

This article was downloaded by: [128.30.101.2]

On: 05 October 2014, At: 17:06

Publisher: Taylor & Francis

Informa Ltd Registered in England and Wales Registered Number: 1072954 Registered office: Mortimer House, 37-41 Mortimer Street, London W1T 3JH, UK

## Journal of Computational and Graphical Statistics

Publication details, including instructions for authors and subscription information:

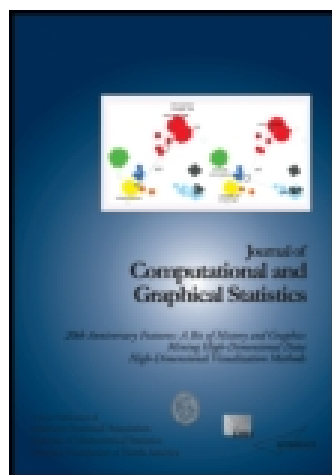
<http://www.tandfonline.com/loi/ucgs20>

### Kernel Logistic Regression and the Import Vector Machine

Ji Zhu<sup>a</sup> & Trevor Hastie<sup>a</sup>

<sup>a</sup> Ji Zhu is Assistant Professor, Department of Statistics, University of Michigan, Ann Arbor, MI 48109-1092 (Email: . Trevor Hastie is Professor, Department of Statistics, Stanford University, Stanford, CA 94305 .

Published online: 01 Jan 2012.



To cite this article: Ji Zhu & Trevor Hastie (2005) Kernel Logistic Regression and the Import Vector Machine, Journal of Computational and Graphical Statistics, 14:1, 185-205, DOI: [10.1198/106186005X25619](https://doi.org/10.1198/106186005X25619)

To link to this article: <http://dx.doi.org/10.1198/106186005X25619>

PLEASE SCROLL DOWN FOR ARTICLE

Taylor & Francis makes every effort to ensure the accuracy of all the information (the "Content") contained in the publications on our platform. However, Taylor & Francis, our agents, and our licensors make no representations or warranties whatsoever as to the accuracy, completeness, or suitability for any purpose of the Content. Any opinions and views expressed in this publication are the opinions and views of the authors, and are not the views of or endorsed by Taylor & Francis. The accuracy of the Content should not be relied upon and should be independently verified with primary sources of information. Taylor and Francis shall not be liable for any losses, actions, claims, proceedings, demands, costs, expenses, damages, and other liabilities whatsoever or howsoever caused arising directly or indirectly in connection with, in relation to or arising out of the use of the Content.

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden. Terms & Conditions of access and use can be found at <http://www.tandfonline.com/page/terms-and-conditions>

# Kernel Logistic Regression and the Import Vector Machine

Ji ZHU and Trevor HASTIE

The support vector machine (SVM) is known for its good performance in two-class classification, but its extension to multiclass classification is still an ongoing research issue. In this article, we propose a new approach for classification, called the import vector machine (IVM), which is built on kernel logistic regression (KLR). We show that the IVM not only performs as well as the SVM in two-class classification, but also can naturally be generalized to the multiclass case. Furthermore, the IVM provides an estimate of the underlying probability. Similar to the support points of the SVM, the IVM model uses only a fraction of the training data to index kernel basis functions, typically a much smaller fraction than the SVM. This gives the IVM a potential computational advantage over the SVM.

**Key Words:** Classification; Kernel methods; Multiclass learning; Radial basis; Reproducing kernel Hilbert space (RKHS); Support vector machines.

## 1. INTRODUCTION

In standard classification problems, we are given a set of training data  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$ , where the input  $\mathbf{x}_i \in \mathcal{R}^p$  and the output  $y_i$  is qualitative and assumes values in a finite set  $\mathcal{C}$ , for example,  $\mathcal{C} = \{1, 2, \dots, C\}$ . We wish to find a classification rule from the training data, so that when given a new input  $\mathbf{x}$ , we can assign a class  $c$  from  $\mathcal{C}$  to it. Usually it is assumed that the training data are an independently and identically distributed sample from an unknown probability distribution  $P(X, Y)$ .

The support vector machine (SVM) works well in two-class classification, that is,  $y \in \{-1, 1\}$ , but its appropriate extension to the multiclass case is still an ongoing research issue (e.g., Vapnik 1998; Weston and Watkins 1999; Bredensteiner and Bennett 1999; Lee, Lin, and Wahba 2002). Another property of the SVM is that it only estimates  $\text{sign}[p(\mathbf{x}) - 1/2]$

---

Ji Zhu is Assistant Professor, Department of Statistics, University of Michigan, Ann Arbor, MI 48109-1092 (E-mail: jizhu@umich.edu). Trevor Hastie is Professor, Department of Statistics, Stanford University, Stanford, CA 94305 (E-mail: hastie@stat.stanford.edu).

©2005 American Statistical Association, Institute of Mathematical Statistics,  
and Interface Foundation of North America

*Journal of Computational and Graphical Statistics*, Volume 14, Number 1, Pages 185–205  
DOI: 10.1198/106186005X25619

(Lin 2002), while the probability  $p(\mathbf{x})$  is often of interest itself, where  $p(\mathbf{x}) = P(Y = 1|X = \mathbf{x})$  is the conditional probability of a point being in class 1 given  $X = \mathbf{x}$ . In this article, we propose a new approach, called the import vector machine (IVM), to address the classification problem. We show that the IVM not only performs as well as the SVM in two-class classification, but also can naturally be generalized to the multiclass case. Furthermore, the IVM provides an estimate of the probability  $p(\mathbf{x})$ . Similar to the *support points* of the SVM, the IVM model uses only a fraction of the training data to index the kernel basis functions. We call these training data *import points*. The computational cost of the SVM is  $O(n^2 n_s)$  (e.g., Kaufman 1999), where  $n_s$  is the number of support points and  $n_s$  usually increases linearly with  $n$ , while the computational cost of the IVM is  $O(n^2 m^2)$ , where  $m$  is the number of import points. Because  $m$  does not tend to increase as  $n$  increases, the IVM can be faster than the SVM. Empirical results show that the number of import points is usually much less than the number of support points.

