

The Genetic Kernel Support Vector Machine: Description and Evaluation

TOM HOWLEY & MICHAEL G. MADDEN*

*Department of Information Technology, National University of Ireland, Galway, Ireland
(*author for correspondence, e-mail: michael.madden@nuigalway.ie)*

Abstract. The Support Vector Machine (SVM) has emerged in recent years as a popular approach to the classification of data. One problem that faces the user of an SVM is how to choose a kernel and the specific parameters for that kernel. Applications of an SVM therefore require a search for the optimum settings for a particular problem. This paper proposes a classification technique, which we call the Genetic Kernel SVM (GK SVM), that uses Genetic Programming to evolve a kernel for a SVM classifier. Results of initial experiments with the proposed technique are presented. These results are compared with those of a standard SVM classifier using the Polynomial, RBF and Sigmoid kernel with various parameter settings.

Keywords: classification, genetic Kernel SVM, genetic programming, Mercer Kernel, model selection, support vector machine

1. Introduction

The SVM is a powerful machine learning tool that is capable of representing non-linear relationships and producing models that generalise well to unseen data. SVMs initially came into prominence in the area of hand-written character recognition (Boser et al., 1992a) and are now being applied to many other areas, e.g. text categorisation (Hearst, 1998; Joachims, 1998) and computer vision (Osuna et al., 1997). An advantage that SVMs have over the widely-used Artificial Neural Network (ANN) is that they typically don't possess the same potential for instability as ANNs do with the effects of different random starting weights (Bleckmann and Meiler, 2003).

Despite this, using an SVM requires a certain amount of model selection. According to Cristianini et al. (1998), "One of the most important design choices for SVMs is the kernel-parameter, which implicitly defines the structure of the high dimensional feature space where a maximal margin hyperplane will be found. Too rich a feature space would cause the system to overfit the data, and conversely the system might not be capable of separating the data if the kernels are

too poor.” However, before this stage is reached in the use of SVMs, the actual kernel must be chosen and, as the experimental results of this paper show, different kernels may exhibit vastly different performance. This paper describes a technique which attempts to alleviate this selection problem by using genetic programming (GP) to evolve a suitable kernel for a particular problem domain. We call our technique the Genetic Kernel SVM (GK SVM).

Section 2 outlines the theory behind SVM classifiers with a particular emphasis on kernel functions. Section 3 gives a very brief overview of genetic programming. Section 4 describes the proposed technique for the evolution of SVM kernels. Experimental results are presented in Section 5. Some related research is described in Section 6. Finally, Section 7 presents the conclusions.

2. Support Vector Machine Classification

The problem of classification can be represented as follows. Given a set of input-output pairs $Z = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_\ell, y_\ell)\}$, construct a classifier function f that maps the input vectors $\mathbf{x} \in X$ onto labels $y \in Y$. In binary classification the set of labels is simply $Y = \{-1, 1\}$. The goal is to find a classifier $f \in F$ which will correctly classify new examples (\mathbf{x}, y) , i.e., $f(\mathbf{x}) = y$ for examples (\mathbf{x}, y) , which were generated under the same probability distribution as the data (Scholkopf, 1998). Binary classification is frequently performed by finding a hyperplane that separates the data, e.g., Linear Discriminant Analysis (LDA) (Hastie et al., 2001). There are two main issues with using a separating hyperplane:

1. The problem of learning this hyperplane is an ill-posed one because there is not a unique solution and many solutions may not generalise well to the unseen examples.
2. The data might not be linearly separable.

SVMs tackle the first problem by finding the hyperplane that realises the maximum margin of separation between the classes (Cristianini and Shawe-Taylor, 2000). A representation of the hyperplane solution used to classify a new sample x_i is:

$$f(x) = \langle w, x_i \rangle + b \quad (1)$$

where $\langle w, x_i \rangle$ is the dot-product of the weight vector w and the input sample, and b is a bias value. The value of each element of w can be

viewed as a measure of the relative importance of each of the sample attributes for the classification of a sample. It has been shown that the optimal hyperplane can be uniquely constructed by solving the following constrained quadratic optimisation problem (Boser et al., 1992b):

$$\text{Minimise } \langle w, w \rangle + C \sum_{i=1}^{\ell} \xi_i \quad (2a)$$

$$\text{subject to } \begin{cases} y_i(\langle w, x_i \rangle + b) \geq 1 - \xi_i, & i = 1, \dots, \ell \\ \xi_i \geq 0, & i = 1, \dots, \ell \end{cases} \quad (2b)$$

