# 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.xxxx

CIP

## Contents

	Preface	vii
1	Introduction to Large Margin Classifiers  Alex J. Smola, Peter Bartlett, Bernhard Schölkopf, and  Dale Schuurmans	1
2	Large Margin Rank Boundaries for Ordinal Regression Ralf Herbrich, Thore Graepel, and Klaus Obermayer	29
	References	46

## 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

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

labels together with corresponding labels

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

The goal is to find some decision function  $g: \mathbb{R}^N \to \{-1,1\}$  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  $f: \mathbb{R}^N \to \mathbb{R}$  whose output is passed through a sign threshold to yield the final classification  $g(\mathbf{x}) = \operatorname{sgn}(f(\mathbf{x}))$ . 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

$$g(\mathbf{x}) := \operatorname{sgn}(f(\mathbf{x})) \text{ where } f(\mathbf{x}) = (\mathbf{x} \cdot \mathbf{w}) + b \text{ for } \mathbf{w} \in \mathbb{R}^N \text{ and } b \in \mathbb{R}.$$
 (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  $\mathrm{sgn}\left((\mathbf{w}\cdot\mathbf{x}_i)+b\right)=y_i$  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.

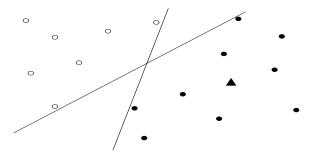


Figure 1.1 A linearly separable classification problem. Note that there may be several possible solutions as depicted by the two lines. The problem becomes non-separable if we replace the triangle by an open circle; in which case no solution  $(\mathbf{w}, b)$  exists.

#### 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

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

Bayes optimal decision function

This function satisfies

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

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  $\Sigma$  centered at  $\mu_+$  and  $\mu_-$  respectively

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

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

By rearranging we find that this condition is equivalent to

$$\mathbf{x}^{\top} \Sigma^{-1} \mathbf{x} - 2\boldsymbol{\mu}_{+}^{\top} \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}_{+}^{\top} \Sigma^{-1} \boldsymbol{\mu}_{+} - \mathbf{x}^{\top} \Sigma^{-1} \mathbf{x} + 2\boldsymbol{\mu}_{-}^{\top} \Sigma^{-1} \mathbf{x} - \boldsymbol{\mu}_{-}^{\top} \Sigma^{-1} \boldsymbol{\mu}_{-} = 0$$

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

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

$$f(\mathbf{x}) = ((\mu_{+} - \mu_{-})^{\top} \Sigma^{-1}) \mathbf{x} - \frac{1}{2} (\mu_{+}^{\top} \Sigma^{-1} \mu_{+} - \mu_{-}^{\top} \Sigma^{-1} \mu_{-}). \tag{1.9}$$

linear discriminant

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  $g(\mathbf{x}) = \operatorname{sgn}(f(\mathbf{x}))$  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
                    Weight vector \mathbf{w} and threshold b.
returns:
function Perceptron (X, Y, \eta)
     initialize \mathbf{w}, b = 0
     repeat
            for all i from i = 1, \ldots, m
                    Compute q(\mathbf{x}_i) = \operatorname{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 \le i \le 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  $g: \mathbf{x} \mapsto \operatorname{sgn}((\mathbf{w} \cdot \mathbf{x}) + b)$ , 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  $\rho$  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

$$y_i\left((\mathbf{w}^* \cdot \mathbf{x}_i) + b^*\right) \ge \rho \text{ for all } 1 \le i \le m.$$
 (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  $R = \max_i ||x_i||^2$ . 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 jth 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,

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

$$= ((\mathbf{w}_j, b_j) \cdot (\mathbf{w}^*, b^*)) + \eta y_i((\mathbf{x}_i, 1) \cdot (\mathbf{w}^*, b^*))$$

$$\geq ((\mathbf{w}_j, b_j) \cdot (\mathbf{w}^*, b^*)) + \eta \rho$$

$$> j\eta\rho.$$

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

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

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  $\rho$  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  $\rho$  is what we shall henceforth call a marqin.

## Definition 1.1 Margin and Margin Errors

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

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

margin

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

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

minimum margin

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  $\rho_f$  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:

- 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  $\operatorname{sgn}(f(\mathbf{x}))$ . 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  $\varepsilon$ . (See, for example, Duda and Hart (1973).)
- Similarly, a slight perturbation to the function f will not affect any of the resulting classifications on the training data  $(\mathbf{x}_1, y_1), ..., (\mathbf{x}_m, y_m)$ . 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,  $y_i f_{\mathbf{w}+\varepsilon}(\mathbf{x}_i)$  will also be positive for small  $\varepsilon$ .

robustness in patterns

robustness in parameters

#### 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

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

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  $f: \mathbb{R}^N \to \mathbb{R}$  by

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

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

 $\begin{array}{c} \text{optimal} \\ \text{hyperplane} \end{array}$ 

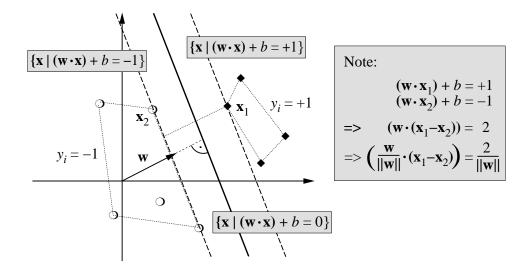


Figure 1.2 A binary classification toy problem: separate balls from diamonds. The optimal hyperplane is orthogonal to the shortest line connecting the convex hulls of the two classes (dotted), and intersects it half-way between the two classes. The problem being separable, there exists a weight vector  $\mathbf{w}$  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  $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) > 1$ . 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  $\mathbf{w}/\|\mathbf{w}\|$ .

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

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

$$(1.15)$$

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

The formulation (1.16) has a simple geometric interpretation:  $-b\mathbf{w}/\|\mathbf{w}\|^2$  is the vector in direction w that ends right on the decision hyperplane (since  $(\mathbf{w} \cdot (-b\mathbf{w}/\|\mathbf{w}\|^2)) = -b)$ , 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  $(\mathbf{w} \cdot \mathbf{x}_i) \mathbf{w} / ||\mathbf{w}||^2 - (-b\mathbf{w} / ||\mathbf{w}||^2)$  appropriatedly 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,  $\rho$ , and maximizing  $\rho$  subject to the constraint that it really is a lower bound:

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

optimization problems

Euclidean

Margin

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

= argmax 
$$\rho$$
 subject to  $\|\mathbf{w}\| = 1$  and  $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \ge \rho$  for  $1 \le i \le m$  (1.18)

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \quad \|\mathbf{w}\|^2 \quad \text{subject to } y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \ge 1 \text{ for } 1 \le i \le m$$
 (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  $\ell_1$  sphere instead of the unit  $\ell_2$  sphere gives the problem of determining

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

$$= \underset{b}{\operatorname{argmax}} \quad \rho \quad \text{subject to } \|\mathbf{w}\|_1 = 1 \text{ and } y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \ge \rho \text{ for } 1 \le i \le m \quad (1.20)$$

 $\ell_{\infty}$  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  $\ell_{\infty}$  distance between the training patterns and the decision hyperplane, in a direct analogue to the original Euclidean formulation (1.15).

Similarly, the constraint that w lie on the unit  $\ell_{\infty}$  sphere yields the problem

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

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

 $\ell_1$  margin

which is also a linear programming problem, but now equivalent to finding the weight vector and threshold that maximize the minimum  $\ell_1$  distance between the training patterns and the decision hyperplane. In general, constraining  $\mathbf{w}$  to lie on the unit  $\ell_p$  sphere yields a convex programming problem

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

$$= \underset{\mathbf{w}, b, \rho}{\operatorname{argmax}} \quad \rho \quad \text{subject to } \|\mathbf{w}\|_p = 1 \text{ and } y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \ge \rho \text{ for } 1 \le i \le m \text{ (1.22)}$$

 $\ell_q$  margin

which is equivalent to finding the weight vector and threshold that maximize the minimum  $\ell_q$  distance between the training patterns and the decision hyperplane, where  $\ell_p$  and  $\ell_q$  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

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  $\rho$  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  $y_i((\mathbf{w} \cdot \mathbf{x}_i) + b) \ge 1$  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

$$\xi_i \ge 0, \text{ for all } i = 1, \dots, m, \tag{1.23}$$

along with relaxed constraints

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

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

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

Expected Risk

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

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

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

Empirical Risk

For thresholded decision functions  $g: \mathbb{R}^N \to \{-1,1\}$  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),

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

In this case, the empirical risk functional is

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

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  $\rho$ , i.e.

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

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  $R_{\rm emp}(g)$  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  $R_{\rm emp}(g) = 0$  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

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.

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

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

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

$$R(g^*) \le \inf_{g \in G} R(g) + \sqrt{\frac{c}{m} \left(h + \ln\left(\frac{1}{\delta}\right)\right)}$$
(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  $\delta$  over m random examples, the error of its hypothesis g satisfies

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

1.2 Theory 11

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  $g(\mathbf{x}) = \operatorname{sgn}(f(\mathbf{x}))$ . Here, the VC dimension is a particularly weak method for measuring the representational capacity of the resulting function class  $G = \operatorname{sgn}(F)$ .

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  $X_0$ ,

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

satisfying  $|f(\mathbf{x})| \ge \rho$  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

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

Note that this can be much smaller than the VC dimension of  $\operatorname{sgn}(F)$  obtained without taking  $\rho$  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.

## 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  $\rho$ -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

$$y_i(f(x_i) - b_i) \ge \rho \text{ for } 1 \le i \le |S|.$$
 (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).) The fat shattering dimension  $\operatorname{fat}_F(\rho)$  of the set F is a function from the positive real numbers to the integers which maps a value  $\rho$  to the size of the largest  $\rho$ -shattered

set, if this is finite, or infinity otherwise.

the VC dimension obtained by thresholding but requiring that outputs are  $\rho$  above the threshold for positive classification and  $\rho$  below for negative.

