

Efficient Cross-Validation of Kernel Fisher Discriminant Classifiers

Gavin C. Cawley and Nicola L. C. Talbot

School of Information Systems
 University of East Anglia
 Norwich, U.K. NR4 7TJ
 gcc@sys.uea.ac.uk

Abstract.

Mika *et al.* [1] introduce a non-linear formulation of the Fisher discriminant based the well-known “kernel trick”, later shown to be equivalent to the Least-Squares Support Vector Machine [2, 3]. In this paper, we show that the cross-validation error can be computed very efficiently for this class of kernel machine, specifically that leave-one-out cross-validation can be performed with a computational complexity of only $\mathcal{O}(\ell^3)$ operations (the same as that of the basic training algorithm), rather than the $\mathcal{O}(\ell^4)$ of a direct implementation. This makes leave-one-out cross-validation a practical proposition for model selection in much larger scale applications of KFD classifiers.

1 Introduction

Assume we are given training data $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_\ell\} = \{\mathcal{X}_1, \mathcal{X}_2\} \subset \mathbb{R}^d$, where $\mathcal{X}_1 = \{\mathbf{x}_1^1, \mathbf{x}_2^1, \dots, \mathbf{x}_{\ell_1}^1\}$ is a set of patterns belonging to class \mathcal{C}_1 and similarly $\mathcal{X}_2 = \{\mathbf{x}_1^2, \mathbf{x}_2^2, \dots, \mathbf{x}_{\ell_2}^2\}$ is a set of patterns belonging to class \mathcal{C}_2 ; Fisher’s linear discriminant (FLD) attempts to find a linear combination of input variables, $\mathbf{w} \cdot \mathbf{x}$, that maximises the average separation of the projections of points belonging to \mathcal{C}_1 and \mathcal{C}_2 , whilst minimising the within class variance of the projections of those points. The Fisher discriminant is given by the vector \mathbf{w} maximising

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}, \quad (1)$$

where \mathbf{S}_B is the between class scatter matrix $\mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$, $\mathbf{m}_j = \ell_j^{-1} \sum_{i=1}^{\ell_j} \mathbf{x}_i^j$ and \mathbf{S}_W the within class scatter matrix

$$\mathbf{S}_W = \sum_{i \in \{1,2\}} \sum_{j=1}^{\ell_i} (\mathbf{x}_j^i - \mathbf{m}_i)(\mathbf{x}_j^i - \mathbf{m}_i)^T.$$

The innovation introduced by Mika *et al.* [1] is to construct Fisher's linear discriminant in a fixed feature space \mathcal{F} ($\phi : \mathcal{X} \rightarrow \mathcal{F}$) induced by a positive definite Mercer kernel $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ defining the inner product $\mathcal{K}(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x}) \cdot \phi(\mathbf{x}')$ (see e.g. Cristianini and Shawe-Taylor [4]). Let the kernel matrices for the entire dataset, \mathbf{K} , and for each class, \mathbf{K}_1 and \mathbf{K}_2 be defined as follows:

$$\mathbf{K} = [k_{ij} = \mathcal{K}(\mathbf{x}_i, \mathbf{x}_j)]_{i,j=1}^{\ell} \quad \text{and} \quad \mathbf{K}_i = [k_{jk}^i = \mathcal{K}(\mathbf{x}_j, \mathbf{x}_k^i)]_{j,k=1}^{j=\ell, k=\ell_i}.$$

The theory of reproducing kernels indicates that \mathbf{w} can then be written as an expansion of the form

$$\mathbf{w} = \sum_{i=1}^{\ell} \alpha_i \phi(\mathbf{x}_i). \quad (2)$$

The objective function (1) can also be written such that the data $\mathbf{x} \in \mathcal{X}$ appear only within inner products, giving

$$J(\boldsymbol{\alpha}) = \frac{\boldsymbol{\alpha}^T \mathbf{M} \boldsymbol{\alpha}}{\boldsymbol{\alpha}^T \mathbf{N} \boldsymbol{\alpha}}, \quad (3)$$

where $\boldsymbol{\alpha} = \{\alpha_i\}_{i=1}^{\ell}$, $\mathbf{M} = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$, $\mathbf{m}_i = \mathbf{K}_i \mathbf{u}_i$, \mathbf{u}_i is a column vector containing ℓ_i elements with a common value of ℓ_i^{-1} and

$$\mathbf{N} = \sum_{i \in \{1,2\}} \mathbf{K}_i (\mathbf{I} - \mathbf{U}_i) \mathbf{K}_i^T,$$

