

ISIS Technical Report

Support Vector Machines for Classification and Regression

Steve Gunn

14 May 1998



Image Speech & Intelligent Systems Group
University of Southampton

Contents

1	Introduction	5
1.1	Statistical Learning Theory	6
1.1.1	VC Dimension	7
1.1.2	Structural Risk Minimisation	7
2	Support Vector Classification	9
2.1	The Optimal Separating Hyperplane	9
2.1.1	Linearly Separable Example	13
2.2	The Generalised Optimal Separating Hyperplane	14
2.2.1	Linearly Non-Separable Example	16
2.3	Generalisation in High Dimensional Feature Space	17
2.3.1	Polynomial Mapping Example	19
2.4	Discussion	19
3	Feature Space	21
3.1	Kernel Functions	21
3.1.1	Polynomial	21
3.1.2	Gaussian Radial Basis Function	21
3.1.3	Exponential Radial Basis Function	22
3.1.4	Multi-Layer Perceptron	22
3.1.5	Fourier Series	22
3.1.6	Splines	22
3.1.7	<i>B</i> splines	23
3.1.8	Additive Kernels	23
3.1.9	Tensor Product Kernels	23
3.2	Implicit vs. Explicit Bias	23
3.3	Data Normalisation	24
3.4	Kernel Selection	24
4	Classification Example: IRIS data	25
4.1	Applications	28
5	Support Vector Regression	31
5.1	Linear Regression	31
	-Insensitive Loss Function	32
5.1.2	Quadratic Loss Function	33
5.1.3	Huber Loss Function	33
5.1.4	Example	34
5.2	Non Linear Regression	35
5.2.1	Examples	35
5.2.2	Comments	38
6	Regression Example: Titanium Data	39
6.1	Applications	42

7 Conclusions	43
References	45
Appendix - Implementation Issues	47
Support Vector Classification.....	47
Support Vector Regression.....	50
MATLAB SVM Toolbox	52

1 Introduction

The problem of empirical data modelling is germane to many engineering applications. In empirical data modelling a process of induction is used to build up a model of the system, from which it is hoped to deduce responses of the system that have yet to be observed. Ultimately the quantity and quality of the observations govern the performance of this empirical model. By its observational nature data obtained is finite and sampled; typically this sampling is non-uniform and due to the high dimensional nature of the problem the data will form only a sparse distribution in the input space. Consequently the problem is nearly always ill posed [16] in the sense of Hadamard [9].

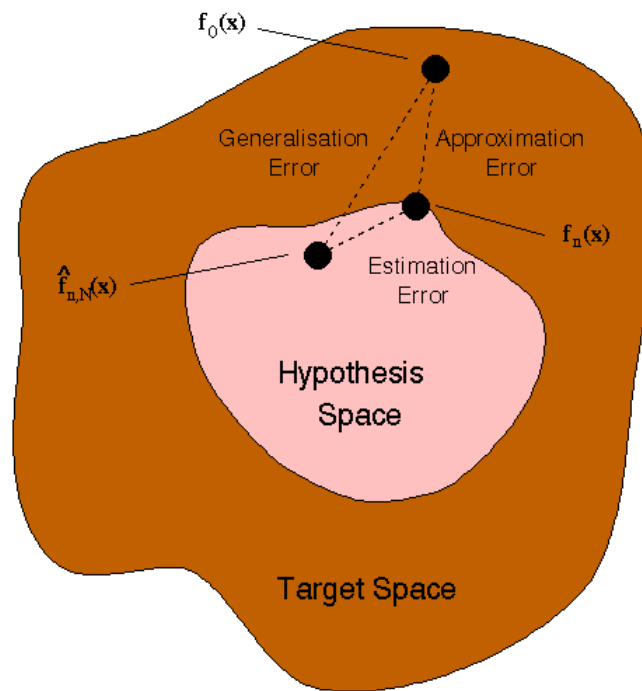
Traditional neural network approaches have suffered difficulties with generalisation, producing models that can overfit the data. This is a consequence of the optimisation algorithms used for parameter selection and the statistical measures used to select the 'best' model.

The foundations of Support Vector Machines (SVM) have been developed by Vapnik [21] and are gaining popularity due to many attractive features, and promising empirical performance. The formulation embodies the Structural Risk Minimisation (SRM) principle, which has been shown to be superior, [8], to traditional Empirical Risk Minimisation (ERM) principle, employed by conventional neural networks. SRM minimises an upper bound on the VC dimension ('generalisation error'), as opposed to ERM that minimises the error on the training data. It is this difference which equips SVM with a greater ability to generalise, which is the goal in statistical learning. SVM were developed to solve the classification problem, but recently they have been extended to the domain of regression problems [20].

In the literature the terminology for SVM can be slightly confusing. The term SVM is typically used to describe classification with support vector methods and support vector regression is used to describe regression with support vector methods. In this report the term SVM will refer to both classification and regression methods, and the terms Support Vector Classification (SVC) and Support Vector Regression (SVR) will be used for specification.

This section continues with an introduction to the structural risk minimisation principle. In section 2 the SVM is introduced in the setting of classification, being both historical and more accessible. This leads onto mapping the input into a higher dimensional feature space by a suitable choice of kernel function. The report then considers the problem of regression. Illustrative examples are given to show the properties of the techniques.

1.1 Statistical Learning Theory



Generalisation error = Estimation error + Approximation error

Ultimately we would like to find the function, f , which minimises the risk,

$$R[f] = \int_{\mathbf{x} \times \mathbf{Y}} (y - f(\mathbf{x}))^2 P(\mathbf{x}, y) d\mathbf{x} dy. \quad (1)$$

However, $P(\mathbf{x}, y)$ is unknown. It is possible to find an approximation according to the empirical risk minimisation principle,

$$R_{emp}[f] = \frac{1}{l} \sum_{i=1}^l (y - f(\mathbf{x}_i))^2 \quad (2)$$

which minimises the empirical risk,

$$\hat{f}_{n,l}(\mathbf{x}) = \arg \min_{f \in H_n} R_{emp}[f] \quad (3)$$

Empirical risk minimisation makes sense only if,

$$\lim_{l \rightarrow \infty} R_{emp}[f] = R[f]$$

which is true from the law of large numbers. However, it must also satisfy,

$$\lim_{l \rightarrow \infty} \min_{f \in H_n} R_{emp}[f] = \min_{f \in H_n} R[f]$$

which is only valid when H_n is 'small' enough. This condition is less intuitive and requires that the minima also converge.

The following bound holds with probability $1 - \delta$,

$$R[f] \leq R_{emp}[f] + \sqrt{\frac{h(\ln \frac{2l}{h} + 1) - \ln(\frac{\delta}{4})}{l}}$$

1.1.1 VC Dimension

The VC dimension is a scalar value that measures the capacity of a set of functions.

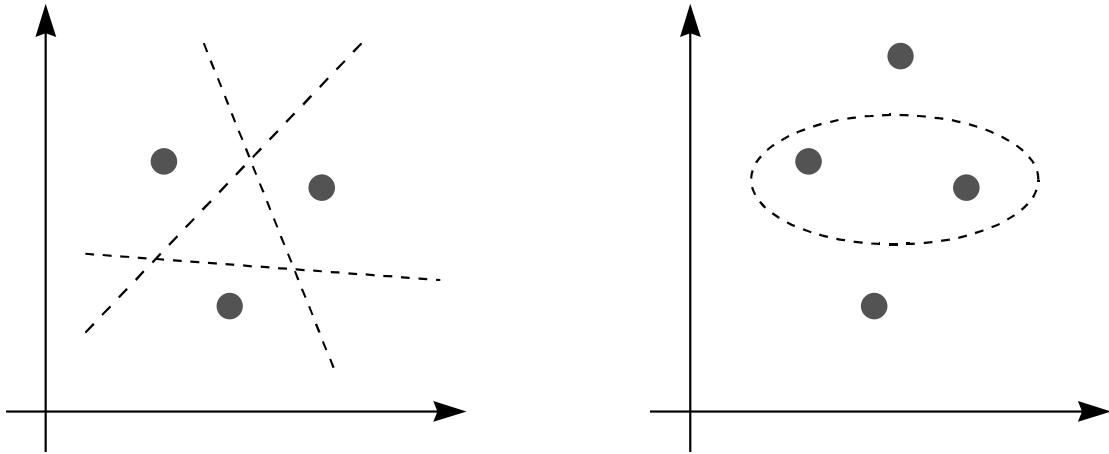


Figure 1 *VC Dimension Illustration*

The set of linear indicator functions,

Equation

has a VC dimension equal to $n+1$. The VC dimension is defined as,

There exists a set of points x^n such that these points can be separated in all 2^n possible configurations, and that no set x^m exists where $m > n$ satisfying this property.

Figure 1 illustrates how three points in the plane can be shattered by the set of linear indicator functions whereas four points cannot. In this case the VC dimension is equal to the number of free parameters, but in general that is not the case. e.g. the function $A \sin(bx)$ has an infinite VC dimension [ref].

SRM principle creates a structure

1.1.2 Structural Risk Minimisation

Create a structure such that S_h is a hypothesis space of VC dimension h ,

$$S_1 \subset S_2 \subset \dots \subset S_l \subset \dots$$

SRM consists in solving the following problem

$$\min_{S_h} \left[R_{emp}[f] + \sqrt{\frac{h(\ln \frac{2l}{h} + 1) - \ln(\frac{\delta}{4})}{l}} \right]$$

If the underlying process being modelled is not deterministic the modelling problem becomes more exacting and consequently this chapter is restricted to deterministic processes. Multiple output problems can usually be reduced to a set of single output

problems that may be considered independent. Hence it is appropriate to consider processes with multiple inputs from which it is desired to predict a single output.

2 Support Vector Classification

The classification problem can be restricted to consideration of the two-class problem without loss of generality. In this problem the goal is to separate the two classes by a function which is induced from available examples. The goal is to produce a classifier that will work well on *unseen* examples, i.e. it generalises well. Consider the example in Figure 2. Here there are many possible linear classifiers that can separate the data, but there is only one that maximises the margin (maximises the distance between it and the nearest data point of each class). This linear classifier is termed the optimal separating hyperplane. Intuitively, we would expect this boundary to generalise well as opposed to the other possible boundaries.

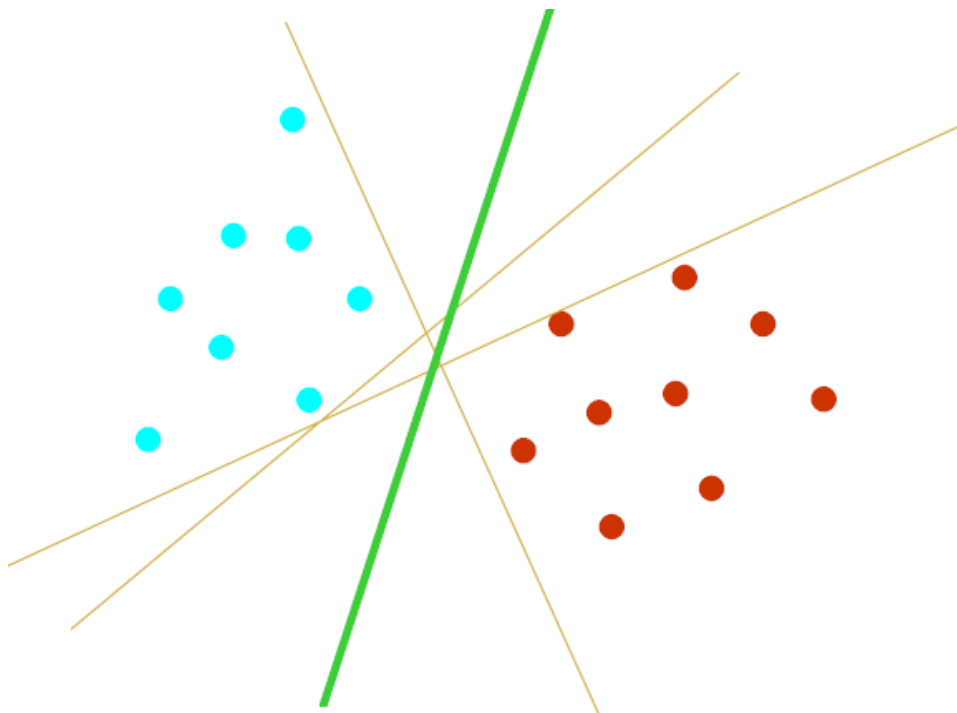


Figure 2 *Optimal Separating Hyperplane*

2.1 The Optimal Separating Hyperplane

Consider the problem of separating the set of training vectors belonging to two separate classes,

$$(y_1, \mathbf{x}_1), \dots, (y_l, \mathbf{x}_l), x \in R^n, y \in \{-1, +1\}, \quad (4)$$

with a hyperplane

$$(\mathbf{w} \cdot \mathbf{x}) + b = 0. \quad (5)$$

The set of vectors is said to be *optimally separated* by the hyperplane if it is separated without error and the distance between the closest vector to the hyperplane is maximal. There is some redundancy in Equation (5), and without loss of generality it is appropriate to consider a canonical hyperplane [21], where the parameters \mathbf{w} , b are constrained by,

$$\min_{\mathbf{x}_i} |\mathbf{x} \cdot \mathbf{w} + b| = 1 \quad (6)$$

This incisive constraint on the parameterisation is preferable to alternatives in simplifying the formulation of the problem. In words it states that: *the norm of the weight vector should be equal to the inverse of the distance, of the nearest point in the data set to the hyperplane.* The idea is illustrated in Figure 3, where the distance from the nearest point to each hyperplane is shown.