This optimisation problem minimises the norm of the vector w which increases the *flatness* (or reduces the complexity) of the resulting model and thereby improves its generalisation ability. With *hard-margin* optimisation the goal is simply to find the minimum $\langle w, w \rangle$ such that the hyperplane $f(x)$ successfully separates all ℓ samples of the training data. The slack variables ξ_i are introduced to allow for finding a hyperplane that misclassifies some of the samples (*soft-margin* optimisation) as many datasets are not linearly separable. The complexity constant $C > 0$ determines the trade-off between the flatness and the amount by which misclassified samples are tolerated. A higher value of C means that more importance is attached to minimising the slack variables than to minimising $\langle w, w \rangle$. Rather than solving this problem in its primal form of (2a) and (2b), it can be more easily solved in its dual formulation (Cristianini and Shawe-Taylor, 2000):

$$\text{Maximise } W(\alpha) = \sum_{i=1}^{\ell} \alpha_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \alpha_i \alpha_j y_i y_j \langle x_i, x_j \rangle \quad (3a)$$

$$\text{subject to } C \geq \alpha_i \geq 0, \quad \sum_{i=1}^{\ell} \alpha_i y_i = 0 \quad (3b)$$

Instead of finding w and b the goal now is find the vector α and bias value b , where each α_i represents the relative importance of a training sample i in the classification of a new sample. To classify a new sample, the quantity $f(x)$ is calculated as:

$$f(x) = \sum_i \alpha_i y_i \langle x, x_i \rangle + b \quad (4)$$

where b is chosen so that $y_i f(x) = 1$ for any i with $C > \alpha_i > 0$. Then, a new sample x_s is classed as negative if $f(x_s)$ is less than zero and positive if $f(x_s)$ is greater than or equal to zero. Samples x_i for which the

corresponding α_i are non-zero are known as *support vectors* since they lie closest to the separating hyperplane. Samples that are not support vectors have no influence on the decision function. In (3b) C places an upper bound (known as the box constraint) on the value that each α_i can take. This limits the influence of outliers, which would otherwise have large α_i values (Cristianini and Shawe-Taylor, 2000).

Training an SVM entails solving the quadratic programming problem of (3a) and (3b). There are many standard techniques that could be applied to SVMs, including the Newton method, conjugate gradient and primal-dual interior-point methods (Cristianini and Shawe-Taylor, 2000). For the experiments reported here the SVM implementation uses the Sequential Minimisation Optimisation (SMO) algorithm of Platt (1999).

2.1. Kernel functions

One key aspect of the SVM model is that the data enters the above expressions (3a and 4) only in the form of the dot product of pairs. This leads to the resolution of the second problem mentioned above, namely that of non-linearly separable data. The basic idea with SVMs is to map the training data into a higher dimensional feature space via some mapping $\phi(x)$ and construct a separating hyperplane with maximum margin there. This yields a non-linear decision boundary in the original input space. By use of a kernel function, $K(x, z) = \langle \phi(x), \phi(z) \rangle$, it is possible to compute the separating hyperplane without explicitly carrying out the mapping into feature space (Scholkopf, 2000). Typical choice for kernels are:

- Linear Kernel: $K(x, z) = \langle x, z \rangle$
- Polynomial Kernel: $K(x, z) = (\langle x, z \rangle)^d$
- RBF Kernel: $K(x, z) = \exp\left(\frac{-\|x-z\|^2}{2\sigma^2}\right)$
- Sigmoid Kernel: $K(x, z) = \tanh(\gamma * \langle x, z \rangle - \theta)$

Each kernel corresponds to some feature space and because no explicit mapping to this feature space occurs, optimal linear separators can be found efficiently in feature spaces with millions of dimensions (Russell and Norvig, 2003). Note that the Linear Kernel is equivalent to a Polynomial Kernel of degree one and corresponds to the original input space. An alternative to using one of the pre-defined kernels is to derive a custom kernel that may be suited to a particular problem, e.g. the string kernel used for text classification by Lodhi

et al. (2002). To ensure that a kernel function actually corresponds to some feature space it must be symmetric, i.e., $K(x, z) = \langle \phi(x), \phi(z) \rangle = \langle \phi(z), \phi(x) \rangle = K(z, x)$. Typically, kernels are also required to satisfy Mercer's theorem, which states that the matrix $K = (K(x_i, x_j))_{i,j=1}^n$ must be positive semi-definite, i.e., it has no non-negative eigenvalues (Cristianini and Shawe-Taylor, 2000). This condition ensures that the solution of (3a) and (3b) produces a global optimum. However, good results have been achieved with non-Mercer kernels, and convergence is expected when the SMO algorithm is used, despite no guarantee of optimality when non-Mercer kernels are used (Bahlmann et al., 2002). Furthermore, despite its wide use, the Sigmoid kernel matrix is not positive semi-definite for certain values of the parameters γ and θ (Lin and Lin, 2003).