where \mathbf{I} is the identity matrix and \mathbf{U}_i is a matrix with all elements equal to ℓ_i^{-1} . The coefficients, $\boldsymbol{\alpha}$, of the expansion (2) are then given by the leading eigenvector of $\mathbf{N}^{-1} \mathbf{M}$. Note that \mathbf{N} is likely to be singular, or at best ill-conditioned, and so a regularised solution is obtained by substituting $\mathbf{N}_{\mu} = \mathbf{N} + \mu \mathbf{I}$, where μ is a regularisation constant. To complete the kernel Fisher discriminant classifier, $f(\mathbf{x}) = \mathbf{w} \cdot \phi(\mathbf{x}) + b$, the bias, b , is given by

$$b = -\boldsymbol{\alpha} \frac{\ell_1 \mathbf{M}_1 + \ell_2 \mathbf{M}_2}{\ell}.$$

Xu *et al.* [3] show that the parameters of the kernel Fisher discriminant classifier are also given by the solution of the following system of linear equations:

$$\begin{bmatrix} \mathbf{K} \mathbf{K} + \mu \mathbf{I} & \mathbf{K} \mathbf{1} \\ (\mathbf{K} \mathbf{1})^T & \ell \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha} \\ b \end{bmatrix} = \begin{bmatrix} \mathbf{K} \\ \mathbf{1} \end{bmatrix} \mathbf{y}, \quad (4)$$

where $\mathbf{1}$ is a column vector of ℓ ones and \mathbf{y} is a column vector with elements $y_i = \ell/\ell_j \forall i : \mathbf{x}_i \in \mathcal{X}_j$. This illustrates the similarities between the kernel Fisher discriminant and the least-squares support vector machine (LS-SVM) [2]. The kernel Fisher discriminant (KFD) classifier has been shown experimentally to demonstrate near state-of-the-art performance on a range of artificial and real world benchmark datasets [1] and so is worthy of consideration for small to medium scale applications. In this paper we present an efficient algorithm for approximate cross-validation of kernel Fisher discriminant models, providing a practical criterion for model selection.

2 Method

The system of linear equations (4) can be written more concisely in the form

$$\begin{bmatrix} \alpha \\ b \end{bmatrix} = \left[\mathbf{R} + \mathbf{Z}^T \mathbf{Z} \right]^{-1} \mathbf{Z}^T \mathbf{y},$$

where $\mathbf{Z} = [\mathbf{K} \ \mathbf{1}]$ and $\mathbf{R} = \text{diag}([\mu \mathbf{1} \ 0])$ (n.b. this is very similar to the set of *normal equations* to be solved in multi-variate linear regression). At each step of the leave-one-out cross-validation procedure, a kernel Fisher discriminant classifier is constructed excluding a single training pattern from the data. The vector of model parameters, $\{\alpha_{(i)}, b_{(i)}\}$ at the i^{th} iteration is then given by the solution of a modified system of linear equations,

$$\begin{bmatrix} \alpha_{(i)} \\ b_{(i)} \end{bmatrix} = \left[\mathbf{R} + \mathbf{Z}_{(i)}^T \mathbf{Z}_{(i)} \right]^{-1} \mathbf{Z}_{(i)}^T \mathbf{y}$$

where $\mathbf{Z}_{(i)}$ is the sub-matrix formed by omitting the i^{th} row of \mathbf{Z} . Normally the most computationally expensive step is the inversion of the matrix $\mathbf{C}_{(i)} = \left[\mathbf{R} + \mathbf{Z}_{(i)}^T \mathbf{Z}_{(i)} \right]$, with a complexity of $\mathcal{O}(\ell^3)$ operations. Fortunately $\mathbf{C}_{(i)}$ can be written as a rank one modification of a matrix \mathbf{C} ,

$$\mathbf{C}_{(i)} = \left[\mathbf{R}_{(i)} + \mathbf{Z}^T \mathbf{Z} - \mathbf{z}_i \mathbf{z}_i^T \right] = \left[\mathbf{C} - \mathbf{z}_i \mathbf{z}_i^T \right], \quad (5)$$

where \mathbf{z}_i is the i^{th} row of \mathbf{Z} . The following matrix inversion lemma then allows $\mathbf{C}_{(i)}^{-1}$ to be found in only $\mathcal{O}(\ell^2)$ operations, given that \mathbf{C}^{-1} is already known:

Lemma 1 (Bartlett Matrix Inversion Formula) *Given an invertible matrix \mathbf{A} and column vectors \mathbf{u} and \mathbf{v} , then assuming $1 - \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq 0$,*

$$(\mathbf{A} + \mathbf{u} \mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}}{1 - \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}}. \quad (6)$$

This is known as the Bartlett matrix inversion formula [5].

