (Blanz et al., 1996; Schölkopf, 1997). Finally, also on large noisy problems SV Machines are very competitive as shown in (Smola, 1998).

1.5.3 Implementation and Available Code

Whilst Boosting can be easily implemented by combining a base learner and following the pseudocode of table 1.2. Hence one only has to provide a base learning algorithm satisfying the properties of a weak learner, which defers all problems to

the underlying algorithm.

<http://www.research.att.com/~yoav/adaboost/>

provides a Java applet demonstrating the basic properties of AdaBoost.

The central problem in Support Vector Machines is a quadratic programming problem. Unfortunately, off-the-shelf packages developed in the context of mathematical programming like MINOS (Murtagh and Saunders, 1993), LOQO (Vanderbei, 1994), OSL (IBM Corporation, 1992), or CPLEX (CPL, 1994) are often prohibitively expensive or unsuitable for optimization problems in more than several thousand variables (whilst the number of variables may be in the tens of thousands in practical applications). Furthermore these programs are often optimized to deal with sparse matrix entries, causing unneeded overhead when solving generic SV optimization problems (which are sparse in the solution, not in the matrix entries).

This situation led to the development of several quadratic optimization algorithms specifically designed to suit the needs of SV machines. Starting from simple subset selection algorithms as initially described by Vapnik (1979) and subsequently implemented in e.g. (Schölkopf et al., 1995), more advanced chunking methods were proposed (Osuna et al., 1997a) (see also (Joachims, 1999) for a detailed description of the algorithm) for splitting up the optimization problem into smaller subproblems that could be easily solved by standard optimization code. Other methods exploit constrained gradient descent techniques (Kaufmann, 1999), or minimize very small subproblems, such as the Sequential Minimal Optimization algorithm (SMO) by Platt (1999). See also chapter ?? for further methods for training a SV classifier. Implementations include SvmLight by Joachims (1999),

http://www-ai.cs.uni-dortmund.de/thorsten/svm_light.html

the Royal Holloway / ATT / GMD Support Vector Machine by Saunders et al. (1998), available at

<http://svm.dcs.rhbnc.ac.uk/>

and the implementation by Steve Gunn which can be downloaded from

<http://www.isis.ecs.soton.ac.uk/resources/svminfo/>.

The first two of these optimizers use the GMD (Smola) implementation of an interior point code along the lines of Vanderbei (1994) as the core optimization engine. It is available as a standalone package at

<http://www.svm.first.gmd.de/software.html>.

This site will also contain pointers to further toolboxes as they become available. Java applets for demonstration purposes can be found at

<http://http://svm.dcs.rhbnc.ac.uk/pagesnew/GPat.shtml>

<http://http://svm.research.bell-labs.com/SVT/SVMsvt.html>.

1.6 Notation

We conclude the introduction with a list of symbols which are used throughout the book, unless stated otherwise.

\mathbb{N}	the set of natural numbers
\mathbb{R}	the set of reals
X	a sample of input patterns
Y	a sample of output labels
\mathcal{X}	an abstract domain
\ln	logarithm to base e
\log_2	logarithm to base 2
$(\mathbf{x} \cdot \mathbf{x}')$	inner product between vectors \mathbf{x} and \mathbf{x}'
$\ \cdot\ $	2-norm (Euclidean distance), $\ \mathbf{x}\ := \sqrt{(\mathbf{x} \cdot \mathbf{x})}$
$\ \cdot\ _p$	p -norm, $\ \mathbf{x}\ _p := \left(\sum_{i=1}^N x_i ^p\right)^{1/p}$
$\ \cdot\ _\infty$	∞ -norm, $\ \mathbf{x}\ _\infty := \max_{i=1}^N x_i $
ℓ_p	ℓ_p metric
$L_2(X)$	space of functions on X square integrable wrt. Borel–Lebesgue measure
$\mathbf{E}(\xi)$	expectation of random variable ξ
$\Pr(\cdot)$	probability of an event
N	dimensionality of input space
m	number of training examples
\mathbf{x}_i	input patterns
y_i	target values, or (in pattern recognition) classes
\mathbf{w}	weight vector
b	constant offset (or threshold)
h	VC dimension
f	a real valued function $f : \mathbb{R}^N \rightarrow \mathbb{R}$ (unthresholded)
F	a family of real valued functions f
g	a decision function $g : \mathbb{R}^N \rightarrow \{-1, 1\}$
F	a family of decision functions g
$\rho_f(\mathbf{x}, y)$	margin of function f on the example (\mathbf{x}, y) , i.e. $y f(\mathbf{x})$
ρ_f	minimum margin, i.e. $\min_{1 \leq i \leq m} \rho_f(\mathbf{x}_i, y_i)$

$c(\mathbf{x}, y, f(\mathbf{x}))$	cost function
$R(g)$	risk of g , i.e. expected fraction of errors
$R_{\text{emp}}(g)$	empirical risk of g , i.e. fraction of training errors
$R(f)$	risk of f
$R_{\text{emp}}(f)$	empirical risk of f
k	Mercer kernel
\mathcal{F}	Feature space induced by a kernel
Φ	map into feature space (induced by k)
α_i	Lagrange multiplier
$\boldsymbol{\alpha}$	vector of all Lagrange multipliers
ξ_i	slack variables
$\boldsymbol{\xi}$	vector of all slack variables
C	regularization constant for SV Machines
λ	regularization constant ($C = \frac{1}{\lambda}$)

Jason Weston

Royal Holloway, University of London
Department of Computer Science,
Egham, Surrey, TW20 OEX, UK
jasonw@dcs.rbnc.ac.uk

Ralf Herbrich

Technical University of Berlin
Department of Computer Science,
Franklinstr. 28/29,
10587 Berlin, Germany
ralfh@cs.tu-berlin.de

In this chapter we present a new learning algorithm, Leave-One-Out (LOO-) SVMs and its generalization Adaptive Margin (AM-) SVMs, inspired by a recent upper bound on the leave-one-out error proved for kernel classifiers by Jaakkola and Haussler. The new approach minimizes the expression given by the bound in an attempt to minimize the leave-one-out error. This gives a convex optimization problem which constructs a sparse linear classifier in *feature space* using the kernel technique. As such the algorithm possesses many of the same properties as SVMs and Linear Programming (LP-) SVMs. These former techniques are based on the minimization of a regularized margin loss, where the margin is treated *equivalently* for each training pattern. We propose a minimization problem such that *adaptive margins* for each training pattern are utilized. Furthermore, we give bounds on the generalization error of the approach which justifies its robustness against outliers. We show experimentally that the generalization error of AM-SVMs is comparable to SVMs and LP-SVMs on benchmark datasets from the UCI repository.

2.1 Introduction

The study of classification learning has shown that algorithms which learn a *real-valued* function for classification can control their generalization error by making use of a quantity known as the *margin* (see Section 1.1.3). Based on these results, learning machines which *directly* control the margin (e.g. SVMs, LP-SVMs) have been proven to be successful in classification learning (Mason and Bartlett, 1998; Vapnik, 1998; Smola, 1998). Moreover, it turned out to be favourable to formulate the decision functions in terms of a symmetric, positive semidefinite, and square integrable function $k(\cdot, \cdot)$ referred to as a *kernel* (see Section 1.3.2). The class of decision functions — also known as *kernel classifiers* (Smola, 1998; Jaakkola and Haussler, 1999) — is then given by¹

$$f(\mathbf{x}) = \sum_{i=1}^m \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) \quad \boldsymbol{\alpha} \geq \mathbf{0}. \quad (2.1)$$

For simplicity we ignore classifiers which use an extra threshold term (cf. Equation (1.74)).