3. Genetic Programming

A GP is an application of the genetic algorithm (GA) approach to derive mathematical equations, logical rules or program functions automatically (Koza, 1992). Rather than representing the solution to a problem as a string of parameters, as in a conventional GA, a GP usually uses a tree structure, the leaves of which represent input variables or numerical constants. Their values are passed to nodes, which perform some numerical or program operation before passing on the result further towards the root of the tree. The GP typically starts off with a random population of individuals, each encoding a function or expression. This population is evolved by selecting better individuals for recombination and using their offspring to create a new population (generation). Mutation is employed to encourage discovery of new individuals. This process is continued until some stopping criteria is met, e.g., homogeneity of the population.

4. Genetic Evolution of Kernels

The approach presented here combines the two techniques of SVMs and GP, using the GP to evolve a kernel for a SVM. The goal is to eliminate the need for testing various kernels and their parameter settings. With this approach it might also be possible to discover new kernels that are particularly useful for the type of data under analysis. An overview of the proposed GK SVM is shown in Figure 1.

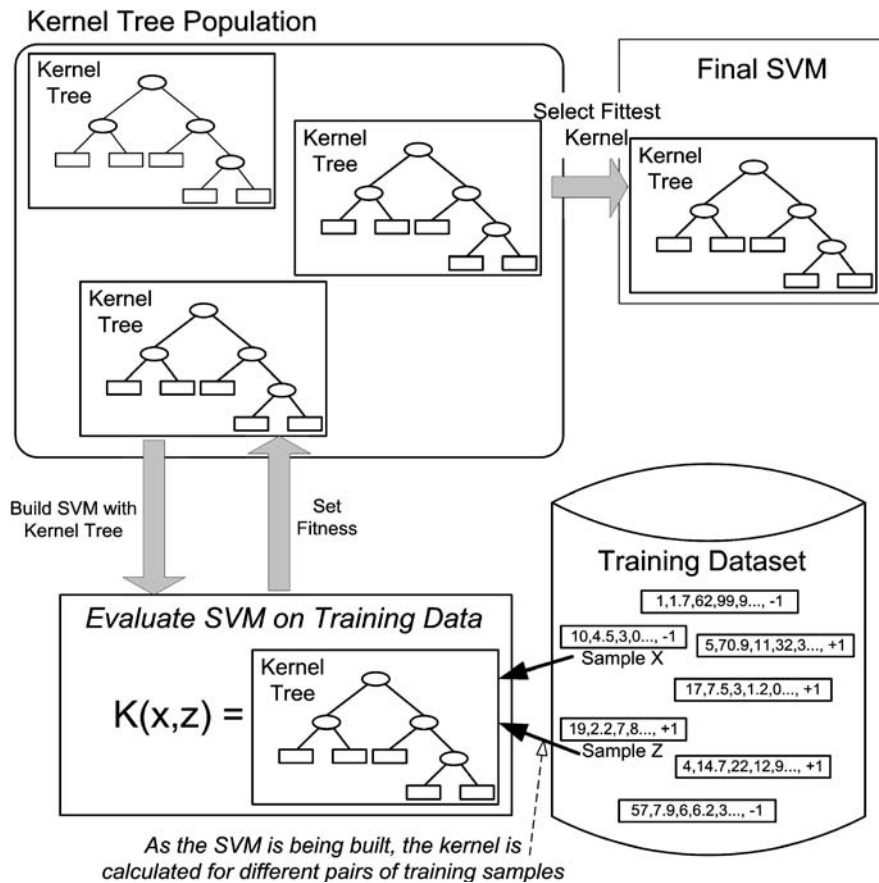


Figure 1. The Genetic Kernel SVM.

The main steps in the building of a GK SVM are:

1. Create a random population of kernel functions, represented as trees—we call these *kernel trees*.
2. Evaluate the fitness of each individual by building an SVM from the kernel tree and test it on the training data.
3. Select the fitter kernel trees as parents for recombination.
4. Perform random mutation on the newly created offspring.
5. Replace the old population with the offspring.
6. Repeat Steps 2 to 5 until the population has converged.
7. Build final SVM using the fittest kernel tree found.