Applying the Bartlett formula to the matrix inversion problem given in (5), we have that

$$\mathbf{C}_{(i)}^{-1} = [\mathbf{C} - \mathbf{z}_i \mathbf{z}_i^T]^{-1} = \mathbf{C} + \frac{\mathbf{C}^{-1} \mathbf{z}_i \mathbf{z}_i^T \mathbf{C}^{-1}}{1 - \mathbf{z}_i^T \mathbf{C}^{-1} \mathbf{z}_i}.$$

The computational complexity of the leave-one-out cross-validation process is thus reduced to only $\mathcal{O}(\ell^3)$ operations, which is the same as that of the basic training algorithm for the kernel Fisher discriminant classifier. In the case of S -fold cross-validation, ℓ/S applications of the Bartlett correction formula (6) are performed in each trial to “erase” one of S disjoint sets of ℓ/S training patterns from the kernel Fisher discriminant classifier trained on the entire dataset, again resulting in a computational complexity of $\mathcal{O}(\ell^3)$ operations.

2.1 A Further Refinement

For model selection purposes, we are not principally concerned with the values of the model parameters themselves, but only statistics such as the leave-one-out error rate

$$E = \frac{1}{\ell} \text{card}\{i : y_i(\mathbf{w}_{(i)} \cdot \boldsymbol{\phi}(\mathbf{x}_i) + b_{(i)}) \leq 0\},$$

or equivalently

$$E = \frac{1}{\ell} \text{card}\{i : \text{sign}(y_i)\{r_{(i)}\}_i \leq -1\},$$

where $\{r_{(i)}\}_i = \text{sign}(y_i) - \mathbf{w}_{(i)} \cdot \boldsymbol{\phi}(\mathbf{x}_i) + b_{(i)}$ is the residual error for the i^{th} training pattern during the i^{th} iteration of the leave-one-out cross-validation procedure. It can be shown that

$$\{r_{(i)}\}_i = \frac{1}{1 - h_{ii}} r_i.$$

where $r_i = \text{sign}(y_i) - \mathbf{w} \cdot \boldsymbol{\phi}(\mathbf{x}_i) + b$ is the residual error for the i^{th} training pattern for a kernel Fisher discriminant classifier trained on the entire dataset, $\mathbf{H} = \mathbf{Z}\mathbf{C}^{-1}\mathbf{Z}^T$ is the *hat* matrix of which h_{ii} is the i^{th} element of the leading diagonal. The leave-one-out error rate can thus be evaluated in closed form without explicit inversion of $\mathbf{C}_{(i)} \forall i \in \{1, 2, \dots, \ell\}$, with a computational complexity of only $\mathcal{O}(\ell^3)$ operations.

3 Results

The proposed approximate leave-one-out cross-validation method is evaluated over a series of randomly generated synthetic datasets. In each case, approximately one quarter of the data belong to class \mathcal{C}_1 and three-quarters to class \mathcal{C}_2 . The patterns comprising class \mathcal{C}_1 are drawn from a bivariate Normal distribution with zero mean and unit variance. The patterns forming class \mathcal{C}_2 form an annulus; the radii of the data are drawn from a normal distribution with a mean of 4 and unit variance, and the angles uniformly distributed. The datasets vary in size between 10 and 500 patterns. Figure 1 shows a graph of run-time as a function of the number of training patterns for direct and fast approximate leave-one-out cross-validation methods. Clearly the fast approximate method is considerably faster and exhibits significantly better scaling properties than the direct implementation. Let the relative approximation error be defined as

$$e = \frac{\|\tilde{\mathbf{r}} - \hat{\mathbf{r}}\|^2}{\|\tilde{\mathbf{r}}\|^2},$$

where $\tilde{\mathbf{r}}$ is a vector of leave-one-out residual errors computed via the direct approach and $\hat{\mathbf{r}}$ is the corresponding vector of residual errors resulting from the proposed approximation. Figure 2 shows a graph of the mean relative approximation error, as a function of the number of training patterns. The approximation error is small for datasets of more than ≈ 30 training patterns.

