
An Improved Training Algorithm for Kernel Fisher Discriminants

Sebastian Mika*
GMD FIRST.IDA
Kekuléstr. 7
12489 Berlin, Germany
mika@first.gmd.de

Alexander Smola*
Australian National University
Canberra 0200 ACT
Australia
alex.smola@anu.edu.au

Bernhard Schölkopf
Microsoft Research
1 Guildhall Street
Cambridge CB2 3NH, UK
bsc@scientist.com

Abstract

We present a fast training algorithm for the kernel Fisher discriminant classifier. It uses a greedy approximation technique and has an empirical scaling behavior which improves upon the state of the art by more than an order of magnitude, thus rendering the kernel Fisher algorithm a viable option also for large datasets.

1 INTRODUCTION

Kernel Fisher Discriminant (KFD) (Mika et al., 1999, 2000) is a nonlinear generalization of Fisher’s Discriminant (Fisher, 1936). The nonlinearity is introduced by the use of kernel functions (Mercer, 1909), in analogy to Support Vector Machines (SVMs) (Boser et al., 1992), Kernel PCA (Schölkopf et al., 1998) and various other techniques.

On a large number of benchmarks, KFD has shown classification accuracies on a par with SVMs. In addition, unlike SVMs, the outputs of KFD lend themselves to a probabilistic interpretation: empirically, the distributions of the two classes projected onto the Fisher direction of discrimination can be approximated very well by Gaussians, which allows the estimation of conditional class probabilities. In this sense, KFD can be considered a “probabilistic variant” of SVMs (cf. also (Tipping, 2000)). Unfortunately, so far, there has been no efficient algorithm for KFD — all the known algorithms effectively scaled like $O(\ell^3)$, where ℓ is the sample size. In the current paper, we propose a much more efficient algorithm, utilizing sparse greedy approximation techniques (Smola and Schölkopf, 2000; Smola and Bartlett, 2001).

*Work done while SM and AS were visiting MSR Cambridge

2 THE KERNEL FISHER DISCRIMINANT REVISITED

For some set \mathcal{X} , let $\{\mathbf{x}_i \in \mathcal{X} \mid i = 1, \dots, \ell\}$ be our training sample and $\mathbf{y} \in \{-1, +1\}^\ell$ be the vector of corresponding class labels. Furthermore define $\mathbf{1} \in \mathbb{R}^\ell$ as the vector of all ones, and let $\mathbf{1}_+, \mathbf{1}_- \in \mathbb{R}^\ell$ be the positive and negative parts of \mathbf{y} , i.e. $\mathbf{1}_+ = \max(\mathbf{y}, 0)$, $\mathbf{1}_- = \max(-\mathbf{y}, 0)$. In the linear case, it is known that Fisher’s discriminant is computed by maximizing the coefficient¹

$$J(\mathbf{w}) = \frac{\mathbf{w}^\top S_B \mathbf{w}}{\mathbf{w}^\top S_W \mathbf{w}} \quad (1)$$

of between and within class variance, i.e. $S_B = (\mathbf{m}_+ - \mathbf{m}_-)(\mathbf{m}_+ - \mathbf{m}_-)^\top$ and $S_W = \sum_{i, y_i=1} (\mathbf{x}_i - \mathbf{m}_+)(\mathbf{x}_i - \mathbf{m}_+)^\top + \sum_{i, y_i=-1} (\mathbf{x}_i - \mathbf{m}_-)(\mathbf{x}_i - \mathbf{m}_-)^\top$, where \mathbf{m}_\pm denotes the sample mean for class ± 1 .