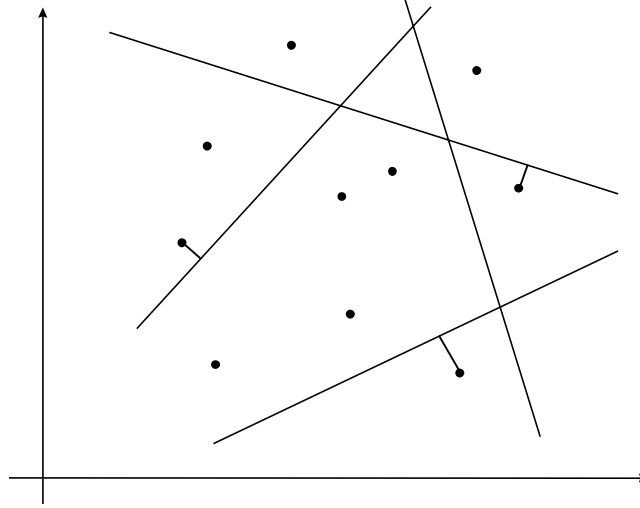


Figure 3 *Canonical Hyperplanes*

A separating hyperplane in canonical form must satisfy the following constraints,

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

The distance $d(\mathbf{w}, b; \mathbf{x})$ of a point \mathbf{x} from the hyperplane (\mathbf{w}, b) is,

$$d(\mathbf{w}, b; \mathbf{x}) = \frac{|\mathbf{w} \cdot \mathbf{x} + b|}{\|\mathbf{w}\|} \quad (8)$$

The optimal hyperplane is given by maximising the margin, $\rho(\mathbf{w}, b)$, subject to the constraints of Equation (7). The margin is given by,

$$\begin{aligned} \rho(\mathbf{w}, b) &= \min_{\{\mathbf{x}_i: y_i=1\}} d(\mathbf{w}, b; \mathbf{x}_i) + \min_{\{\mathbf{x}_j: y_j=-1\}} d(\mathbf{w}, b; \mathbf{x}_j) \\ &= \min_{\{\mathbf{x}_i: y_i=1\}} \frac{|\mathbf{w} \cdot \mathbf{x}_i + b|}{\|\mathbf{w}\|} + \min_{\{\mathbf{x}_j: y_j=-1\}} \frac{|\mathbf{w} \cdot \mathbf{x}_j + b|}{\|\mathbf{w}\|} \\ &= \frac{1}{\|\mathbf{w}\|} \left(\min_{\{\mathbf{x}_i: y_i=1\}} |\mathbf{w} \cdot \mathbf{x}_i + b| + \min_{\{\mathbf{x}_j: y_j=-1\}} |\mathbf{w} \cdot \mathbf{x}_j + b| \right) \\ &= \frac{2}{\|\mathbf{w}\|} \end{aligned} \quad (9)$$

Hence the hyperplane that optimally separates the data is the one that minimises

$$\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2. \quad (10)$$

It is independent of b because provided Equation (7) is satisfied (i.e. it is a separating hyperplane) changing b will move it in the normal direction to itself. Accordingly the margin remains unchanged but the hyperplane is no longer optimal in that it will be nearer to one class than the other.

To consider how minimising Equation (10) is equivalent to implementing the *SRM* principle, suppose that the following bound holds,

$$\|\mathbf{w}\| \leq A. \quad (11)$$

Then from Equation (7) and (8),

$$d(\mathbf{w}, b; \mathbf{x}) \geq \frac{1}{A}. \quad (12)$$

Accordingly the hyperplanes cannot be nearer than $1/A$ to any of the data points and intuitively it can be seen in Figure 4 how this reduces the possible hyperplanes, and hence the capacity.

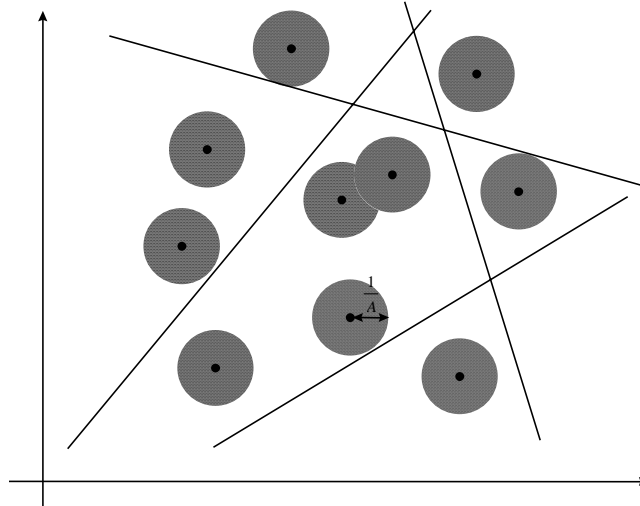


Figure 4 *Constraining the Canonical Hyperplanes*

The VC dimension, h , of the set of canonical hyperplanes in n dimensional space is,

$$h \leq \min[R^2 A^2, n] + 1, \quad (13)$$

where R is the radius of a hypersphere enclosing all the data points. Hence minimising Equation (10) is equivalent to minimising an upper bound on the VC dimension.

The solution to the optimisation problem of Equation (10) under the constraints of Equation (7) is given by the saddle point of the Lagrange functional (Lagrangian) [11],

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^l \alpha_i \{[(\mathbf{x}_i \cdot \mathbf{w}) + b] y_i - 1\}. \quad (14)$$

where α_i are the Lagrange multipliers. The Lagrangian has to be minimised with respect to \mathbf{w} , b and maximised with respect to $\alpha_i \geq 0$. Classical Lagrangian duality enables the *primal* problem, Equation (14), to be transformed to its *dual* problem, which is easier to solve. The *dual* problem is given by,

$$\max_{\alpha} W(\alpha) = \max_{\alpha} \left\{ \min_{\mathbf{w}, b} L(\mathbf{w}, b, \alpha) \right\} \quad (15)$$

The minimum with respect to \mathbf{w} and b of the Lagrangian, L , is given by,

$$\begin{aligned} \frac{\partial L}{\partial b} = 0 & \quad \Rightarrow \quad \sum_{i=1}^l \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \mathbf{w}} = 0 & \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^l \alpha_i \mathbf{x}_i y_i \end{aligned} \quad (16)$$

Hence from Equations (14), (15) and (16), the *dual* problem is,

$$\max_{\alpha} W(\alpha) = \max_{\alpha} -\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) + \sum_{i=1}^l \alpha_i \quad (17)$$

and hence the solution to the problem is given by,

$$\bar{\alpha} = \arg \min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^l \alpha_i \quad (18)$$

with constraints,

$$\begin{aligned} \alpha_i & \geq 0, \quad i = 1, \dots, l \\ \sum_{i=1}^l \alpha_i y_i & = 0 \end{aligned} \quad (19)$$

Solving Equation (18) with constraints Equation (19) determines the Lagrange multipliers, and the optimal separating hyperplane is given by,

$$\bar{\mathbf{w}} = \sum_{i=1}^l \bar{\alpha}_i \mathbf{x}_i y_i \quad (20)$$

$$\bar{b} = -\frac{1}{2} \bar{\mathbf{w}} \cdot [\mathbf{x}_r + \mathbf{x}_s]$$

where \mathbf{x}_r and \mathbf{x}_s are any support vector from each class satisfying,

$$\bar{\alpha}_r, \bar{\alpha}_s > 0, \quad y_r = 1, \quad y_s = -1. \quad (21)$$

The hard classifier is then,

$$f(\mathbf{x}) = \text{sign}(\bar{\mathbf{w}} \cdot \mathbf{x} + \bar{b}). \quad (22)$$

Alternatively, a soft classifier may be used which linearly interpolates the margin,

$$f(\mathbf{x}) = h(\bar{\mathbf{w}} \cdot \mathbf{x} + \bar{b}) \quad \text{where} \quad h(x) = \begin{cases} -1 & x < -1 \\ x & -1 \leq x \leq 1 \\ 1 & 1 \end{cases} \quad (23)$$

This may be more appropriate than the hard classifier of Equation (22), because it produces a real valued output between -1 and 1 when the classifier is queried within the margin, where no training data resides.

From the Kuhn-Tucker conditions,

$$\bar{\alpha}_i [y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) - 1] = 0, \quad (24)$$

and hence only the points \mathbf{x}_i which satisfy,

$$y_i (\bar{\mathbf{w}} \cdot \mathbf{x}_i + \bar{b}) = 1, \quad (25)$$

will have non-zero Lagrange multipliers. These points are termed *Support Vectors* (SV). If the data is linearly separable all the SV will lie on the margin and hence the number of SV can be very small. Consequently the hyperplane is determined by a small subset of the training set; the other points could be removed from the training set and recalculating the hyperplane would produce the same answer. Hence SVM can be used to summarise the information contained in a data set by the SV produced. If the data is linearly separable the following equality will hold,

$$\|\bar{\mathbf{w}}\|^2 = \sum_{i=1}^l \bar{\alpha}_i = \sum_{SVs} \bar{\alpha}_i = \sum_{SVs} \sum_{SVs} \bar{\alpha}_i \bar{\alpha}_j (\mathbf{x}_i \cdot \mathbf{x}_j) y_i y_j. \quad (26)$$

Hence from Equation (13) the VC dimension of the classifier is bounded by,

$$h \leq \min \left[R^2 \sum_{SVs} \bar{\alpha}_i, n \right] + 1, \quad (27)$$

and if the training data, \mathbf{x} , is normalised to lie in the unit hypersphere,

$$h \leq 1 + n \min \left[\sum_{SVs} \bar{\alpha}_i, 1 \right]. \quad (28)$$

2.1.1 Linearly Separable Example

To illustrate the method consider the training set in Table 1.

X_1	X_2	Class
1	1	-1
3	3	1
1	3	1
3	1	-1
2	2.5	1
3	2.5	-1
4	3	-1

Table 1 *Linearly Separable Classification Data*

The SVC solution is shown in Figure 4, where the dotted lines describe the locus of the margin and the circled data points represent the SV, which all lie on the margin.

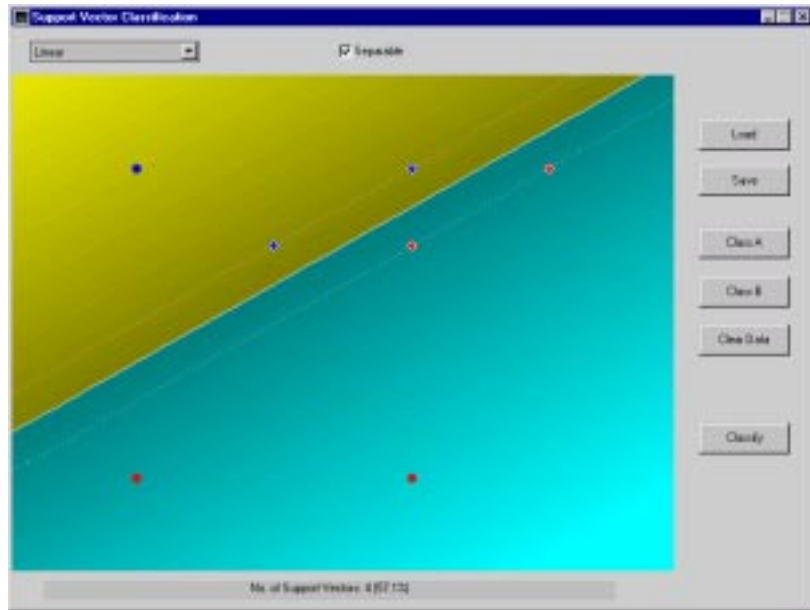


Figure 5 *Optimal Separating Hyperplane*

2.2 The Generalised Optimal Separating Hyperplane

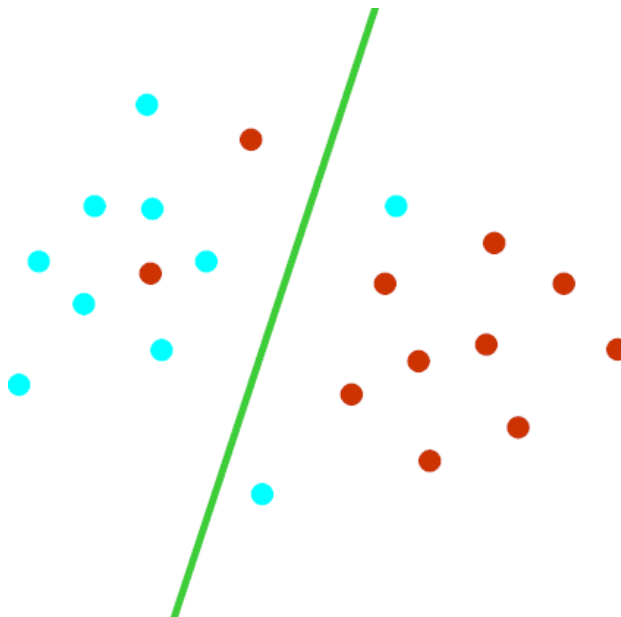


Figure 6 *Generalised Optimal Separating Hyperplane*

So far the discussion has been restricted to the case where the training data is linearly separable. However, in general this will not be the case, Figure 6. There are two approaches to generalising the problem, which are dependent upon prior knowledge of the problem and an estimate of the noise on the data. In the case where it is expected (or possibly even known) that a hyperplane can correctly separate the data, a method of introducing an additional cost function associated with misclassification is appropriate. Alternatively a more complex function can be used to describe the boundary, as discussed in section 2.3. To enable the optimal separating hyperplane method to be generalised, Cortes [5] introduced non-negative variables $\xi_i \geq 0$ and a penalty function,

$$F_{\sigma}(\xi) = \sum_{i=1}^l \xi_i^{\sigma}, \quad \sigma > 0,$$

where the ξ are a measure of the misclassification error. The optimisation problem is now posed so as to minimise the classification error as well as minimising the bound on the VC dimension of the classifier. The constraints of Equation (7) are modified for the non-separable case to,

$$y_i [(\mathbf{w} \cdot \mathbf{x}_i) + b] \geq 1 - \xi_i, \quad i = 1, \dots, l \quad (29)$$

where $\xi_i \geq 0$. The generalised optimal separating hyperplane is determined by the vector \mathbf{w} , that minimises the functional,

$$\Phi(\mathbf{w}, \xi) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^l \xi_i, \quad (30)$$

(where C is a given value) subject to the constraints of Equation (29).