Whilst the algorithms proposed so far are restricted to a *fixed margin* (the same constant value) at each training pattern (\mathbf{x}_i, y_i) , we show that *adaptive* margins can successfully be used. Moreover, it turns out that adaptive margins effectively control the complexity of the model. The chapter is structured as follows: In Section 2.2 we describe the LOO-SVM algorithm. The generalization of LOO-SVMs to control the margin adaptively, which gives AM-SVMs, is then presented in Section 2.3 and their relation to SVMs and LP-SVMs is revealed in Section 2.4. In Section 2.5 we give bounds on the generalization error of AM-SVMs which justify the use of adaptive margins as a regularizer. In Section 2.6 results of a comparison of AM-SVMs with SVMs on artificial and benchmark datasets from the UCI repository² are presented. Finally, in Section 2.7 we summarize the chapter and discuss further directions.

2.2 Leave-One-Out Support Vector Machines

Support Vector Machines obtain sparse solutions that yield a direct assessment of generalization: the leave-one-out error is bounded by the expected ratio of the number of non-zero coefficients α_i to the number m of training examples (Vapnik, 1995). In Jaakkola and Haussler (1999) a bound on this error is derived for a class of classifiers which includes SVMs but can be applied to non-sparse solutions. In

1. Although this class of functions is dependent on the training set, the restrictions put on $k(\cdot, \cdot)$ automatically ensure that the influence of each *new* basis function $k(\mathbf{x}_i, \cdot)$ decreases rapidly for increasing training set sizes m . Thus we can assume the existence of a *fixed* feature space (see e.g. Graepel et al. (1999)).

2. <http://www.ics.uci.edu/mllearn/MLRepository.html>.

leave-one-out
bound

order to motivate our reasoning we restate their result which is given by (1.84) in a more concise form.

Theorem 1

For any training set of examples $\mathbf{x}_i \in R^N$ and labels $y_i \in \{\pm 1\}$, for an SVM the leave-one-out error estimate of the classifier is bounded by

$$\frac{1}{m} \sum_{i=1}^m \theta \left(-y_i \sum_{j \neq i} y_j \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \right). \quad (2.2)$$

where $\theta(\cdot)$ is the step function.

This bound is slightly tighter than the classical SVM leave-one-out bound. This is easy to see when one considers that all training points that have $\alpha_i = 0$ cannot be leave-one-out errors in either bound. Vapnik's bound assumes all support vectors (all training points with $\alpha_i > 0$) are errors, whereas they only contribute as errors in Equation (2.2) if

$$y_i \sum_{j \neq i} \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) \leq 0. \quad (2.3)$$

In practice this means the bound is tighter for less sparse solutions.

Theorem 1 motivates the following algorithm (Weston, 1999): directly minimize the expression in the bound. In order to achieve this, one introduces slack variables following the standard approach in Cortes and Vapnik (1995) to give the following optimization problem:

$$\text{minimize} \quad \sum_{i=1}^m \xi_i^\delta \quad (2.4)$$

$$\text{subject to} \quad y_i \sum_{j \neq i} \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 1 - \xi_i, \quad \text{for all } i = 1, \dots, m \quad (2.5)$$

$$\boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\xi} \geq \mathbf{0}. \quad (2.6)$$

where one chooses a fixed constant for the margin to ensure non-zero solutions.

To make the optimization problem tractable, the smallest value for δ for which we obtain a convex objective function is $\delta = 1$. Noting also that $y_i \sum_{j \neq i} \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) = y_i f(\mathbf{x}_i) - \alpha_i k(\mathbf{x}_i, \mathbf{x}_i)$ we obtain the equivalent linear program:

Leave-one-out
SVM

$$\text{minimize} \quad \sum_{i=1}^m \xi_i \quad (2.7)$$

$$\text{subject to} \quad y_i f(\mathbf{x}_i) \geq 1 - \xi_i + \alpha_i k(\mathbf{x}_i, \mathbf{x}_i), \quad \text{for all } i = 1, \dots, m \quad (2.8)$$

$$\boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\xi} \geq \mathbf{0}. \quad (2.9)$$

As in other kernel classifiers, one uses the decision rule given in Equation (2.1). Note that Theorem 1 is no longer valid for this learning algorithm. Nevertheless, let us study the resulting method which we call a Leave-One-Out Support Vector Machine (LOO-SVM).

regularization

Firstly, the technique appears to have no free regularization parameter ³. This

should be compared with Support Vector Machines which control the amount of regularization with the free parameter C (see Section 1.3). For SVMs, in the case of $C = \infty$ one obtains a *hard margin* classifier with no training errors. In the case of noisy or linear inseparable datasets⁴ (through noise, outliers, or class overlap) one must admit some training errors (by constructing a so called *soft margin* – see Section 1.1.4). To find the best choice of training error/margin tradeoff one has to choose the appropriate value of C . In LOO-SVMs a soft margin is automatically constructed. This happens because the algorithm does not attempt to minimize the number of training errors – it minimizes the number of training points that are classified incorrectly even when they are removed from the linear combination that forms the decision rule. However, if one can classify a training point correctly when it is removed from the linear combination then it will always be classified correctly when it is placed back into the rule. This can be seen as $\alpha_i y_i k(\mathbf{x}_i, \mathbf{x}_i)$ has always the same sign as y_i , any training point is pushed further from the decision boundary by its own component of the linear combination. Note also that summing for all $j \neq i$ in the constraint 2.5 is equivalent to setting the diagonal of the kernel matrix to zero and instead summing for all j . Thus the regularization employed by LOO-SVMs disregards the values $k(\mathbf{x}_i, \mathbf{x}_i) = 0$ for all i .

sparsity

Secondly, like Support Vector machines, the solutions can be sparse; that is, only some of the coefficients α_i , $i = 1, \dots, m$ are non-zero (see Section 2.6.2 for computer simulations confirming this). As the coefficient of a training point does not contribute to its leave-one-out error in constraint (2.5) the algorithm does not assign a non-zero value to the coefficient of a training point in order to correctly classify it. A training point has to be classified correctly by the training points of the same label that are close to it (in *feature space*), but the training point itself makes no contribution to its own classification.

In the next Section we show how this method does in fact have an implicit regularization parameter and generalize the method to control the regularization on the set of decision functions.

2.3 Adaptive Margin SVMs

In the setting of the optimization problem (2.7)–(2.9) it is easy to see that a training point \mathbf{x}_i is linearly penalized for failing to obtain a margin of $\rho_f(\mathbf{x}_i, y_i) \geq 1 + \alpha_i k(\mathbf{x}_i, \mathbf{x}_i)$. That is, the larger the contribution the training point has to the decision rule (the larger the value of α_i), the larger its margin must be. Thus, the algorithm controls the margin for each training point *adaptively*. From this

3. As we shall see later there is an implicit regularization parameter, but it is fixed. The generalization of this problem which allows one to control this parameter gives Adaptive Margin SVMs.

4. Here we refer to linearly inseparability in *feature space*. Both SVMs and LOO-SVM Machines are essentially linear classifiers.

Adaptive Margin SVM

formulation one can generalize the algorithm to control regularization through the margin loss. To make the margin at each training point a controlling variable we propose the following learning algorithm:

$$\text{minimize} \quad \sum_{i=1}^m \xi_i \tag{2.10}$$

$$\text{subject to} \quad y_i f(\mathbf{x}_i) \geq 1 - \xi_i + \lambda \alpha_i k(\mathbf{x}_i, \mathbf{x}_i), \quad \text{for all } i = 1, \dots, m. \tag{2.11}$$

$$\boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\xi} \geq \mathbf{0}. \tag{2.12}$$

This algorithm can then be viewed in the following way (see Figure 2.1): Suppose the data lives on the surface of a hypersphere in \mathcal{F} , i.e. $k(\cdot, \cdot)$ is an RBF kernel given by Equation (1.73). Then $k(\mathbf{x}_i, \mathbf{x}_j)$ is the cosine of the angle between $\Phi(\mathbf{x}_i)$ and $\Phi(\mathbf{x}_j)$. As soon as a point $\Phi(\mathbf{x}_k)$ is an outlier (the cosine of the angles to points in its class are small and to points in the other class are large) α_k in Equation (2.11) has to be large in order to classify $\Phi(\mathbf{x}_k)$ correctly. Whilst SVMs and LP-SVMs use the same margin for such an outlier, they attempt to classify $\Phi(\mathbf{x}_k)$ correctly. In AM-SVMs the margin is automatically increased to $1 + \lambda \alpha_k k(\mathbf{x}_k, \mathbf{x}_k)$ for $\Phi(\mathbf{x}_k)$ and thus less attempt is made to change the decision function.