In Section 2, we briefly review some results of the SVM for two-class classification and compare it with kernel logistic regression (KLR). In Section 3, we propose our IVM algorithm. In Section 4, we show some numerical results. In Section 5, we generalize the IVM to the multiclass case.

## 2. SUPPORT VECTOR MACHINES AND KERNEL LOGISTIC REGRESSION

The standard SVM produces a nonlinear classification boundary in the original input space by constructing a linear boundary in a transformed version of the original input space. The dimension of the transformed space can be very large, even infinite in some cases. This seemingly prohibitive computation is achieved through a positive definite reproducing kernel  $K(\cdot, \cdot)$ , which gives the inner product in the transformed space.

Many people have noted the relationship between the SVM and regularized function estimation in the reproducing kernel Hilbert spaces (RKHS). An overview can be found in Burges (1998), Evgeniou, Pontil, and Poggio (1999), Wahba (1999), and Hastie, Tibshirani, and Friedman (2001). Fitting an SVM is equivalent to

$$\min_{f \in \mathcal{H}_K} \frac{1}{n} \sum_{i=1}^n [1 - y_i f(\mathbf{x}_i)]_+ + \frac{\lambda}{2} \|f\|_{\mathcal{H}_K}^2, \quad (2.1)$$

where  $\mathcal{H}_K$  is the RKHS generated by the kernel  $K(\cdot, \cdot)$ . The classification rule is given by  $\text{sign}[f(\mathbf{x})]$ . For the purpose of simple notation, we omit the constant term in  $f(\mathbf{x})$ .

By the representer theorem (Kimeldorf and Wahba 1971), the optimal  $f(\mathbf{x})$  has the form:

$$f(\mathbf{x}) = \sum_{i=1}^n a_i K(\mathbf{x}, \mathbf{x}_i). \quad (2.2)$$

It often happens that a sizeable fraction of the  $n$  values of  $a_i$  can be zero. This is a consequence of the truncation property of the first part of criterion (2.1). This seems to be an

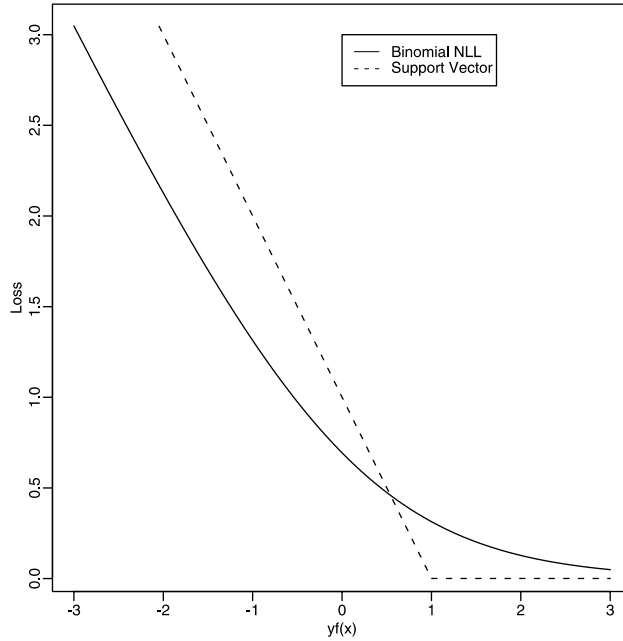


Figure 1. Two loss functions,  $y \in \{-1, 1\}$ .

attractive property, because only the points on the wrong side of the classification boundary, and those on the right side but near the boundary have an influence in determining the position of the boundary, and hence have nonzero  $a_i$ 's. The corresponding  $\mathbf{x}_i$ 's are called *support points*.

Notice that (2.1) has the form *loss + penalty*. The loss function  $(1 - yf)_+$  is plotted in Figure 1, along with the negative log-likelihood (NLL) of the binomial distribution. As we can see, the NLL of the binomial distribution has a similar shape to that of the SVM: both increase linearly as  $yf$  gets very small (negative) and both encourage  $y$  and  $f$  to have the same sign. If we replace  $(1 - yf)_+$  in (2.1) with  $\ln(1 + e^{-yf})$ , the NLL of the binomial distribution, the problem becomes a kernel logistic regression (KLR) problem:

$$\min_{f \in \mathcal{H}_K} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y_i f(\mathbf{x}_i)}) + \frac{\lambda}{2} \|f\|_{\mathcal{H}_K}^2. \quad (2.3)$$

Because of the similarity between the two loss functions, we expect that the fitted function performs similarly to the SVM for two-class classification.

There are two immediate advantages of making such a replacement: (a) Besides giving a classification rule, KLR also offers a natural estimate of the probability  $p(\mathbf{x}) = e^{f(\mathbf{x})} / (1 + e^{f(\mathbf{x})})$ , while the SVM only estimates  $\text{sign}[p(\mathbf{x}) - 1/2]$  (Lin 2002); (b) KLR can naturally be generalized to the multiclass case through kernel multi-logit regression. However, because KLR compromises the hinge loss function of the SVM, it no longer has the support points property; in other words, all the  $a_i$ 's in (2.2) are nonzero.