The solution to the optimisation problem of Equation (30) under the constraints of Equation (29) is given by the saddle point of the Lagrangian [11],

$$L(\mathbf{w}, b, \xi, \alpha, \beta) = \frac{1}{2} (\mathbf{w} \cdot \mathbf{w}) + C \sum_{i=1}^l \xi_i - \sum_{i=1}^l \alpha_i \{[(\mathbf{x}_i \cdot \mathbf{w}) + b] y_i - 1 + \xi_i\} - \sum_{i=1}^l \beta_i \xi_i, \quad (31)$$

where α_i, β_i are the Lagrange multipliers. The Lagrangian has to be minimised with respect to \mathbf{w} , b , ξ and maximised with respect to $\alpha_i, \beta_i \geq 0$. As before, classical Lagrangian duality enables the *primal* problem, Equation (31), to be transformed to its *dual* problem. The *dual* problem is given by,

$$\max_{\alpha, \beta} W(\alpha, \beta) = \max_{\alpha, \beta} \left\{ \min_{\mathbf{w}, b, \xi} L(\mathbf{w}, b, \xi, \alpha, \beta) \right\} \quad (32)$$

The minimum with respect to \mathbf{w} , b and ξ_i of the Lagrangian, L , is given by,

$$\begin{aligned} \frac{\partial L}{\partial b} = 0 & \quad \Rightarrow \quad \sum_{i=1}^l \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \mathbf{w}} = 0 & \quad \Rightarrow \quad \mathbf{w} = \sum_{i=1}^l \alpha_i \mathbf{x}_i y_i. \\ \frac{\partial L}{\partial \xi_i} = 0 & \quad \Rightarrow \quad \alpha_i + \beta_i = C \end{aligned} \quad (33)$$

Hence from Equations (31), (32) and (33), the *dual* problem is,

$$\max_{\alpha} W(\alpha) = \max_{\alpha} - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) + \sum_{i=1}^l \alpha_i \quad (34)$$

and hence the solution to the problem is given by,

$$\bar{\alpha} = \arg \min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^l \alpha_i \quad (35)$$

with constraints,

$$0 \leq \alpha_i \leq C, \quad i = 1, \dots, l$$

$$\sum_{i=1}^l \alpha_i y_i = 0 \quad . \quad (36)$$

The solution to this minimisation problem is identical to the separable case except for a modification of the bounds of the Lagrange multipliers. The uncertain part of Cortes's approach is that the coefficient C has to be determined. This parameter introduces additional capacity control within the classifier. In some circumstances C can be directly related to a regularisation parameter [7,17]. Blanz [4] uses a value of $C=5$, but ultimately C must be chosen to reflect the knowledge of the noise on the data. This warrants further work, but a more practical discussion is given in section 4.

2.2.1 Linearly Non-Separable Example

Two additional data points are added to the separable data of Table 1 to produce a linearly non-separable data set, Table 2.

X_1	X_2	Class
1	1	-1
3	3	1
1	3	1
3	1	-1
2	2.5	1
3	2.5	-1
4	3	-1
1.5	1.5	1
1	2	-1

Table 2 *Linearly Non-Separable Classification Data*

The resulting SVC is shown in Figure 7, for $C=10$. The SV are no longer required to lie on the margin, as in Figure 5, and the orientation of the hyperplane and the width of the margin are different.

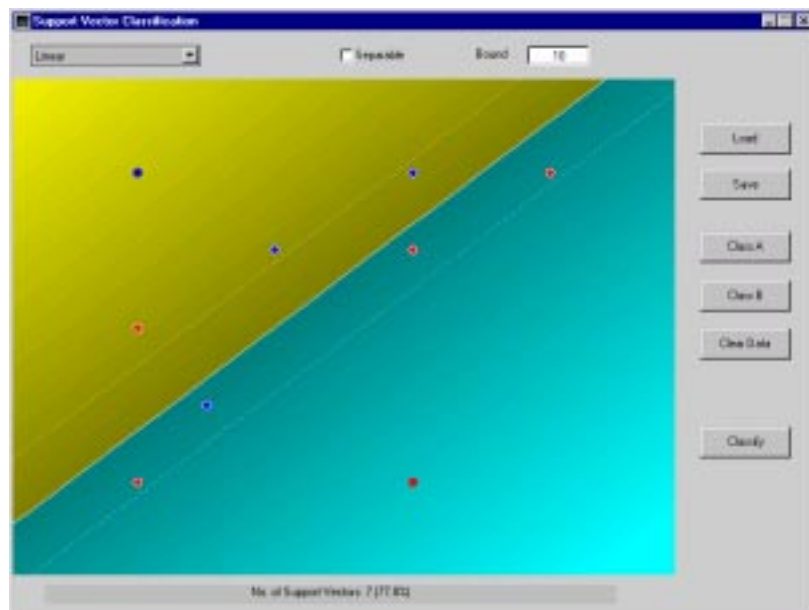


Figure 7 *Generalised Optimal Separating Hyperplane Example ($C=10$)*

As $C \rightarrow \infty$ the solution converges towards the solution of obtained by the optimal separating hyperplane (on this non-separable data), Figure 7.

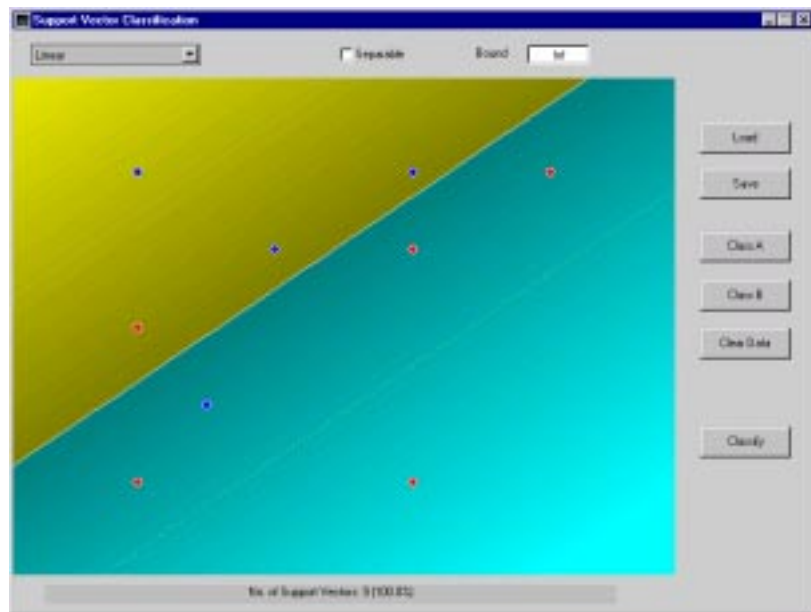


Figure 8 *Generalised Optimal Separating Hyperplane Example ($C=\infty$)*

In the limit as $C \rightarrow 0$ the solution converges to one which minimises the error associated with the slack variables, Figure 9. There is now no emphasis on maximising the margin, but purely on minimising the misclassification error, producing a zero width margin. Consequently as C decreases so does the width of the margin.

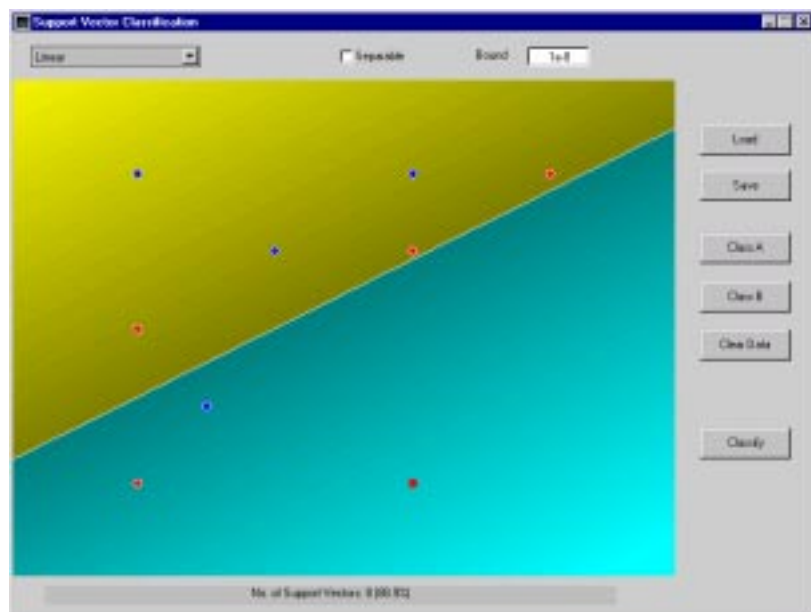


Figure 9 *Generalised Optimal Separating Hyperplane Example ($C=10^{-8}$)*

2.3 Generalisation in High Dimensional Feature Space

In the case where a linear boundary is inappropriate the *SVM* can map the input vector, \mathbf{x} , into a high dimensional feature space, \mathbf{z} . By choosing a non-linear mapping

a priori, the SVM constructs an optimal separating hyperplane in this higher dimensional space, Figure 10. The idea exploits the method of [1] which, enables the curse of dimensionality [3] to be addressed.

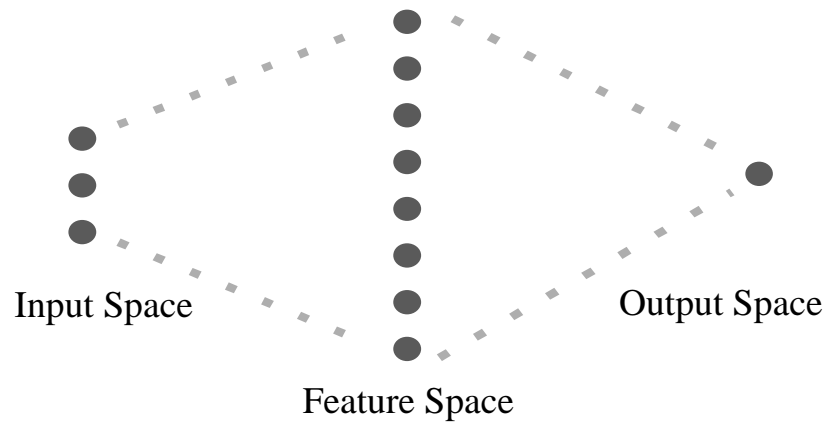


Figure 10 Mapping the Input Space into a High Dimensional Feature Space

There are some restrictions on the non-linear mapping that can be employed, see Chapter 3, but it turns out, surprisingly, that most commonly employed functions are acceptable. Among acceptable mappings are polynomials, radial basis functions and certain sigmoid functions. The optimisation problem of Equation (35) becomes,

$$\bar{\alpha} = \arg \min_{\alpha} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) - \sum_{i=1}^l \alpha_i \quad (37)$$

where $K(\mathbf{x}, \mathbf{y})$ is the kernel function performing the non-linear mapping into feature space, and the constraints are unchanged,

$$\begin{aligned} \alpha_i &\geq 0, \quad i = 1, \dots, l \\ \sum_{i=1}^l \alpha_i y_i &= 0 \end{aligned} \quad (38)$$

Solving Equation (37) with constraints Equation (38) determines the Lagrange multipliers, and a hard classifier implementing the optimal separating hyperplane in the feature space is given by,

$$f(\mathbf{x}) = \text{sign} \left(\sum_{SVs} \bar{\alpha}_i y_i K(\mathbf{x}_i, \mathbf{x}) + \bar{b} \right) \quad (39)$$

where

$$\begin{aligned} \bar{\mathbf{w}} \cdot \mathbf{x} &= \sum_{SVs} \bar{\alpha}_i y_i K(\mathbf{x}_i, \mathbf{x}) \\ \bar{b} &= -\frac{1}{2} \sum_{SVs} \bar{\alpha}_i y_i [K(\mathbf{x}_r, \mathbf{x}_i) + K(\mathbf{x}_s, \mathbf{x}_i)] \end{aligned} \quad (40)$$

The bias is computed here using two support vectors, but can be computed using all the SV on the margin for stability [20]. If the Kernel contains a bias term, the bias can be accommodated within the Kernel, and hence the classifier is simply,

$$f(\mathbf{x}) = \text{sign}\left(\sum_{SVs} \bar{\alpha}_i y_i K(\mathbf{x}_i, \mathbf{x})\right). \quad (41)$$

Many employed kernels have a bias term and any finite Kernel can be made to have one [7]. This simplifies the optimisation problem by removing the equality constraint of Equation (38). Chapter 3 discusses the necessary conditions that must be satisfied by valid kernel functions.

2.3.1 Polynomial Mapping Example

Consider a polynomial kernel of the form,

$$K(\mathbf{x}, \mathbf{y}) = [(\mathbf{x} \cdot \mathbf{y}) + 1]^2,$$

which maps a two dimensional input vector into a six dimensional feature space. Applying the non-linear SVC to the linearly non-separable training data of Table 2, produces the classification illustrated in Figure 11 ($C=\infty$). The margin is no longer of constant width due to the non-linear projection into the input space. The solution is in contrast to Figure 8, in that the training data is now classified correctly. However, even though *SVMs* implement the *SRM* principle and hence can generalise well, a careful choice of the kernel function is necessary to produce a classification boundary that is topologically appropriate. It is always possible to map the input space into a dimension greater than the number of training points and produce a classifier with no classification errors on the training set. However, this will generalise badly.