The *Grow* method (Banzhaf et al., 1998) is used to initialise the population of trees, each tree being grown until no more leaves could be expanded (i.e., all leaves are terminals) or until a preset initial

maximum depth (2 for the experiments reported here) is reached. Rank-based selection is employed with a crossover probability of 0.9. Mutation with probability 0.2 is carried out on offspring by randomly replacing a sub-tree with a newly generated (via Grow method) tree. To prevent the proliferation of massive tree structures, pruning is carried out on trees after crossover and mutation, maintaining a maximum depth of 12. In the experiments reported here, five populations are evolved in parallel and the best individual over all populations is selected after all populations have converged. This reduces the likelihood of the procedure converging on a poor solution.

4.1. Terminal & function set

In the construction of kernel trees the approach adopted was to use the entire sample vector as input. An example of a kernel tree is shown in Figure 2 (Section 5). Since a kernel function only operates on two samples the resulting terminal set comprises only two vector elements: x and z . The evaluation of a kernel on a pair of samples is:

$$K(x, z) = \langle \text{treeEval}(x, z), \text{treeEval}(z, x) \rangle \tag{5}$$

The kernel is first evaluated on the two samples x and z . These samples are swapped and the kernel is evaluated again. The dot-product

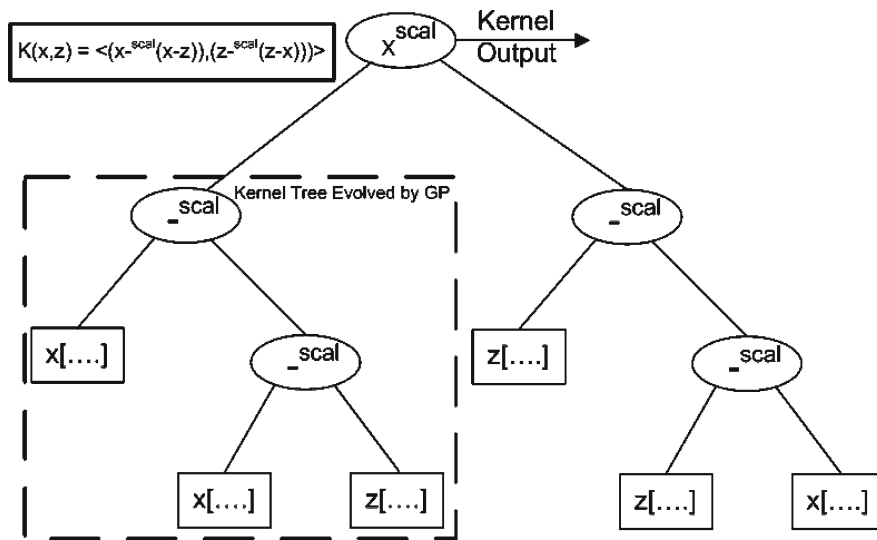


Figure 2. Example of a Kernel found on the Wisconsin Breast Cancer Dataset.

of these two evaluations is returned as the kernel output. This current approach produces symmetric kernels, but does not guarantee that they obey Mercer's theorem. Ensuring that such a condition is met would add considerable time to kernel fitness evaluation and, as stated earlier, using a non-Mercer kernel does not preclude finding a good solution.

The use of vector inputs requires corresponding vector operators to be used as functions in the kernel tree. The design employed uses two versions of the $+$, $-$ and \times mathematical functions: *scalar* and *vector*. Scalar functions return a single scalar value regardless of the operand's type, e.g., $x *^{\text{scal}} z$ calculates the dot-product of the two vectors. For the two other operators ($+$ and $-$) the operation is performed on each pair of elements and the *magnitude* of the resulting vector is returned as the output. Vector functions return a vector provided at least one of the inputs is a vector. For the vector versions of addition and subtraction (e.g., $x +^{\text{vect}} z$) the operation is performed on each pair of elements as with the scalar function, but in this case the resulting vector is returned as the output. No multiplication operator that returns a vector is used. If two inputs to a vector function are scalar (as could happen in the random generation of a kernel tree) then it behaves as the scalar operator. If only one input is scalar then that input is treated as a vector of the same length as the other vector operand with each element set to the same original scalar value.

4.2. Fitness function

Another key element to this approach (and to any evolutionary approach) is the choice of fitness function. An obvious choice for the fitness estimate is the classification error on the training set, but there is a danger that this estimate might produce SVM kernel tree models that are overfitted to the training data. One alternative is to base the fitness on a cross-validation test (e.g., leave-one-out cross-validation) in order to give a better estimation of a kernel tree's ability to produce a model that generalises well to unseen data. However, this would obviously increase computational effort greatly. Therefore, our solution (after experimenting with a number of alternatives) is to use a tiebreaker to limit overfitting. The fitness function used is:

$$\text{fitness}(\text{tree}) = \text{Error, with tiebreaker: fitness} = \sum \alpha_i * R^2 \quad (6)$$