We may think of the fat-shattering dimension of a set of real-valued functions as

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

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  $\varepsilon$  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

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

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

Covering a class of functions F with an  $\varepsilon$ -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

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

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).)

fat shattering

covering number

## Theorem 1.4 Bounds on 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  $m \ge \operatorname{fat}_F(\varepsilon/4)$ . Then

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

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.  $\epsilon$ , 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 N and $\rho$

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

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

Furthermore,

Pr ("some 
$$f$$
 in  $F$  has  $R(f) \ge R_{\rho}(f) + \varepsilon$ ")  $\le 2 \mathcal{N}\left(\frac{\rho}{2}, F, 2m\right) e^{-\varepsilon^2 m/8}$ . (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).

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  $\mathcal{X} \times [0,1]$ ,  $\varepsilon$  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

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

anatomy of a uniform convergence bound Comparing with (1.41), notice that the scale of the covering number depends on the desired accuracy  $\varepsilon$ , whereas in (1.41) it depends on the scale  $\rho$  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  $\ell_2$  ball of radius R in  $\mathbb{R}^n$ , centered at the origin, and consider the set

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

Then

$$fat_F(\rho) \le \left(\frac{R}{\rho}\right)^2$$
. (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  $\rho$  on all training examples (i.e.  $R_{\rho}(f) = 0$ ) satisfies

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

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

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

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

$$F_R = \{ \mathbf{x} \mapsto \mathbf{w}^T \mathbf{x} \mid \|\mathbf{x}\|_{\infty} < 1, \|\mathbf{w}\|_1 < R \}$$

$$(1.47)$$

we have

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

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

$$co(G) = \left\{ \sum_{i} \alpha_{i} g_{i} \mid \alpha_{i} > 0, \sum_{i} \alpha_{i} = 1, g_{i} \in G \right\}.$$

$$(1.49)$$

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  $\rho > 0$ .

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

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

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

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

Lagrangian

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

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*  $\alpha_i$  (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  $\alpha_i$ . 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.

 $\begin{array}{c} {\rm KKT} \\ {\rm conditions} \end{array}$ 

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  $\alpha_i$  must be 0: this is the value of  $\alpha_i$  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,

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

leads to

$$\sum_{i=1}^{m} \alpha_i y_i = 0 \tag{1.55}$$

and

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

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  $\alpha_i$  is non-zero. By the Karush-Kuhn-Tucker complementarity conditions these training patterns are the ones for which

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

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).

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  $\alpha_i$  which

dual optimization problem

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

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

The hyperplane decision function can thus be written as

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

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

$$\Phi: \mathbb{R}^N \to \mathcal{F},\tag{1.61}$$

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  $j_1,\ldots,j_d\in\{1,\ldots,N\}$ . 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  $16 \times 16$  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

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

polynomial kernel

For instance, the polynomial kernel

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

can be shown to correspond to a map  $\Phi$  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)

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

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  $\Phi$  such that (1.62) holds (Mercer, 1909; Aizerman et al., 1964; Boser et al., 1992; Dunford and Schwartz, 1963):

## Theorem 1.11 Mercer

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

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

with

$$\int_{\mathcal{X}\times\mathcal{X}} k(\mathbf{x}, \mathbf{x}') f(\mathbf{x}) f(\mathbf{x}') d\mathbf{x} d\mathbf{x}' \ge 0$$
(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  $\psi_j$  and positive eigenvalues  $\lambda_j$ ,

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

where  $N_{\mathcal{F}} < \infty$  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  $\Phi$  into a potentially infinite-dimensional  $l_2$  space which satisfies (1.62). For instance, we may use

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

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  $f(.) = k(\mathbf{x}_i, .)$  ( $\mathbf{x}_i \in \mathcal{X}$ ). To ensure that the map  $\Phi : \mathcal{X} \to \mathcal{H}_k$ , which in this case is defined as

$$\Phi(\mathbf{x}) = k(\mathbf{x}, .),\tag{1.69}$$

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

$$\langle k(\mathbf{x},.), k(\mathbf{x}',.) \rangle = k(\mathbf{x}, \mathbf{x}'), \tag{1.70}$$

Figure 1.3 The idea of SV machines: map the training data nonlinearly into a higher-dimensional feature space via  $\Phi$ , 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  $\psi_i$   $(i = 1, ..., N_F)$ 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  $\delta_{jn}$ . From this, we can construct  $\langle .,. \rangle$ such that

$$\langle \sqrt{\lambda_i} \psi_i, \sqrt{\lambda_n} \psi_n \rangle = \delta_{in}. \tag{1.71}$$

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

Besides (1.63), SV practictioners use sigmoid kernels

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

kernel

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

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

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  $\Phi$  (1.69), only single bumps  $k(\mathbf{x}, \cdot)$  do have pre-images under  $\Phi$ .

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

sigmoid

kernel

Gaussian RBF

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

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

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

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

subject to 
$$\alpha_i \ge 0$$
,  $i = 1, \dots, m$ , and  $\sum_{i=1}^m \alpha_i y_i = 0$ . (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

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

The only difference from the separable case (1.76) is the upper bound C on the Lagrange multipliers  $\alpha_i$ . 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  $\xi_i$  is zero (this again follows from the Karush-Kuhn-Tucker complementarity conditions), and hence

$$\sum_{j=1}^{m} y_j \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) + b = y_i.$$
(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  $0 \le \nu \le 1$  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.

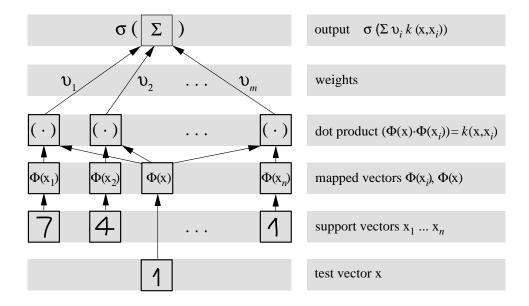


Figure 1.4 Architecture of SV machines. The input  $\mathbf{x}$  and the Support Vectors  $\mathbf{x}_i$  are nonlinearly mapped (by  $\Phi$ ) into a feature space  $\mathcal{F}$ , where dot products are computed. By the use of the kernel k, these two layers are in practice computed in one single step. 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  $\sigma$  (in pattern recognition,  $\sigma(x) = \operatorname{sgn}(x + b)$ ; 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

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

regularized risk

(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  $R_{\rm reg}(f)$  on the set of functions

$$f(\mathbf{x}) = \sum_{i} \alpha_{i} k(\mathbf{x}_{i}, \mathbf{x}) + b, \tag{1.80}$$

provided that k and P are interrelated by

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

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  $(k(\mathbf{x}_i,.)\cdot(P^*Pk)(\mathbf{x}_j,.))=(k(\mathbf{x}_i,.)\cdot\delta_{\mathbf{x}_j}(.))=k(\mathbf{x}_i,\mathbf{x}_j)$ . 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.

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.

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

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

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

regularization networks

Gaussian processes

1.4 Boosting 23

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

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

$$R_{\text{out}}(f) \leq \frac{1}{m} \sum_{i=1}^{m} 1_{\left\{ y_{i} \sum_{j \neq i} \alpha_{j} y_{j} k(\mathbf{x}_{j}, \mathbf{x}_{i}) + y_{i} b > 0 \right\}}$$

$$= \frac{1}{m} \sum_{i=1}^{m} 1_{\left\{ y_{i} f(\mathbf{x}_{i}) - \alpha_{i} k(\mathbf{x}_{i}, \mathbf{x}_{i})) > 0 \right\}}.$$
(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 ?? and ?? for futher 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