KLR is a well-studied problem; see Green and Yandell (1985), Hastie and Tibshirani

(1990), Wahba, Gu, Wang, and Chappell (1995) and the references therein; however, they are all under the smoothing spline analysis of variance scheme.

We use a simulation example to illustrate the similar performances between KLR and the SVM. The data in each class are simulated from a mixture of Gaussian distribution (Hastie et al. 2001): first we generate 10 means  $\mu_k$  from a bivariate Gaussian distribution  $N((1, 0)^T, \mathbf{I})$  and label this class +1. Similarly, 10 more are drawn from  $N((0, 1)^T, \mathbf{I})$  and labeled class -1. Then for each class, we generate 100 observations as follows: for each observation, we pick an  $\mu_k$  at random with probability 1/10, and then generate a  $N(\mu_k, \mathbf{I}/5)$ , thus leading to a mixture of Gaussian clusters for each class.

We use the radial basis kernel

$$K(\mathbf{x}_i, \mathbf{x}_{i'}) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_{i'}\|^2}{2\sigma^2}}. \quad (2.4)$$

The regularization parameter  $\lambda$  is chosen to achieve good misclassification error. The results are shown in Figure 2. The radial basis kernel produces a boundary quite close to the Bayes optimal boundary for this simulation. We see that the fitted model of KLR is quite similar in classification performance to that of the SVM. In addition to a classification boundary, since KLR estimates the log-odds of class probabilities, it can also produce probability contours (Figure 2).

## 2.1 KLR AS A MARGIN MAXIMIZER

The SVM was initiated as a method to maximize the margin, that is,  $\min_i y_i f(\mathbf{x}_i)$ , of the training data; KLR is motivated by the similarity in shape between the NLL of the binomial distribution and the hinge loss of the SVM. Then a natural question is: what does KLR do with the margin?

Suppose the dictionary of the basis functions of the transformed feature space is

$$\{h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_q(\mathbf{x})\},$$

where  $q$  is the dimension of the transformed feature space. Note if  $q = p$  and  $h_j(\mathbf{x})$  is the  $j$ th component of  $\mathbf{x}$ , the transformed feature space is reduced to the original input space. The classification boundary, a hyperplane in the transformed feature space, is given by

$$\{\mathbf{x} : f(\mathbf{x}) = \beta_0 + \mathbf{h}(\mathbf{x})^T \boldsymbol{\beta} = 0\}.$$

Suppose the transformed feature space is so rich that the training data are separable, then the margin-maximizing SVM can be written as:

$$\max_{\beta_0, \boldsymbol{\beta}, \|\boldsymbol{\beta}\|_2=1} D \quad (2.5)$$

$$\text{subject to } y_i (\beta_0 + \mathbf{h}(\mathbf{x}_i)^T \boldsymbol{\beta}) \geq D, \quad i = 1, \dots, n \quad (2.6)$$

where  $D$  is the shortest distance from the training data to the separating hyperplane and is defined as the *margin* (Burges 1998).

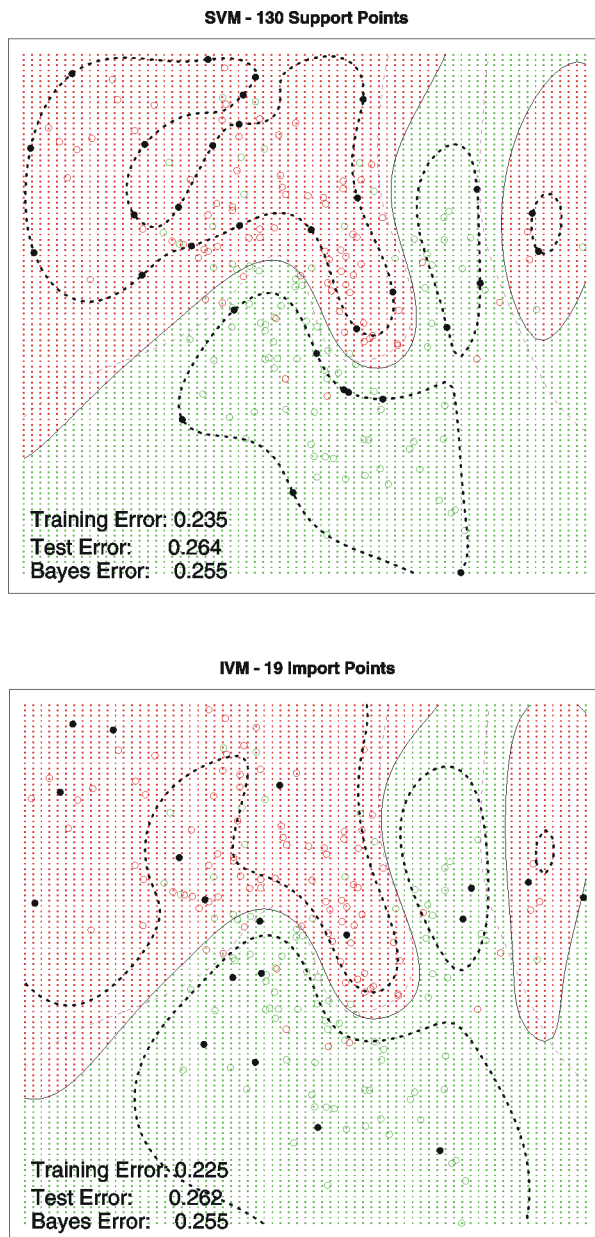


Figure 2. The solid black lines are classification boundaries; the dashed purple lines are Bayes optimal boundaries. For the SVM, the dotted black lines are the edges of the margins and the black points are the points exactly on the edges of the margin. For the IVM, the dotted black lines are the  $p_1(\mathbf{x}) = .25$  and  $.75$  lines and the black points are the import points. Because the classification boundaries of KLR and the IVM are almost identical, we omit the picture of KLR here.