This firstly differentiates between kernel trees based on their training error. For kernel trees of equal training error, a second evaluation is used as a tiebreaker. This is based on the sum of the support vector values, $\sum \alpha_i$ ($\alpha_i = 0$ for non-support vectors). The rationale behind this fitness estimate is based on the following definition of the geometric margin of a hyperplane, γ (Cristianini and Shawe-Taylor, 2000):

$$\gamma = \left(\sum_{i \in sv} \alpha_i \right)^{-\frac{1}{2}} \quad (7)$$

Therefore, the smaller the sum of the α_i s, the bigger the margin and the smaller the chance of overfitting to the training data. The fitness function also incorporates a penalty corresponding to R , the radius of the smallest hypersphere, centred at the origin, that encloses the training data in feature space. R is computed as (Cristianini and Shawe-Taylor, 2000):

$$R = \max_{1 \leq i \leq \ell} (K(x_i, x_i)) \quad (8)$$

where ℓ is the number of samples in the training dataset. This fitness function therefore favours a kernel tree that produces a SVM with a large margin relative to the radius of its feature space.

5. Experimental Results

Table 1 shows the performance of the GK SVM classifier compared with three commonly used SVM kernels, Polynomial, RBF and Sigmoid, on a number of datasets. (These are the only datasets with which the GK SVM has been evaluated to date.) The first four datasets contain the Raman spectra for 24 sample mixtures, made up of different combinations of the following four solvents: Acetone, Cyclohexanol, Acetonitrile and Toluene; see Hennessy et al. (2005) for a description of the dataset. The classification task considered here is to identify the presence or absence of one of these solvents in a mixture. For each solvent, the dataset was divided into a training set of 14 samples and a validation set of 10. The validation set in each case contained five positive and five negative samples. The final two datasets, Wisconsin Breast Cancer Prognosis (WBCP) and Glass2, are readily available from the UCI machine learning database repository

Table 1. Percentage classification accuracy of GK SVM compared to that of Polynomial, RBF and Sigmoid Kernel SVM on six datasets

Classifier	Dataset					
	Acetone	Cyclohexanol	Acetonitrile	Toluene	WBCP	Glass2
Kernel - degree d						
1	100.00	100.00	100.00	90.00	78.00	62.00
2	90.00	90.00	100.00	90.00	77.00	70.91
3	50.00	90.00	100.00	60.00	86.00	78.18
4	50.00	50.00	50.00	50.00	87.00	74.55
5	50.00	50.00	50.00	50.00	84.00	76.36
RBF kernel σ						
0.0001	50.00	50.00	50.00	50.00	78.00	58.18
0.001	50.00	90.00	50.00	50.00	78.00	58.18
0.01	60.00	80.00	50.00	60.00	78.00	59.64
0.1	50.00	50.00	50.00	50.00	78.00	63.64
1	50.00	50.00	50.00	50.00	81.00	70.91
10	50.00	50.00	50.00	50.00	94.44	83.64
100	50.00	50.00	50.00	50.00	94.44	81.82
Sigmoid kernel	90.00	90.00	100.00	90.00	75.76	70.91
GK SVM	100.00	100.00	100.00	80.00	93.43	87.27

(Blake and Merz, 1998). The results for WBCP dataset show the average classification accuracy based on a three-fold cross validation test on the whole dataset. Experiments on the Glass2 dataset use a training set of 108 instances and a validation set of 55 instances.

For all SVM classifiers the complexity parameter, C , was set to 1. An initial population of 100 randomly generated kernel trees was used for the WBCP and Glass2 datasets and a population of 30 was used for finding a model for the Raman spectra datasets. The behaviour of the GP search differed for each dataset. For the spectral datasets, the search quickly converged to the simple solution after an average of only five generations, whereas the WBCP and Glass2 datasets required an average of 17 and 31 generations, respectively. (As stated earlier, five populations are evolved in parallel and the best individual chosen.)

The results clearly demonstrate both the large variation in accuracy between the Polynomial, RBF and Sigmoid kernels as well as

Table 2. Best Parameter Settings for the Sigmoid kernel

Dataset	Ranges tested		Best Accuracy	Best parameter settings (γ , θ)
	γ	θ		
Acetone	1–10	0–1000	90	(3,500), (5,400),(6,700)
Cyclohexanol	1–10	0–1000	90	(1,100), (2,200), (3,300), (3,400), (5,600), (6,600–900), (7,700)
Acetonitrile	1–10	0–1000	100	(7,200),(8,300),(8,400),(9,400),(10,600)
Toluene	1–10	0–1000	90	(1,800),(3,900),(4,400),(6,1000)
WBCP	0–1	1–10	75.76	(0.4,0)
Glass2	0–1	1–10	70.91	(0.9,6)

the variation between the performance of models using the same kernel but with different parameter settings: degree d for the Polynomial kernel, σ for the RBF kernel, γ and θ for the Sigmoid kernel. For the Polynomial and RBF kernel, the accuracy for different settings is shown. As there are two parameters to set for the Sigmoid kernel, only the best accuracy, over all combinations of parameters tested, for each dataset is shown. The actual values of γ and θ used to get this accuracy on each dataset is shown in Table 2.