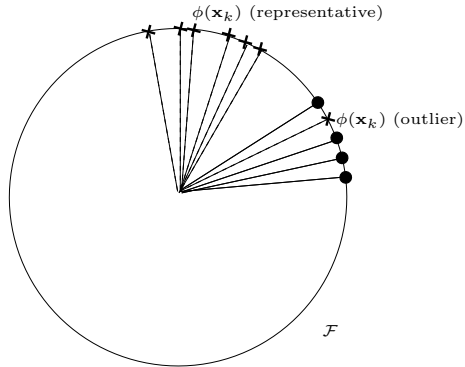


Figure 2.1 Adaptation of margins at each training pattern depending on the distance $k(\mathbf{x}_i, \mathbf{x}_j)$ in feature space \mathcal{F} . Note that $k(\mathbf{x}_i, \mathbf{x}_j)$ is large if the enclosed angle between data points is small. See the text for explanation.

Cluster centres

Moreover, in AM-SVMs the points $\Phi(\mathbf{x}_k)$ which are representatives of clusters (centres) in feature space \mathcal{F} , i.e. those which have large values of the cosine of the angles to points from its class, will have non-zero α_k . In order to see this we consider two points k and k' of the same class. Let us assume that k having $\xi_k > 0$ is the centre of a cluster (in the metric induced by Φ) and k' (having $\xi_{k'} > 0$) lies

at the boundary of the cluster. Hence we subdivide the set of all points into

$$\begin{aligned} i \in C^+ & & \xi_i = 0, y_i = y_k, i \neq k, i \neq k' \\ i \in C^- & & \xi_i = 0, y_i \neq y_k \\ i \in I^+ & & \xi_i > 0, y_i = y_k, i \neq k, i \neq k' \\ i \in I^- & & \xi_i > 0, y_i \neq y_k \end{aligned}$$

We consider the change in ξ if we increase α_k by $\Delta > 0$ (giving ξ') and simultaneously decrease $\alpha_{k'}$ by Δ (giving ξ''). From Equation (2.10)-(2.12) we know that

$$\begin{aligned} i \in C^+ & & \xi'_i &= \xi_i & & \xi''_i &\leq \Delta k(\mathbf{x}_i, \mathbf{x}_{k'}) \\ i \in C^- & & \xi'_i &\leq \Delta k(\mathbf{x}_i, \mathbf{x}_k) & & \xi''_i &= \xi_i \\ i \in I^+ & & \xi'_i &\geq \xi_i - \Delta k(\mathbf{x}_i, \mathbf{x}_k) & & \xi''_i &= \xi_i + \Delta k(\mathbf{x}_i, \mathbf{x}_{k'}) \\ i \in I^- & & \xi'_i &= \xi_i + \Delta k(\mathbf{x}_i, \mathbf{x}_k) & & \xi''_i &\geq \xi_i - \Delta k(\mathbf{x}_i, \mathbf{x}_{k'}) \\ i = k & & \xi'_k &\geq \xi_k - \Delta(1 - \lambda)k(\mathbf{x}_k, \mathbf{x}_k) & & \xi''_k &= \xi_k + \Delta k(\mathbf{x}_k, \mathbf{x}_{k'}) \\ i = k' & & \xi'_{k'} &\geq \xi_{k'} - \Delta k(\mathbf{x}_{k'}, \mathbf{x}_k) & & \xi''_{k'} &\geq \xi_{k'} + (1 - \lambda)\Delta k(\mathbf{x}_{k'}, \mathbf{x}_{k'}) \end{aligned}$$

Now we choose the biggest Δ such that all inequalities for $i \in \{I^+, I^-, k, k'\}$ become equalities and the r.h.s for $i \in \{C^+, C^-\}$ equals zero⁵. Hence, the relative change in the objective function is given by

$$\frac{1}{\Delta} \sum_{i=1}^{\ell} (\xi'_i - \xi''_i) = \underbrace{\sum_{i \in I^+} (k(\mathbf{x}_i, \mathbf{x}_{k'}) - k(\mathbf{x}_i, \mathbf{x}_k))}_{\text{change of intra-class distance}} - \underbrace{\sum_{i \in I^-} (k(\mathbf{x}_i, \mathbf{x}_{k'}) - k(\mathbf{x}_i, \mathbf{x}_k))}_{\text{change of inter-class distance}},$$

where we assumed that $k(\mathbf{x}_k, \mathbf{x}_k) = k(\mathbf{x}_{k'}, \mathbf{x}_{k'})$ and $k(\mathbf{x}_k, \mathbf{x}_{k'}) = k(\mathbf{x}_{k'}, \mathbf{x}_k)$. Since the cluster centres in feature space \mathcal{F} minimize the intra-class distance whilst maximizing the inter-class distances it becomes apparent that their α_k will be higher. Taking into account that the maximal Δ to be considerable for this analysis is decreasing as λ increases we see that for suitable small λ AM-SVMs tends to give non-zero α 's only to cluster centres in feature space \mathcal{F} (see also Section 2.6 and Figure 2.4).

It is worthwhile to study the influence of λ :

- If $\lambda = 0$ no adaptation of the margins is performed. This is equivalent to minimizing training error with no regularization, i.e. approximating the expected risk $R(f)$ (1.26) with the empirical risk (1.27) (see Section 1.2).
- If $\lambda \rightarrow \infty$ the margin at each point tends to infinity ($1 + \lambda \alpha_i k(\mathbf{x}_i, \mathbf{x}_i)$) and the solution is thus to set all α 's to an equal and small value. This corresponds to paying *no* attention to $R_{\text{emp}}(f)$ and is equivalent to density estimation on each

5. As for all kernels $k(\mathbf{x}_i, \mathbf{x}_i) \geq k(\mathbf{x}_i, \mathbf{x}_j)$ the inequalities $\xi'_k \geq \xi_k - \Delta(1 - \lambda)k(\mathbf{x}_k, \mathbf{x}_k)$ and $\xi''_{k'} \geq \xi_{k'} + (1 - \lambda)\Delta k(\mathbf{x}_{k'}, \mathbf{x}_{k'})$ basically determine the maximal Δ to be considered.

class (Parzen windows) (Parzen, 1962).

- If $\lambda = 1$ the resulting algorithm is equivalent to LOO-SVMs.

2.4 Relationship of AM-SVMs to other SVMs

Using the soft margin loss

$$c(\mathbf{x}, y, f(\mathbf{x})) = \max(1 - yf(\mathbf{x}), 0) \quad (2.13)$$

one can derive SVMs and LP-SVMs by choosing different regularizers. If we use the quadratic regularization functional

$$Q_{\text{QP}}(f) = \|\mathbf{w}\|_2^2, \quad (2.14)$$

SVMs

we directly obtain the well known class of SVMs (see Section 1.3), i.e.

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^m \xi_i + \lambda \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ \text{subject to} \quad & y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \quad \text{for all } i = 1, \dots, m \\ & \boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\xi} \geq \mathbf{0}. \end{aligned} \quad (2.15)$$

Here we used

$$\mathbf{w} = \sum_{j=1}^m \alpha_j y_j \Phi(\mathbf{x}), \quad (2.16)$$

LP-SVMs

where $\Phi(\cdot)$ maps into a feature space \mathcal{F} such that $(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')) = k(\mathbf{x}, \mathbf{x}')$. It is known that $Q_{\text{QP}}(f)$ controls the covering number $\mathcal{N}(\cdot, F)$ of the induced loss-function class (Theorem 1.2) (Shawe-Taylor et al., 1998; Smola, 1998). This choice of regularizer favours flat functions in feature space.