Now consider an equivalent setup of KLR:

$$\min_{\beta_0, \beta} \sum_{i=1}^n \ln \left( 1 + e^{-y_i f(\mathbf{x}_i)} \right) \quad (2.7)$$

$$\text{subject to} \quad \|\beta\|_2^2 \leq s \quad (2.8)$$

$$f(\mathbf{x}_i) = \beta_0 + \mathbf{h}(\mathbf{x}_i)^T \beta, \quad i = 1, \dots, n. \quad (2.9)$$

Then we have Theorem 1.

**Theorem 1.** *Suppose the training data are separable, that is,  $\exists \beta_0, \beta$ , s.t.  $y_i(\beta_0 + \mathbf{h}(\mathbf{x}_i)^T \beta) > 0$ ,  $\forall i$ . Let the solution of (2.7)–(2.9) be denoted by  $\hat{\beta}(s)$ , then*

$$\frac{\hat{\beta}(s)}{s} \rightarrow \beta^* \quad \text{as } s \rightarrow \infty,$$

where  $\beta^*$  is the solution of the margin-maximizing SVM (2.5)–(2.6), if  $\beta^*$  is unique.

If  $\beta^*$  is not unique, then  $\frac{\hat{\beta}(s)}{s}$  may have multiple convergence points, but they will all represent margin-maximizing separating hyperplanes.

The proof of the theorem appears in the Appendix. Theorem 1 implies that KLR, similar to the SVM, can also be considered as a margin maximizer. We have also proved a more general theorem relating loss functions and margin maximizers in Rosset, Zhu, and Hastie (2004).

## 2.2 COMPUTATIONAL CONSIDERATIONS

Because (2.3) is convex, it is natural to use the Newton-Raphson method to fit KLR. In order to guarantee convergence, suitable bisection steps can be combined with the Newton-Raphson iterations. The drawback of the Newton-Raphson method is that in each iteration, an  $n \times n$  matrix needs to be inverted. Therefore the computational cost of KLR is  $O(n^3)$ . Recently Keerthi, Duan, Shevade, and Poo (2002) proposed a dual algorithm for KLR which avoids inverting huge matrices. It follows the spirit of the popular sequential minimal optimization (SMO) algorithm (Platt 1999). Preliminary computational experiments show that the algorithm is robust and fast. Keerthi et al. (2002) described the algorithm for two-class classification; we have generalized it to the multiclass case (Zhu and Hastie 2004).

Although the sequential minimal optimization method helps reduce the computational cost of KLR, in the fitted model (2.2), all the  $a_i$ 's are nonzero. Hence, unlike the SVM, KLR does not allow for data compression and does not have the advantage of less storage and quicker evaluation.

In this article, we propose an import vector machine (IVM) model that finds a submodel to approximate the full model (2.2) given by KLR. The submodel has the form:

$$f(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathcal{S}} a_i K(\mathbf{x}, \mathbf{x}_i), \quad (2.10)$$

where  $\mathcal{S}$  is a subset of the training data  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , and the data in  $\mathcal{S}$  are called *import points*. The advantage of this submodel is that the computational cost is reduced, especially

for large training datasets, while not jeopardizing the performance in classification; and since only a subset of the training data are used to index the fitted model, data compression is achieved.

Several other researchers have also investigated techniques in selecting the subset  $\mathcal{S}$ . Lin et al. (2000) divided the training data into several clusters, then randomly selected a representative from each cluster to make up  $\mathcal{S}$ . Smola and Schölkopf (2000) developed a greedy technique to sequentially select  $m$  columns of the kernel matrix  $[K(\mathbf{x}_i, \mathbf{x}_{i'})]_{n \times n}$ , such that the span of these  $m$  columns approximates the span of  $[K(\mathbf{x}_i, \mathbf{x}_{i'})]_{n \times n}$  well in the Frobenius norm. Williams and Seeger (2001) proposed randomly selecting  $m$  points of the training data, then using the Nystrom method to approximate the eigen-decomposition of the kernel matrix  $[K(\mathbf{x}_i, \mathbf{x}_{i'})]_{n \times n}$ , and expanding the results back up to  $n$  dimensions. None of these methods uses the output  $y_i$  in selecting the subset  $\mathcal{S}$  (i.e., the procedure only involves  $\mathbf{x}_i$ ). The IVM algorithm uses both the output  $y_i$  and the input  $\mathbf{x}_i$  to select the subset  $\mathcal{S}$ , in such a way that the resulting fit approximates the full model well. The idea is similar to that used in Luo and Wahba (1997), which also used the output  $y_i$  to select a subset of the training data, but under the regression scheme.