It can be shown that this approach results in the optimal (in this case linear) decision for two Gaussian distributions with equal covariance structure. In spite of the fact that Fisher’s discriminant often yields useful results even when this assumption is violated, its basic limitation is that the discriminating direction is linear. To overcome this limitation, Mika et al. (1999) proposed to use the same approach as in Support Vector Machines (Boser et al., 1992) or Kernel PCA (Schölkopf et al., 1998): kernel functions. In a nutshell, the idea is to first apply a nonlinear mapping $\Phi : \mathcal{X} \rightarrow \mathcal{F}$ to the data and then to perform the same linear algorithm on the mapped data. If \mathcal{F} is sufficiently rich, this increases the chance of finding a good linear, separating direction in the mapped space. This linear direction in the *feature* space \mathcal{F} then implicitly yields a nonlinear direction in the input space. To avoid having to work in \mathcal{F} explicitly, we use the *kernel trick*, i.e. we choose a feature space whose dot product can efficiently be computed by a

¹In this case \mathcal{X} has to be a subset of a metric space, e.g. $\mathcal{X} \subseteq \mathbb{R}^N$.

kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$,

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)),$$

and formulate everything in terms of dot products in \mathcal{F} , i.e. in kernels on $\mathcal{X} \times \mathcal{X}$.

3 FISHER'S DISCRIMINANT IN FEATURE SPACE

To solve Fisher's problem in a kernel feature space \mathcal{F} one needs a formulation which makes use of the training samples only in terms of dot products. One can prove (Mika et al., 1999) that the solution $\mathbf{w} \in \mathcal{F}$ of (1) can be expanded as

$$\mathbf{w} = \sum_{i=1}^{\ell} \alpha_i \Phi(\mathbf{x}_i), \quad \alpha_i \in \mathbb{R}. \quad (2)$$

It is straightforward to find an expression similar to (1) for the coefficients α (Mika et al., 1999). However, here we will use a different formulation for finding α , building on the following observation: the goal of Fisher's discriminant is to find a one dimensional projection on which the class means are far apart while the within class variance is small.

In (Mika et al., 2001) it was shown that KFD can be cast in a slightly more general form as the following convex, quadratic optimization problem:

$$\min_{\alpha, b, \xi} \frac{1}{2} \|\xi\|^2 + \frac{C}{2} P(\alpha) \quad (3)$$

subject to:

$$K\alpha + \mathbf{1}b = \mathbf{y} + \xi \quad (3a)$$

$$\mathbf{1}_+^\top \xi = 0, \quad \mathbf{1}_-^\top \xi = 0. \quad (3b)$$

Here, C is a regularization constant, and P a regularization functional (Mika et al., 1999) which we assume to be quadratic in the following, e.g. $P(\alpha) = \|\alpha\|^2$ or $P(\alpha) = \alpha^\top K \alpha$. For $C = 0$ one obtains Fisher's original algorithm in the feature space \mathcal{F} . However, as pointed out e.g. by Friedman (1989) working in the space \mathcal{F} which has an effective dimension of ℓ with only ℓ training samples is an ill posed problem. Thus, introducing some form of regularization is mandatory. The projection of a test point onto the discriminant is computed by $(\mathbf{w} \cdot \Phi(\mathbf{x})) = \sum_i \alpha_i k(\mathbf{x}_i, \mathbf{x})$. The program essentially states that the output for each training sample should be close to its label, where we penalize the squared error of the deviation (constraint (3a)), and that the average deviation from the label should be zero, separately for each class (constraints (3b)).

However, for large data sets, solving (3) is expensive in terms of time and memory. Contrary to SVMs, the

solutions are not sparse and deriving efficient decomposition techniques for the programming problem is difficult. In (Mika et al., 2001), it was proposed to use a ℓ_1 regularizer in (3) as an approximation to a ℓ_0 regularizer which would just count the number of non zero elements in α . While this solved the problem of non-sparsity and was a promising candidate to use chunking techniques there was no really efficient algorithm for large sample sizes yet.

Using a ℓ_0 regularizer or to add a nonlinear constraint of the form: *Find a solution α with at most m non-zero elements*, would in principle be optimal. Unfortunately, such a constraint is impossible to deal with analytically. Finding the true optimal solution would require to search the space of all possible solutions which make use of m possible α_i , i.e. one had to solve $\binom{\ell}{m}$ problems.

We will presently derive an algorithm which might be viewed as a greedy approximation to such a solution. Along the lines of (Smola and Bartlett, 2001; Smola and Schölkopf, 2000), we will iteratively approximate the solution to (3) with as few non-zero α_i as possible.