The RBF kernel performs poorly on the spectral datasets but then outperforms the Polynomial kernel on the Wisconsin Breast Cancer Prognosis and Glass2 datasets. The Sigmoid kernel performs much better than the RBF kernel on the spectral datasets, and slightly worse than the Polynomial kernel on these datasets. However, on the WBCP and Glass2 datasets, it performs much worse than the other kernels. For the first three spectral datasets, the GK SVM achieves 100% accuracy, each time finding the same simple linear kernel as the best kernel tree:

$$K(x, z) = \langle x, z \rangle \quad (9)$$

For the Toluene dataset, the GK SVM manages to find a kernel of higher fitness (according to the fitness function detailed in Section 4.2) than the linear kernel, but which happens to perform worse on the test dataset. One drawback with the use of these spectral datasets is that the small number of samples is not very suitable for a complex search procedure such as used in GK SVM. A small training dataset increases the danger of an evolutionary technique, such as GP, finding a model that fits the training set well but performs poorly on the test data.

On the Wisconsin Breast Cancer Prognosis dataset, the GK SVM performs better than the best Polynomial kernel ($d=4$). The best kernel tree found during the final fold of the three-fold cross-validation test is shown in Figure 2. This tree represents the following kernel function:

$$K(x, z) = \langle (x -^{\text{scal}} (x -^{\text{scal}} z)), (z -^{\text{scal}} (z -^{\text{scal}} x)) \rangle \quad (10)$$

The performance of the GK SVM on this dataset demonstrates its potential to find new non-linear kernels for the classification of data. The GK SVM does, however, perform marginally worse than the RBF kernel on this dataset. This may be due to the fact that the kernel trees are constructed using only three basic mathematical operators and therefore cannot find a solution to compete with the exponential function of the RBF kernel. Despite this apparent disadvantage, the GK SVM clearly outperforms all kernels on the Glass2 dataset.

Table 2 details the settings for γ and θ of the Sigmoid kernel that resulted in the best accuracy for the SVM on each dataset. For example, with $\gamma = 5$ and $\theta = 400$, an accuracy of 90% was achieved in classifying the Acetone test dataset. Note that these results show the best accuracy over a range of settings for (γ, θ) . The range for each parameter was divided into ten partitions, including the starting and end value, i.e., 110 different pairs were tested on the spectral datasets. A different range of values was required to find the best accuracy on the WBCP and Glass2 datasets. This highlights further the problem of finding the best setting for a kernel, especially when there is more than one parameter involved.

Overall, these results show the ability of the GK SVM to automatically find kernel functions that perform competitively in comparison with the widely used Polynomial, RBF and Sigmoid kernels, but without requiring a manual parameter search to achieve optimum performance.

6. Related Research

6.1. SVM model selection

Research on the tuning of kernel parameters or model selection is of particular relevance to the work presented here, which is attempting to automate kernel selection. A common approach is to use a

grid-search of the parameters, e.g., complexity parameter C and width of RBF kernel, σ (Hsu et al., 2003). In this case, pairs of (C, σ) are tried and the one with best cross-validation accuracy is picked. A similar algorithm for the selection of SVM parameters is presented in Staelin (2002). That algorithm starts with a very coarse grid covering the whole search space and iteratively refines both grid resolution and search boundaries, keeping the number of samples at each iteration roughly constant. It is based on a search method from the design of experiments (DOE) field. Those techniques still require selection of a suitable kernel in addition to knowledge of a suitable starting range for the kernel parameters being optimised. The same can be said for the model selection technique proposed in Cristianini et al. (1998), in which an on-line gradient ascent method is used to find the optimal σ for an RBF kernel.

6.2. Application of evolutionary techniques with SVM classifiers

Some research has been carried out on the use of evolutionary approaches in tandem with SVMs. Fröhlich et al. (2003) use GAs for feature selection and train SVMs on the reduced data. The novelty of this approach is in its use of a fitness function based on the calculation of the theoretical bounds on the generalisation error of the SVM. This approach was found to achieve better results than when a fitness function based on cross-validation error was used. A RBF kernel was used in all reported experiments.