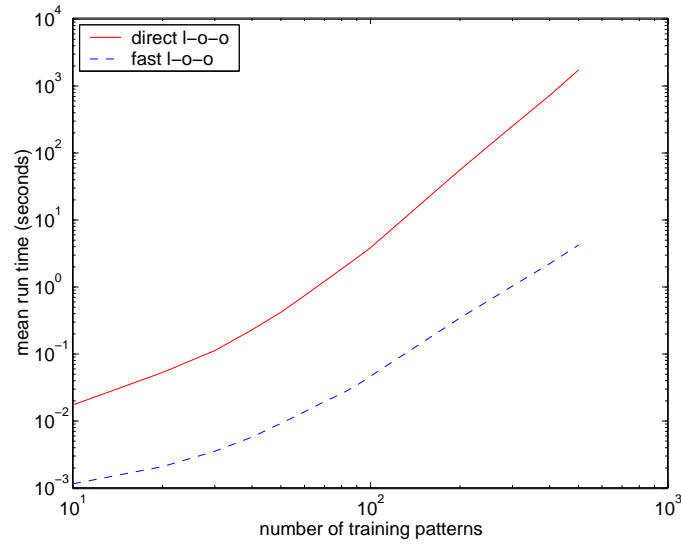


Figure 1: Graph of run-time as a function of the number of training patterns for leave-one-out cross-validation of kernel Fisher discriminant classifiers via direct and fast approximate methods (mean of 20 trials).

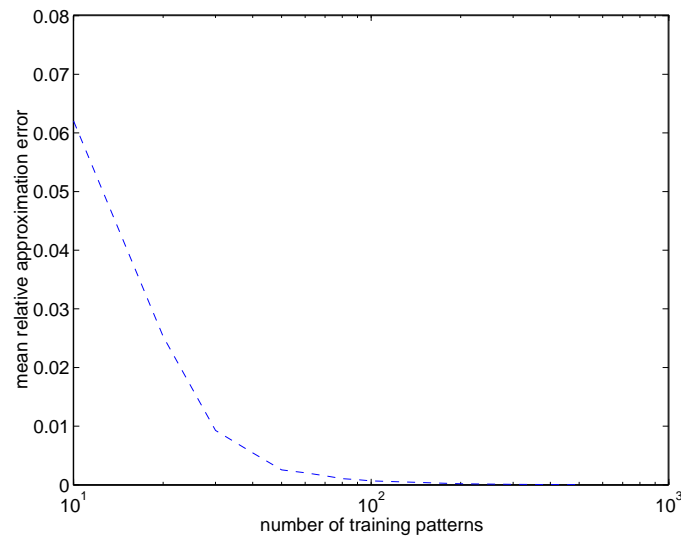


Figure 2: Graph of the mean relative approximation error as a function of the number of training patterns for the proposed fast approximate leave-one-out cross-validation method (mean of 20 trials).

4 Summary

In this paper we have generalised an existing algorithm for leave-one-out cross-validation of multi-variate linear regression models (see e.g. [6]) to provide an approximation to the leave-one-out error rate of kernel Fisher discriminant classifiers. The proposed algorithm allows approximate leave-one-out cross-validation of this class of model with a computational complexity of only $\mathcal{O}(\ell^3)$ operations, instead of the $\mathcal{O}(\ell^4)$ of a direct approach. Furthermore, profiling information reveals that, providing \mathbf{C}^{-1} is cached, the time taken to estimate the leave-one-out error rate is considerably *less* than the time taken to train the KFD classifier. As a result leave-one-out cross-validation becomes a practical model selection criterion in far larger scale applications of KFD models.

5 Acknowledgements

The authors would like to thank Rob Foxall for his helpful comments on previous drafts of this manuscript. This work was supported by the Royal Society (grant number RSRG-22270).

References

- [1] S. Mika, G. Rätsch, J. Weston, B. Schölkopf, and K.-R. Müller. Fisher discriminant analysis with kernels. In *Neural Networks for Signal Processing*, volume IX, pages 41–48. IEEE Press, New York, 1999.
- [2] J. A. K. Suykens and J. Vandewalle. Least squares support vector machine classifiers. *Neural Processing Letters*, 9(3):293–300, June 1990.
- [3] J. Xu, X. Zhang, and Y. Li. Kernel MSE algorithm: A unified framework for KFD, LS-SVM and KRR. In *Proc. IJCNN*, pages 1486–1491, Washington, DC, July 2001.
- [4] N. Cristianini and J. Shawe-Taylor. *An Introduction to Support Vector Machines (and other kernel-based learning methods)*. Cambridge University Press, Cambridge, U.K., 2000.
- [5] M. S. Bartlett. An inverse matrix adjustment arising in discriminant analysis. *Annals of Mathematical Statistics*, 22(1):107–111, 1951.
- [6] S. Weisberg. *Applied linear regression*. John Wiley and Sons, New York, second edition, 1985.