4 THE ALGORITHM

To proceed, let us rewrite (3). Define

$$\mathbf{a} = \begin{bmatrix} b \\ \alpha \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} \ell_+ - \ell_- \\ K^\top \mathbf{y} \end{bmatrix}, \quad \mathbf{A}_\pm = \begin{bmatrix} \ell_\pm \\ K^\top \mathbf{1}_\pm \end{bmatrix}$$

$$H = \begin{bmatrix} \ell & \mathbf{1}^\top K \\ K^\top \mathbf{1} & K^\top K + C P \end{bmatrix}. \quad (4)$$

Here, ℓ_\pm denotes the number of samples in class ± 1 . Then the problem (3) can equivalently be rewritten as:

$$\min_{\mathbf{a}} \frac{1}{2} \mathbf{a}^\top H \mathbf{a} - \mathbf{c}^\top \mathbf{a} + \frac{\ell}{2} \quad (5)$$

subject to:

$$\mathbf{A}_+^\top \mathbf{a} - \ell_+ = 0 \quad (5a)$$

$$\mathbf{A}_-^\top \mathbf{a} + \ell_- = 0. \quad (5b)$$

Forming the Lagrangian of (5) with multipliers λ_\pm

$$L(\mathbf{a}, \lambda_+, \lambda_-) = \frac{1}{2} \mathbf{a}^\top H \mathbf{a} - \mathbf{c}^\top \mathbf{a}$$

$$+ \lambda_+ (\mathbf{A}_+^\top \mathbf{a} - \ell_+) + \lambda_- (\mathbf{A}_-^\top \mathbf{a} + \ell_-) + \frac{\ell}{2}, \quad (6)$$

and taking derivatives with respect to the primal variables \mathbf{a} one obtains the dual

$$\max_{\mathbf{a}, \lambda_+, \lambda_-} -\frac{1}{2} \mathbf{a}^\top H \mathbf{a} - \lambda_+ \ell_+ + \lambda_- \ell_- + \frac{\ell}{2} \quad (7)$$

subject to:

$$H \mathbf{a} - \mathbf{c} + (\lambda_+ \mathbf{A}_+ + \lambda_- \mathbf{A}_-) = 0. \quad (8a)$$

Now we use the dual constraint (8a) to solve for \mathbf{a} , i.e.

$$\mathbf{a} = H^{-1}(\mathbf{c} - (\lambda_+ \mathbf{A}_+ + \lambda_- \mathbf{A}_-)). \quad (8)$$

This equation is well defined if H has full rank. If not we can still perform this step as we will approximate H^{-1} instead of computing it directly. Resubstituting (8) into the dual problem (which has no constraints left) yields the following problem in the two variables λ_+ and λ_- :

$$\begin{aligned} \max_{\lambda_+, \lambda_-} & -\frac{1}{2} \begin{bmatrix} \lambda_+ \\ \lambda_- \end{bmatrix}^\top \begin{bmatrix} \mathbf{A}_+^\top H^{-1} \mathbf{A}_+ & \mathbf{A}_+^\top H^{-1} \mathbf{A}_- \\ \mathbf{A}_-^\top H^{-1} \mathbf{A}_+ & \mathbf{A}_-^\top H^{-1} \mathbf{A}_- \end{bmatrix} \begin{bmatrix} \lambda_+ \\ \lambda_- \end{bmatrix} \\ & + \begin{bmatrix} -\ell_+ + \mathbf{c}^\top H^{-1} \mathbf{A}_+ \\ \ell_- + \mathbf{c}^\top H^{-1} \mathbf{A}_- \end{bmatrix}^\top \begin{bmatrix} \lambda_+ \\ \lambda_- \end{bmatrix} - \frac{1}{2} \mathbf{c}^\top H^{-1} \mathbf{c} + \frac{\ell}{2}. \end{aligned} \quad (9)$$

This problem can be solved analytically, yielding optimal values for λ_+ and λ_- which substituted into (8) yield values for \mathbf{a} or $\boldsymbol{\alpha}$ and b , respectively.

4.1 A SPARSE GREEDY APPROXIMATION

Of course, this problem is no easier to solve than the original one nor does it yield a sparse solution: H^{-1} is an $(\ell + 1) \times (\ell + 1)$ matrix and for large datasets its inversion is not feasible, neither in terms of time nor memory cost. Now, the idea is to use the following, greedy approximation scheme (cf. (Smola and Bartlett, 2001)). Instead of trying to find a full set of ℓ α_i 's for the solution (2), we approximate the optimal solution by a shorter expansion containing only $m \ll \ell$ terms.