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

```
argument: Training sample, X=\{\mathbf{x}_1,\ldots,\mathbf{x}_m\}\subset\mathcal{X}, Y=\{y_1,\ldots,y_m\}\subset\{\pm 1\}
                      Number of iterations, {\cal T}
                       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, \ldots, m
             D_1(i) := 1/m
     endfor
     for all t from \{1, \ldots, T\}
             q_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)}
             D_{t+1}(i) := \left\{ egin{array}{ll} D_t(i)e^{-lpha_t}/Z_t & 	ext{if } y_i = g_t(x_i) \ D_t(i)e^{lpha_t}/Z_t & 	ext{otherwise,} \end{array} 
ight.
             endfor
     endfor
     \mathbf{return} \ f = \frac{\sum_{t=1}^{T} \alpha_t g_t}{\sum_{t=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  $\rho$  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

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

In particular, if 
$$\varepsilon_t \leq 1/2 - 2\rho$$
, then 
$$R_{\rho}(f) < (1 - \rho^2)^{T/2},$$
 and this is less than  $\varepsilon$  for  $T \geq (2/\rho^2) \ln(1/\varepsilon)$ . (1.87)

## 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 27

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

$c(\mathbf{x}, y, f(\mathbf{x}))$	cost function
R(g)	risk of $g$ , i.e. expected fraction of errors
$R_{\mathrm{emp}}(g)$	empirical risk of $g$ , i.e. fraction of training errors
R(f)	risk of $f$
$R_{ m emp}(f)$	empirical risk of $f$
k	Mercer kernel
$\mathcal F$	Feature space induced by a kernel
$\Phi$	map into feature space (induced by $k$ )
$\alpha_i$	Lagrange multiplier
lpha	vector of all Lagrange multipliers
$\xi_i$	slack variables
ξ	vector of all slack variables
C	regularization constant for SV Machines
$\lambda$	regularization constant $(C = \frac{1}{\lambda})$

# 2 Large Margin Rank Boundaries for Ordinal Regression

Ralf Herbrich, Klaus Obermayer, and Thore Graepel
Technical University of Berlin
Department of Computer Science
Franklinstr. 28/29,
10587 Berlin,
Germany
ralfh|oby|graepel2@cs.tu-berlin.de

In contrast to the standard machine learning tasks of classification and metric regression we investigate the problem of predicting variables of ordinal scale, a setting referred to as ordinal regression. This problem arises frequently in the social sciences and in information retrieval where human preferences play a major role. Whilst approaches proposed in statistics rely on a probability model of a latent (unobserved) variable we present a distribution independent risk formulation of ordinal regression which allows us to derive a uniform convergence bound. Applying this bound we present a large margin algorithm that is based on a mapping from objects to scalar utility values thus classifying pairs of objects. We give experimental results for an information retrieval task which show that our algorithm outperforms more naive approaches to ordinal regression such as Support Vector Classification and Support Vector Regression in the case of more than two ranks.

#### 2.1 Introduction

Let us shortly recall the model presented in Chapter 1. Given an iid sample (X,Y), and a set F of mappings  $f: \mathcal{X} \mapsto \mathcal{Y}$ , a learning procedure aims at finding  $f^*$  such that — using a predefined loss  $c: \mathcal{X} \times \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$  — the risk functional (1.26) is minimized. Using the principle of Empirical Risk Minimization (ERM), one chooses the function  $f_{\text{emp}}$  which minimizes the mean of the loss  $R_{\text{emp}}(f)$  (Equation 1.27)

classification and regression

given the sample (X,Y). Introducing a quantity which characterizes the capacity of F, bounds for the deviation  $|R(f_{\rm emp}) - \inf_{f \in F} R(f)|$  can be derived (see Theorems 1.2, 1.3, 1.5, and 1.6). Two main scenarios were considered in the past: (i) If  $\mathcal{Y}$  is a finite unordered set (nominal scale), the task is referred to as classification learning. Since  $\mathcal{Y}$  is unordered, the 0-1 loss, i.e.,  $c_{\rm class}(\mathbf{x},y,f(\mathbf{x}))=1_{f(\mathbf{x})\neq y}$ , is adequate to capture the loss at each point  $(\mathbf{x},y)$ . (ii) If  $\mathcal{Y}$  is a metric space, e.g., the set of real numbers, the task is referred to as regression estimation. In this case the loss function can take into account the full metric structure. Different metric loss functions have been proposed which are optimal under given probability models  $\P(y|\mathbf{x})$  (c.f Huber (1981)). Usually, optimality is measured in terms of the mean squared error of  $f_{\rm emp}$ .

Here, we consider a problem which shares properties of both cases (i) and (ii). Like in (i)  $\mathcal{Y}$  is a finite set and like in (ii) there exists an ordering among the elements of  $\mathcal{Y}$ . In contrast to regression estimation we have to deal with the fact that  $\mathcal{Y}$  is a non-metric space. A variable of the above type exhibits an ordinal scale and can be considered as the result of a coarsely measured continuous variable (Anderson and Philips, 1981). The ordinal scale leads to problems in defining an appropriate loss function for our task (see also McCullagh (1980) and Anderson (1984)): On the one hand, there exists no metric in the space  $\mathcal{Y}$ , i.e., the distance (y-y') of two elements is not defined. On the other hand, the simple 0-1 loss does not reflect the ordering in  $\mathcal{Y}$ . Since no loss function  $c(\mathbf{x}, y, f(\mathbf{x}))$  can be found that acts on true ranks y and predicted ranks  $f(\mathbf{x})$ , we suggest to exploit the ordinal nature of the elements of  $\mathcal{Y}$  by considering the order on the space  $\mathcal{X}$  induced by each mapping  $f: \mathcal{X} \mapsto \mathcal{Y}$ . Thus our loss function  $c_{\text{pref}}(\mathbf{x}_1, \mathbf{x}_2, y_1, y_2, f(\mathbf{x}_1), f(\mathbf{x}_2))$  acts on pairs of true ranks  $(y_1, y_2)$  and predicted ranks  $(f(\mathbf{x}_1), f(\mathbf{x}_2))$ . Such an approach makes it possible to formulate a distribution independent theory of ordinal regression and to give uniform bounds for the risk functional. Roughly speaking, the proposed risk functional measures the probability of misclassification of a randomly drawn pair  $(\mathbf{x}_1, \mathbf{x}_2)$  of observations, where the two classes are  $\mathbf{x}_1 \succ_{\mathcal{X}} \mathbf{x}_2$  and  $\mathbf{x}_2 \succ_{\mathcal{X}} \mathbf{x}_1$  (see Section 2.3). Problems of ordinal regression arise in many fields, e.g., in information retrieval (Wong et al., 1988; Herbrich et al., 1998), in econometric models (Tangian and Gruber, 1995; Herbrich et al., 1999), and in classical statistics (McCullagh, 1980; Fahrmeir and Tutz, 1994; Anderson, 1984; de Moraes and Dunsmore, 1995; Keener and Waldman, 1985).

distribution independent theory of ordinal regression

by intervals on the real line. Then the task is to find a latent utility function that maps objects to scalar values. Due to the ordering of ranks, the function is restricted to be transitive and asymmetric, because these are the defining properties of a preference relation. The resulting learning task is also referred to as learning of preference relations (see Herbrich et al. (1998)). One might think that learning of preference relations reduces to a standard classification problem if pairs of objects are considered. This, however, is not true in general because the properties of transitivity and asymmetry may be violated by traditional Bayesian approaches

which may violate stochastic transitivity (Suppes et al., 1989). Considering pairs of

As an application of the above—mentioned theory, we suggest to model ranks

preference relation

large margin

objects, the task of learning reduces to finding a utility function that best reflects the preferences induced by the unknown distribution  $p(\mathbf{x}, y)$ . Our learning procedure on pairs of objects is an application of the large margin idea known from data–dependent Structural Risk Minimization (Shawe-Taylor et al., 1998). The resulting algorithm is similar to Support Vector Machines (see Section 1.3). Since during learning and application of SVMs only inner products of object representations  $\mathbf{x}_i$  and  $\mathbf{x}_j$  have to be computed, the method of potential functions can be applied (see Aizerman et al. (1964) or Section 1.3.2).

In Section 2.2 we introduce the setting of ordinal regression and shortly present well known results and models from the field of statistics. In Section 2.3 we introduce our model for ordinal regression and give a bound for the proposed loss function. In the following section we present an algorithm for ordinal regression based on large margin techniques. In Section 2.5 we give learning curves of our approach in a controlled experiment and in a real—world experiment on data from information retrieval.

### 2.2 Classical Models for Ordinal Regression

In this section we shortly recall the well–known cumulative or threshold model for ordinal regression (McCullagh and Nelder, 1983).

In contrast to Equation (1.2) we assume that there is an outcome space  $\mathcal{Y} = \{r_1, \ldots, r_q\}$  with ordered ranks  $r_q \succ_{\mathcal{Y}} r_{q-1} \succ_{\mathcal{Y}} \cdots \succ_{\mathcal{Y}} r_1$ . The symbol  $\succ_{\mathcal{Y}}$  denotes the ordering between different ranks and can be interpreted as "is preferred to". Since  $\mathcal{Y}$  contains only a finite number of ranks,  $P(y = r_i | \mathbf{x})$  is a multinomial distribution.