### 3. IMPORT VECTOR MACHINES

In KLR, we want to minimize

$$H = \frac{1}{n} \sum_{i=1}^n \ln \left( 1 + e^{-y_i f(\mathbf{x}_i)} \right) + \frac{\lambda}{2} \|f\|_{\mathcal{H}_K}^2. \quad (3.1)$$

Let

$$p_i = \frac{1}{1 + e^{-y_i f(\mathbf{x}_i)}}, \quad i = 1, \dots, n \quad (3.2)$$

$$\mathbf{a} = (a_1, \dots, a_n)^T \quad (3.3)$$

$$\mathbf{p} = (p_1, \dots, p_n)^T \quad (3.4)$$

$$\mathbf{y} = (y_1, \dots, y_n)^T \quad (3.5)$$

$$\mathbf{K}_1 = (K(\mathbf{x}_i, \mathbf{x}_{i'}))_{i,i'=1}^n \quad (3.6)$$

$$\mathbf{K}_2 = \mathbf{K}_1 \quad (3.7)$$

$$\mathbf{W} = \text{diag}(p_1(1-p_1), \dots, p_n(1-p_n)). \quad (3.8)$$

With some abuse of notation, using (2.2), (3.1) can be written in a finite dimensional form:

$$H = \frac{1}{n} \mathbf{1}^T \ln \left( \mathbf{1} + e^{-\mathbf{y} \cdot (\mathbf{K}_1 \mathbf{a})} \right) + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K}_2 \mathbf{a}, \quad (3.9)$$

where “ $\cdot$ ” denotes element-wise multiplication. To find  $\mathbf{a}$ , we set the derivative of  $H$  with respect to  $\mathbf{a}$  equal to 0, and use the Newton-Raphson method to iteratively solve the score equation. It can be shown that the Newton-Raphson step is a weighted least squares step:

$$\mathbf{a}^{(k)} = \left( \frac{1}{n} \mathbf{K}_1^T \mathbf{W} \mathbf{K}_1 + \lambda \mathbf{K}_2 \right)^{-1} \mathbf{K}_1^T \mathbf{W} \mathbf{z}, \quad (3.10)$$



where  $\mathbf{a}^{(k)}$  is the value of  $\mathbf{a}$  in the  $k$ th step, and

$$\mathbf{z} = \frac{1}{n} \left( \mathbf{K}_1 \mathbf{a}^{(k-1)} + \mathbf{W}^{-1}(\mathbf{y} \cdot \mathbf{p}) \right). \quad (3.11)$$

### 3.1 BASIC ALGORITHM

As mentioned in Section 2, we want to find a subset  $\mathcal{S}$  of  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , such that the submodel (2.10) is a good approximation of the full model (2.2). Because searching for every subset  $\mathcal{S}$  is a combinatorial problem and computationally prohibitive, we use the following greedy forward strategy: we start with the null model, that is,  $\mathcal{S} = \emptyset$ , then iteratively build up  $\mathcal{S}$  one element at a time. Basically, we look for a data point among  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \setminus \mathcal{S}$ , so that if it is added into the current  $\mathcal{S}$ , the new submodel will decrease the regularized negative log-likelihood the most:

**Algorithm 1:** *Basic IVM Algorithm.*

1. Let  $\mathcal{S} = \emptyset$ ,  $\mathcal{L} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $k = 1$ .
2. For each  $\mathbf{x}_l \in \mathcal{L}$ , let

$$f_l(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathcal{S} \cup \{\mathbf{x}_l\}} a_i K(\mathbf{x}, \mathbf{x}_i).$$

Use the Newton-Raphson method to find  $\mathbf{a}$  to minimize

$$H(\mathbf{x}_l) = \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i f_l(\mathbf{x}_i))) + \frac{\lambda}{2} \|f_l(\mathbf{x})\|_{\mathcal{H}_K}^2 \quad (3.12)$$

$$= \frac{1}{n} \mathbf{1}^T \ln(\mathbf{1} + \exp(-\mathbf{y} \cdot (\mathbf{K}_1^l \mathbf{a}))) + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K}_2^l \mathbf{a}, \quad (3.13)$$

where the regressor matrix

$$\mathbf{K}_1^l = (K(\mathbf{x}_i, \mathbf{x}_{i'}))_{n \times k}, \quad \mathbf{x}_i \in \{\mathbf{x}_1, \dots, \mathbf{x}_n\}, \mathbf{x}_{i'} \in \mathcal{S} \cup \{\mathbf{x}_l\};$$

the regularization matrix

$$\mathbf{K}_2^l = (K(\mathbf{x}_i, \mathbf{x}_{i'}))_{k \times k}, \quad \mathbf{x}_i, \mathbf{x}_{i'} \in \mathcal{S} \cup \{\mathbf{x}_l\};$$

and  $k = |\mathcal{S}| + 1$ .

3. Find

$$\mathbf{x}_{l^*} = \operatorname{argmin}_{\mathbf{x}_l \in \mathcal{L}} H(\mathbf{x}_l).$$

Let  $\mathcal{S} = \mathcal{S} \cup \{\mathbf{x}_{l^*}\}$ ,  $\mathcal{L} = \mathcal{L} \setminus \{\mathbf{x}_{l^*}\}$ ,  $H_k = H(\mathbf{x}_{l^*})$ ,  $k = k + 1$ .

4. Repeat Steps 2 and 3 until  $H_k$  converges.

The points in  $\mathcal{S}$  are the import points.

### 3.2 REVISED ALGORITHM

The above algorithm is computationally feasible, but in Step 2 we need to use the Newton-Raphson method to find  $\mathbf{a}$  iteratively. When the number of import points  $k$  becomes large, the Newton-Raphson computation can be expensive. To reduce this computation, we use a further approximation.