Starting with an empty expansion $m = 0$, one selects in each iteration a new sample \mathbf{x}_i (or an index i) and resolves the problem for the expansion (2) containing this new index and all previously picked indices; we stop as soon as a suitable criterion is satisfied. This approach would still be infeasible in terms of computational cost if we had to solve the quadratic program (5) anew in each iteration or invert H in (8) and (9). But with the derivation made before it is possible to find a close approximation to the optimal solution in each iteration at a cost of $\mathcal{O}(\kappa \ell m^2)$ where κ is a user defined value (see below).

Writing down the quadratic program (3) for KFD when the expansion for the solution is restricted to an m element subset $\mathcal{I} \subset [\ell]$

$$\mathbf{w}_{\mathcal{I}} = \sum_{i \in \mathcal{I}} \alpha_i \Phi(\mathbf{x}_i) \quad (10)$$

of the training patterns amounts to replacing the $\ell \times \ell$ matrix K by the $\ell \times m$ matrix K^m , where $K_{ij}^m = k(\mathbf{x}_i, \mathbf{x}_j)$, $i = 1, \dots, \ell$ and $j \in \mathcal{I}$. Analogously, we can derive the formulation (5) using the matrix K^m

in (4). The problem is of order $m \times m$ now. Assume we already know the optimal solution (and inverse of H) using m kernel-functions. Then H^{-1} for $m + 1$ samples can be obtained by a rank one update of the previous H^{-1} using only m basis functions: The following Lemma (e.g. (Golub and van Loan, 1996)) tells us how to obtain the new H^{-1} .

Lemma 1 (Sherman–Woodbury–Formula). *The inverse of a symmetric, positive matrix can be computed as:*

$$\begin{bmatrix} H & B \\ B^\top & C \end{bmatrix}^{-1} = \begin{bmatrix} H^{-1} + (H^{-1}B)\gamma(H^{-1}B)^\top & -\gamma(H^{-1}B) \\ -(\gamma(H^{-1}B))^\top & \gamma \end{bmatrix},$$

where $\gamma = (C - B^\top H^{-1}B)^{-1}$.

Note that for our case B is a vector and C a scalar. This is an operation of cost $\mathcal{O}(m^2)$ as we already know the inverse of the smaller system.

4.1.1 Selection Rules

The last major problem is to pick an index i in each iteration. Choosing one index at a time can be considered as a coordinate wise descent method with the difference that we update all coordinates which were already chosen in each iteration. As we are dealing with a convex, quadratic optimization problem any selection rule will finally lead to the (optimal) solution of the full problem. As we are only adding kernel functions to the expansion, it is easy to see that this convergence is monotonic with respect to the primal objective. The interesting question is how one has to choose the sequence of indices such that this greedy approach will get close to the full solution as soon as possible. Here we propose two heuristic selection criteria which we believe to be reasonable:

Minimizing the objective

One possibility is to choose the i for which we get the biggest decrease in the primal-objective (or equivalently as they are identical for the optimal coefficients \mathbf{a} , the dual-objective (7)). This would correspond to a steepest descent method in a restricted space of the gradients.

Minimizing the dual infeasibility

Another, slightly more sophisticated way can be derived as follows. For quadratic programs there are three optimality criteria: (i) primal feasibility, (ii) dual feasibility, and (iii) the duality gap, i.e. the difference between the primal and dual objective value. For an optimal solution it is necessary and also sufficient to

fulfill these three conditions. If we inspect the sequence of solutions \mathbf{a} generated by our approach we find that by construction (i) and (iii) are always fulfilled, i.e. each intermediate solution is primal feasible and the primal and dual objectives are identical². But the dual constraints (8a) are only fulfilled for those indices, i.e. rows, which have already been chosen. Thus another possible selection criterion is to reduce the dual infeasibility as much as possible by picking the index i which corresponds to the biggest (absolute) violation. But, contrary to the value of the objective, this violation is not guaranteed to decrease monotonically. Such an approach exhibits strong connections to column generation algorithms as used e.g. in Boosting (Bennett et al., 2000; Rätsch et al., 2000).