stochastic ordering Let us make the assumption of stochastic ordering of the related space  $\mathcal{X}$ , i.e., for all different  $\mathbf{x}_1$  and  $\mathbf{x}_2$  either

$$\Pr(y \le r_i | \mathbf{x}_1) \ge \Pr(y \le r_i | \mathbf{x}_2)$$
 for all  $r_i \in \mathcal{Y}$ , (2.1)

or

$$\Pr(y \le r_i | \mathbf{x}_1) \le \Pr(y \le r_i | \mathbf{x}_2)$$
 for all  $r_i \in \mathcal{Y}$ . (2.2)

Stochastic ordering is satisfied by a model of the form

$$l^{-1}(\Pr(y < r_i | \mathbf{x})) = \theta(r_i) - (\mathbf{w} \cdot \mathbf{x}), \tag{2.3}$$

where  $l^{-1}:[0,1] \mapsto (-\infty,+\infty)$  is a monotonic function often referred to as the inverse link function and  $\theta: \mathcal{Y} \mapsto \mathbb{R}$  is increasing for increasing ranks. The stochastic ordering follows from the fact that

$$\Pr(y \le r_i | \mathbf{x}_1) \ge \Pr(y \le r_i | \mathbf{x}_2) \Leftrightarrow \Pr(y \le r_i | \mathbf{x}_1) - \Pr(y \le r_i | \mathbf{x}_2) \ge 0$$
$$\Leftrightarrow l^{-1}(\Pr(y \le r_i | \mathbf{x}_1)) - l^{-1}(\Pr(y \le r_i | \mathbf{x}_2)) \ge 0$$
$$\Leftrightarrow (\mathbf{w} \cdot (\mathbf{x}_2 - \mathbf{x}_1)) > 0,$$

which no longer depends on  $r_i$  (the same applies to  $\Pr(y \le r_i | \mathbf{x}_1) \le \Pr(y \le r_i | \mathbf{x}_2)$ ).

model	inverse link function $P_{\epsilon}^{-1}(\Delta)$	distribution $P_{\epsilon}(\epsilon = \eta)$
logit	$\ln \frac{\Delta}{1-\Delta}$	$\frac{\exp(\eta)}{(1+\exp(\eta))^2}$
probit	$N^{-1}(\Delta)$	$\frac{(1+\exp(\eta))^2}{\frac{1}{\sqrt{2\pi}}}\exp\left\{-\frac{\eta^2}{2}\right\}$
${\color{red} {\rm complementary\ log-log}}$	$\ln(-\ln(1-\Delta))$	$\exp\{\eta - \exp(\eta)\}$

**Table 2.1** Inverse link functions for different models for ordinal regression (taken from McCullagh and Nelder (1983)). Here,  $N^{-1}$  denotes the inverse normal function.

cumulative model

Such a model is called a cumulative or threshold model and can be motivated by the following argument: Let us assume that the ordinal response is a coarsely measured *latent* continuous variable  $U(\mathbf{x})$ . Thus, we observe rank  $r_i$  in the training set iff

$$y = r_i \Leftrightarrow U(\mathbf{x}) \in [\theta(r_{i-1}), \theta(r_i)], \tag{2.4}$$

where the function U (latent utility) and  $\boldsymbol{\theta} = (\theta(r_0), \dots, \theta(r_q))^T$  are to be determined from the data. By definition  $\theta(r_0) = -\infty$  and  $\theta(r_q) = +\infty$ . We see that the real line is divided into q consecutive intervals, where each interval corresponds to a rank  $r_i$ . Let us make a linear model of the latent variable  $U(\mathbf{x})$ 

linear utility model

$$U(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + \epsilon, \tag{2.5}$$

where  $\epsilon$  is the random component of zero expectation,  $\mathbf{E}_{\epsilon}(\epsilon) = 0$ , and distributed according to  $P_{\epsilon}$ . It follows from Equation (2.4) that

$$\Pr(y \le r_i | \mathbf{x}) = \sum_{j=1}^{i} \Pr(y = r_j | \mathbf{x}) = \sum_{j=1}^{i} \Pr(U(\mathbf{x}) \in [\theta(r_{j-1}), \theta(r_j)])$$

$$= \Pr(U(\mathbf{x}) \in [-\infty, \theta(r_i)]) = \Pr((\mathbf{w} \cdot \mathbf{x}) + \epsilon \le \theta(r_i))$$

$$= P(\epsilon \le \underbrace{\theta(r_i) - (\mathbf{w} \cdot \mathbf{x})}_{\eta}) = P_{\epsilon}(\theta(r_i) - (\mathbf{w} \cdot \mathbf{x})).$$

If we now make a distributional assumption  $P_{\epsilon}$  for  $\epsilon$  we obtain the cumulative model by choosing as the inverse link function  $l^{-1}$  the inverse distribution function  $P_{\epsilon}^{-1}$  (quantile function). Note that each quantile function  $P_{\epsilon}^{-1}:[0,1]\mapsto (-\infty,+\infty)$  is a monotonic function. Different distributional assumptions for  $\epsilon$  yield the logit, probit, or complementary log-log model (see Table 2.1).

In order to estimate **w** and  $\theta$  from model (2.3), for the observation ( $\mathbf{x}_i, y$ ) we see

$$\underbrace{\begin{pmatrix}
o_1(\mathbf{x}_i) \\
o_2(\mathbf{x}_i) \\
\vdots \\
o_{q-2}(\mathbf{x}_i) \\
o_{q-1}(\mathbf{x}_i)
\end{pmatrix}}_{\mathbf{o}(\mathbf{x}_i)} = \underbrace{\begin{pmatrix}
-\mathbf{x}_i & 1 & 0 & \cdots & 0 & 0 \\
-\mathbf{x}_i & 0 & 1 & \cdots & 0 & 0 \\
\vdots & \vdots & \vdots & \ddots & \vdots \\
-\mathbf{x}_i & 0 & 0 & \cdots & 1 & 0 \\
-\mathbf{x}_i & 0 & 0 & \cdots & 0 & 1
\end{pmatrix}}_{\mathbf{Z}(\mathbf{x}_i)} \underbrace{\begin{pmatrix}
\mathbf{w} \\
\theta(r_1) \\
\theta(r_2) \\
\vdots \\
\theta(r_{q-2}) \\
\theta(r_{q-1})
\end{pmatrix}}_{\mathbf{w}_{GLM}},$$

design matrix

where  $o_j(\mathbf{x}_i) = P_{\epsilon}^{-1}(\Pr(y \leq r_j|\mathbf{x}_i))$  is the transformed probability of ranks less than or equal to  $r_j$  given  $\mathbf{x}_i$ , which will be estimated from the sample by the transformed frequencies of that event. Note that the complexity of the model is determined by the linearity assumption (2.5) and by  $P_{\epsilon}^{-1}$  which can be thought of as a regularizer in the resulting likelihood equation. For the complete training set we obtain

$$\underbrace{\begin{pmatrix} \mathbf{o}(\mathbf{x}_{1}) \\ \vdots \\ \mathbf{o}(\mathbf{x}_{\ell}) \end{pmatrix}}_{l^{-1}(\mathbf{y}) \text{ (random)}} = \underbrace{\begin{pmatrix} \mathbf{Z}(\mathbf{x}_{1}) & \dots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \dots & \mathbf{Z}(\mathbf{x}_{\ell}) \end{pmatrix}}_{\mathbf{Z} \text{ (random)}} \underbrace{\begin{pmatrix} \mathbf{w}_{GLM} \\ \vdots \\ \mathbf{w}_{GLM} \end{pmatrix}}_{\mathbf{W}_{GLM} \text{ (parameters)}}.$$
(2.6)

The last equation is called the design matrix of a multivariate generalized linear model (GLM). A generalized linear model  $\mathbf{y} = l(\mathbf{Z}\mathbf{W}_{\text{GLM}})$  is mainly determined by the design matrix  $\mathbf{Z}$  and the link function  $l(\cdot) = P_{\epsilon}(\cdot)$ . Then given a sample (X,Y) and a link function — which coincides with a distributional assumption about the data — methods for calculating the maximum likelihood estimate  $\mathbf{W}_{\text{GLM}}$  exist (see McCullagh and Nelder (1983) or Fahrmeir and Tutz (1994) for a detailed discussion). The main difficulty in maximizing the likelihood is introduced by the nonlinear link function.

maximum likelihood estimate

To conclude this review of classical statistical methods we want to highlight the two main assumptions made for ordinal regression: (i) the assumption of stochastic ordering of the space  $\mathcal{X}$  (ii) and a distributional assumption on the unobservable latent variable.

#### 2.3 A Risk Formulation for Ordinal Regression

Instead of the distributional assumptions made in the last section, we now consider a parameterized model space G of mappings from objects to ranks. Each such function g induces an ordering  $\succ_{\mathcal{X}}$  on the elements of the input space by the following rule

$$\mathbf{x}_i \succ_{\mathcal{X}} \mathbf{x}_i \Leftrightarrow g(\mathbf{x}_i) \succ_{\mathcal{V}} g(\mathbf{x}_i) \,. \tag{2.7}$$

If we neglect the ordering of the space  $\mathcal{Y}$ , it was already shown in Section 1.1.1 that the Bayes–optimal function  $g_{\text{class}}^*$  given by Equation (1.5) is known to minimize

$$R_{\text{class}}(g) = \mathbf{E}_{\mathbf{x},y} \left( \mathbf{1}_{g(\mathbf{x}) \neq y} \right) = \mathbf{E}_{\mathbf{x},y} \left( c_{\text{class}}(\mathbf{x}, y, g(\mathbf{x})) \right). \tag{2.8}$$

Let us rewrite  $R_{\rm class}(g)$  by

$$R_{\text{class}}(g) = \int_{\mathcal{X}} \mathcal{Q}_{\text{class}}(\mathbf{x}, g) \ p(\mathbf{x}) d\mathbf{x},$$
where
$$\mathcal{Q}_{\text{class}}(\mathbf{x}, g) = \sum_{i=1}^{q} \Pr(r_i | \mathbf{x}) - \Pr(g(\mathbf{x}) | \mathbf{x}) = 1 - \Pr(g(\mathbf{x}) | \mathbf{x}).$$
(2.9)

A closer look at Equation (2.9) shows that a sufficient condition for two mappings  $g_1$  and  $g_2$  to incur equal risks  $R_{\text{class}}(g_1)$  and  $R_{\text{class}}(g_2)$  is given by  $\Pr(g_1(\mathbf{x})|\mathbf{x}) = \Pr(g_2(\mathbf{x})|\mathbf{x})$  for every  $\mathbf{x}$ . Assuming that  $\Pr(r_i|\mathbf{x})$  is one for every  $\mathbf{x}$  at a certain rank  $r_k$  the risks are equal — independently of how "far away" (in terms of rank difference) the mappings  $g_1(\mathbf{x})$  and  $g_2(\mathbf{x})$  are from the optimal rank argmax  $r_i \in \mathcal{Y}$   $\Pr(r_i|\mathbf{x})$ . This evidently shows that  $c_{\text{class}}$  is inappropriate for the case where a natural ordering is defined on the elements of  $\mathcal{Y}$ .

Since the only available information given by the ranks is the induced ordering of the input space  $\mathcal{X}$  (see Equation (2.7)) we argue that a distribution independent model of ordinal regression has to single out that function  $g_{\text{pref}}^*$  which induces the ordering of the space  $\mathcal{X}$  that incurs the smallest number of inversions on pairs  $(\mathbf{x}_1, \mathbf{x}_2)$  of objects (for a similar reasoning see Sobel (1993)). To model this property we note that due to the ordering of the space  $\mathcal{Y}$ , each mapping g induces an ordering on the space  $\mathcal{X}$  by Equation (2.7). Let use define the rank difference  $\Theta: \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{Z}$  by

$$r_i \ominus r_j := i - j. \tag{2.10}$$