Similarly using a linear regularization functional

$$Q_{\text{LP}}(f) = \sum \alpha_i \quad (2.17)$$

we obtain LP-SVMs. The corresponding minimization problem is given by⁶

$$\begin{aligned} \text{minimize} \quad & \sum_{i=1}^m \xi_i + \lambda \sum_{i=1}^m \alpha_i \\ \text{subject to} \quad & y_i f(\mathbf{x}_i) \geq 1 - \xi_i, \quad \text{for all } i = 1, \dots, m \\ & \boldsymbol{\alpha} \geq \mathbf{0}, \boldsymbol{\xi} \geq \mathbf{0}. \end{aligned} \quad (2.18)$$

Recently it was shown that also $Q_{\text{LP}}(f)$ can also be used to control the covering number of $c(\cdot, \cdot, f(\cdot))$ (Smola, 1998). In contrast to the quadratic regularizer, $Q_{\text{LP}}(f)$ favours non-smooth functions by strongly penalizing basis functions $\Phi_j(\cdot)$ with a

6. Note, that we require $\boldsymbol{\alpha} \geq \mathbf{0}$ which allows us to omit the absolute values on the α_i 's.

small eigenvalue (Smola, 1998).

Comparing these algorithms to AV-SVMs, one can see all three algorithms produce a sparse kernel classifier. It is easy to see that for $\lambda = 0$ and $\lambda \rightarrow \infty$ all three algorithms revert to the same learnt function. It is only how λ stratifies the set of decision functions to form the type regularization that differentiates the three algorithms.

2.5 Theoretical Analysis

To obtain margin distribution bounds for Adaptive Margin Machines we apply the following theorem to be found in Shawe-Taylor and Cristianini (1998):

Theorem 2.1

Consider a fixed but unknown probability distribution on the input space \mathcal{X} with support in the ball of radius R about the origin. Then with probability $1 - \delta$ over randomly drawn training sets (X, Y) of size m for all $\rho > 0$ such that $d((\mathbf{x}, y), \mathbf{w}, \rho) = 0$, for some $(\mathbf{x}, y) \in (X, Y)$, the generalization of a linear classifier \mathbf{w} on \mathcal{X} satisfying $\|\mathbf{w}\|_{\mathcal{X}} \leq 1$ is bounded from above by

$$\epsilon = \frac{2}{m} \left(\kappa \log_2 \left(\frac{8em}{\kappa} \right) \log_2(32m) + \log_2 \left(\frac{2m(28 + \log_2(m))}{\delta} \right) \right), \quad (2.19)$$

where

$$\kappa = \left\lfloor \frac{65[(R + D)^2 + 2.25RD]}{\rho^2} \right\rfloor, \quad (2.20)$$

$$D = D(S, \mathbf{w}, \rho) = \sqrt{\sum_{i=1}^m d_i^2}$$

$$d_i = d((\mathbf{x}_i, y), \mathbf{w}, \rho) = \max\{0, \rho - y(\mathbf{w} \cdot \mathbf{x}_i)\}$$

and provided $m \geq \max\{2/\epsilon, 6\}$ and $\kappa \leq em$.

Applying the bound to AM-SVMs we can give the following theorem.

Theorem 2.2

Consider a fixed but unknown probability distribution on the feature space \mathcal{F} with support in the ball of radius R about the origin. Then with probability $1 - \delta$ over randomly drawn training sets (X, Y) of size m for $\boldsymbol{\alpha} \geq \mathbf{0}$ and $\boldsymbol{\xi} \geq \mathbf{0}$ which are feasible solutions of AM-SVMs such that $d((\mathbf{x}, y), \mathbf{w}, 1) = 0$ for some $(\mathbf{x}, y) \in (X, Y)$, the generalization error $R(f)$ is bounded by

$$\epsilon = \frac{2}{m} \left(\kappa \log_2 \left(\frac{8em}{\kappa} \right) \log_2(32m) + \log_2 \left(\frac{2m(28 + \log_2(m))}{\delta} \right) \right), \quad (2.21)$$

where

$$\begin{aligned}\kappa &\leq \lceil 65[(WR + 3D)^2] \rceil, \\ D &= \sqrt{\sum_{i=1}^m [\max\{0, \xi_i - \lambda\alpha_i k(\mathbf{x}_i, \mathbf{x}_i)\}]^2}, \\ W^2 &= \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j),\end{aligned}$$

provided $m \geq \max\{2/\epsilon, 6\}$ and $\kappa \leq em$.

Proof Firstly, AM-SVMs are linear classifiers $f(\mathbf{x}) = (\mathbf{w} \cdot \Phi(\mathbf{x}))$ where \mathbf{w} is defined by Equation (2.16). We wish to redefine the measure of margin error $d((\mathbf{x}, y), \mathbf{w}, \rho) = \rho - y_i f(\mathbf{x}_i)$ in Theorem 2.1 in terms of ξ_i and $\lambda\alpha_i k(\mathbf{x}_i, \mathbf{x}_i)$ to capture the adaptive margin of a training point \mathbf{x}_i . Then we know from the assumption of a feasible solution $\boldsymbol{\alpha}, \boldsymbol{\xi}$ that

$$\max\{0, \rho - y_i f(\mathbf{x}_i)\} \leq \max\{0, \rho - 1 + \xi_i - \lambda\alpha_i k(\mathbf{x}_i, \mathbf{x}_i)\}. \quad (2.22)$$

In order to apply Theorem 2.1 for *any* vector \mathbf{w} we have to normalize ρ , D , and $\boldsymbol{\alpha}$ by the norm of $\|\mathbf{w}\|_{\mathcal{F}} = W$ given by (2.16). This results in

$$\kappa = \left\lceil \frac{65[(R + \frac{1}{W}D)^2 + 2.25\frac{1}{W}RD]}{\rho^2} W^2 \right\rceil. \quad (2.23)$$

Now we fix $\rho = 1$ as done by AM-SVMs. This gives for Equation (2.22)

$$\max\{0, \rho - y_i f(\mathbf{x}_i)\} \leq \max\{0, \xi_i - \lambda\alpha_i k(\mathbf{x}_i, \mathbf{x}_i)\}. \quad (2.24)$$

Making use of

$$\left[\left(R + \frac{1}{W}D \right)^2 + 2.25\frac{1}{W}RD \right] W^2 \leq [(WR + 3D)^2], \quad (2.25)$$

the theorem is proven. ■

From the theorem, one can gain the following insights. Our goal to minimize the generalization error is achieved by minimizing κ , the minimum of which is a tradeoff between minimizing W (the margin) and D (the loss with adaptive margin). We require a small value of both but small values of one term automatically gives a large value of the other. By minimizing $\sum_{i=1}^m \xi_i$ AM-SVMs effectively control the tradeoff between the two terms through the parameter λ . For small values of λ , the resulting D is small and W can take any value as it is not minimized (it can be forced to very large values). For large λ the increased margin in D acts a regularizer, penalizing large values of α . This results in small values of W (a smooth function) but large values of D (large training error). This bound motivates the objective function of AM-SVMs which at first appears to only minimize error and have no regularization. In fact, as we have seen, the regularization comes from the adaptive margin in the constraints controlled by λ .