Instead of iteratively computing  $\mathbf{a}^{(k)}$  until it converges, we can just do a one-step iteration, and use it as an approximation to the converged one. This is equivalent to approximating the negative binomial log-likelihood with a different weighted quadratic loss function at each iteration. To get a good approximation, we take advantage of the fitted result from the current “optimal”  $\mathcal{S}$ , that is, the submodel when  $|\mathcal{S}| = k - 1$ , and use it to compute  $\mathbf{z}$  in (3.11). This one-step update is similar to the score test in generalized linear models (GLM); but the latter does not have a penalty term. The updating formula allows the weighted regression (3.10) to be computed in  $O(nm)$  time.

Hence, we have the revised Steps 1 and 2 for the basic algorithm:

**Algorithm 2:** *Revised Steps 1 and 2.*

- 1\* Let  $\mathcal{S} = \emptyset$ ,  $\mathcal{L} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $k = 1$ . Let  $\mathbf{a}^{(0)} = \mathbf{0}$ , hence  $\mathbf{z} = 2\mathbf{y}/n$ .
- 2\* For each  $\mathbf{x}_l \in \mathcal{L}$ , correspondingly augment  $\mathbf{K}_1$  with a column, and  $\mathbf{K}_2$  with a column and a row. Use the current submodel from iteration  $(k - 1)$  to compute  $\mathbf{z}$  in (3.11) and use the updating formula (3.10) to find  $\mathbf{a}$ . Compute (3.13).

### 3.3 STOPPING RULE FOR ADDING POINT TO $\mathcal{S}$

In Step 4 of the basic algorithm, we need to decide whether  $H_k$  has converged. A natural stopping rule is to look at the regularized NLL. Let  $H_1, H_2, \dots$  be the sequence of regularized NLL's obtained in Step 3. At each step  $k$ , we compare  $H_k$  with  $H_{k-\Delta k}$ , where  $\Delta k$  is a prechosen small integer, for example  $\Delta k = 1$ . If the ratio  $\frac{|H_k - H_{k-\Delta k}|}{|H_k|}$  is less than some prechosen small number  $\epsilon$ , for example,  $\epsilon = .001$ , we stop adding new import points to  $\mathcal{S}$ .

### 3.4 CHOOSING THE REGULARIZATION PARAMETER $\lambda$

So far, we have assumed that the regularization parameter  $\lambda$  is fixed. In practice, we also need to choose an “optimal”  $\lambda$ . We can randomly split all the data into a training set and a tuning set, and use the misclassification error on the tuning set as a criterion for choosing  $\lambda$ . To reduce the computation, we take advantage of the fact that the regularized NLL converges faster for a larger  $\lambda$ . Thus, instead of running the entire revised algorithm for each  $\lambda$ , we propose the following procedure, which combines both adding import points to  $\mathcal{S}$  and choosing the optimal  $\lambda$ .

**Algorithm 3:** *Simultaneous Selection of  $\mathcal{S}$  and  $\lambda$ .*

1. Start with a large regularization parameter  $\lambda$ .

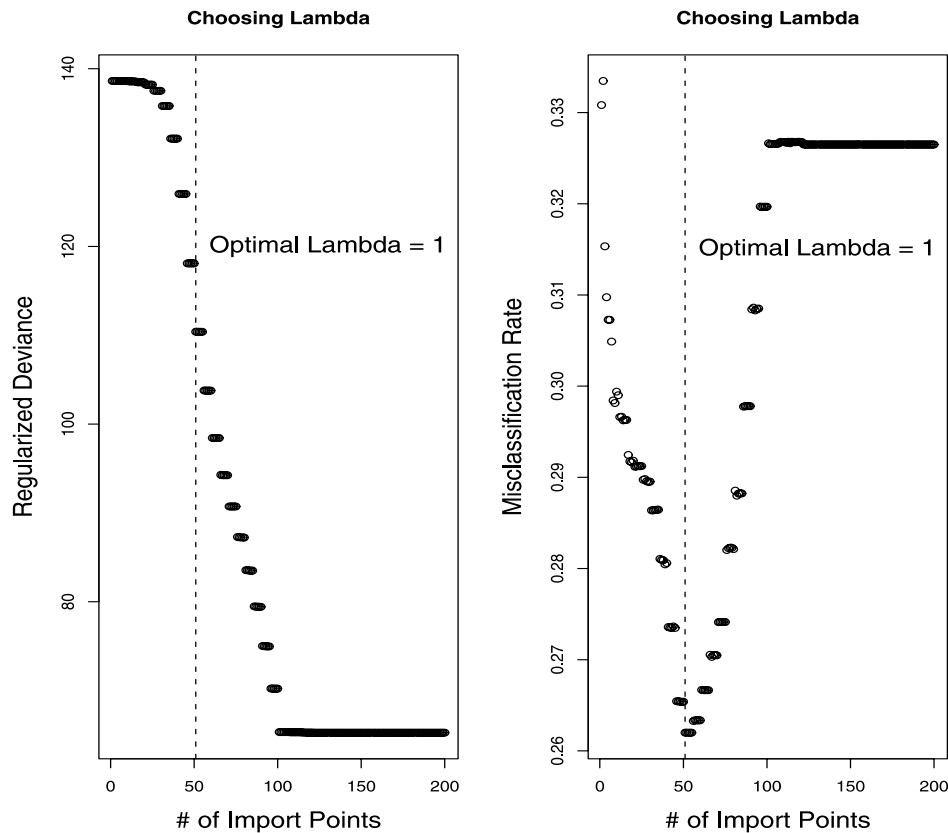


Figure 3. Radial kernel is used.  $n = 200$ ,  $\sigma^2 = .7$ ,  $\Delta k = 3$ ,  $\epsilon = .001$ ,  $\lambda$  decreases from  $e^{10}$  to  $e^{-10}$ . The minimum misclassification rate .262 is found to correspond to  $\lambda = 1$ .