Now given a pair  $(\mathbf{x}_1, y_1)$  and  $(\mathbf{x}_2, y_2)$  of objects we distinguish between two different events:  $y_1 \ominus y_2 > 0$  and  $y_1 \ominus y_2 < 0$ . According to Equation (2.7) a function g violates the ordering if  $y_1 \ominus y_2 > 0$  and  $g(\mathbf{x}_1) \ominus g(\mathbf{x}_2) \leq 0$ , or  $y_1 \ominus y_2 < 0$  and  $g(\mathbf{x}_1) \ominus g(\mathbf{x}_2) \geq 0$ . Additionally taking into account that each weak order  $\succ_{\mathcal{Y}}$  induces an equivalence  $\sim_{\mathcal{Y}}$  (Fishburn, 1985) the case  $y_1 \ominus y_2 = 0$  is automatically taken care of. Thus, an appropriate loss function is given by

loss function for ordinal regression

$$c_{\text{pref}}(\mathbf{x}_{1}, \mathbf{x}_{2}, y_{1}, y_{2}, g(\mathbf{x}_{1}), g(\mathbf{x}_{2})) = \begin{cases} 1 & y_{1} \ominus y_{2} > 0 \land g(\mathbf{x}_{1}) \ominus g(\mathbf{x}_{2}) \leq 0 \\ 1 & y_{2} \ominus y_{1} > 0 \land g(\mathbf{x}_{2}) \ominus g(\mathbf{x}_{1}) \leq 0 \end{cases}$$
(2.11)  
$$0 \quad \text{else}$$

Note, that we can obtain  $m^2$  samples drawn according to  $p(\mathbf{x}_1, \mathbf{x}_2, y_1, y_2)$ . It is important that these samples do not provide  $m^2$  iid samples of the function  $c_{\text{pref}}(\mathbf{x}_1, \mathbf{x}_2, y_1, y_2, g(\mathbf{x}_1), g(\mathbf{x}_2))$  for any g. Furthermore, if we define

$$c_g(\mathbf{x}_1, y_1, g(\mathbf{x}_1)) = \mathbf{E}_{\mathbf{x}, y} \left[ c_{\text{pref}}(\mathbf{x}_1, \mathbf{x}, y_1, y, g(\mathbf{x}_1), g(\mathbf{x})) \right], \tag{2.12}$$

risk functional for ordinal regression the risk functional to be minimized is given by

$$R_{\text{pref}}(g) = \mathbf{E}_{\mathbf{x}_1, y_1, \mathbf{x}_2, y_2} \left( c_{\text{pref}}(\mathbf{x}_1, \mathbf{x}_2, y_1, y_2, g(\mathbf{x}_1), g(\mathbf{x}_2)) \right)$$

$$= \mathbf{E}_{\mathbf{x}_1, y_1} \left( c_g(\mathbf{x}_1, y_1, g(\mathbf{x}_1)) \right). \tag{2.13}$$

Although Equation (2.13) shows great similarity to the classification learning risk functional (2.8) we see that due to the loss function  $c_g$ , which exploits the ordinal nature of  $\mathcal{Y}$ , we have a different pointwise loss function for each g. Thus we have found a risk functional which can be used for ordinal regression and takes into account the ordering as proposed by McCullagh and Nelder (1983).

In order to relate  $R_{\text{pref}}(g)$  to a simple classification risk we slightly redefine the empirical risk based on  $c_{\text{pref}}$  and the training data (X,Y). For notational simplification let us define the space  $\mathcal{E}$  of events of pairs  $\mathbf{x}$  and y with unequal ranks by

$$\mathcal{E} := \{ (\mathbf{z}, t) \mid \mathbf{z} = (\mathbf{x}_i, \mathbf{x}_i) \in \mathcal{X} \times \mathcal{X}, t = \Omega(y_k, y_l), y_k \in \mathcal{Y}, y_l \in \mathcal{Y}, |y_k \ominus y_l| > 0 \}$$

Furthermore, using the shorthand notation  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$  to denote the first and second object of a pair a new training set (X', Y') can be derived from (X, Y) if we use all 2–sets in  $\mathcal{E}$  derivable from (X, Y), i.e.

$$\forall 0 < |y_i^{(1)} - y_i^{(2)}| \qquad (X', Y') = \left\{ \left( \left( \mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)} \right), \Omega \left( y_i^{(1)}, y_i^{(2)} \right) \right) \right\}_{i=1}^{m'}$$

$$\Omega(y_1, y_2) := \operatorname{sgn} (y_1 \oplus y_2),$$
(2.14)

where  $\Omega$  is an indicator function for rank differences and m' is the cardinality of (X', Y').

preference learning  $\Leftrightarrow$  classification

#### Theorem 2.1 Equivalence of risk functionals

Assume an unknown probability measure  $p(\mathbf{x}, y)$  on  $\mathcal{X} \times \mathcal{Y}$  is given. Then for each  $q: \mathcal{X} \mapsto \mathcal{Y}$  the following equalities hold true

$$R_{\text{pref}}(g) = \mathbf{E}_{y_1, y_2} \left( |\Omega(y_1, y_2)| \right) \mathbf{E}_{\mathbf{z}, t} \left( c_{\text{class}}(\mathbf{z}, t, \Omega(g(\mathbf{x}_1), g(\mathbf{x}_2))) \right) , \qquad (2.16)$$

$$R_{\text{emp}}(g) = \frac{m'}{m^2} \sum_{i=1}^{m'} c_{\text{class}} \left( \left( \mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)} \right), \Omega\left( y_i^{(1)}, y_i^{(2)} \right), \Omega\left( g\left( \mathbf{x}_i^{(1)} \right), g\left( \mathbf{x}_i^{(2)} \right) \right) \right) .$$

**Proof** Let us derive the probability  $p(\mathbf{z},t)$  on  $\mathcal{E}$  derived from  $p(\mathbf{x}_1,\mathbf{x}_2,y_1,y_2)$ :

$$p(\mathbf{z},t) = \begin{cases} 0 & t=0 \\ p(\mathbf{x}_1, \mathbf{x}_2, y_1, y_2)/\Delta & t \neq 0 \end{cases}$$
,

$$\Delta = \mathbf{E}_{y_1,y_2}(|\Omega(y_1,y_2)|) = \Pr(|y_1 \ominus y_2| > 0).$$

Now exploiting the definition (2.11) of  $c_{\text{pref}}$  we see

$$\forall \mathbf{x}_1, \mathbf{x}_2, y_1, y_2, g : t = c_{\text{pref}}(\mathbf{x}_1, \mathbf{x}_2, y_1, y_2, g(\mathbf{x}_1), g(\mathbf{x}_2))$$
.

The first statement is proven. The second statement follows by setting  $\mathcal{X} = X, \mathcal{Y} = Y$  and assigning constant mass of  $1/m^2$  at each point  $(\mathbf{x}_1, \mathbf{x}_2, y_1, y_2)$ .

Taking into account that each function  $g \in G$  defines a function  $p_g : \mathcal{X} \times \mathcal{X} \mapsto \{-1, 0, +1\}$  by

$$p_g(\mathbf{x}_1, \mathbf{x}_2) := \Omega(g(\mathbf{x}_1), g(\mathbf{x}_2)), \tag{2.17}$$

reduction to classification problem

Theorem 2.1 states that the empirical risk of a certain mapping g on a sample (X,Y) is equivalent to the  $c_{\text{class}}$  loss of the related mapping  $p_g$  on the sample (X',Y') up to a constant factor  $m'/m^2$  which depends neither on g nor on  $p_g$ . Thus, the problem of distribution independent ordinal regression can be reduced to a classification problem on pairs of objects. It is important to emphasize the chain of argument that lead to this equivalence. The original problem was to find a function g that maps objects to ranks given a sample (X,Y). Taking the ordinal nature

of ranks into account leads to the equivalent formulation of finding a function  $p_g$  that maps pairs of objects to the three classes  $\succ_{\mathcal{Y}}$ ,  $\prec_{\mathcal{Y}}$ , and  $\sim_{\mathcal{Y}}$ . Reverting the chain of argumentation may lead to difficulties by observing that only those  $p_g$  are admissible — in the sense that there is a function g that fulfills Equation (2.17) — which define an asymmetric, transitive relation on  $\mathcal{X}$ . Therefore we also call this the problem of preference learning. It was shown that the Bayes optimal decision function given by (1.5) on pairs of objects can result in a function  $p_g$  which is no longer transitive on  $\mathcal{X}$  (Herbrich et al., 1998). This is also known as the problem of stochastic transitivity (Suppes et al., 1989). Note also that the conditions of transitivity and asymmetry effectively reduce the space of admissible classification functions  $p_g$  acting on pairs of objects.

uniform convergence bounds

However, the above formulation is — in the form presented — not amenable to the straightforward application of classical results from learning theory. The reason is that the constructed samples of pairs of objects violate the iid assumption. In order to still be able to give upper bounds on a risk for preference learning we have to reduce our sample such that the resulting realization of the loss (2.11) is distributed iid. Under this condition it is then possible to bound the deviation of the expected risk from the empirical risk. Let  $\sigma$  be any permutation of the numbers  $1, \ldots, m$ . Furthermore, for notational convenience let  $\mathcal{C}_g(i,j)$  abbreviate  $c_{\text{pref}}(\mathbf{x}_i, \mathbf{x}_j, y_i, y_j, g(\mathbf{x}_i), g(\mathbf{x}_j))$ . Then we see that for any  $g \in G$ 

$$\Pr(\mathcal{C}_g(\sigma(1), \sigma(2)), \mathcal{C}_g(\sigma(2), \sigma(3)), \dots, \mathcal{C}_g(\sigma(m-1), \sigma(m))) = \\ \Pr(\mathcal{C}_g(\sigma(1), \sigma(2))) \cdot \Pr(\mathcal{C}_g(\sigma(2), \sigma(3))) \cdot \dots \cdot \Pr(\mathcal{C}_g(\sigma(m-1), \sigma(m))).$$
 (2.18)

Clearly, m-1 is the maximum number of pairs of objects that still fulfil the iid assumption. In order to see this consider that by transitivity the ordering  $g(\mathbf{x}_1) \prec_{\mathcal{Y}} g(\mathbf{x}_2)$  and  $g(\mathbf{x}_2) \prec_{\mathcal{Y}} g(\mathbf{x}_3)$  implies  $g(\mathbf{x}_1) \prec_{\mathcal{Y}} g(\mathbf{x}_3)$  (and vice versa for  $\succ_{\mathcal{Y}}$  and  $\sim_{\mathcal{Y}}$ ). Now we can give the following theorem.

#### Theorem 2.2 A Margin Bound for Ordinal Regression

Let p be a probability measure on  $\mathcal{X} \times \{r_1, \ldots, r_q\}$ , let (X, Y) be a sample of size m drawn iid from p. Let  $\sigma$  be any permutation of the numbers  $1, \ldots, m$ . For each function  $g: \mathcal{X} \mapsto \{r_1, \ldots, r_q\}$  there exists a function  $f \in F$  and a vector  $\boldsymbol{\theta}$  such that 1

$$g(\mathbf{x}) = r_i \Leftrightarrow f(\mathbf{x}) \in [\theta(r_{i-1}), \theta(r_i)]. \tag{2.19}$$

Let the fat–shattering dimension of the set of functions F be bounded above by the function  $\operatorname{afat}_F:\mathbb{R}\to\mathbb{N}$ . Then for each function g with zero training error, i.e.  $\sum_{i=1}^{m-1}Q_g(\sigma(i),\sigma(i+1))=0$  and

$$\rho_f = \min_{i=1,\dots,m-1} \Omega\left(y_{\sigma(i)}, y_{\sigma(i+1)}\right) |f(\mathbf{x}_{\sigma(i)}) - f(\mathbf{x}_{\sigma(i+1)})|$$

<sup>1.</sup> Note the close relationship to the cumulative model presented in Section 2.2.

with probability  $1 - \delta$ 

$$\begin{split} R_{\mathrm{pref}}(g) &\leq \frac{2}{m-1} \left( k \log_2 \left( \frac{8e(m-1)}{k} \right) \log_2(32(m-1)) + \log_2 \left( \frac{8(m-1)}{\delta} \right) \right) \,, \\ \text{where } k &= \mathrm{afat}_F(\rho_f/8) \leq e(m-1). \end{split}$$