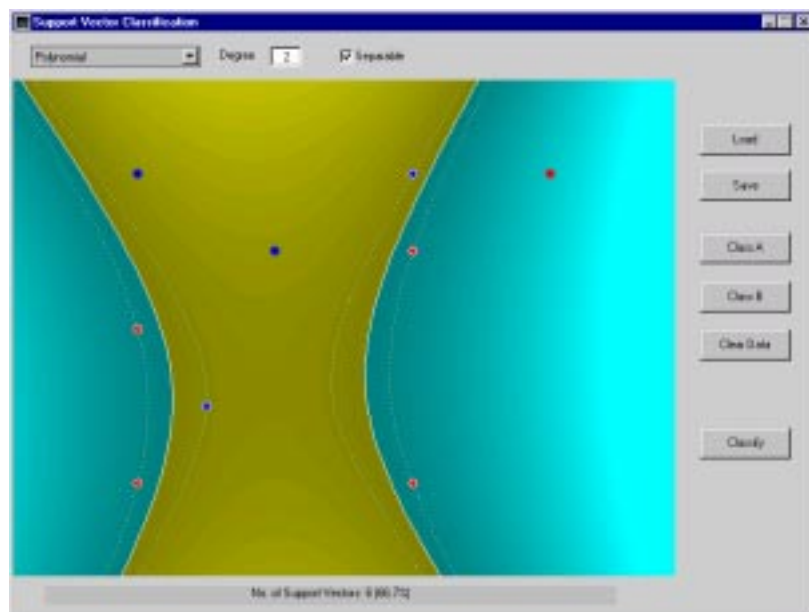


Figure 11 Mapping input space into Polynomial Feature Space

2.4 Discussion

Typically the data will only be linearly separable in some, possibly very high dimensional feature space. It may not make sense to try and separate the data exactly, particularly when only a finite amount of training data is available which is potentially corrupted by noise. Hence in practice it will be necessary to employ the non-separable approach which places an upper bound on the Lagrange multipliers.

This raises the question of how to determine the parameter C . It is similar to the problem in regularisation where the regularisation coefficient has to be determined, and it has been shown that the parameter C can be directly related to a regularisation parameter for certain kernels [17]. A process of cross-validation can be used to determine this parameter, although more efficient and potentially better methods are sought after. In removing the training patterns that are not support vectors, the solution is unchanged and hence a fast method for validation may be available when the support vectors are sparse.

3 Feature Space

This chapter discusses the method that can be used to construct a mapping into a high dimensional feature space by the use of reproducing kernels. The idea of the kernel function is to enable operations to be performed in the input space rather than the potentially high dimensional feature space. Hence the inner product does not need to be evaluated in the feature space. This provides a way of addressing the curse of dimensionality. However, the computation is still critically dependent upon the number of training patterns and to provide a good data distribution for a high dimensional problem will generally require a large training set.

3.1 Kernel Functions

The following theory is based upon Reproducing Kernel Hilbert Spaces (RKHS) [2, 22, 7, 10]. An inner product in feature space has an equivalent kernel in input space,

$$K(\mathbf{x}, \mathbf{y}) = k(\mathbf{x}) \cdot k(\mathbf{y}), \quad (42)$$

provided certain conditions hold. If K is a symmetric positive definite function, which satisfies Mercer's Conditions,

$$K(\mathbf{x}, \mathbf{y}) = \sum_{m=1}^{\infty} \alpha_m \psi(\mathbf{x}) \psi(\mathbf{y}), \quad \alpha_m \geq 0, \quad (43)$$

$$\iint K(\mathbf{x}, \mathbf{y}) g(\mathbf{x}) g(\mathbf{y}) d\mathbf{x} d\mathbf{y} > 0, \quad \int g^2(\mathbf{x}) d\mathbf{x} < \infty$$

then the kernel represents a legitimate inner product in feature space. Valid functions that satisfy Mercer's conditions are now given, which unless stated are valid for all real \mathbf{x} and \mathbf{y} .

3.1.1 Polynomial

A polynomial mapping is a popular method for non-linear modelling,

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x} \cdot \mathbf{y})^d, \quad d = 1, \dots \quad (44)$$

$$K(\mathbf{x}, \mathbf{y}) = ((\mathbf{x} \cdot \mathbf{y}) + 1)^d$$

The second kernel is usually preferable as it avoids problems with the hessian becoming zero.

3.1.2 Gaussian Radial Basis Function

Radial basis functions have received significant attention, most commonly with a Gaussian of the form,

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{(\mathbf{x} - \mathbf{y})^2}{2\sigma^2}\right) \quad (45)$$

Classical techniques utilising radial basis functions employ some method of determining a subset of centres. Typically a method of clustering is first employed to select a subset of centres. An attractive feature of the SVM is that this selection is implicit, with each support vectors contributing one local Gaussian function, centred at that data point. By further considerations it is possible to select the global basis function width, σ , using the *SRM* principle [21].

3.1.3 Exponential Radial Basis Function

A radial basis function of the form,

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{|\mathbf{x} - \mathbf{y}|}{2\sigma^2}\right) \quad (46)$$

produces a piecewise linear solution which can be attractive when discontinuities are acceptable.

3.1.4 Multi-Layer Perceptron

The long established MLP, with a single hidden layer, also has a valid kernel representation,

$$K(\mathbf{x}, \mathbf{y}) = \tanh(\text{scale} \cdot (\mathbf{x} \cdot \mathbf{y}) - \text{offset}) \quad (47)$$

for certain values of the *scale* and *offset* paramters. Here the SV correspond to the first layer and the Lagrange multipliers to the weights.

3.1.5 Fourier Series

A Fourier series can be considered an expansion in the following $2N+1$ dimensional feature space. The kernel is defined on the interval $[-\frac{\pi}{2}, \frac{\pi}{2}]$,

$$K(\mathbf{x}, \mathbf{y}) = \frac{\sin(N + \frac{1}{2})(\mathbf{x} - \mathbf{y})}{\sin(\frac{1}{2}(\mathbf{x} - \mathbf{y}))}. \quad (48)$$

However, this kernel is probably not a good choice because its regularisation capability is poor, which is evident by consideration of its Fourier transform [17].

3.1.6 Splines

Splines are a popular choice for modelling due to their flexibility. A finite spline, of order p , with knots located at τ_s is given by,

$$K(x, y) = \sum_{r=0}^p x^r y^r + \sum_{s=1}^N (x - \tau_s)_+^p (y - \tau_s)_+^p \quad (49)$$

An infinite spline is defined on the interval $[0,1)$ by,

$$K(x, y) = \sum_{r=0}^p x^r y^r + \int_0^1 (x - \tau)_+^p (y - \tau)_+^p d\tau. \quad (50)$$

In the case when $p=1$ (S_1^∞) the kernel is given by,

$$K(x, y) = 1 + xy + xy \min(x, y) - \frac{(x + y)}{2} (\min(x, y))^2 + \frac{1}{3} (\min(x, y))^3 \quad (51)$$

where the solution is a piece-wise cubic.

3.1.7 Bsplines

Bsplines are another popular spline formulation. The kernel is defined on the interval $[-1, 1]$, and has an attractive closed form,

$$K(x, y) = B_{2n+1}(x - y). \quad (52)$$

3.1.8 Additive Kernels

More complicated kernels can be obtained by forming summing kernels, since the sum of two positive definite functions is positive definite.

$$K(\mathbf{x}, \mathbf{y}) = \sum_i K_i(\mathbf{x}, \mathbf{y}). \quad (53)$$

3.1.9 Tensor Product Kernels

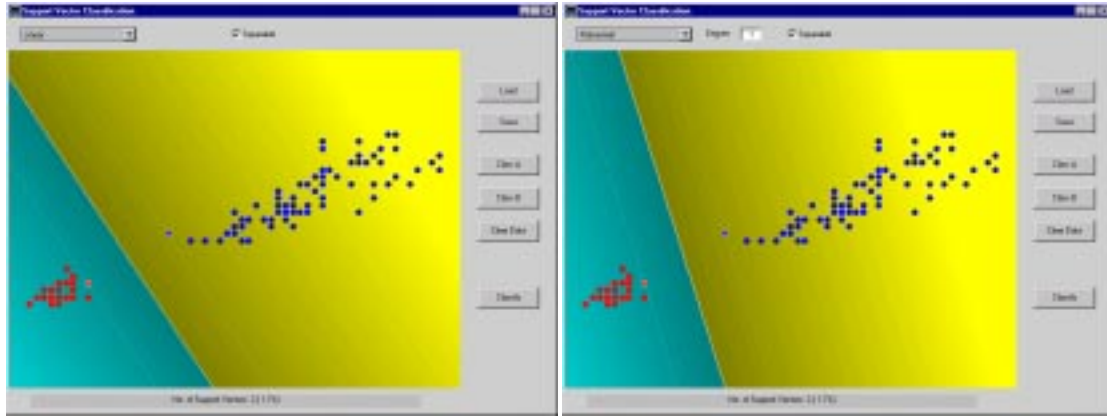
Multidimensional kernels can be obtained by forming tensor products of kernels [2],

$$K(\mathbf{x}, \mathbf{y}) = \prod_{m=1}^n K_m(\mathbf{x}_m, \mathbf{y}_m). \quad (54)$$

This is particularly useful in the construction of multidimensional spline kernels, which are simply obtained from the product of the univariate kernels.

3.2 Implicit vs. Explicit Bias

It was remarked in the previous chapter that kernels may or may not contain an implicit bias. The inclusion of a bias within the kernel function can lead to a slightly more efficient method of implementation. However, the solutions obtained with an implicit and explicit bias are not the same, which may initially come as a surprise. This difference helps to highlight the difficulties with the interpretation of generalisation in high dimensional feature spaces. Figure 12 compares a linear kernel with explicit bias against polynomial of degree 1 with implicit bias. It is evident that the solutions are different, although both solutions would seem to offer good generalisation.



(a) Explicit (linear)

(b) Implicit (polynomial degree 1)

Figure 12 *Comparison between Implicit and Explicit bias for a linear kernel.*

3.3 Data Normalisation

Data normalisation is required for particular kernels due to their restricted domain, and may also be advantageous for unrestricted kernels. To determine if normalisation (isotropic or non-isotropic) of the data is necessary requires the consideration of the input features. Additionally, normalisation will improve the condition number of the hessian in the optimisation problem.

3.4 Kernel Selection

The obvious question that arises is that with so many different mappings to choose from, which is the best for a particular problem? This is not a new question, but with the inclusion of many mappings within one framework it is easier to make a comparison. The upper bound on the VC dimension, Equation (13), is a potential avenue to provide a means of comparing the kernels. However, it requires the estimation of the radius of the hypersphere enclosing the data in the non-linear feature space. As a final caution, even if a strong theoretical method for selecting a kernel is developed, unless this can be validated using independent test sets on a large number of problems, methods such as bootstrapping and cross-validation will remain the preferred method for kernel selection.

4 Classification Example: IRIS data

The iris data set is an established data set used for demonstrating the performance of classification algorithms. The data set contains four attributes of an iris, and the goal is to classify the class of iris based on these four attributes. To visualise the problem we restrict ourselves to the two features that contain the most information about the class, namely the petal length and the petal width. The distribution of the data is illustrated in Figure 13.

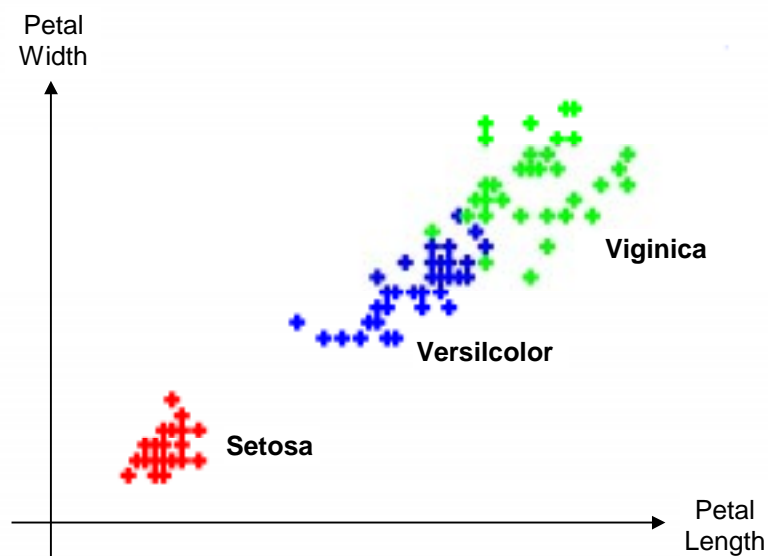


Figure 13 *Iris data set*

The Setosa and Versicolor classes are easily separated with a linear boundary and the SVC solution using an inner product kernel is illustrated in Figure 14, with the two support vectors circled.

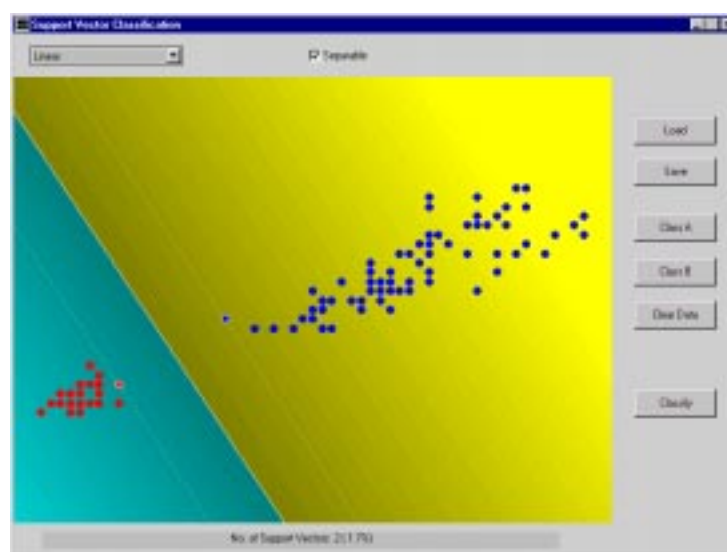


Figure 14 *Separating Setosa with a linear SVC ($C=\infty$)*

The two support vectors contain the important information about the classification boundary and hence illustrate the potential of SVC for data selection. The separation of the class *Viginica* from the other two classes is not so trivial. In fact, two of the examples are identical in petal length and width, but correspond to different classes. Figure 15 illustrates the SVC solution obtained using a degree 2 polynomial and it is clear that the area of input space where there is little data is classified as *Viginica*.