An example of GPs and SVMs is found in Eads et al. (2002), which reports on the use of SVMs for identification of lightning types based on time series data. However, in this case the GP was used to extract a set of features for each time series sample in the dataset. This derived dataset was then used as the training data for building the SVM which mapped each feature set or vector onto a lightning category. A GA was then used to evolve a chromosome of multiple GP trees (each tree was used to generate one element of the feature vector) and the fitness of a single chromosome was based on the cross validation error of an SVM using the set of features it encoded. With this approach the SVM kernel (along with σ) still had to be selected, in this case the RBF kernel was used.

Some more recent work has been carried out in the use of an evolutionary strategy (ES) for SVM model selection. ES is an evolutionary approach which is generally applied to real-valued representations of optimisation problems, and which tends to emphasise mutation

over crossover (Whitley, 2001). Runarsson and Sigurdsson (2004) use an ES to evolve an optimal value for C and σ for an RBF kernel of an SVM. Four different criteria are used for evaluating a particular set of parameters. Two of these criteria are based on the kernel radius and $\sum \alpha_i$ measures (discussed in section 4.2). The fourth criterion used is simply the count of support vectors used in the SVM model, with a lower count indicating a better model. The best overall performance appears to be obtained using the following evaluation

$$f(x) = \left(R^2 + \frac{1}{C} \right) \sum_{i=1}^{\ell} \alpha_i \quad (11)$$

where $f(x)$ is the fitness function used to evaluate a particular SVM kernel. This paper also reports the usage of ES and SVM to classify a dataset of chromosomes, which are represented by variable-length strings. In this case, an RBF kernel is used with the Euclidean distance replaced by the string edit (or Levenshtein) distance (another example of a custom kernel). The ES is used to evolve a set of costs for each of the symbols used to describe a chromosome, where the costs are required to calculate the distance between two chromosome strings. They found that minimising the number of support vectors resulted in overfitting to the training data and conclude that this criterion is not suitable when dealing with small training sets.

Another example of the use of ES methods to tune an RBF kernel is presented in Friedrichs and Igel (2004), which involves the adaption of not only the scaling, but also the orientation of the kernel. Three optimisation methods are reported in this work:

1. The σ of the RBF kernel is adapted.
2. Independent scalings of the components of the input vector are adapted.
3. Both the scaling and the rotation of the input vector is adapted.

The fitness of a kernel variation was based on its error on a separate test set. The performance of this ES approach was compared with a SVM that was tuned using a grid search. Better results were achieved with both the scaled kernel and the scaled and rotated kernel. However, it must be noted that the results of the grid search were used as initial values for the ES approach, i.e., used to initialise σ .

Again, the focus in these last two examples of research in this area is on the tuning of parameters for an RBF kernel. This appears to be the most popular kernel (particularly when model selection is considered), but as the results presented in Section 5 show, does not always

achieve the best performance. The goal of our research is to devise a method that overcomes this problem, and produces the best kernel for a given dataset.

7. Conclusions

This paper has proposed a novel approach to tackle the problem of kernel selection for SVM classifiers. The proposed GK SVM uses a GP to evolve a suitable kernel for a particular problem. The initial experimental results show that the GK SVM is capable of matching or beating the best performance of the standard SVM kernels on the majority of the datasets tested. These experiments also demonstrate the potential for this technique to discover new kernels for a particular problem domain. Future work will involve testing the GK SVM on more datasets and possibly finding more kernels to compare its performance with. The effect of restricting the GP search to Mercer kernels will be investigated. In order to help the GK SVM find better solutions, further experimentation is required with increasing the range of functions available for construction of kernel trees, e.g. to include the exponential or *tanh* function. In addition, future work will investigate alternative fitness evaluations for the ranking of kernels, e.g. including the support vector count in the fitness estimate.

Acknowledgements

This research has been funded by Enterprise Ireland's Basic Research Grant Programme. The authors are grateful to Dr. Alan Ryder and Jennifer Conroy for providing the spectral datasets.

References

- Bahlmann, C., Haasdonk, B. & Burkhart, H. (2002). On-line Handwriting Recognition with Support Vector Machines – A Kernel Approach. In *Proceedings of the 8th International Workshop on Frontiers in Handwriting Recognition*, 49–54.
- Banzhaf, W., Nordin, P., Keller, R. & Francone, F. (1998). *Genetic Programming – An Introduction*. Morgan Kaufmann.
- Blake, C. & Merz, C. (1998). UCI Repository of machine learning databases, <http://www.ics.uci.edu/~mlearn/MLRepository.html>. University of California, Irvine, Department of Information and Computer Sciences.