**Proof** Let us recall the following theorem based on Theorem 1.5.

#### Theorem 2.3 (Shawe-Taylor et al., 1998)

Consider a real valued function class F having fat shattering function bounded above by the a function  $\operatorname{afat}_F:\mathbb{R}\to\mathbb{N}$  which is continuous from the right. Fix  $\theta\in\mathbb{R}$ . Then with probability  $1-\delta$  a learner that correctly classifies m iid generated examples  $(\mathbf{x}_1,y_1),\ldots,(\mathbf{x}_m,y_m)$  with  $h=T_\theta(f)\in T_\theta(F)$  such that  $h(\mathbf{x}_i)=y_i, i=1,\ldots,m$  and  $\rho_f=\min_i y_i (|f(\mathbf{x}_i)-\theta|)$  will have error of h bounded from above by

$$\frac{2}{m} \left( k \log_2 \left( \frac{8em}{k} \right) \log_2(32m) + \log_2 \left( \frac{8m}{\delta} \right) \right) , \tag{2.20}$$

where  $k = \operatorname{afat}_F(\rho_f/8) \leq em$ .

Taking into account that by construction we got m-1 iid examples and that the classification of a pair is carried out by a decision based on the difference  $f(\mathbf{x}_{\sigma(i)}) - f(\mathbf{x}_{\sigma(i+1)})$  we can upper bound  $R_{\text{pref}}(g)$  by replacing each m with m-1 and using  $\theta = 0$ .

The afat<sub>F</sub>( $\rho$ )—shattering dimension of F can be thought of as the maximum number of objects that can be arranged in any order using functions from F and a minimum margin  $\min \Omega(y_1,y_2)|f(\mathbf{x}_1)-f(\mathbf{x}_2)|$  of  $\rho$  (utilizing Equation (2.7) together with (2.19)). Note, that the zero training error condition for the above bound is automatically satisfied for any  $\sigma$  if  $R_{\rm emp}(g)=0$ . Even though this empirical risk was not based on an iid sample its minimization allows the application of the above bound. In the following section we will present an algorithm which aims at minimizing exactly that empirical risk while at the same time enforcing large margin rank boundaries.

#### 2.4 An Algorithm for Ordinal Regression

Based on the results of Theorem 2.2 we suggest to model ranks as intervals on the real line. In equivalence with the classical cumulative model used in ordinal regression, let us introduce a (latent) linear function  $f: \mathcal{X} \to \mathbb{R}$  for each function g

$$f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}), \tag{2.21}$$

which are related by (2.19). In order to apply the given theorem we see that we have to find a function  $f^*$  which incurs no training error on (X', Y') while controlling the generalization error by maximizing the margin  $\rho_f$ . Note, that

ranks as intervals on the real line

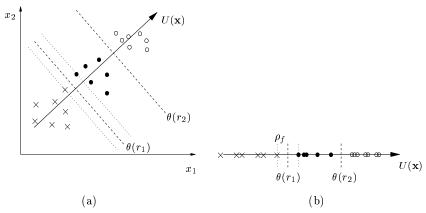


Figure 2.1 (a) Mapping of objects from rank  $r_1$  (×), rank  $r_2$  (•), and rank  $r_3$  (o) to the axis  $f(\mathbf{x})$ , where  $\mathbf{x} = (x_1, x_2)^T$ . Note that by  $\theta(r_1)$  and  $\theta(r_2)$  two coupled hyperplanes are defined. (b) The margin of the coupled hyperplanes  $\rho_f = \min_{(X',Y')} \Omega(y_i^{(1)}, y_i^{(2)}) |f(\mathbf{x}_i^{(1)}) - f(\mathbf{x}_i^{(2)})|$  is this time defined at the rank boundaries  $\theta(r_i)$ .

$$f(\mathbf{x}_i) - f(\mathbf{x}_i) = (\mathbf{w} \cdot (\mathbf{x}_i - \mathbf{x}_i))$$

which makes apparent that each pair  $(\mathbf{x}_i, \mathbf{x}_j) \in X'$  is represented by its difference vector  $(\mathbf{x}_i - \mathbf{x}_j)$  assuming a linear model of f. This allows the straightforward application of the large margin algorithm given by Equation (1.51) and (1.52) replacing each  $\mathbf{x}_i$  by  $(\mathbf{x}_i^{(1)} - \mathbf{x}_i^{(2)})$ . Hence, the maximization of the margin takes place at the rank boundaries  $\theta(r_i)$  (see Equation (2.19) and Figure 2.1). In practice it is preferable to use the soft margin extension of the large margin algorithm (see Equation (1.25)). Furthermore due to the KKT conditions (see Equation (1.54))  $\mathbf{w}^*$  can be written in terms of the training data. This gives

$$\mathbf{w}^* = \sum_{i=1}^{m'} \alpha_i^* t_i \left( \mathbf{x}_i^{(1)} - \mathbf{x}_i^{(2)} \right) , \qquad (2.22)$$

soft margin

where  $\alpha^*$  is given by

$$\boldsymbol{\alpha}^* = \underset{\substack{\sigma : 1 \geq \boldsymbol{\alpha} \geq \mathbf{o} \\ (\boldsymbol{\alpha} \cdot \mathbf{t}) = 0}}{\operatorname{argmax}} \left[ \sum_{i=1}^{m'} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{m'} \alpha_i \alpha_j t_i t_j \left( (\mathbf{x}_i^{(1)} - \mathbf{x}_i^{(2)}) \cdot (\mathbf{x}_j^{(1)} - \mathbf{x}_j^{(2)}) \right) \right], \quad (2.23)$$

and  $\mathbf{t} = (\Omega(y_1^{(1)}, y_1^{(2)}), \dots, \Omega(y_{m'}^{(1)}, y_{m'}^{(2)}))$ . Note, however, that due to the expansion of the last term in (2.23),

$$\left( (\mathbf{x}_i^{(1)} - \mathbf{x}_i^{(2)}) \cdot (\mathbf{x}_j^{(1)} - \mathbf{x}_j^{(2)}) \right) = (\mathbf{x}_i^{(1)} \cdot \mathbf{x}_j^{(1)}) - (\mathbf{x}_i^{(1)} \cdot \mathbf{x}_j^{(2)}) - (\mathbf{x}_i^{(2)} \cdot \mathbf{x}_j^{(1)}) + (\mathbf{x}_i^{(2)} \cdot \mathbf{x}_j^{(2)}) \,,$$

the solution  $\alpha^*$  to this problem can be calculated solely in terms of the inner products between the feature vectors without reference to the feature vectors themselves. Hence, the idea of (implicitly) mapping the data X via a nonlinear

kernel trick

mapping  $\Phi: \mathcal{X} \mapsto \mathcal{F}$  into a feature space  $\mathcal{F}$  can successfully applied (for further details see Section 1.3.2). Replacing each occurrence of  $\mathbf{x}$  by  $\Phi(\mathbf{x})$  gives

$$\boldsymbol{\alpha}^* = \underset{\substack{\alpha 1 \geq \boldsymbol{\alpha} \geq \mathbf{0} \\ (\boldsymbol{\alpha}^*, \mathbf{t}) = 0}}{\operatorname{argmax}} \left[ \sum_{i=1}^t \alpha_i - \frac{1}{2} \sum_{i,j=1}^t \alpha_i \alpha_j t_i t_j \mathcal{K}\left(\mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)}, \mathbf{x}_j^{(1)}, \mathbf{x}_j^{(2)}\right) \right]. \tag{2.24}$$

where K is for a given function k defined by

$$\mathcal{K}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4) = k(\mathbf{x}_1, \mathbf{x}_3) - k(\mathbf{x}_1, \mathbf{x}_4) - k(\mathbf{x}_2, \mathbf{x}_3) + k(\mathbf{x}_2, \mathbf{x}_4). \tag{2.25}$$

Here,  $k: \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}$  is a Mercer kernel and for a fixed mapping  $\Phi$  is defined by

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

Some kernels k to be used in learning are given by Equations (1.63) and (1.73). Note that the usage of kernels instead of explicitly performing the mapping  $\Phi$  allows us to deal with nonlinear functions f without running into computational difficulties. Moreover, as stated in Theorem 2.2 the bound on the risk  $R_{\text{pref}}(\mathbf{w})$  does not depend on the dimension of  $\mathcal{F}$  but on the margin  $\rho_f$ .

In order to estimate the rank boundaries we note that due to Equations (1.52) the difference in  $f^*$  is greater or equal to one for all training examples which constitute a correctly classified pair. These can easily be obtained by checking  $0 < \alpha_i^* < C$ , i.e. training patterns which do not meet the box constraint (see Section 1.1.4). Thus if  $\Theta(k) \subset (X',Y')$  is the fraction of objects from the training set with  $0 < \alpha_i^* < C$  and rank difference  $\ominus$  exactly one starting from rank  $r_k$ , i.e.

rank boundaries

$$\Theta(k) = \left\{ \left( \mathbf{x}_i^{(1)}, \mathbf{x}_i^{(2)} \right) \middle| y_i^{(1)} = r_k \land y_i^{(2)} = r_{k+1} \land 0 < \alpha_i^* < C \right\}$$
 (2.26)

then the estimation of  $\theta(r_k)$  is given by

$$\theta^*(r_k) = \frac{f^*(\mathbf{x}_1) + f^*(\mathbf{x}_2)}{2}, \tag{2.27}$$

where

$$(\mathbf{x}_1, \mathbf{x}_2) = \operatorname*{argmin}_{(\mathbf{x}_i, \mathbf{x}_j) \in \Theta(k)} [f^*(\mathbf{x}_i) - f^*(\mathbf{x}_j)]. \tag{2.28}$$

In other words, the optimal threshold  $\theta^*(r_k)$  for rank  $r_k$  lies in the middle of the utilities of the closest (in the sense of their utility) objects of rank  $r_k$  and  $r_{k+1}$ . After the estimation of the rank boundaries  $\theta(r_k)$  a new object is assigned to a rank according to Equation (2.19).

coupled hyperplanes We want to emphasize that taking the difference vector as a representation of a pair of objects effectively couples all hyperplanes  $f(\mathbf{x}) = \theta(r_k)$  thus resulting in a standard QP problem. Furthermore, the effective coupling is retained if we use general  $\ell_q$ -margins (see Section 1.1.4). It is the reduction of the hypothesis space which makes the presented algorithm suited for the task of ordinal regression. Note, that also the kernel  $\mathcal{K}$  derived from k acts only in  $\mathcal{F}$  and thus avoids considering too large a hypothesis space. All properties are consequences of the modeling of ranks as intervals on the real line and of the prior knowledge of the ordering of  $\mathcal{Y}$ .

#### 2.5 Experimental Results

In this section we present some experimental results for the algorithm presented in Section 2.4. We start by giving results for artificial data which allows us to analyze our algorithm in a controlled setting. Then we give learning curves for an example from the field of information retrieval.