4.1.2 A probabilistic speed-up

Whichever selection rule one uses, testing it for all $\ell - m$ indices which are unused so far is — again, too expensive. One possible solution lies in a second approximation. Instead of choosing the *best* possible index it is usually sufficient to find an index for which with high probability we achieve something close to the optimal achievement. It turns out (Smola and Schölkopf, 2000) that it can be enough to consider 59 randomly chosen indices from the remaining ones:

Lemma 2 (Maximum of Random Variables).

Denote by ρ_1, \dots, ρ_m identically distributed independent random variables with a common cumulative distribution function F . Then the cumulative distribution function of $\rho = \max_{i \in [m]} \rho_i$ is F^m .

This means that e.g. for the uniform distribution on $[0, 1]$ $\max_{i \in [m]} \rho_i$ is distributed according to ρ^m . Thus, to obtain an estimate that is with probability 0.95 among the best 0.05 of all estimates, a random sample of size $\kappa := (\log 0.05 / \log 0.95) = 59$ is enough.

4.1.3 Termination

Still open is the question when to stop. If one wanted to compute the full solution this approach would not be very efficient as it would take $\mathcal{O}(\kappa \ell^3)$ which is worse than the original problem. A principled stopping rule would be to measure the distance of $\mathbf{w}_{\mathcal{I}}$ to the solution of the full problem and then to stop when this falls below a certain threshold. Unfortunately the full solution is, for obvious reasons, not available. Instead one could try to bound the difference of the objective (3) for the current solution to the optimal value obtained for the full problem as done in (Smola and Bartlett,

²Strictly speaking this only holds true for $m > 1$. The prove is straight forward and can be done by using the relation $\mathbf{c} = \mathbf{A}_+ - \mathbf{A}_-$, the condition (8), and the optimality of λ_{\pm} for (9)

2001). But in our case an efficient way to bound this difference is not available. Instead we have chosen two very simple heuristic which turned out to work well in the experiments: Using the selection criterion which minimizes the objective function, stop when the average improvement in the dual objective (7) over the last p iterations is less than some threshold θ . If using the selection criterion which minimizes the maximal dual infeasibility, terminate if this maximum is smaller than some threshold, again averaged over the last p iterations. The longer the averaging process, the more confident we are that the current solution is not at a plateau. The smaller the threshold, the closer we are to the original solution (indeed, setting the threshold to zero forces the algorithm in both cases to take all training samples into account).

The complete algorithm for a sparse greedy solution to the KFD problem is schematized in Figure 1. It is easy to implement using a linear algebra package like BLAS and has the potential to be easily parallelized (the matrix update) and distributed. Furthermore, in multi-class problems (if one is using a one against the rest scheme) it is possible to consider all two-class problems simultaneously. Testing the same subset of indices for each classifier would result in a reasonable speedup, as the computation of the rows of the kernel matrix for each picked index is among the most expensive parts of this algorithm and must now only be done once.

5 OBTAINING PROBABILITIES AS OUTPUTS

One of the advantages of KFD over e.g. SVM is that the outputs of KFD can be interpreted in a probabilistic manner. If one is interested in probabilities these are straightforward to obtain. Implicitly the optimization problem for KFD assumes Gaussian distribution for the likelihood functions, an assumption which is empirically supported by examining the output histograms on different datasets: they exhibit a strong Gaussianity (Mika et al., 2001; Bradshaw et al., 2000). If we estimate mean and variance of these Gaussians using the training data it turns out that the class label is the mean (i.e. $\mu_+ = 1$, $\mu_- = -1$), due to the constraints (3b), and that the variances of the class conditional densities are given by $\sigma_+^2 = \frac{1}{\ell_+ - 1} \sum_{y_i=1} \xi_i^2$ and $\sigma_-^2 = \frac{1}{\ell_- - 1} \sum_{y_i=-1} \xi_i^2$. In some cases it might be advantageous to estimate μ_{\pm} and σ_{\pm} from a separate validation set in which case they can simply be computed from the outputs