- Bleckmann, A. & Meiler, J. (2003). Epothilones: Quantitative Structure Activity Relations Studied by Support Vector Machines and Artificial Neural Networks. *QSAR & Combinatorial Science* **22**, 722–728.
- Boser, B., Guyon, I. & Vapnik, V. (1992a). A Training Algorithm for Optimal Margin Classifiers. In Haussler, D. (ed.) *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, 144–152, ACM Press.
- Boser, B., Guyon, I. & Vapnik, V. (1992b). A Training Algorithm for Optimal Margin Classifiers. In Haussler, D. (ed.) *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, 144–152, ACM Press.
- Cristianini, N., Campbell, C. & Shawe-Taylor, J. (1998). Dynamically Adapting Kernels in Support Vector Machines. Technical Report NC2-TR-1998-017, NeuroCOLT2.
- Cristianini, N. & Shawe-Taylor, J. (2000). *An Introduction to Support Vector Machines*. Cambridge University Press.
- Eads, D., Hill, D., Davis, S., Perkins, S., Ma, J., Porter, R. & Theiler, J. (2002). Genetic Algorithms and Support Vector Machines for Times Series Classification. In *Proceedings of the Fifth Conference on the Applications and Science of Neural Networks, Fuzzy Systems, and Evolutionary Computation. Symposium on Optical Science and Technology of the 2002 SPIE Annual Meeting*, 74–85.
- Friedrichs, F. & Igel, C. (2004). Evolutionary Tuning of Multiple SVM Parameters. In *Proceedings of the 12th European Symposium on Artificial Neural Network*, 519–524.
- Frolich, H., Chapelle, O. & Scholkopf, B. (2003). Feature Selection for Support Vector Machines by Means of Genetic Algorithms. In *Proceedings of the International IEEE Conference on Tools with AI*, 142–148.
- Hastie, T., Tibshirani, R. & Friedman, J. (2001). *The Elements of Statistical Learning*. Springer.
- Hearst, M. (1998). Using SVMs for Text Categorisation. *IEEE Intelligent Systems* **13**(4), 18–28.
- Hennessy, K., Madden, M., Conroy, J. & Ryder, A. (2005). An Improved Genetic Programming Technique for the Classification of Raman Spectra. *Knowledge Based Systems*.
- Hsu, C., Chang, C. & Lin, C. (2003). A Practical Guide to Support Vector Classification, <http://www.csie.ntu.edu.tw/~cjlin/guide/guide.pdf>. Department of Computer Science and Information Engineering, National Taiwan University.
- Joachims, T. (1998). Text categorisation with support vector machines. In *Proceedings of European Conference on Machine Learning (ECML)*.
- Koza, J. (1992). *Genetic Programming*. MIT Press.
- Lin, H. & Lin, C. (2003). A Study on Sigmoid Kernels for SVM and the Training of non-PSD Kernels by SMO-type Methods. Technical report, Department of Computer Science and Information Engineering, National Taiwan University.
- Lodhi, H., Saunders, C., Shawe-Taylor, J., Cristianini, N. & Watkins, C. (2002). Text Classification Using String Kernels. *Journal of Machine Learning Research* **2**, 419–444.
- Osuna, E., Freund, R. & Girosi, F. (1997). Training Support Vector Machines: An Application to Face Detection. In *Proceedings of Computer Vision and Pattern Recognition*, 130–136.

- Platt, J. (1999). Using Analytic QP and Sparseness to Speed Training of Support Vector Machines. In *Proceedings of Neural Information Processing Systems (NIPS)*, 557–563.
- Runarsson, T. & Sigurdsson, S. (2004). Asynchronous Parallel Evolutionary Model Selection for Support Vector Machines. *Neural Information Processing – Letters and Reviews* **3**, 59–67.
- Russell, S. & Norvig, P. (2003). *Artificial Intelligence A Modern Approach*. Prentice-Hall.
- Scholkopf, B. (1998). Support Vector Machines – A Practical Consequence of Learning Theory. *IEEE Intelligent Systems* **13**(4), 18–28.
- Scholkopf, B. (2000). Statistical Learning and Kernel Methods. Technical Report MSR-TR-2000-23, Microsoft Research, Microsoft Corporation.
- Staelin, C. (2002). Parameter selection for support vector machines. Technical Report HPL-2002-354, HP Laboratories, Israel.
- Whitley, D. (2001). An Overview of Evolutionary Algorithms: Practical Issues and Common Pitfalls. *Information and Software Technology* **43**(14), 817–831.