#### 2.5.1 Learning Curves for Ordinal Regression

multi-class SVM and support vector regression In this experiment we want to compare the generalization behavior of our algorithm with the multi-class SVM (Weston and Watkins, 1998) and Support Vector regression (SVR) (c.f. Smola (1998)) — the methods of choice, if one does not pay attention to the ordinal nature of  $\mathcal Y$  and instead treats ranks as classes (classification) or continuous response values (regression estimation). Another reason for choosing those algorithms is their similar regularizer  $\|\mathbf w\|^2$  and hypothesis space F which make them as comparable as possible. We generated 1000 observations  $\mathbf x = (x_1, x_2)$  in the unit square  $[0, 1] \times [0, 1] \subset \mathbb{R}^2$  according to a uniform distribution. We assigned to each observation  $\mathbf x$  a value y according to

$$y = i \Leftrightarrow \underbrace{10((x_1 - 0.5) \cdot (x_2 - 0.5))}_{f(\mathbf{x})} + \epsilon \in [\theta(r_{i-1}), \theta(r_i)], \qquad (2.29)$$

example utility function

where  $\epsilon$  was normally distributed, i.e.  $\epsilon \sim N(0,0.125)$ , and  $\boldsymbol{\theta} = (-\infty,-1,-0.1,0.25,1,+\infty)$  is the vector of predefined thresholds. In Figure 2.2 (a) the points  $\mathbf{x}_i$  which are assigned to a different rank after the addition of the normally distributed quantity  $\epsilon_i$  are shown. If we treat the whole task as a classification problem, we would call them incorrectly classified training examples. The solid lines in Figure 2.2 (a) indicate the "true" rank boundaries  $\boldsymbol{\theta}$  on  $f(\mathbf{x})$ .

comparison to other methods

In order to compare the three different algorithms we randomly drew 100 training samples (X,Y) of training set sizes m ranging from 5 to 45, thereby making sure that at least one representative of each rank was within the drawn training set. Classification with multi-class SVMs was carried out by computing the pairwise  $5 \cdot 4/2 = 10$  hyperplanes. For all algorithms, i.e. multi-class SVMs, SVR, and the algorithm presented in Section 2.4, we chose the kernel  $k(\mathbf{x}_i, \mathbf{x}_j) = ((\mathbf{x}_i \cdot \mathbf{x}_j) + 1)^2$  and a trade-off parameter C = 1000000. In the particular case of Support Vector regression we used a value of  $\varepsilon = 0.5$  for the  $\varepsilon$ -insensitive loss function (see (Vapnik, 1995) for the definition of this loss function) and thresholds  $\boldsymbol{\theta} = (0.5, 1.5, 2.5, 3.5, 4.5)$  to transform real valued predictions into ranks.

learning curves

In order to estimate the risk  $R_{\text{pref}}(g^*)/\mathbf{E}_{y_1,y_2}(|\Omega(y_1,y_2)|)$  from the remaining 995 to 955 data points we averaged over all 100 results for a given training set size. Thus we obtained the three learning curves shown in Figure 2.2 (b). Note that we used the scaled  $R_{\text{pref}}$  — which is larger by a constant factor. It can be seen that the algorithm proposed for ordinal regression generalizes much faster by exploiting the ordinal nature underlying  $\mathcal{Y}$  compared to classification. This can be

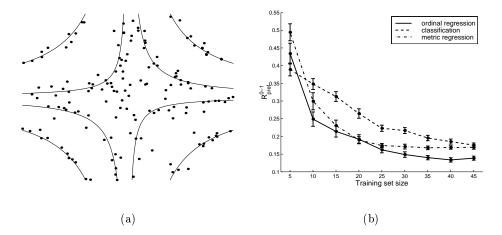


Figure 2.2 (a) Scatter plot of data points  $\mathbf{x}$  which  $f(\mathbf{x})$  maps to a different interval than  $f(\mathbf{x}) + \epsilon$  (see Equation (2.29)). (b) Learning curves for multi-class SVM (dashed lines), SV regression (dashed-dotted line) and the algorithm for ordinal regression (solid line) if we measure  $R_{\text{pref}}$ . The error bars indicate the 95% confidence intervals of the estimated risk  $R_{\text{pref}}$ .

explained by the fact that due to the model of a latent utility all "hyperplanes"  $f(\mathbf{x}) = \theta(r_k)$  are coupled (see Figure 2.1) which does not hold true for the case of multi-class SVMs. Furthermore, the learning curves for SVR and the proposed ordinal regression algorithm are very close which can be explained by the fact that the predefined thresholds  $\theta(r_k)$  are defined in such a way that their pairwise difference is about 0.5 — the size of the  $\varepsilon$ -tube chosen beforehand. Thus the utility and the continuous ranks estimated by the regression algorithm are of the same magnitude which results in the same generalization behavior.

In Figure 2.3 we plotted the assignments of the unit square to ranks  $r_1$  (black areas) to ranks  $r_5$  (white areas) for the functions  $g^*(\mathbf{x})$  learned from randomly drawn training sets ranging from size m=5 (top row) to m=25 (bottom row). We used the same parameters as for the computation of the learning curves. In the rightmost column (e) the true assignment, i.e.,  $y=r_i \Leftrightarrow f(\mathbf{x}) \in [\theta(r_{i-1}), \theta(r_i)]$  is shown. In the first column (a) we can see how the algorithm presented in Section 2.4 performs for varying training set sizes. As expected, for the training set size m=25, the method found a utility function together with a set of thresholds which represent the true ranking very well. The second column (b) shows the results of the abovementioned multi-class SVM on the task. Here the pairwise hyperplanes are not coupled since the ordinal nature of  $\mathcal{Y}$  is not taken into account. This results in a worse generalization, especially in regions, where no training points were given. The third column (c) gives the assignments made by the SVR algorithm if we represent each rank  $r_i$  by i. Similar to the good results seen in the learning curve, the generalization behavior is comparable to the ordinal regression method (first

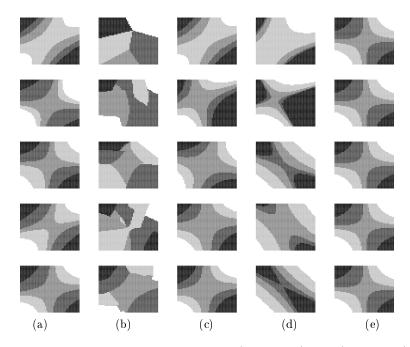


Figure 2.3 Assignments of points to ranks  $r_1$  (black area) to  $r_5$  (white area) by the learned function  $g^*(\mathbf{x})$  based on randomly drawn training samples of size 5, 10, 15, 20, and 25 (top row to bottom row). (a) Results of the algorithm presented in Section 2.4. (b) Results of multi-class SVM if we treat each rank as a class. (c) Results of SVR if we assign rank  $r_i$  to number i. (d) Results of SVR if we assign rank  $r_i$  to real number  $\exp(i)$ . (e) Underlying assignment uncorrupted by noise.

representation of ranks

column). The deficiency of SVR for this task becomes apparent when we change the representation of ranks. In the fourth column (d) we applied the same SVR algorithm, this time on the representation  $\exp(i)$  for rank  $r_i$ . As can be seen, this dramatically changes the generalization behavior of the SVR method. We conclude that the crucial task for application of metric regression estimation methods to the task of ordinal regression is the definition of the representation of ranks. This is automatically — although more time—consuming — solved by the proposed algorithm.

#### 2.5.2 An Application to Information Retrieval

information retrieval In this experiment we make the following assumption: After an initial (textual) query a user makes to an IR system, the system returns a bundle of documents to the user. Now the user assigns ranks to a small fraction of the returned documents and the task for the learning algorithm is to assign ranks to the remaining unranked documents in order to rank the remaining documents. We assume that the quantity of interest is the percentage of inversions incurred by the ranking induced by the

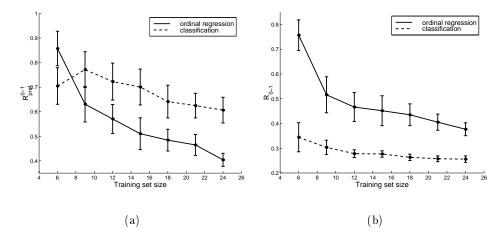


Figure 2.4 Learning curves for multi-class SVM (dashed lines) and the algorithm for ordinal regression (solid line) for the OHSUMED dataset query 1 if we measure (a)  $R_{\text{pref}}$  and (b)  $R_{\text{class}}$ . Error bars indicate the 95% confidence intervals.

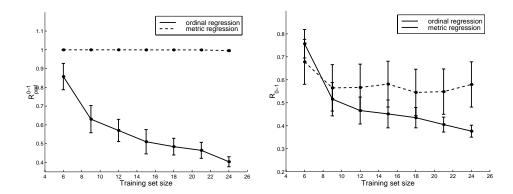
learning algorithm. This quantity is captured by  $R_{\rm emp}(g)/m'$  (m' = |(X', Y')|,see Equation (2.14) for an exact definition) and thus after using m=6 up to m=24 documents and their respective ranking we measure this value on the remaining documents. For this experiment we used the same parameters as in the previous experiment. The investigated dataset was the OHSUMED dataset collected by William Hersh<sup>2</sup>, which consists of 348 566 documents and 106 queries with their respective ranked results. There are three ranks: "document is relevant", "document is partially relevant", and "irrelevant document" wrt. the given textual query. For our experiments we used the results of query 1 ("Are there adverse effects on lipids when progesterone is given with estrogen replacement therapy?") which consists of 107 documents taken from the whole database. In order to apply our algorithm we used the bag-of-words representation (Salton, 1968), i.e., we computed for every document the vector of "term-frequencies-inverse-documentfrequencies" (TFIDF) components. The TFIDF is a weighting scheme for the bagof-words representation which gives higher weights to terms which occur very rarely in all documents. We restricted ourselves to terms that appear at least three times in the whole database. This results in  $\approx 1700$  terms which leads for a certain document to a very high-dimensional but sparse vector. We normalized the length of each document vector to unity (see Joachims (1998)).

bag-of-words representation

Figure 2.4 (a) shows the learning curves for multi-class SVMs and our algorithm for ordinal regression measured in terms of the number of incurred inversions. As can be seen from the plot, the proposed algorithm shows very good generalization behavior compared to the algorithm which treats each rank as a separate class.

<sup>2.</sup> This dataset is publicly available at ftp://medir.ohsu.edu/pub/ohsumed/.