2.6 Experiments

2.6.1 Artificial Data

2.6.1.1 LOO-SVMs

We first describe some two dimensional examples to illustrate how the new technique works. Let us first consider AM-SVMs with regularization parameter $\lambda = 1$ (this corresponds to LOO-SVMs, see Section 2.2). Figures 2.2 and 2.3 show two artificially constructed training problems with various solutions. We fixed $k(\cdot, \cdot)$ to be a radial basis function (RBF) kernel

$$k(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / (2\sigma^2)), \quad (2.26)$$

and then found the solution to the problems with LOO-SVM, which has no other free parameters, and with SVMs, for which one controls the *soft margin* with the free parameter $C = \frac{1}{\lambda}$. The first solution (left) for both training problems is the LOO-SVM solution and the other two solutions for each problem are SVMs with different choices of *soft margin* using parameter $C = 1$ (middle) and $C = 100$ (right).

In the first problem (Figure 2.2) the two classes (represented by crosses and dots) are almost linearly separable apart from a single outlier. The automatic *soft margin* control of LOO-SVM constructs a classifier which incorrectly classifies the far right dot, assuming that it is an outlier. The Support Vector solutions both classify the outlier correctly resulting in non-smooth decision rules. In the second problem (Figure 2.3) the two classes occupy opposite sides (horizontally) of the picture, but slightly overlap. In this case the data is only separable with a highly nonlinear decision rule, as reflected in the solution by an SVM with parameter $C = 100$ (right). Both problems highlight the difficulty of choosing the parameter C in SVMs, whereas LOO-SVM (AM-SVM with $\lambda = 1$) appears to produce robust⁷, natural decision rules.

2.6.1.2 AM-SVMs

In order to demonstrate how the regularization parameter λ in AM-SVMs (rather than being fixed to $\lambda = 1$ as in LOO-SVMs) affects the generated decision rule we give a comparison on the same toy problem as SVMs and LP-SVMs. We generated another two class problem in \mathbb{R}^2 (represented by crosses and dots) and trained an AM-SVM using RBF-kernels ($\sigma = 0.5$) with $\lambda = 1, 2, 5, 10$ (see Figure 2.4). As can be seen increasing λ allows AM-SVM to widen the margin for points far away

7. As there is no unique definition of robustness (see e.g. (Huber, 1981)) we call a classification learning algorithm *robust* if a few pattern far apart from the remaining ones (in the metric induced by Φ) have no influence on the resulting decision function.

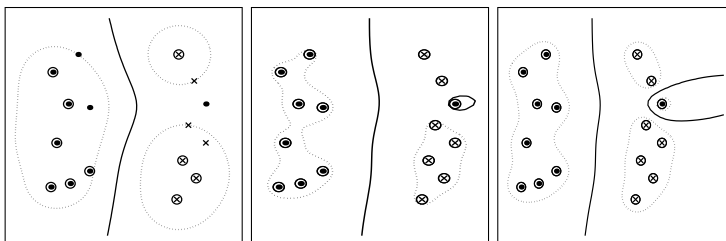


Figure 2.2 A simple two dimensional problem with one outlier solved by LOO-SVMs (left) and SVMs with $C = 1$ (middle) and $C = 100$ (right). LOO-SVMs soft margin regularization appears to perform better than the choices of parameter for SVMs.

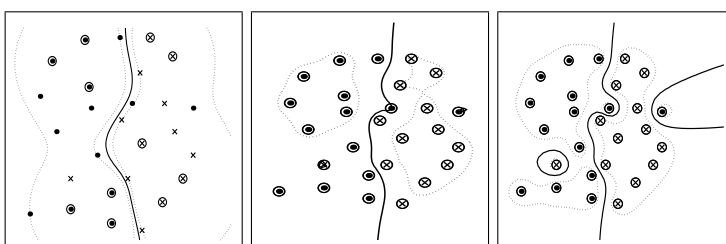


Figure 2.3 A simple two dimensional problem of two overlapping classes solved by LOO-SVMs (left) and SVMs with $C = 1$ (middle) and $C = 100$ (right). LOO-SVMs soft margin regularization appears to perform better than the choices of parameter for SVMs.

from the decision surface. Consequently, the algorithm is more robust to outliers which results in very smooth decision functions. In Figure 2.5 we used the same dataset and trained ν LP-SVMs (Graepel et al., 1999). ν LP-SVMs are obtained by reparameterizing Equation (2.18) where ν upper-bounds the number of margin errors. Varying $\nu = 0.0, 0.1, 0.2, 0.5$ shows that *margin* errors are sacrificed in order to lower the complexity of the decision function f measured in the one-norm (see Equation (2.17) where λ can be replaced by a fixed function of ν). As already mentioned this leads to non-smooth functions. Furthermore it should be noted that the outlier (dot) on the far left side leads to very rugged decision functions. Similar conclusions can be drawn for ν SVMs (Schölkopf et al., 1998d) (see Figure 2.6) though the decision functions are smoother. Thus, AM-SVMs turn out to provide robust solutions (through control of the regularization parameter) which provide a new approach when compared to the solutions of SVMs and LP-SVMs. In these toy examples AM-SVMs appear to provide decision functions which are less influenced by single points (outliers).

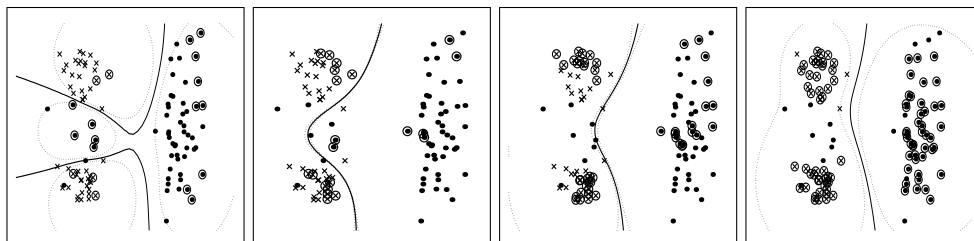


Figure 2.4 Decision functions (solid lines) obtained by AM-SVMs with different choices of the regularization parameter λ . The dashed line represents the minimal margin over all training points. (a) $\lambda = 1$ is equivalent to LOO-SVMs (b) $\lambda = 2$, (c) $\lambda = 5$, and (d) $\lambda = 10$ widens the amount to which margin errors at each point are accepted and thus results in very flat functions. Note, that less attention is paid to the outlier (dot) at the left hand side.

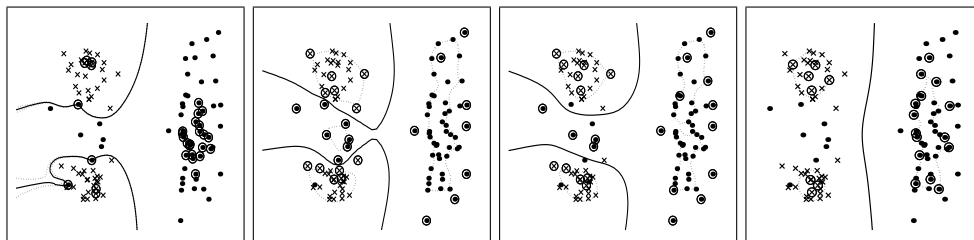


Figure 2.5 Decision functions (solid lines) obtained by ν LP-SVMs with different choices of the assumed noise level ν . The dashed line represents the margin. (a) $\nu = 0.0$ leads to very non-smooth and overfitted decision functions. (b) $\nu = 0.1$, (c) $\nu = 0.2$, and (d) $\nu = 0.5$ smooth the decision function.

2.6.2 Benchmark Datasets

We conducted computer simulations using 6 artificial and real world datasets from the UCI benchmark repositories, following the same experimental setup as in Rätsch et al. (1998). The authors of this article also provide a website to obtain the data⁸. Briefly, the setup is as follows: the performance of a classifier is measured by its average error over one hundred partitions of the datasets into training and testing sets. Free parameter(s) in the learning algorithm are chosen as the median value of the best model chosen by cross validation of the first five training datasets.