$$q(\mathbf{x}) := (\mathbf{w} \cdot \mathbf{x}) + b = \sum_i \alpha_i k(\mathbf{x}, \mathbf{x}_i) + b$$

```

arguments: Sample  $X = \{\mathbf{x}_1, \dots, \mathbf{x}_\ell\}, \mathbf{y} = \{y_1, \dots, y_\ell\}$ 
             Maximum number of coefficients
             or parameters of other stopping criterion: OPTS
             Regularization constant  $C$ 
              $\kappa$  and kernel  $k$ 
returns:  Set of indices  $I$  and corresponding  $\alpha$ 's.
             Threshold  $b$ .
function  SG-KFD( $X, \mathbf{y}, C, \kappa, k, \text{OPTS}$ )
     $m \leftarrow 0$ 
     $I \leftarrow \emptyset$ 
    while termination criterion not satisfied do
       $S \leftarrow (\kappa \text{ elements from } [\ell] \setminus I)$ 
       $\text{obj}_{\max} \leftarrow \infty$ 
      for  $i \in S$  do
        Compute column  $i$  of kernel matrix
        Update inverse adding the  $i$ -th kernel, compute optimal  $\mathbf{a}$ 
        Compute new dual objective
        if dual objective  $<$   $\text{obj}_{\max}$  do
           $i_{\text{opt}} \leftarrow i$ 
           $\text{obj}_{\max} \leftarrow$  dual objective
        endif
      endfor
      Update inverse  $H$  and solution  $\mathbf{a}$  with kernel  $i_{\text{opt}}$ 
       $I \leftarrow I \cup \{i_{\text{opt}}\}$ 
      Check termination criterion
    endwhile

```

Figure 1: The Sparse Greedy Kernel Fisher Algorithm (for the selection criterion which minimizes the objective)

obtained by projecting this data set onto the direction \mathbf{w} . No matter which way the parameters are estimated, one obtains the following class-conditional densities:

$$p(\mathbf{x}|y = \pm 1) = p(q(\mathbf{x})|y = \pm 1) = (2\pi\sigma_{\pm}^2)^{-1/2} \exp\left(-\frac{(q(\mathbf{x}) - \mu_{\pm})^2}{2\sigma_{\pm}^2}\right)$$

The prior, if it is unknown for the problem at hand, can be estimated from the training data, i.e. $P(y = \pm 1) = \ell_{\pm}/\ell$, ℓ_{\pm} denoting the number of samples from the respective class. Using Bayes' theorem the conditional probabilities are then given by:

$$P(y = \pm 1|\mathbf{x}) = \frac{P(y = \pm 1|q) p(q|\mathbf{x})}{p(q|y = 1)P(y = 1) + p(q|y = -1)P(y = -1)},$$

where $q = q(\mathbf{x})$ is defined as above.

6 EXPERIMENTAL EVALUATION

We now report some experiments to illustrate that the run-time behavior of our new algorithm improves significantly over the full quadratic optimization of (3).

Furthermore, we show that the approximation does not significantly degrade the quality of the solutions. All timing experiments were carried out on a Pentium III, 500 MHz and 512 MB memory running Linux. To implement the new approach we used a single threaded, optimized BLAS. Timings were measured with the system command `clock()`. We compared this to an implementation of the quadratic program given by (3), which is partly in matlab and calls external C functions for all time-consuming operations, i.e. solving the quadratic program. Here timings were measured by the matlab command `cputime()`. The quadratic optimizer used was `loqo` which implements an interior point algorithm (Vanderbei, 1997).

6.1 TIMING

First we compare the runtime of the new algorithm to the previous implementation. We used a one-against-the-rest task constructed from the USPS handwritten digit data set. The data are $N = 256$ dimensional and the set contains 7291 samples, categorized into ten classes. All experiments were done with a Gaussian kernel $\exp(-\|\mathbf{x} - \mathbf{y}\|^2 / (0.3 \cdot N))$ and using a regularization constant $C = 1$. We compared against the program given by (3) with the regularizer $P = \|\boldsymbol{\alpha}\|^2$.