Figure 2.4 (b) shows the learning curves for both algorithms if we measure the number of misclassifications — treating the ranks as classes. As expected, the multiclass SVMs perform much better than our algorithm. It is important to note again, that minimizing the zero—one loss  $R_{\rm class}$  does not automatically lead to a minimal number of inversions and thus to an optimal ordering.



**Figure 2.5** Learning curves for SVR (dashed lines) and the algorithm for ordinal regression (solid line) for the OHSUMED dataset query 1 if we measure (a)  $R_{\text{pref}}$  and (b)  $R_{0-1}$ . Error bars indicate the 95% confidence intervals.

Figure 2.5 (a) shows the learning curves for SVR and for our algorithm for ordinal regression, measured the number of incurred inversions. While the former performs quite well on the artificial dataset, in the real world dataset the SVR algorithm fails to find a ranking which minimizes the number of inversions. This can be explained by fact that for the real—world example the equidistance in the assumed utility may no longer hold — especially taking into account that the data space is very sparse for this type of problem. Similarly, Figure 2.5 (b) shows the learning curves for both algorithms if we measure the number of misclassifications. As expected from the curves on the right the SVR algorithm is worse even on that measure. Note that the SVR algorithm minimizes neither  $R_{\rm pref}$  nor  $R_{0-1}$  which may explain its bad generalization behavior. Also note that we made no adaptation of the parameter  $\varepsilon$  — the size of the tube. The reason is that in this particular task there would not be enough training examples available to set aside a reasonable portion of them for validation purposes.

## 2.6 Discussion and Conclusion

In this chapter we considered the task of ordinal regression which is mainly characterized by the ordinal nature of the outcome space  $\mathcal{Y}$ . All known approaches to this problem (see Section 2.2) make distributional assumptions on an underlying

continuous random variable. In contrast, we proposed a loss function which allows for application of distribution independent methods to solve ordinal regression problems. By exploiting the fact that the induced loss function class is a set of indicator functions we could give a distribution independent bound on our proposed risk. Moreover, we could show that to each ordinal regression problem there exists a corresponding preference learning problem on pairs of objects. This result built the link between ordinal regression and classification methods — this time on pairs of objects. For the representation of ranks by intervals on the real line, we could give margin bounds on our proposed risk — this time applied at the rank boundaries. Based on this result we presented an algorithm which is very similar to the well known Support Vector algorithm but effectively couples the hyperplanes used for rank determination.

learning of equivalence relation

Noting that our presented loss involves pairs of objects we see that the problem of multi-class classification can also be reformulated on pairs of objects which leads to the problem of learning an equivalence relation. Usually, in order to extend a binary classification method to multiple classes, one—against—one or one—against—all techniques are devised (Hastie and Tibshirani, 1996; Weston and Watkins, 1998). Such techniques increase the size of the hypothesis space quadratically or linearly in the number of classes, respectively. Recent work (Phillips, 1998) has shown that learning equivalence relations can increase the generalization behavior of binary—class methods when extended to multiple classes.

Further investigations will include the following question: Does the application of the GLM methods presented in Section 2.2 lead automatically to large margins (see Theorem 2.2)? The answer to such a question would finally close the gap between methods extensively used in the past to theories developed currently in the field of Machine Learning.

#### Acknowledgments

First of all we are indebted to our collaborator Peter Bollmann-Sdorra who first stimulated research on this topic. We also thank Nello Cristianini, Ulrich Kockelkorn, Gerhard Tutz, and Jason Weston for fruitful discussions. Finally, we would like to thank our anonymous reviewers for very helpful comments on the uniform convergence results.

# 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.
- J.A. Anderson. Regression and ordered categorical variables (with discussion).

  Journal of the Royal Statistical Society Series B, 46:1–30, 1984.
- J.A. Anderson and P.R. Philips. Regression, discriminantion and measurement models for ordered categorial variables. *Applied Statistics*, 30:22–31, 1981.
- 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.
- Aipore de Moraes and Ian R. Dunsmore. Predictive comparisons in ordinal models. Communications in Statistics Theory and Methods, 24(8):2145–2164, 1995.
- 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.
- Ludwig Fahrmeir and Gerhard Tutz. Multivariate Statistical Modelling Based on Generalized Linear Models. Springer-Verlag, 1994.
- Peter C. Fishburn. Interval Orders and Interval Graphs. Jon Wiley and Sons, 1985.
- 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.
- L. Gurvits. A note on a scale-sensitive dimension of linear bounded functionals in banach spaces. In *Proceedings of Algorithm Learning Theory*, ALT-97, pages 352–363. Springer Verlag, 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.
- Trevor Hastie and Robert Tibshirani. Classification by pairwise coupling. Technical report, Department of Public Health Sciences and Statistics, University of Toronto, Canada, 1996.
- Ralf Herbrich, Thore Graepel, Peter Bollmann-Sdorra, and Klaus Obermayer. Learning a preference relation in IR. In *Proceedings Workshop Text Categorization and Machine Learning*, International Conference on Machine Learning 1998, pages 80–84, 1998.
- Ralf Herbrich, Max Keilbach, Thore Graepel, Peter Bollmann-Sdorra, and Klaus Obermayer. Neural networks in economics: Background, applications, and new developments. *Advances in Computational Economics*, 11:169–196, 1999.
- W. Hoeffding. Probability inequalities for sums of bounded random variables. Journal of the American Statistical Association, 58:13–30, 1963.
- 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. Text categorization with support vector machines. In European Conference on Machine Learning (ECML), 1998.
- 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.

- Robert W. Keener and Donald M. Waldman. Maximum liklihood regression of rank-censored data. *Journal of the American Statistical Association*, 80:385–392, 1985.
- 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. 2<sup>nd</sup> Berkeley Symposium on Mathematical Statistics and Probabilistics, pages 481–492, Berkeley, 1951. University of California Press.
- P.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.
- P. McCullagh and J. A. Nelder. Generalized Linear Models. Chapman & Hall, London, 1983.
- Peter McCullagh. Regression models for ordinal data (with discussion). *Journal of the Royal Statistical Society Series B*, 42:109–142, 1980.
- 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.
- P. Jonathon Phillips. Support Vector Machines applied to face recognition. In *Proceedings of the Neural Information Processing Conference*, Denver, USA, 1998. in press.
- 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.

51

- 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.
- S. Saitoh. Theory of Reproducing Kernels and its Applications. Longman Scientific & Technical, Harlow, England, 1988.
- Gerard Salton. Automatic Information Organization and Retrieval. McGraw-Hill, New York, 1968.
- 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.rhbnc.ac.uk/research/compint/areas/comp\_learn/sv/pub/ report98-03.ps; SVM available at http://svm.dcs.rhbnc.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 Techology, 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 Informa*tion Theory, 44(5):1926–1940, 1998.
- 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.
- Marc Sobel. Bayes and empirical bayes procedures for comparing parameters. Journal of the American Statistical Association, 88:687–693, 1993.
- Patrick Suppes, David H. Krantz, R. Duncan Luce, and Amos Tversky. Foundations of Measurement Vol. II. Academic Press Inc., San Diego, 1989.
- M. Talagrand. Sharper bounds for gaussian and empirical processes. Annals of Probability, 22:28–76, 1994.

REFERENCES

Andranik Tangian and Josef Gruber. Constructing quadratic and polynomial objective functions. In *Proceedings of the 3rd International Conference on Econometric Decision Models*, pages 166–194, Schwerte, Germany, 1995. Springer.

53

- 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 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.
- J. Weston and C. Watkins. Multi-class support vector machines. Technical Report CSD-TR-98-04, Department of Computer Science, Royal Holloway, University of London, Egham, TW20 0EX, UK, 1998.
- C. K. I. Williams. Prediction with gaussian processes: From linear regression to linear prediction and beyond. In M. I. Jordan, editor, *Learning and Infer*ence 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.
- S. K. M. Wong, Y. Y. Yao, and Peter Bollmann. Linear structure in information retrieval. In *Proceedings of the 11th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 219–232, 1988.