Table 2.1 compares percentage test error of LOO-SVMs to AdaBoost (AB), Regularized AdaBoost (AB_R) and SVMs which are all known to be excellent

8. <http://svm.first.gmd.de/~raetsch/data/benchmarks.htm>. The datasets have been pre-processed to have zero mean and standard deviation one, and the exact one hundred splits of training and testing sets used in the author's experiments can be obtained.

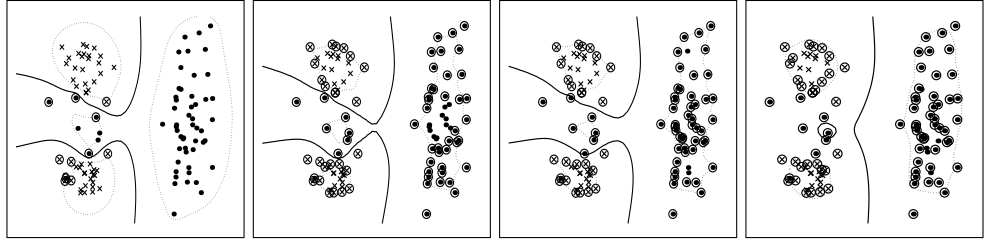


Figure 2.6 Decision functions (solid lines) obtained by ν SVMs with different choices of the assumed noise level ν . The dashed line represents the margin. (a) $\nu = 0.0$ leads to an overfitted decision functions (note the captured outlier in the lower left region). (b) $\nu = 0.1$, (c) $\nu = 0.2$, and (d) $\nu = 0.5$ allow for much flatter functions though regularizing differently to AM-SVMs.

	AB	AB _R	SVM	LOO-SVM
Banana	12.3	10.9	11.5	10.6
B. Cancer	30.4	26.5	26.0	26.3
Diabetes	26.5	23.9	23.5	23.4
Heart	20.3	16.6	16.0	16.1
Thyroid	4.4	4.4	4.8	5.0
Titanic	22.6	22.6	22.4	22.7

Table 2.1 Comparison of percentage test error of AdaBoost (AB), Regularized AdaBoost (AB_R), Support Vector Machines (SVMs) and Leave-One-Out Machines (LOO-SVMs) on 6 datasets.

classifiers⁹. The competitiveness of LOO-SVMs to SVMs and AB_R (which both have a *soft margin* control parameter) is remarkable considering LOO-SVMs have no free parameter. This indicates that the *soft margin* automatically selected by LOO-SVMs is close to optimal. AdaBoost loses out to the three other algorithms, being essentially an algorithm designed to deal with noise-free data.

To give more insight into the behaviour of the algorithm we give two plots in Figure 2.7. The left graph shows the fraction of training points that have non-zero coefficients (*SVs*) plotted against $\log(\sigma)$ (*RBF width*) on the thyroid dataset. Here, one can see the sparsity of the decision rule, the sparseness of which depends on the chosen value of σ . The right graph shows the percentage training and test error (*train err* and *test err*), the value of $\sum_{i=1}^m \xi_i$ (*slacks*) and the value of the bound given in Theorem 1 (*l-o-o bound*). One can see the training and test error (and the bound) closely match. The minimum of all four plots is roughly at $\log(\sigma) = -1$, indicating one could perform model selection using one of the known expressions.

9. The results for AB, AB_R and SVMs were taken from (Rätsch et al., 1998)

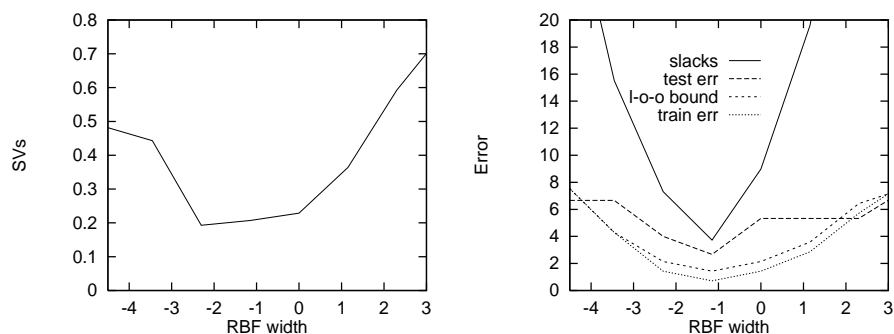


Figure 2.7 The fraction of training patterns that are Support Vectors (top) and various error rates (bottom) both plotted against RBF kernel width for Leave-One-Out Machines on the thyroid dataset.

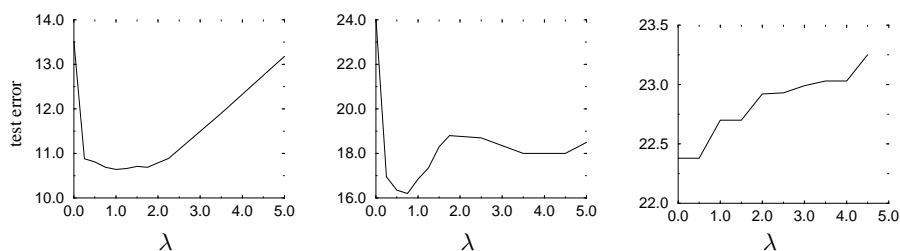


Figure 2.8 Test error plotted against the regularization parameter λ in AM-SVMs. The three plots from left to right are (a) the banana dataset, (b) heart dataset and (c) titanic dataset. Note how $\lambda = 1$ is close to the optimum of the bowl in the first two plots, but in the third plot the plot is not a bowl at all – the best choice of regularization is to choose no regularization ($\lambda = 0$).

Note also that for a reasonable range of σ the test error is roughly the same, indicating the *soft margin* control overcomes overfitting problems.

Finally, we conducted experiments to assess the effect in generalization performance by controlling the regularization parameter λ in AM-SVMs. Figure 2.8 plots λ against test error for three of the datasets averaged over 10 runs for the first two, and over all 100 runs for the last. The banana dataset (left) and the heart dataset (middle) gave bowl-shaped graphs with the minimum exactly (banana) or almost (heart) at $\lambda = 1$. The optimum choice of λ for the titanic dataset, on the hand, is at $\lambda = 0$. In this case the best choice of the regularization parameter λ is to have no regularization at all – the training points give enough information about the unknown decision function. Note this error rate for $\lambda = 0$ is as good as the best SVM solution (see Table 2.1). The first two plots and the results in Table 2.1 justify the choice of $\lambda = 1$ in LOO-SVMs. The last plot in Figure 2.8 justifies AM-SVMs.

2.7 Discussion

In this chapter we presented a new learning algorithm for kernel classifiers. Motivated by minimizing a bound on leave-one-out error we obtained LOO-SVMs and generalizing this approach to control regularization through the margin loss we obtained AM-SVMs. This approach introduced a novel method of capacity control via margin maximization by allowing adaptive rather than fixed margins at each training pattern. We have shown experimentally that this reformulation results in an algorithm which is robust against outliers. Nevertheless, our algorithm has a parameter λ which needs to be optimized for a given learning problem. Further investigations will be made in the derivation of bounds on the leave-one-out error of this algorithm which allows for efficient model order selection. Finally, we note that penalization of the diagonal of the kernel matrix is a well known technique in regression estimation known as Ridge Regression (Hoerl and Kennard, 1970).

Acknowledgements

The authors would like to thank Alex Gammerman, Thore Graepel, Tom Melliush, and Craig Saunders for discussions. In particular, we are indebted to both John Shawe-Taylor and Vladimir Vapnik for their help with this work. Ralf Herbrich would like to thank the Department of Computer Science at Royal Holloway for the warm hospitality during his research stay. Jason Weston thanks the ESPRC for providing financial support through grant GR/L35812.