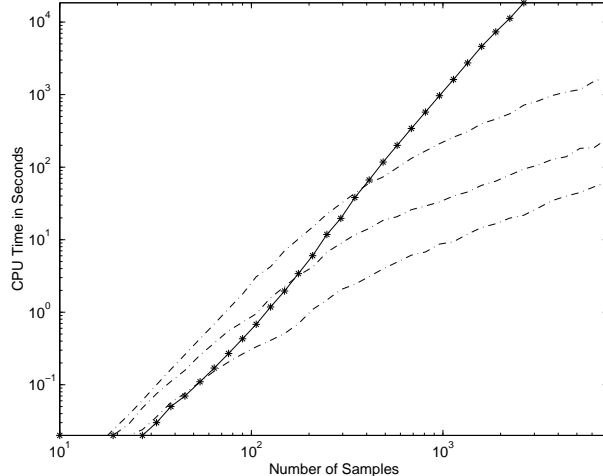


Figure 2: Runtime of sparse greedy KFD training. Depicted is the number of samples in the training set versus the CPU time of the proposed algorithm (dash dotted lines; using the selection criterion minimizing the objective (cf. Section 4.1.1)) and the QP formulation (3) (solid line). The estimates are averages over ten trials, one for each of the ten one-against-the-rest problems in the USPS database. The three lines for KFD are generated by requiring different accuracies on the dual error function in the stopping criterion, namely 10^{-a} , $a = 1, \dots, 3$ relative to the function value (in that order from bottom to top). There is a speed-accuracy tradeoff in that for large a , the algorithm converges more slowly. In the log-log plot it can be seen that the QP algorithm roughly scales cubic in the number of samples while the new algorithm scales with an exponent of about $\frac{3}{2}$ for large sample sizes.

The results of our findings are given in Figure 2. It can be seen that the new training algorithm halves the scaling exponent of the training time for large sample sizes. In addition, it is important to keep in mind that the sparse greedy approach only needs to store at most an $m \times m$ matrix, where m is the maximal number of kernel functions chosen before termination. In contrast, previous approaches needed to store $\ell \times \ell$ matrices.

6.2 PERFORMANCE

As our new approach is an approximation to the original (theoretically exact) algorithm, the question arises how good the quality of this approximation is. To this end, we repeated the above experiment on the USPS database for different regularization constants $C = 10^{-3}, 10^{-4}, 10^{-5}$ and different kernel widths $c = 0.3 \cdot N, 0.4 \cdot N, 0.5 \cdot N$. Using the objective-based selection criterion, the algorithm was terminated when the average achievement in the dual objective over the last five iterations was less than $10^{-1}, 10^{-2}, 5 \cdot 10^{-3}, 10^{-3}$, respectively, relative to the objective or when a maximum of 500 coefficients was found. As the purpose of this experiment is to show that our new approach is capable of producing results comparably to the full system, no model selection was performed and just the best results on the *test set* are reported (cf. Table 1). A small improvement in the test error can be achieved

by using an optimized threshold b rather than the one given by the algorithm itself. This optimized b is found by training a linear support vector machine on the one dimensional outputs of the training data, i.e. we try to find a threshold which maximizes the smallest distance of the projections to the decision boundary (details e.g. in (Mika et al., 1999)). So far the best result for KFD on the USPS dataset (without using prior knowledge) was 3.7% (Mika et al., 2000), using an expansion restricted to the first 3000 training patterns and the optimized threshold. From Table 1 it can be seen that our new approach produces results close to the QP solution, however, using a significantly smaller number of kernel functions (less than 500 vs. 3000). It can be observed that the chosen precision for the termination criterion is an important parameter. Still, although the high precision of 10^{-3} takes longer to train, the runtime of our new approach is more than ten times smaller than solving the QP with 3000 patterns.