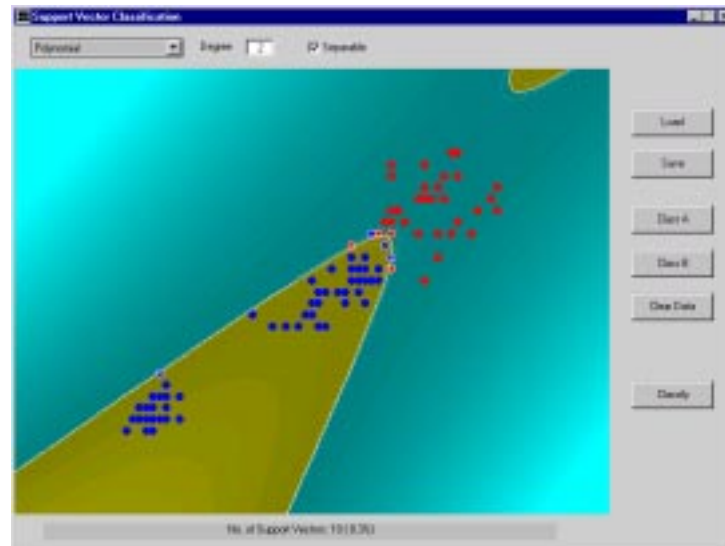


Figure 15 *Separating Viginica with a polynomial SVM (degree 2, $C=\infty$)*

Figure 16 illustrates the use of a higher order polynomial to separate the *Viginica*, with no additional capacity control. This SVC determines a hyperplane in a 55 dimensional feature space. There is evidence of overfitting due to the high dimensional nature of the kernel function, which is emphasised by the disjoint region in the top of the illustration.



Figure 16 *Separating Viginica with a polynomial SVM (degree 10, $C=\infty$)*

Figure 17 illustrates a Gaussian radial basis function SVC using a pre-specified variance. The result is similar to that of the degree 2 polynomial.

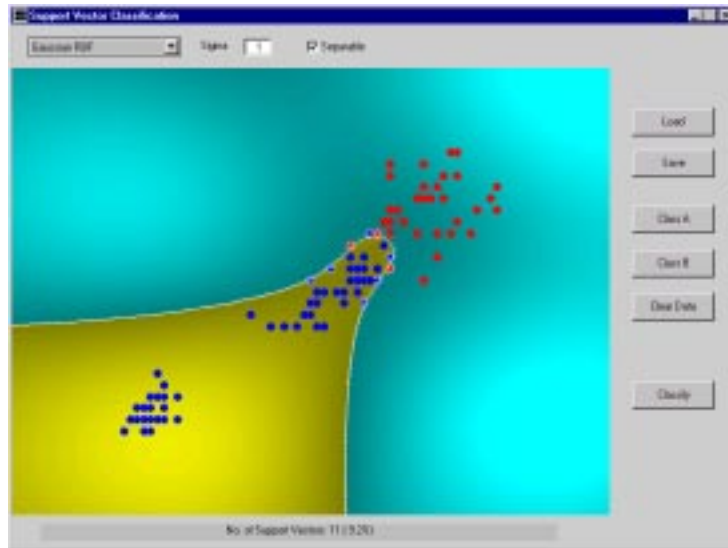


Figure 17 *Separating Viginica with a Radial Basis Function SVM ($\sigma=1.0, C=\infty$)*

Figure 18 illustrates the SVC solution obtained using the degree 2 polynomial with some tolerance to misclassification errors ($C=10$). This can be seen to produce a solution with good expected generalisation, emphasising the importance of tolerating misclassification errors in this example. This is necessary due to the non-separable nature of the data using just two input features.

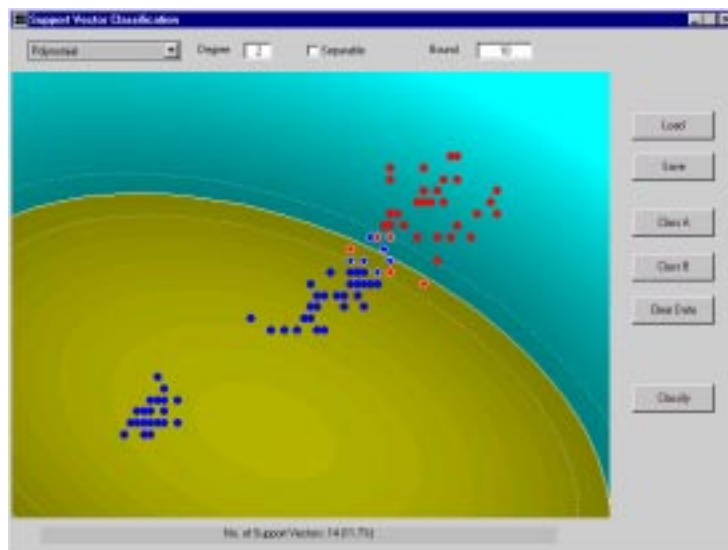


Figure 18 *Separating Viginica with a polynomial SVM (degree 2, $C=10$)*

To visualise the effect of the tolerance to misclassification errors on the topology of the classifier boundary, Figure 19 shows the results of a linear spline SVC for various degrees of misclassification tolerance.

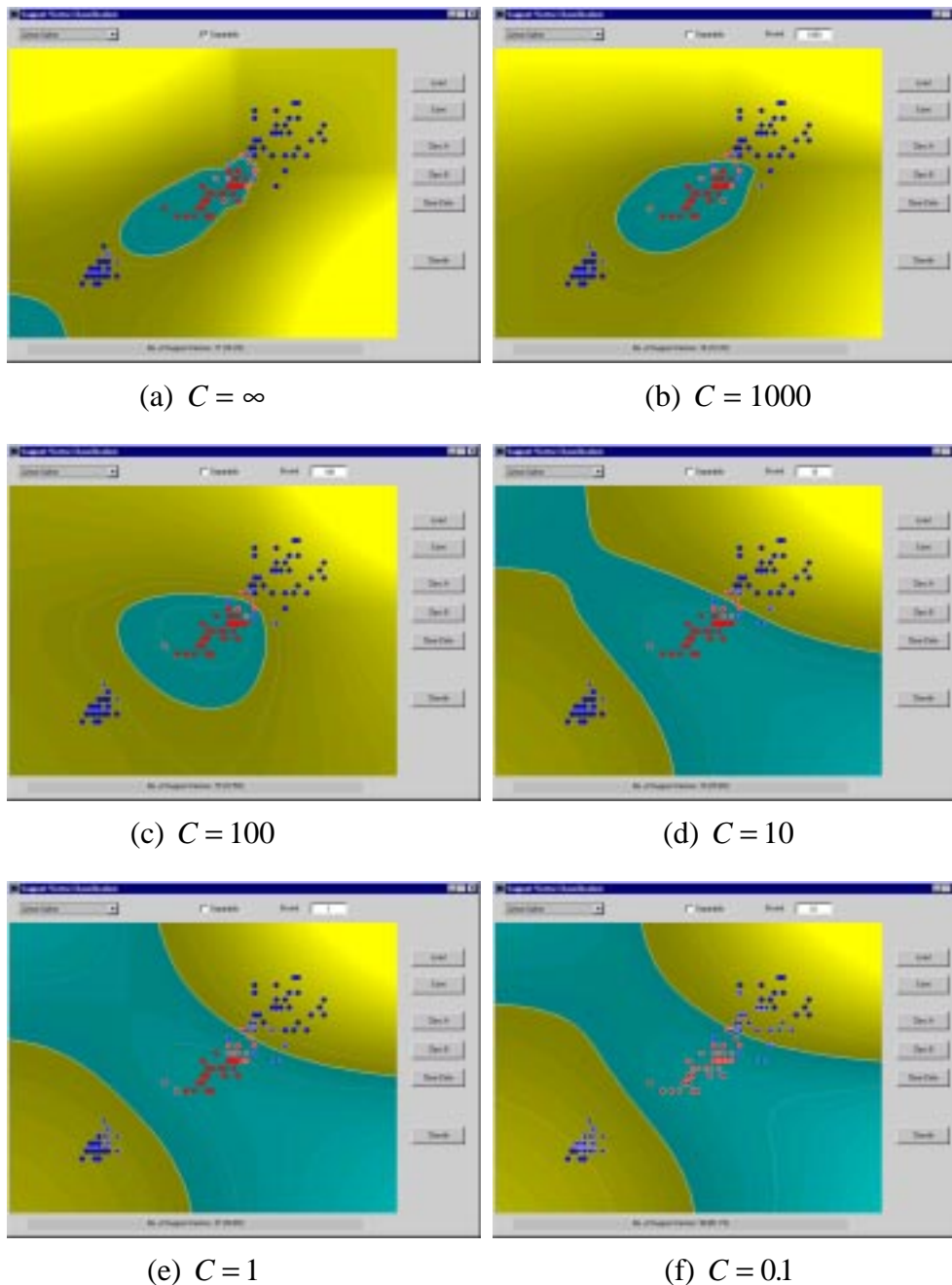


Figure 19 The effect of C on the separation of Versicolor with a linear spline SVM

Interestingly, the values of $C=1$ and $C=100$ seem to offer good solutions, depending upon whether an open boundary, Figure 19(a) or a closed boundary, Figure 19(c) is more appropriate. This demonstrates that the parameter C may have more than one optimal value and prior knowledge about the problem under consideration may be required to select the final solution.

4.1 Applications

Larger and more complex classification problems have been attacked with SVC. Notably, Osuna [14] has applied SVC to the exacting problem of face recognition, with encouraging results. In conclusion, SVC provides a robust method for pattern classification by minimising overfitting problems by adopting the SRM principle. Use

of a kernel function enables the curse of dimensionality to be addressed, and the solution implicitly contains support vectors that provide a description of the significant data for classification.

5 Support Vector Regression

SVMs can also be applied to regression problems by the introduction of an alternative loss function. The loss function must be modified to include a distance measure. Figure 20 illustrates four possible loss functions.

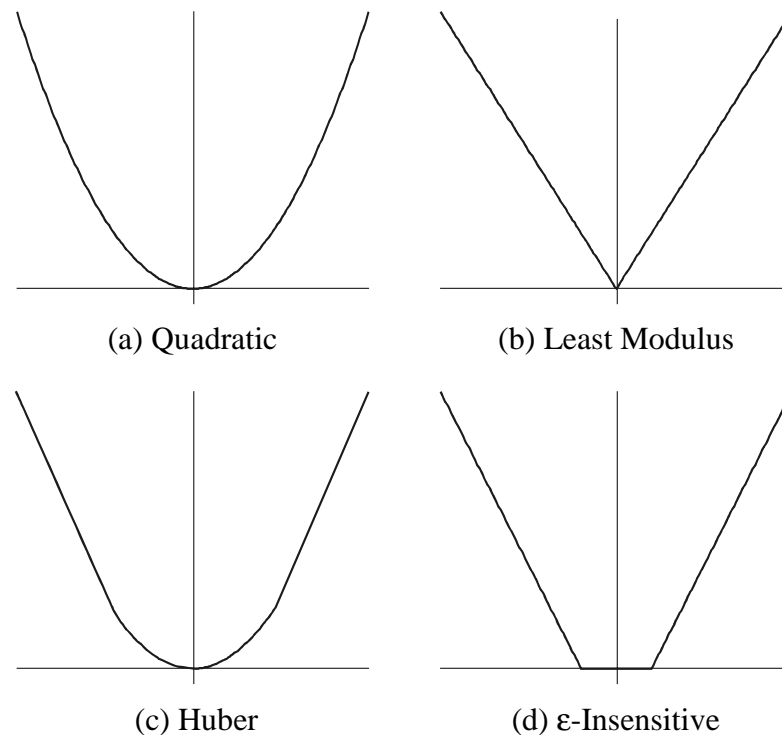


Figure 20 *Loss Functions*

The loss function in Figure 20(a) corresponds to the conventional least squares error criterion. The loss function in Figure 20(b) is a Laplacian loss function that is less sensitive to outliers than the quadratic loss function. Huber proposed the loss function in Figure 20(c) as a robust loss function that has optimal properties when the underlying distribution of the data is unknown. These three loss functions will produce no sparseness in the support vectors. To address this issue Vapnik proposed the loss function in Figure 20(d) as an approximation to Huber's loss function that enables a sparse set of support vectors to be obtained.

5.1 Linear Regression

Consider the problem of approximating the set of data,

$$(y_1, \mathbf{x}_1), \dots, (y_l, \mathbf{x}_l), x \in R^n, y \in R, \quad (55)$$

with a linear function,

$$f(\mathbf{x}) = (\mathbf{w} \cdot \mathbf{x}) + b. \quad (56)$$

the optimal regression function is given by the minimum of the functional,

$$\Phi(\mathbf{w}, \xi^*, \xi) = \frac{1}{2} \|\mathbf{w}\|^2 + C \left(\sum_{i=1}^l \xi_i + \sum_{i=1}^l \xi_i^* \right), \quad (57)$$

where C is a pre-specified value, and ξ, ξ^* are slack variables representing upper and lower constraints on the outputs of the system.