References

- M. Aizerman, E. Braverman, and L. Rozonoer. Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, 25:821 – 837, 1964.
- K. S. Alexander. Probability inequalities for empirical processes and a law of the iterated logarithm. *Annals of Probability*, 12:1041–1067, 1984.
- N. Alon, S. Ben-David, N. Cesa-Bianchi, and D. Haussler. Scale-sensitive Dimensions, Uniform Convergence, and Learnability. *Journal of the ACM*, 44(4):615–631, 1997.
- M. Anthony and P. Bartlett. *Neural Network Learning: Theoretical Foundations*. Cambridge University Press, 1999. (to appear).
- N. Aronszajn. Theory of reproducing kernels. *Transactions of the American Mathematical Society*, 68:337 – 404, 1950.
- P. Bartlett and J. Shawe-Taylor. Generalization performance of support vector machines and other pattern classifiers. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods — Support Vector Learning*, pages 43–54, Cambridge, MA, 1999. MIT Press.
- P. L. Bartlett. 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, 1998.
- P. L. Bartlett, S. R. Kulkarni, and S. E. Posner. Covering numbers for real-valued function classes. *IEEE Transactions on Information Theory*, 43(5):1721–1724, 1997.
- D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, MA, 1995.
- V. Blanz, B. Schölkopf, H. Bülthoff, C. Burges, V. Vapnik, and T. Vetter. Comparison of view-based object recognition algorithms using realistic 3D models. In C. von der Malsburg, W. von Seelen, J. C. Vorbrüggen, and B. Sendhoff, editors, *Artificial Neural Networks — ICANN’96*, pages 251 – 256, Berlin, 1996. Springer Lecture Notes in Computer Science, Vol. 1112.
- B. E. Boser, I. M. Guyon, and V. N. Vapnik. A training algorithm for optimal margin classifiers. In D. Haussler, editor, *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152, Pittsburgh, PA, July 1992. ACM Press.

- C. J. C. Burges and B. Schölkopf. Improving the accuracy and speed of support vector learning machines. In M. Mozer, M. Jordan, and T. Petsche, editors, *Advances in Neural Information Processing Systems 9*, pages 375–381, Cambridge, MA, 1997. MIT Press.
- B. Carl and I. Stephani. *Entropy, compactness, and the approximation of operators*. Cambridge University Press, Cambridge, UK, 1990.
- S. Chen, D. Donoho, and M. Saunders. Atomic decomposition by basis pursuit. Technical Report 479, Department of Statistics, Stanford University, May 1995.
- C. Cortes and V. Vapnik. Support vector networks. *Machine Learning*, 20:273 – 297, 1995.
- D. Cox and F. O’Sullivan. Asymptotic analysis of penalized likelihood and related estimators. *Ann. Statist.*, 18:1676–1695, 1990.
- CPLEX Optimization Incorporated, Incline Village, Nevada. *Using the CPLEX Callable Library*, 1994.
- Luc Devroye and Gábor Lugosi. Lower bounds in pattern recognition and learning. *Pattern Recognition*, 28, 1995.
- R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. John Wiley & Sons, 1973.
- N. Dunford and J. T. Schwartz. *Linear Operators Part II: Spectral Theory, Self Adjoint Operators in Hilbert Space*. Number VII in Pure and Applied Mathematics. John Wiley & Sons, New York, 1963.
- A. Ehrenfeucht, D. Haussler, M. Kearns, and L. Valiant. A general lower bound on the number of examples needed for learning. *Information and Computation*, 82: 247–261, 1989.
- Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. In *Computational Learning Theory: Eurocolt ’95*, pages 23–37. Springer-Verlag, 1995.
- Y. Freund and R. E. Schapire. Experiments with a new boosting algorithm. In *Proc. 13th International Conference on Machine Learning*, pages 148–146. Morgan Kaufmann, 1996.
- F. Girosi. An equivalence between sparse approximation and support vector machines. *Neural Computation*, 10(6):1455–1480, 1998.
- F. Girosi, M. Jones, and T. Poggio. Priors, stabilizers and basis functions: From regularization to radial, tensor and additive splines. A.I. Memo No. 1430, MIT, 1993.
- H. Goldstein. *Classical Mechanics*. Addison-Wesley, Reading, MA, 1986.
- Thore Graepel, Ralf Herbrich, Bernhard Schölkopf, Alex Smola, Peter Bartlett, Klaus Robert-Müller, Klaus Obermayer, and Bob Williamson. Classification on proximity data with LP-machines. In *Proceedings of the International Conference of Artificial Neural Networks*, 1999. in press.

- L. Gurvits. A note on a scale-sensitive dimension of linear bounded functionals in banach spaces. In *Proceedings of Algorithm Learning Theory, ALT-97*, 1997. Also: NECI Technical Report, 1997.
- I. Guyon, B. Boser, and V. Vapnik. Automatic capacity tuning of very large VC-dimension classifiers. In S. J. Hanson, J. D. Cowan, and C. L. Giles, editors, *Advances in Neural Information Processing Systems*, volume 5, pages 147–155. Morgan Kaufmann, San Mateo, CA, 1993.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58:13–30, 1963.
- A.E. Hoerl and R. W. Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- P. J. Huber. *Robust Statistics*. John Wiley and Sons, New York, 1981.
- IBM Corporation. IBM optimization subroutine library guide and reference. *IBM Systems Journal*, 31, 1992. SC23-0519.
- T. S. Jaakkola and D. Haussler. Probabilistic kernel regression models. In *Proceedings of the 1999 Conference on AI and Statistics*, 1999.
- T. Joachims. Making large-scale SVM learning practical. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods — Support Vector Learning*, pages 169–184, Cambridge, MA, 1999. MIT Press.
- W. Karush. Minima of functions of several variables with inequalities as side constraints. Master’s thesis, Dept. of Mathematics, Univ. of Chicago, 1939.
- L. Kaufmann. Solving the quadratic programming problem arising in support vector classification. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods — Support Vector Learning*, pages 147–168, Cambridge, MA, 1999. MIT Press.
- G. Kimeldorf and G. Wahba. Some results on Tchebycheffian spline functions. *J. Math. Anal. Applic.*, 33:82–95, 1971.
- H. W. Kuhn and A. W. Tucker. Nonlinear programming. In *Proc. 2nd Berkeley Symposium on Mathematical Statistics and Probabilistics*, pages 481–492, Berkeley, 1951. University of California Press.
- Philip M. Long. The complexity of learning according to two models of a drifting environment. In *Proceedings of the 11th Annual Conference on Computational Learning Theory*, pages 116–125. ACM Press, 1998.
- D. G. Luenberger. *Introduction to Linear and Nonlinear Programming*. Addison-Wesley, Reading, MA, 1973.
- O. L. Mangasarian. Multi-surface method of pattern separation. *IEEE Transactions on Information Theory*, IT-14:801–807, 1968.
- O. L. Mangasarian. Mathematical programming in data mining. *Data Mining and Knowledge Discovery*, 42(1):183–201, 1997.
- Lew Mason and Peter Bartlett. Direct optimization of margins. In *Advances in*