Preliminary experiments with the second selection criterion, i.e. selecting the kernel which has the largest dual infeasibility have shown no significant advantage yet. However, when terminating the algorithm if the average constraint violation over the last iterations is less than some threshold (analogously to the criterion used for the first selection criterion) this threshold parameter is slightly more well behaved. For the first selection criterion this parameter is difficult to adjust in

Table 1: Minimal 10-class test error on the USPS dataset using the parameters described in the text. Shown is the threshold on the improvement in the dual objective used to terminate the algorithm (Tolerance), the test error using the threshold given by the algorithm itself, and the test error using an extra, optimized threshold b (see text). The best result of 3.8% is almost identical to the result of 3.7% obtained on the same dataset using an expansion fixed to the first 3000 training samples (Mika et al., 2000). Note, moreover, that for our new algorithm the number of samples in the expansion (2) is less than 500 in each single classifier.

| Tolerance | 10^{-1} | 10^{-2} | $5 \cdot 10^{-3}$ | 10^{-3} |
|-------------------------------------|-----------|-----------|-------------------|-----------|
| test error with QP threshold | 10.4% | 6.4% | 5.3% | 3.9% |
| test error with optimized threshold | 10.3% | 6.3% | 5.3% | 3.8% |

so far that small changes can lead to a large change in the number of iterations performed. The performance on the USPS set, adjusting the parameters such that an approximately equal number of kernels is chosen, is slightly worse for the infeasibility criterion (4.2% for the optimal parameters).

7 CONCLUSION

We presented a new algorithm for kernel Fisher discriminants. These algorithmic advances are crucial for the possibility of applying the KFD algorithms to problems that had previously been beyond its reach. First, since it trains significantly faster, and second, since it requires less memory.

Acknowledgments

SM was supported by the DFG under contract Ja 379/9-1. Thanks to Gunnar Rätsch for fruitful discussions.

References

K.P. Bennett, A. Demiriz, and J. Shawe-Taylor. A column generation algorithm for boosting. In P. Langley, editor, *Proceedings, 17th ICML*, pages 65–72, San Francisco, 2000. Morgan Kaufmann.

B.E. Boser, I.M. Guyon, and V.N. Vapnik. A training algorithm for optimal margin classifiers. In D. Haussler, editor, *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152, 1992.

B. Bradshaw, B. Schölkopf, and J. Platt. Kernel methods for extracting local image semantics. unpublished manuscript, private communication, 2000.

R.A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7:179–188, 1936.

J.H. Friedman. Regularized discriminant analysis. *Journal of the American Statistical Association*, 84(405):165–175, 1989.

G.H. Golub and C.F. van Loan. *Matrix Computations*. John Hopkins University Press, Baltimore, London, 3rd edition, 1996.

J. Mercer. Functions of positive and negative type and their connection with the theory of integral equations. *Philos. Trans. Roy. Soc. London*, A 209:415–446, 1909.

S. Mika, G. Rätsch, and K.-R. Müller. A mathematical programming approach to the Kernel Fisher algorithm. In *Advances in Neural Information Processing Systems 13*, 2001. to appear.

S. Mika, G. Rätsch, J. Weston, B. Schölkopf, and K.-R. Müller. Fisher discriminant analysis with kernels. In Y.-H. Hu, J. Larsen, E. Wilson, and S. Douglas, editors, *Neural Networks for Signal Processing IX*, pages 41–48. IEEE, 1999.

S. Mika, G. Rätsch, J. Weston, B. Schölkopf, A.J. Smola, and K.-R. Müller. Invariant feature extraction and classification in kernel spaces. In S.A. Solla, T.K. Leen, and K.-R. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 526–532. MIT Press, 2000.

G. Rätsch, A. Demiriz, and K. Bennett. Sparse regression ensembles in infinite and finite hypothesis spaces. Neurocolt techreport, Royal Holloway College, 2000. submitted to Special Issue of Machine Learning.

B. Schölkopf, A.J. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10:1299–1319, 1998.

A.J. Smola and P.L. Bartlett. Sparse greedy gaussian process regression. In *Advances in Neural Information Processing Systems 13*, 2001. to appear.

A.J. Smola and B. Schölkopf. Sparse greedy matrix approximation for machine learning. In P. Langley, editor, *Proc. ICML’00*, pages 911–918, San Francisco, 2000. Morgan Kaufmann.

M.E. Tipping. The relevance vector machine. In S.A. Solla, T.K. Leen, and K.-R. Müller, editors, *Advances in Neural Information Processing Systems 12*, pages 652–658. MIT Press, 2000.

R.J. Vanderbei. LOQO user’s manual – version 3.10. Technical Report SOR-97-08, Princeton University, Statistics and Operations Research, 1997. Code available at <http://www.princeton.edu/~rvdb/>.