5.1.1 ε -Insensitive Loss Function

Using an ε -insensitive loss function, Figure 20(d),

$$L_\varepsilon(y) = \begin{cases} 0 & \text{for } |f(\mathbf{x}) - y| < \varepsilon \\ |f(\mathbf{x}) - y| - \varepsilon & \text{otherwise} \end{cases}. \quad (58)$$

the solution is given by,

$$\max_{\alpha, \alpha^*} W(\alpha, \alpha^*) = \max_{\alpha, \alpha^*} \left\{ \begin{aligned} & -\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*)(\mathbf{x}_i \cdot \mathbf{x}_j) \\ & + \sum_{i=1}^l \alpha_i (y_i - \varepsilon) - \alpha_i^* (y_i + \varepsilon) \end{aligned} \right\}, \quad (59)$$

or alternatively,

$$\bar{\alpha}, \bar{\alpha}^* = \arg \min_{\alpha, \alpha^*} \left\{ \begin{aligned} & \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*)(\mathbf{x}_i \cdot \mathbf{x}_j) \\ & - \sum_{i=1}^l (\alpha_i - \alpha_i^*) y_i + \sum_{i=1}^l (\alpha_i + \alpha_i^*) \varepsilon \end{aligned} \right\} \quad (60)$$

with constraints,

$$\begin{aligned} 0 &\leq \alpha_i \leq C, \quad i = 1, \dots, l \\ 0 &\leq \alpha_i^* \leq C, \quad i = 1, \dots, l \\ \sum_{i=1}^l (\alpha_i - \alpha_i^*) &= 0 \end{aligned}. \quad (61)$$

Solving Equation (59) with constraints Equation (61) determines the Lagrange multipliers, α_i, α_i^* , and the regression function is given by Equation (56), where

$$\begin{aligned} \bar{\mathbf{w}} &= \sum_{i=1}^l (\alpha_i - \alpha_i^*) \mathbf{x}_i \\ \bar{b} &= -\frac{1}{2} \bar{\mathbf{w}} \cdot [\mathbf{x}_r + \mathbf{x}_s] \end{aligned} \quad (62)$$

The Karush-Kuhn-Tucker (KKT) conditions that are satisfied by the solution are,

$$\bar{\alpha}_i \bar{\alpha}_i^* = 0, \quad i = 1, \dots, l. \quad (63)$$

Therefore the support vectors are points where exactly one of the Lagrange multipliers is greater than zero.

When $\varepsilon = 0$, we get the L_1 loss function and the optimisation problem is simplified,

$$\min_{\beta} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \beta_i \beta_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^l \beta_i y_i \quad (64)$$

with constraints,

$$\begin{aligned} -C \leq \beta_i \leq C, \quad i=1, \dots, l \\ \sum_{i=1}^l \beta_i = 0 \end{aligned} \quad , \quad (65)$$

and the regression function is given by Equation (56), where

$$\begin{aligned} \bar{\mathbf{w}} &= \sum_{i=1}^l \bar{\beta}_i \mathbf{x}_i \\ \bar{b} &= -\frac{1}{2} \bar{\mathbf{w}} \cdot [\mathbf{x}_r + \mathbf{x}_s] \end{aligned} \quad (66)$$

5.1.2 Quadratic Loss Function

Using a quadratic loss function, Figure 20(a),

$$L_{quad}(f(\mathbf{x}) - y) = (f(\mathbf{x}) - y)^2 \quad (67)$$

the solution is given by,

$$\max_{\alpha, \alpha^*} W(\alpha, \alpha^*) = \max_{\alpha, \alpha^*} \left\{ \begin{aligned} &-\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*)(\mathbf{x}_i \cdot \mathbf{x}_j) \\ &+ \sum_{i=1}^l (\alpha_i - \alpha_i^*) y_i - \frac{1}{2C} \sum_{i=1}^l (\alpha_i^2 + \alpha_i^{*2}) \end{aligned} \right\} \quad (68)$$

The corresponding optimisation can be simplified by exploiting the KKT conditions, Equation (63) and noting that these imply $\beta_i^* = |\beta_i|$. The resultant optimisation problems is,

$$\min_{\beta} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \beta_i \beta_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^l \beta_i y_i + \frac{1}{2C} \sum_{i=1}^l \beta_i^2 \quad (69)$$

with constraints,

$$\sum_{i=1}^l \beta_i = 0 \quad (70)$$

and the regression function is given by Equations (56) and (66).

5.1.3 Huber Loss Function

Using a Huber loss function, Figure 20(c),

$$L_{huber}(f(\mathbf{x}) - y) = \begin{cases} \frac{1}{2} (f(\mathbf{x}) - y)^2 & \text{for } |f(\mathbf{x}) - y| < \mu \\ \mu |f(\mathbf{x}) - y| - \frac{\mu^2}{2} & \text{otherwise} \end{cases} \quad (71)$$

the solution is given by,

$$\max_{\alpha, \alpha^*} W(\alpha, \alpha^*) = \max_{\alpha, \alpha^*} \left\{ \begin{aligned} & -\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*)(\mathbf{x}_i \cdot \mathbf{x}_j) \\ & + \sum_{i=1}^l (\alpha_i - \alpha_i^*) y_i - \frac{1}{2C} \sum_{i=1}^l (\alpha_i^2 + \alpha_i^{*2}) \mu \end{aligned} \right\}, \quad (72)$$

The resultant optimisation problems is,

$$\min_{\beta} \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \beta_i \beta_j (\mathbf{x}_i \cdot \mathbf{x}_j) - \sum_{i=1}^l \beta_i y_i + \frac{1}{2C} \sum_{i=1}^l \beta_i^2 \mu \quad (73)$$

with constraints,

$$\begin{aligned} -C \leq \beta_i \leq C, \quad i=1, \dots, l \\ \sum_{i=1}^l \beta_i = 0 \end{aligned} \quad (74)$$

and the regression function is given by Equations (56) and (66).

5.1.4 Example

Consider the example data set in Table 3.

X	Y
1.0	-1.6
3.0	-1.8
4.0	-1.0
5.6	1.2
7.8	2.2
10.2	6.8
11.0	10.0
11.5	10.0
12.7	10.0

Table 3 Regression Data

The SVR solution for an L_1 loss function with no additional capacity control is shown in Figure 21.

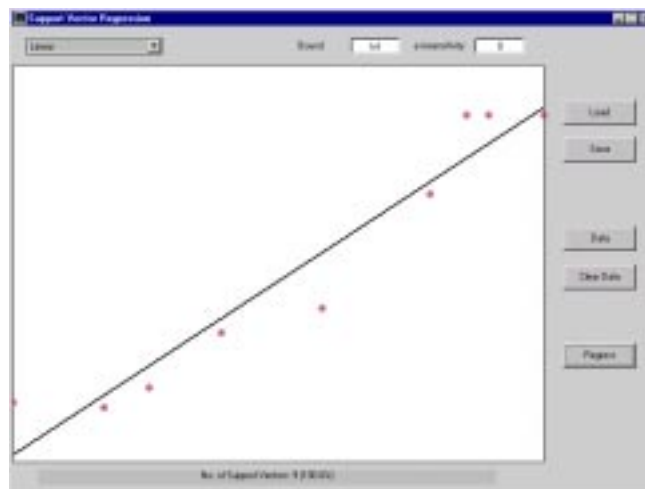


Figure 21 Linear regression

5.2 Non Linear Regression

Similarly to classification problems, a non-linear model is usually required to adequately model data. In the same manner as the non-linear SVC approach, a non-linear mapping can be used to map the data into a high dimensional feature space where linear regression is performed. The kernel approach is again employed to address the curse of dimensionality. The non-linear SVR solution, using an ε -insensitive loss function, Figure 20(c), is given by,

$$\max_{\alpha, \alpha^*} W(\alpha, \alpha^*) = \max_{\alpha, \alpha^*} \left\{ \begin{array}{l} \sum_{i=1}^l \alpha_i^* (y_i - \varepsilon) - \alpha_i (y_i + \varepsilon) \\ -\frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) K(\mathbf{x}_i, \mathbf{x}_j) \end{array} \right\} \quad (75)$$

with constraints,

$$\begin{aligned} 0 \leq \alpha_i \leq C, \quad i = 1, \dots, l \\ 0 \leq \alpha_i^* \leq C, \quad i = 1, \dots, l. \end{aligned} \quad (76)$$

$$\sum_{i=1}^l (\alpha_i^* - \alpha_i) = 0$$

Solving Equation (75) with constraints Equation (76) determines the Lagrange multipliers, α_i, α_i^* , and the regression function is given by,

$$f(\mathbf{x}) = \sum_{SVs} (\bar{\alpha}_i - \bar{\alpha}_i^*) K(\mathbf{x}_i, \mathbf{x}) + \bar{b} \quad (77)$$

where

$$\begin{aligned} \bar{\mathbf{w}} \cdot \mathbf{x} &= \sum_{SVs} (\bar{\alpha}_i - \bar{\alpha}_i^*) K(\mathbf{x}_i, \mathbf{x}) \\ \bar{b} &= -\frac{1}{2} \sum_{SVs} (\bar{\alpha}_i - \bar{\alpha}_i^*) [K(\mathbf{x}_r, \mathbf{x}_i) + K(\mathbf{x}_s, \mathbf{x}_i)] \end{aligned} \quad (78)$$

As with the SVC the equality constraint may be dropped if the Kernel contains a bias term, b being accommodated within the Kernel function, and the regression function is given by,

$$f(\mathbf{x}) = \sum_{SVs} (\bar{\alpha}_i - \bar{\alpha}_i^*) K(\mathbf{x}_i, \mathbf{x}). \quad (79)$$

The optimisation criteria for the other loss functions of section 5.1 are similarly obtained by replacing the dot product with a kernel function. The ε insensitive loss function is attractive because unlike the quadratic and Huber cost functions, where all the data points will be support vectors, the SV solution can be sparse. The quadratic loss function produces a solution which is equivalent to ridge regression, or zeroth order regularisation, where the regularisation parameter $\lambda = \frac{1}{2C}$.

5.2.1 Examples

To illustrate some of the non-linear SVR solutions, various kernel functions were used to model the regression data in Table 3, with an ε -insensitive loss function

($\epsilon=0.5$) and no additional capacity control. Figure 22 shows the SVR solution for a degree 2 polynomial, with the SV circled as before. The dotted line describes the ϵ -insensitive region around the solution (N.B. if all the data points lie within this region there will be zero error associated with the loss function). The result demonstrates that there are no support vectors within the ϵ -insensitive region.

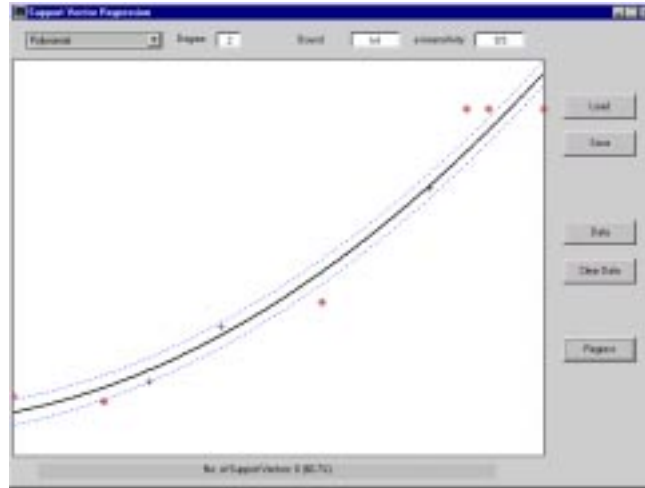


Figure 22 *Polynomial Regression*

Figure 23 illustrates the SVR solution for a radial basis function with $\sigma = 1.0$. In this example the model is flexible enough to model the function with zero error associated with the loss function, as is verified by the fact that all the data points lie on, or within, the ϵ -insensitive zone.

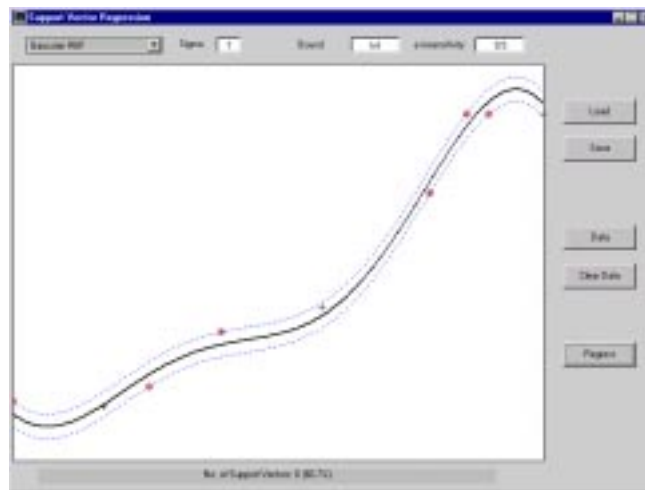


Figure 23 *Radial Basis Function Regression*

Figure 24 shows the SVR solution for a linear spline kernel. The resulting model is a piecewise cubic spline, and again due to the high capacity of this function it is able to model the data with zero loss function error, but notice that overfitting is controlled.

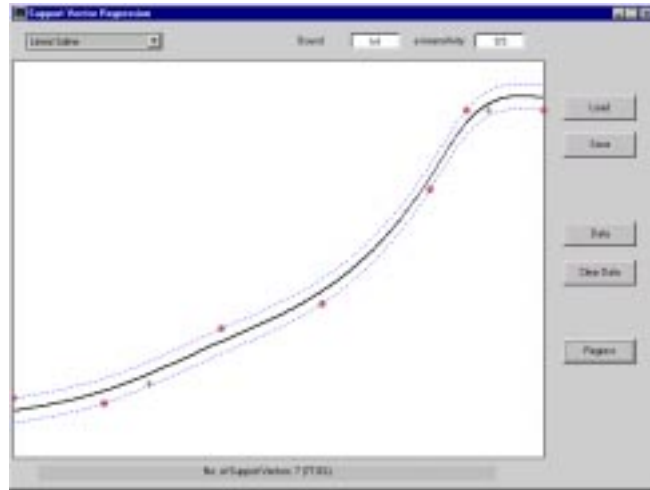


Figure 24 *Spline Regression*

Figure 25 shows the SVR solution for an infinite B-spline kernel, which has a similar solution to the spline kernel except for the endpoints.

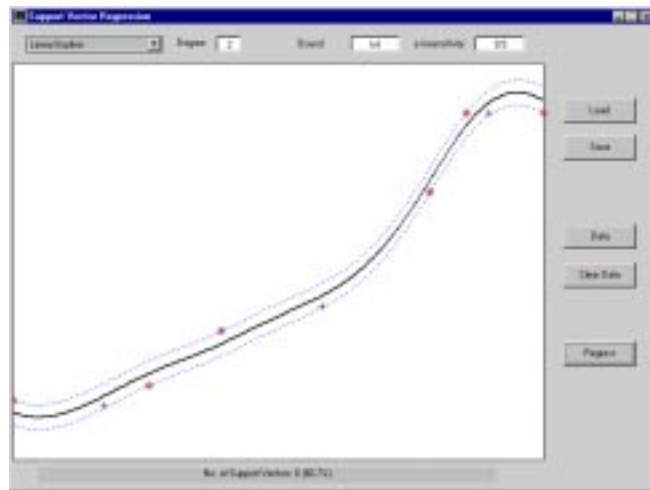


Figure 25 *B-spline regression*

Figure 26 shows the solution for an exponential RBF kernel, which is a piecewise linear spline. Although this model has a high capacity it shows sensible behaviour in the extremity regions.