- Neural Information Processing Systems*, page in press, San Mateo, CA, 1998. Morgan Kaufmann.
- J. Mercer. Functions of positive and negative type and their connection with the theory of integral equations. *Philos. Trans. Roy. Soc. London*, A 209:415–446, 1909.
- B. A. Murtagh and M. A. Saunders. MINOS 5.4 user’s guide. Technical Report SOL 83.20, Stanford University, 1993.
- M. Oren, C. Papageorgiou, P. Sinha, E. Osuna, and T. Poggio. Pedestrian detection using wavelet templates. In *Proc. Computer Vision and Pattern Recognition*, pages 193–199, Puerto Rico, June 16–20 1997.
- E. Osuna, R. Freund, and F. Girosi. An improved training algorithm for support vector machines. In J. Principe, L. Gile, N. Morgan, and E. Wilson, editors, *Neural Networks for Signal Processing VII — Proceedings of the 1997 IEEE Workshop*, pages 276 – 285, New York, 1997a. IEEE.
- E. Osuna, R. Freund, and F. Girosi. Training support vector machines: An application to face detection. In *Proc. Computer Vision and Pattern Recognition ’97*, pages 130–136, 1997b.
- E. Parzen. An approach to time series analysis. *Ann. Math. Statist.*, 32:951–989, 1962.
- J. Platt. Fast training of support vector machines using sequential minimal optimization. In B. Schölkopf, C. J. C. Burges, and A. J. Smola, editors, *Advances in Kernel Methods — Support Vector Learning*, pages 185–208, Cambridge, MA, 1999. MIT Press.
- T. Poggio. On optimal nonlinear associative recall. *Biological Cybernetics*, 19: 201–209, 1975.
- J. R. Quinlan. Bagging, boosting, and C4.5. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence and the Eighth Innovative Applications of Artificial Intelligence Conference*, pages 725–730, Menlo Park, 1996. AAAI Press / MIT Press.
- G. Rätsch. Ensemble learning for classification. Master’s thesis, University of Potsdam, 1998. in German.
- Gunnar Rätsch, T. Onoda, and Klaus-Robert Müller. Soft margins for adaboost. Technical report, Royal Holloway, University of London, 1998. TR–98–21.
- S. Saitoh. *Theory of Reproducing Kernels and its Applications*. Longman Scientific & Technical, Harlow, England, 1988.
- C. Saunders, M. O. Stitson, J. Weston, L. Bottou, B. Schölkopf, and A. Smola. Support vector machine - reference manual. Technical Report CSD-TR-98-03, Department of Computer Science, Royal Holloway, University of London, Egham, TW20 0EX, UK, 1998. TR available as http://www.dcs.rhnc.ac.uk/research/compint/areas/comp_learn/sv/pub/report98-03.ps; SVM available at <http://svm.dcs.rhnc.ac.uk/>.

- R. Schapire, Y. Freund, P. Bartlett, and W. Sun Lee. Boosting the margin: A new explanation for the effectiveness of voting methods. *Annals of Statistics*, 1998. (To appear. An earlier version appeared in: D.H. Fisher, Jr. (ed.), *Proceedings ICML97*, Morgan Kaufmann.).
- B. Schölkopf, C. Burges, and V. Vapnik. Extracting support data for a given task. In U. M. Fayyad and R. Uthurusamy, editors, *Proceedings, First International Conference on Knowledge Discovery & Data Mining*. AAAI Press, Menlo Park, CA, 1995.
- B. Schölkopf. *Support Vector Learning*. R. Oldenbourg Verlag, Munich, 1997.
- B. Schölkopf, S. Mika, A. Smola, G. Rätsch, and K.-R. Müller. Kernel PCA pattern reconstruction *via* approximate pre-images. In L. Niklasson, M. Bodén, and T. Ziemke, editors, *Proceedings of the 8th International Conference on Artificial Neural Networks*, Perspectives in Neural Computing, pages 147 – 152, Berlin, 1998a. Springer Verlag.
- B. Schölkopf, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson. Generalization bounds via eigenvalues of the Gram matrix. Submitted to COLT99, February 1999.
- B. Schölkopf, P. Simard, A. Smola, and V. Vapnik. Prior knowledge in support vector kernels. In M. Jordan, M. Kearns, and S. Solla, editors, *Advances in Neural Information Processing Systems 10*, pages 640 – 646, Cambridge, MA, 1998b. MIT Press.
- B. Schölkopf, A. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10:1299–1319, 1998c.
- B. Schölkopf, A. Smola, R. Williamson, and P. Bartlett. New support vector algorithms. NeuroCOLT Technical Report NC-TR-98-031, Royal Holloway College, University of London, UK, 1998d.
- B. Schölkopf, K. Sung, C. Burges, F. Girosi, P. Niyogi, T. Poggio, and V. Vapnik. Comparing support vector machines with gaussian kernels to radial basis function classifiers. A.I. Memo No. 1599, Massachusetts Institute of Technology, 1996.
- B. Schölkopf, K. Sung, C. Burges, F. Girosi, P. Niyogi, T. Poggio, and V. Vapnik. Comparing support vector machines with gaussian kernels to radial basis function classifiers. *IEEE Trans. Sign. Processing*, 45:2758 – 2765, 1997.
- H. Schwenk and Y. Bengio. Training methods for adaptive boosting of neural networks. In Michael I. Jordan, Michael J. Kearns, and Sara A. Solla, editors, *Advances in Neural Information Processing Systems*, volume 10. The MIT Press, 1998.
- J. Shawe-Taylor, P. L. Bartlett, R. C. Williamson, and M. Anthony. Structural risk minimization over data-dependent hierarchies. *IEEE Transactions on Information Theory*, 1998. To appear. Also: NeuroCOLT Technical Report NC-TR-96-053, 1996, ftp://ftp.dcs.rhnc.ac.uk/pub/neurocolt/tech_reports.
- John Shawe-Taylor and Nello Cristianini. Margin distribution bounds on general-

- ization. Technical report, Royal Holloway, University of London, 1998. NC2-TR-1998-020.
- Hans Ulrich Simon. General bounds on the number of examples needed for learning probabilistic concepts. *J. of Comput. Syst. Sci.*, 52(2):239–254, 1996. Earlier version in 6th COLT, 1993.
- A. Smola and B. Schölkopf. From regularization operators to support vector kernels. In M. Jordan, M. Kearns, and S. Solla, editors, *Advances in Neural Information Processing Systems 10*, pages 343 – 349, Cambridge, MA, 1998a. MIT Press.
- A. Smola and B. Schölkopf. On a kernel-based method for pattern recognition, regression, approximation and operator inversion. *Algorithmica*, 22:211 – 231, 1998b.
- A. Smola, B. Schölkopf, and K.-R. Müller. General cost functions for support vector regression. In T. Downs, M. Frean, and M. Gallagher, editors, *Proc. of the Ninth Australian Conf. on Neural Networks*, pages 79 – 83, Brisbane, Australia, 1998. University of Queensland.
- A. J. Smola. *Learning with Kernels*. PhD thesis, Technische Universität Berlin, 1998.
- M. Talagrand. Sharper bounds for gaussian and empirical processes. *Annals of Probability*, 22:28–76, 1994.
- R. Vanderbei. LOQO: An interior point code for quadratic programming. Technical Report SOR 94-15, Princeton University, 1994.
- R. J. Vanderbei. LOQO user’s manual – version 3.10. Technical Report SOR-97-08, Princeton University, Statistics and Operations Research, 1997. Code available at <http://www.princeton.edu/~rvdb/>.
- V. Vapnik. *Estimation of Dependences Based on Empirical Data [in Russian]*. Nauka, Moscow, 1979. (English translation: Springer Verlag, New York, 1982).
- V. Vapnik. *The Nature of Statistical Learning Theory*. Springer Verlag, New York, 1995.
- V. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998. forthcoming.
- V. Vapnik and A. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2): 264–280, 1971.
- G. Wahba. Convergence rates of certain approximate solutions to Fredholm integral equations of the first kind. *Journal of Approximation Theory*, 7:167 – 185, 1973.
- Jason Weston. Leave-one-out support vector machines. In *Proceedings of the International Joint Conference on Artificial Intelligence*, Sweden, 1999.
- C. K. I. Williams. Prediction with gaussian processes: From linear regression to linear prediction and beyond. In M. I. Jordan, editor, *Learning and Inference in Graphical Models*. Kluwer, 1998. To appear. Also: Technical Report NCRG/97/012, Aston University.

- R. Williamson, A. Smola, and B. Schölkopf. Entropy numbers, operators and support vector kernels. In *Advances in Neural Information Processing Systems 11*, 1998a. Submitted.
- R. C. Williamson, A. J. Smola, and B. Schölkopf. Generalization performance of regularization networks and support vector machines via entropy numbers of compact operators. NeuroCOLT Technical Report NC-TR-98-019, Royal Holloway College, University of London, UK, 1998b.