2. Let  $\mathcal{S} = \emptyset$ ,  $\mathcal{L} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ ,  $k = 1$ . Let  $\mathbf{a}^{(0)} = \mathbf{0}$ , hence  $\mathbf{z} = 2\mathbf{y}/n$ .
3. Run Steps 2\*, 3, and 4 of the revised Algorithm 2, until the stopping criterion is satisfied at  $\mathcal{S} = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}\}$ . Along the way, also compute the misclassification error on the tuning set.
4. Decrease  $\lambda$  to a smaller value.
5. Repeat Steps 3 and 4, starting with  $\mathcal{S} = \{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}\}$ .

We choose the optimal  $\lambda$  as the one that corresponds to the minimum misclassification error on the tuning set.

## 4. NUMERICAL RESULTS

In this section, we use both simulation and real data to illustrate the IVM method.

### 4.1 SIMULATION RESULTS

The data are generated in the same way as Figure 2. The simulation results are shown in Figures 3–5.

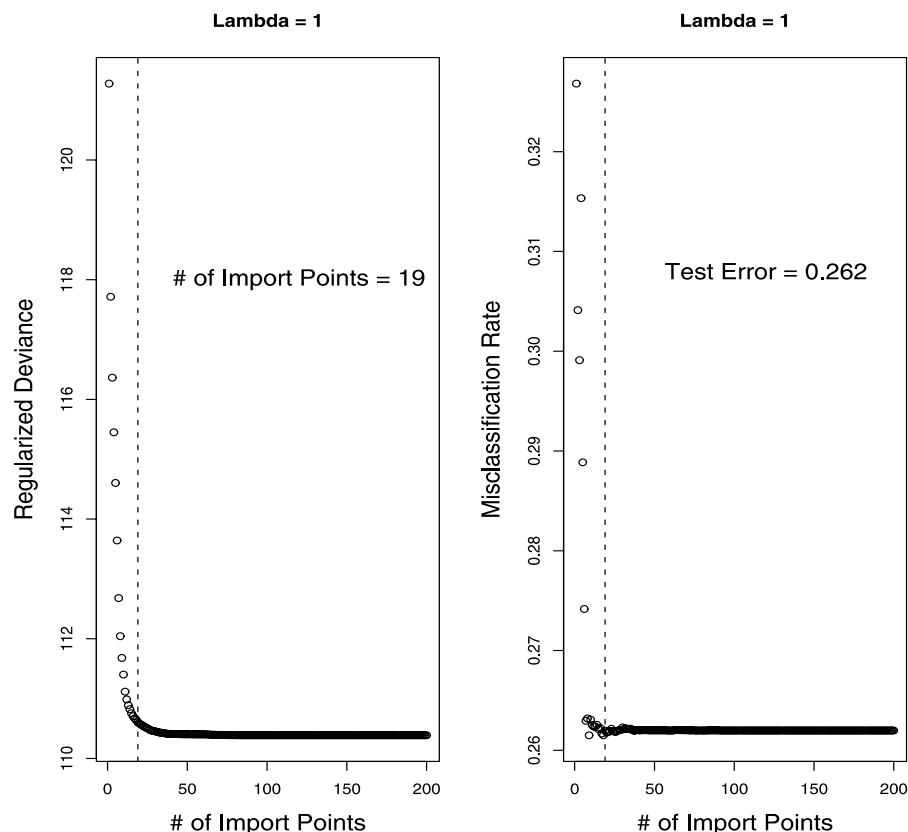


Figure 4. Radial kernel is used.  $n = 200$ ,  $\sigma^2 = .7$ ,  $\Delta k = 3$ ,  $\epsilon = .001$ ,  $\lambda = 1$ . The stopping criterion is satisfied when  $|\mathcal{S}| = 19$ .

Figure 3 shows how the tuning parameter  $\lambda$  is selected. The optimal  $\lambda$  is found to be equal to 1 and corresponds to a misclassification rate .262. Figure 4 fixes the tuning parameter to  $\lambda = 1$  and finds 19 import points. Figure 2 compares the results of the SVM and the IVM: the SVM has 130 support points, and the IVM uses 19 import points; they give similar classification boundaries. Figure 5 is for the same simulation but different sizes of training data:  $n = 200, 400, 600, 800$ . We see that as the size of training data  $n$  increases, the number of import points does not tend to increase.