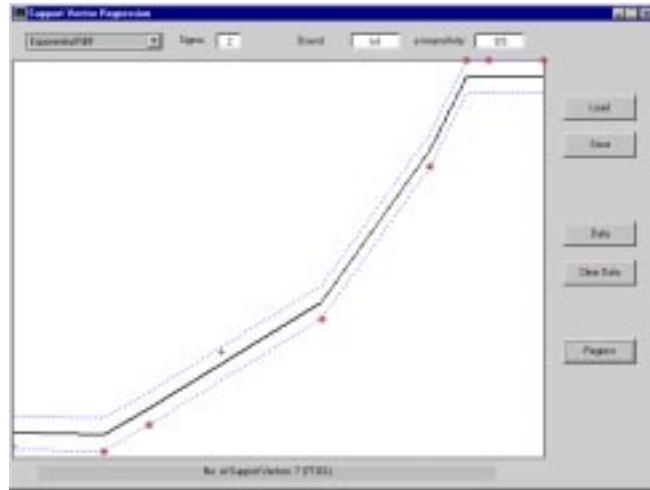


Figure 26 *Exponential RBF*

5.2.2 Comments

In the regression method it is necessary to select both a representative loss function and any additional capacity control that may be required. These considerations must be based on prior knowledge of the problem and the distribution of the noise. In the absence of such information Huber's robust loss function, Figure 20(d), has been shown to be the best alternative [21]. Vapnik developed the ϵ -insensitive loss function as a trade-off between the robust loss function of Huber and one that enables sparsity within the SVs. However, its implementation is more computationally expensive and the ϵ -insensitive region can have drawbacks, as will be demonstrated in the next section.

6 Regression Example: Titanium Data

The example given here considers the titanium data [6] as an illustrative example of a one dimensional non-linear regression problem. There are three methods for controlling the regression model, the loss function, the kernel, and additional capacity control, C . The results shown in this chapter were obtained using an ϵ -insensitive loss function ($\epsilon=0.05$), with different kernels and different degrees of capacity control.

Figure 27 illustrates the solution for a linear spline kernel and no additional capacity control. It is evident that the solution lies within the ϵ -insensitive region.

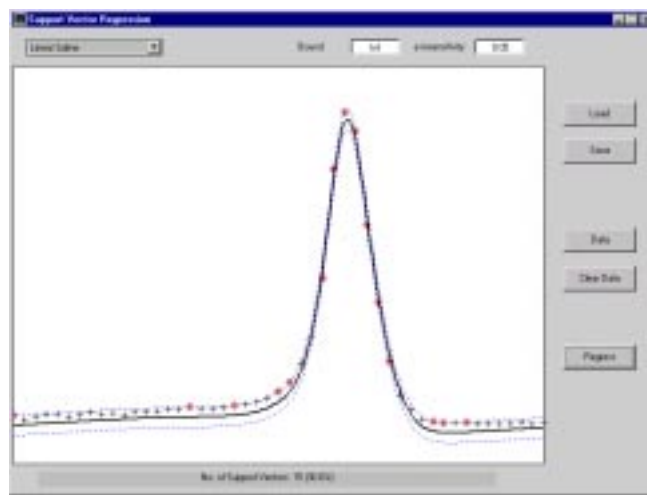


Figure 27 *Linear Spline Regression* ($\epsilon=0.05$, $C=\infty$)

Figure 28 illustrates the solution for a B-spline kernel with no additional capacity control. This particular B-spline kernel would appear to be prone to oscillation when an ϵ -insensitive region is used, and hence the linear spline kernel, or an alternative loss function is to be preferred.

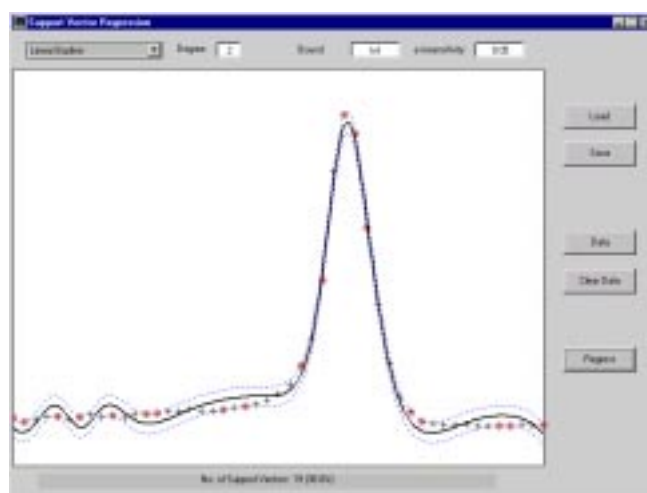


Figure 28 *B-Spline Regression* ($\epsilon=0.05$, $C=\infty$)

Figure 29 illustrates the solution for a Gaussian RBF kernel ($\sigma=1.0$) with no additional capacity control. It can be seen that the RBF is too wide to accurately model the data.

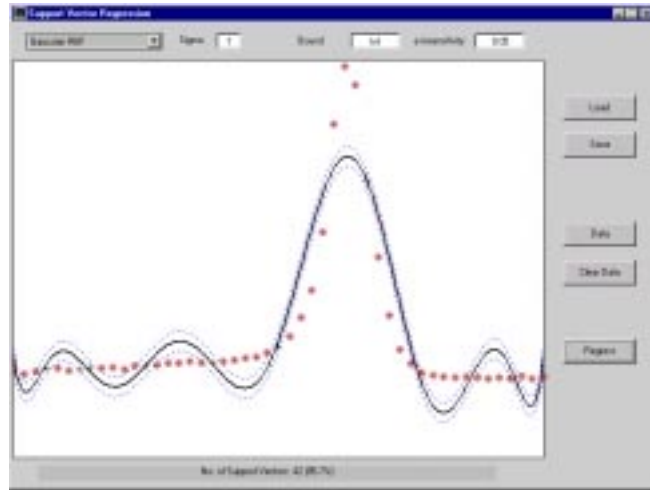


Figure 29 *Gaussian RBF Regression* ($\epsilon=0.05$, $\sigma=1.0$, $C=\infty$)

Figure 30 illustrates the solution for a Gaussian RBF kernel ($\sigma=0.3$) with no additional capacity control. It can be seen that the RBF is now able to accurately model the data. However, this is at the expense of the introduction of oscillation, which is not penalised in the ϵ -insensitive region.

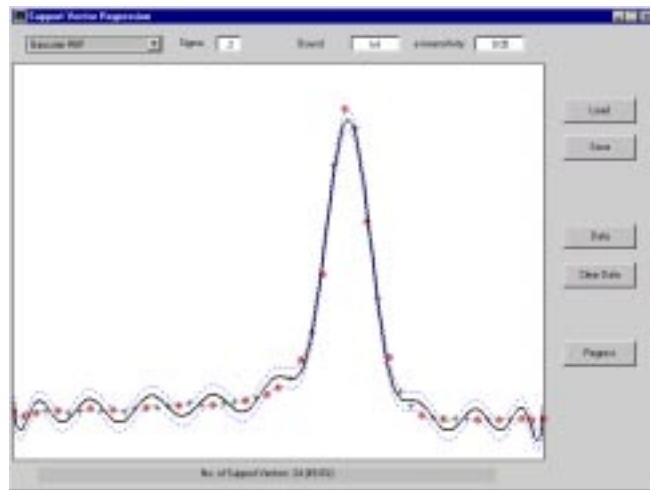


Figure 30 *Gaussian RBF Regression* ($\epsilon=0.05$, $\sigma=0.3$, $C=\infty$)

Figure 31 illustrates the solution for an exponential RBF kernel ($\sigma=0.3$) with no additional capacity control. The corresponding solution is a piece-wise linear function and consequently oscillation is avoided.

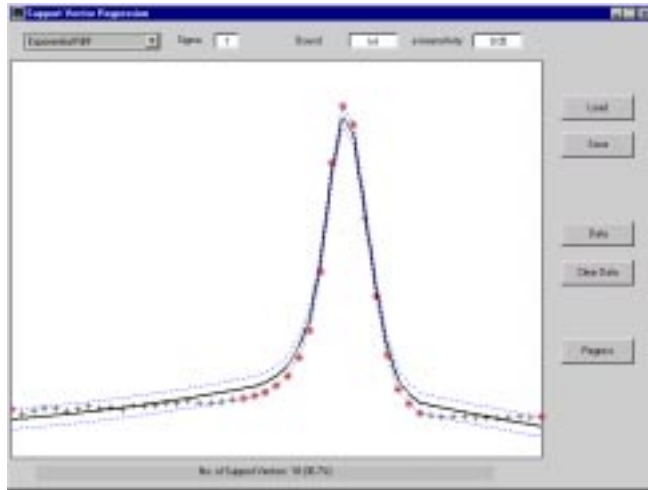


Figure 31 *Exponential RBF Regression* ($\epsilon=0.05$, $\sigma=1.0$)

Figure 32 illustrates the solution for a degree 3 Fourier kernel with no additional capacity control. The solution suffers similar problems to the wide Gaussian RBF kernel in that the kernel cannot accurately model the data.

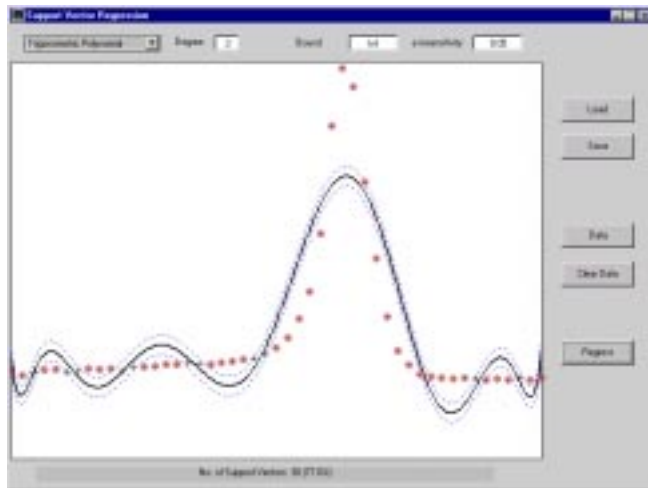


Figure 32 *Fourier Regression* ($\epsilon=0.05$, degree 3, $C=\infty$)

Figure 33 illustrates the solution for a linear spline kernel, with additional capacity control, ($C=10$). The extra capacity control renders the solution incapable of accurately modelling the peak in the data, in contrast to Figure 27.

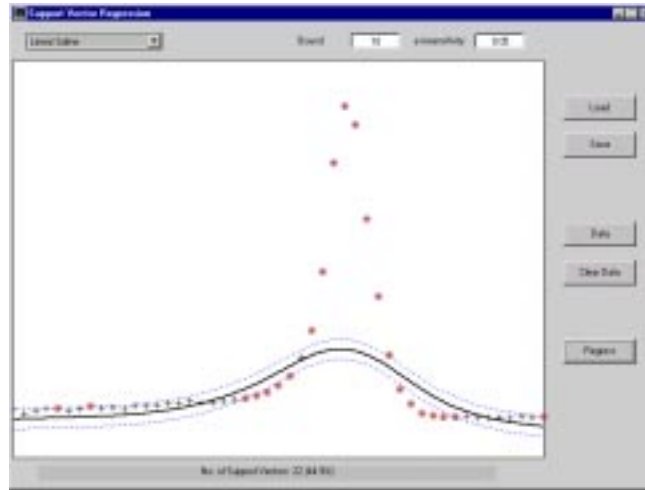


Figure 33 *Linear Spline Regression* ($\epsilon=0.05$, $C=10$)

Figure 34 illustrates the solution for a B-spline kernel, with additional capacity control, ($C=10$). The extra capacity control renders the solution incapable of accurately modelling the peak in the data, in contrast to Figure 28.

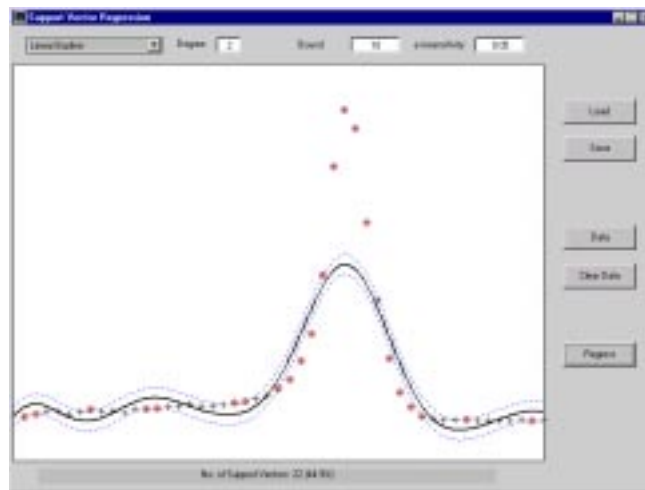


Figure 34 *B-Spline Regression* ($\epsilon=0.05$, $C=10$)

The examples that have been shown here are not a representative set. The ϵ -insensitive region has been exaggerated for the purposes of illustration, and typically careful selection of additional capacity control with methods such as cross validation will be required. The ϵ -insensitive loss function may be an inappropriate choice for particular kernels causing the solution to oscillate within the ϵ -insensitive region.

6.1 Applications

SVR has been applied to some time series modelling problems [12]. Notably, [13] has achieved excellent results in applying SVR to one of the data sets from the Santa Fe time series competition.

7 Conclusions

Support Vector Machines are an attractive approach to data modelling. They combine generalisation control with a technique to address the curse of dimensionality. The formulation results in a global quadratic optimisation problem with box constraints, which is readily solved by interior point methods. The kernel mapping provides a unifying framework for most of the commonly employed model architectures, enabling comparisons to be performed. In classification problems generalisation control is obtained by maximising the margin, which corresponds to minimisation of the weight vector in a canonical framework. The solution is obtained as a set of support vectors that can be sparse. These lie on the boundary and as such summarise the information required to separate the data. The minimisation of the weight vector can be used as a criterion in regression problems, with a modified loss function.

Future directions include: A technique for choosing the kernel function and additional capacity control; Development of kernels with invariances.

References

- [1] 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.
- [2] N. Aronszajn. Theory of Reproducing Kernels. *Trans. Amer. Math. Soc.*, **68**:337-404, 1950.
- [3] R. Bellman. *Dynamic Programming*. Princeton University Press, Princeton, NJ, 1957.
- [4] V. Blanz, B. Schölkopf, H. Bülthoff, C. Burges, V. Vapnik, and T. Vetter. 1996. 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 (eds.): *Artificial Neural Networks - ICANN'96*. Springer Lecture Notes in Computer Science Vol. 1112, Berlin, 251-256.
- [5] C. Cortes, and V. Vapnik. 1995. Support Vector Networks. *Machine Learning* 20:273-297.
- [6] P. Dierckx. *Curve and Surface Fitting with Splines*. Monographs on Numerical Analysis, Clarendon Press, Oxford, 1993.
- [7] F. Girosi. An Equivalence Between Sparse Approximation and Support Vector Machines. Technical Report AIM-1606, Artificial Intelligence Laboratory, Massachusetts Institute of Technology (MIT), Cambridge, Massachusetts, 1997.
- [8] S. R. Gunn, M. Brown and K. M. Bossley. Network Performance Assessment for Neurofuzzy Data Modelling. *Lecture Notes in Computer Science*, **1280**:313-323, 1997.
- [9] J. Hadamard. *Lectures on the Cauchy Problem in Linear Partial Differential Equations*. Yale University Press, 1923.
- [10] N. E. Heckman. The Theory and Application of Penalized Least Squares Methods or Reproducing Kernel Hilbert Spaces Made Easy, 1997. <ftp://newton.stat.ubc.ca/pub/nancy/PLS.ps>
- [11] M. Minoux. *Mathematical Programming: Theory and Algorithms*. John Wiley and Sons, 1986.
- [12] S. Mukherjee, E. Osuna, and F. Girosi. Nonlinear Prediction of Chaotic Time Series using Support Vector Machines. To appear in *Proc. of IEEE NNSP'97*, Amelia Island, FL, 24-26 Sep., 1997.

- [13] K. R. Müller, A. Smola, G. Rätsch, B. Schölkopf, J. Kohlmorgen and V. Vapnik. 1997. Predicting Time Series with Support Vector Machines. Submitted to: *ICANN'97*.
- [14] E. Osuna, R. Freund and F. Girosi. An Improved Training Algorithm for Support Vector Machines. To appear in *Proc. of IEEE NNSP'97*, Amelia Island, FL, 24-26 Sep., 1997.
- [15] E. Osuna, R. Freund and F. Girosi. 1997. Improved Training Algorithm for Support Vector Machines. To appear in: *NNSP'97*.
- [16] T. Poggio, V. Torre and C. Koch. Computational Vision and Regularization Theory. *Nature*, **317**(26):314-319, 1985.
- [17] A. Smola and B. Schölkopf. 1997. On a Kernel-based Method for Pattern Recognition, Regression, Approximation and Operator Inversion. GMD Technical Report No. 1064. (To appear in *Algorithmica*, special issue on Machine Learning)
- [18] M. O. Stitson and J. A. E. Weston. Implementational Issues of Support Vector Machines. Technical Report CSD-TR-96-18, Computational Intelligence Group, Royal Holloway, University of London, 1996.
- [19] M. O. Stitson, J. A. E. Weston, A. Gammernan, V. Vovk and V. Vapnik. Theory of Support Vector Machines. Technical Report CSD-TR-96-17, Computational Intelligence Group, Royal Holloway, University of London, 1996.
- [20] V. Vapnik, S. Golowich and A. Smola. 1997. Support Vector Method for Function Approximation, Regression Estimation, and Signal Processing. In: M. Mozer, M. Jordan, and T. Petsche (eds.): *Neural Information Processing Systems*, Vol. 9. MIT Press, Cambridge, MA, 1997.
- [21] V. Vapnik. 1995. *The Nature of Statistical Learning Theory*. Springer-Verlag, New York.
- [22] G. Wahba. *Spline Models for Observational Data*. Series in Applied Mathematics, Vol. 59, SIAM, Philadelphia, 1990.

Appendix - Implementation Issues

The resulting optimisation problems are dependent upon the number of training examples. As such, when this data is large methods have been proposed for speeding up the algorithm by decomposing the problem into smaller ones. Approximate [18] and exact [15] methods have been proposed.

MATLAB implementations of the main support vector routines are shown below. Note that these routines are not optimised in any sense! Typically the quadratic matrix, H , is badly conditioned which can render quadratic program optimisers incapable of producing an accurate solution. To address this, a quick fix of using zero order regularisation can be used, but too large a value will perturb the solution significantly. (Note that this is the capacity control used in some of the SVR routines.)

Support Vector Classification

The optimisation problem can be expressed in matrix notation as,

$$\min_x \frac{1}{2} \alpha^T H \alpha + c^T \alpha \quad (80)$$

where

$$H = ZZ^T, \quad c^T = (-1, \dots, -1) \quad (81)$$

with constraints

$$\alpha^T Y = 0, \quad \alpha_i \geq 0, i = 1, \dots, l. \quad (82)$$

where

$$Z = \begin{bmatrix} y_1 \mathbf{x}_1 \\ \vdots \\ y_l \mathbf{x}_l \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_l \end{bmatrix} \quad (83)$$

The MATLAB implementation is given below:

```
function [nsv, alpha, b0] = svc(X,Y,ker,C)
%SVC Support Vector Classification
%
% Usage: [nsv alpha bias] = svc(X,Y,ker,C)
%
% Parameters: X      - Training inputs
%             Y      - Training targets
%             ker     - kernel function
%             C      - upper bound (non-separable case)
%             nsv    - number of support vectors
%             alpha  - Lagrange Multipliers
%             b0     - bias term
%
% Author: Steve Gunn (srg@ecs.soton.ac.uk)

if (nargin < 2 | nargin > 4) % check correct number of arguments
```

```

help svc
else

n = size(X,1);
if (margin<4) C=Inf; end
if (margin<3) ker='linear'; end
epsilon = 1e-10;

% Construct the H matrix and c vector

H = zeros(n,n);
for i=1:n
    for j=1:n
        H(i,j) = Y(i)*Y(j)*svkernel(ker,X(i,:),X(j,:));
    end
end
c = -ones(n,1);

% Add small amount of zero order regularisation to
% avoid problems when Hessian is badly conditioned.

if (abs(cond(H)) > 1e+10)
    fprintf('Hessian badly conditioned, regularising ...\n');
    fprintf('    Old condition number: %4.2g\n',cond(H));
    H = H+0.00000001*eye(size(H));
    fprintf('    New condition number: %4.2g\n',cond(H));
end

% Set up the parameters for the Optimisation problem

vlb = zeros(n,1);      % Set the bounds: alphas >= 0
vub = C*ones(n,1);    %                alphas <= C
x0 = [ ];             % The starting point is [0 0 0  0]
neqcstr = nobias(ker); % Set the number of equality constraints (1 or 0)
if neqcstr
    A = Y'; b = 0;    % Set the constraint Ax = b
else
    A = []; b = [];
end

% Solve the Optimisation Problem

st = cputime;

if ( vlb == zeros(size(vlb)) & min(vub) == Inf & neqcstr == 0 )
    % Separable problem with Implicit Bias term
    % Use Non Negative Least Squares
    alpha = fnnls(H,-c);
else
    % Otherwise
    % Use Quadratic Programming
    alpha = qp(H, c, A, b, vlb, vub, x0, neqcstr, -1);
end

fprintf('Execution time: %4.1f seconds\n',cputime - st);
fprintf('|w0|^2      : %f\n',alpha'*H*alpha);
fprintf('Sum alpha : %f\n',sum(alpha));

% Compute the number of Support Vectors
svi = find( abs(alpha) > epsilon);
nsv = length(svi);

if neqcstr == 0
    % Implicit bias, b0
    b0 = 0;
else
    % Explicit bias, b0;
    % find b0 from pair of support vectors, one from each class
    classAsvi = find( abs(alpha) > epsilon & Y == 1);
    classBsvi = find( abs(alpha) > epsilon & Y == -1);
    nAsv = length( classAsvi );
    nBsv = length( classBsvi );
    if ( nAsv > 0 & nBsv > 0 )
        svpair = [classAsvi(1) classBsvi(1)];
        b0 = -(1/2)*sum(Y(svpair)'*H(svpair)*alpha(svi));
    end
end

```



```
    else
      b0 = 0;
    end
  end
end
end
```

Support Vector Regression

The optimisation problem for an ε -insensitive loss function can be expressed in matrix notation as,

$$\min_x \frac{1}{2} x^T H x + c^T x \quad (84)$$

where

$$H = \begin{bmatrix} XX^T & -XX^T \\ -XX^T & XX^T \end{bmatrix}, \quad c = \begin{bmatrix} \varepsilon + Y \\ \varepsilon - Y \end{bmatrix}, \quad x = \begin{bmatrix} \alpha \\ \alpha^* \end{bmatrix} \quad (85)$$

with constraints

$$x \cdot (1, \dots, 1, -1, \dots, -1) = 0, \quad \alpha_i, \alpha_i^* \geq 0, i = 1, \dots, l. \quad (86)$$

where

$$X = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_l \end{bmatrix}, \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_l \end{bmatrix} \quad (87)$$

The MATLAB implementation is given below:

```
function [nsv, beta, b0] = svr(X,Y,ker,e,C)
%SVR Support Vector Regression
%
% Usage: alpha = svr(X,Y,ker,e,C)
%
% Parameters: X      - Training inputs
%              Y      - Training targets
%              ker    - kernel function
%              e      - insensitivity
%              C      - upper bound (non-separable case)
%              nsv    - number of support vectors
%              beta   - Difference of Lagrange Multipliers
%              b0     - bias term
%
% Author: Steve Gunn (srg@ecs.soton.ac.uk)

if (nargin < 3 | nargin > 5) % check correct number of arguments
    help svr
else

    n = size(X,1);
    if (nargin < 5) C=Inf; end
    if (nargin < 4) e=0.05; end
    if (nargin < 3) ker='linear'; end
    epsilon = 1e-10; % tolerance for Support Vector Detection

    % Construct the H matrix and c vector

    H = zeros(n,n);
    for i=1:n
        for j=1:n
            H(i,j) = svkernel(ker,X(i,:),X(j,:));
        end
    end
    end
    Hb = [H -H; -H H];
```

```

c = [(e*ones(n,1) + Y); (e*ones(n,1) - Y)];

% Add small amount of zero order regularisation to
% avoid problems when Hessian is badly conditioned.
% Rank is always less than or equal to n.
% Note that adding to much reg will perturb solution

if (abs(cond(Hb)) > 1e+10)
    fprintf('Hessian badly conditioned, regularising ....\n');
    fprintf('    Old condition number: %4.2g\n',cond(Hb));
    Hb = Hb+0.000000000001*eye(size(Hb));
    fprintf('    New condition number: %4.2g\n',cond(Hb));
end

% Set up the parameters for the Optimisation problem

vlb = zeros(2*n,1); % Set the bounds: alphas >= 0
vub = C*ones(2*n,1); %                    alphas <= C
x0 = [ ]; % The starting point is [0 0 0 0]
neqcstr = nobias(ker); % Set the number of equality constraints (1 or 0)
if neqcstr
    A = [ones(1,n) -ones(1,n)]; b = 0; % Set the constraint Ax = b
else
    A = []; b = [];
end

% Solve the Optimisation Problem

st = cputime;

if ( vlb == zeros(size(vlb)) & min(vub) == Inf & neqcstr == 0 )
    % Separable problem with Implicit Bias term
    % Use Non Negative Least Squares
    alpha = fnnlsl(Hb,-c);
else
    % Otherwise
    % Use Quadratic Programming
    alpha = qp(Hb, c, A, b, vlb, vub, x0, neqcstr, -1);
end

fprintf('Execution time: %4.1f seconds\n',cputime - st);
fprintf('|w0|^2 : %f\n',alpha'*Hb*alpha);
fprintf('Sum alpha : %f\n',sum(alpha));

% Compute the number of Support Vectors
beta = alpha(1:n) - alpha(n+1:2*n);
svi = find( abs(beta) > epsilon );
nsv = length( svi );

if neqcstr == 0
    % Implicit bias, b0
    b0 = 0;
else
    % Explicit bias, b0;
    % compute using robust method of Smola
    % find b0 from average of support vectors with interpolation error e
    svbi = find( abs(beta) > epsilon & abs(beta) < C );
    nsvb = length(svbi);
    if nsvb > 0
        b0 = (1/nsvb)*sum(Y(svbi) + e*sign(beta(svbi)) - H(svbi,svi)*beta(svi));
    else
        b0 = (max(Y)+min(Y))/2;
    end
end
end
end

```

MATLAB SVM Toolbox

A MATLAB toolbox implementing SVM is freely available for academic purposes:

- 1) Download it from: <http://www.isis.ecs.soton.ac.uk/research/svm/>
- 2) Extract the tar file *svm.tar* under the matlab toolbox directory
- 3) Add/matlab/toolbox/svm to your MATLAB path.
- 4) Type *help svm* at the MATLAB prompt for help.

The two main user interfaces are for 2D classification and 1D regression (*uiclass* and *uiregress* respectively).