*Remarks.* The support points of the SVM are those which are close to the classification boundary or misclassified and usually have large weights  $p(\mathbf{x})(1 - p(\mathbf{x}))$ . The import points of the IVM are those that decrease the regularized NLL the most, and can be either close to or far from the classification boundary. This difference is natural, because the SVM is concerned only with the classification sign  $[p(\mathbf{x}) - 1/2]$ , while the IVM also focuses on the unknown probability  $p(\mathbf{x})$ . Though points away from the classification boundary do not contribute to determining the position of the classification boundary, they may contribute to estimating the unknown probability  $p(\mathbf{x})$ . The total computational cost of the SVM

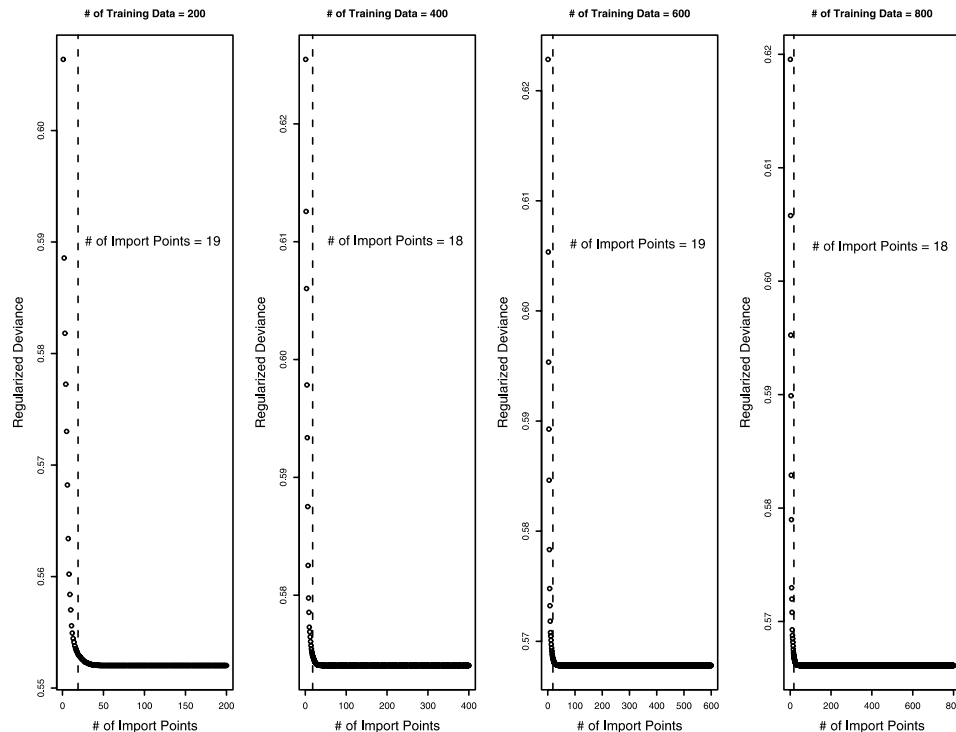


Figure 5. The data are generated in the same way as Figures 2–4. Radial kernel is used.  $\sigma^2 = .7$ ,  $\lambda = 1$ ,  $\Delta k = 3$ ,  $\epsilon = .001$ . The sizes of training data are  $n = 200, 400, 600, 800$ , and the corresponding numbers of import points are 19, 18, 19, 18.

is  $O(n^2 n_s)$  (e.g., Kaufman 1999), where  $n_s$  is the number of support points, while the computational cost of the IVM method is  $O(n^2 m^2)$ , where  $m$  is the number of import points. Since  $m$  does not tend to increase as  $n$  increases, as illustrated in Figure 5, the computational cost of the IVM can be smaller than that of the SVM.

## 4.2 REAL DATA RESULTS

In this section, we compare the performances of the IVM and the SVM on some real datasets. Ten benchmark datasets are used for this purpose: Banana, Breast-cancer, Flare-solar, German, Heart, Image, Ringnorm, Splice, Thyroid, Titanic, Twonorm and Waveform. Detailed information about these datasets can be found in Rätsch, Onoda, and Müller (2000).

Table 1 contains a summary of these datasets. Radial kernel (2.4) is used throughout these datasets. The parameters  $\sigma$  and  $\lambda$  are fixed at specific values that are optimal for the SVMs generalization performance (Rätsch et al. 2000). Each dataset has 20 realizations of the training and test data. The results are in Table 2 and Table 3. The number outside each bracket is the mean over 20 realizations of the training and test data, and the number in each bracket is the standard error. From Table 2, we can see that the IVM performs as well as the SVM in classification on these benchmark datasets. From Table 3, we can see that

Table 1. Summary of the Ten Benchmark Datasets.  $n$  is the size of the training data,  $p$  is the dimension of the original input,  $\sigma^2$  is the parameter of the radial kernel,  $\lambda$  is the tuning parameter, and  $N$  is the size of the test data.

<i>Dataset</i>	$n$	$p$	$\sigma^2$	$\lambda$	$N$
Banana	400	2	1	$3.16 \times 10^{-3}$	4900
Breast-cancer	200	9	50	$6.58 \times 10^{-2}$	77
Flare-solar	666	9	30	.978	400
German	700	20	55	.316	300
Heart	170	13	120	.316	100
Image	1300	18	3	.002	1010
Ringnorm	400	20	10	$10^{-9}$	7000
Thyroid	140	5	3	.1	75
Titanic	150	3	2	$10^{-5}$	2051
Twonorm	400	20	40	.316	7000
Waveform	400	21	20	1	4600

the IVM typically uses a much smaller fraction of the training data than the SVM to index kernel basis functions. This may give the IVM a computational advantage over the SVM.

## 5. MULTICLASS CASE

In this section, we briefly describe a generalization of the IVM to multiclass classification. Suppose there are  $C$  classes. The conditional probability of a point being in class  $c$  given  $X = \mathbf{x}$  is denoted as  $p_c(\mathbf{x}) = P(Y = c | X = \mathbf{x})$ . Hence the Bayes classification rule is given by:

$$c(\mathbf{x}) = \operatorname{argmax}_{c \in \{1, \dots, C\}} p_c(\mathbf{x})$$

The model has the form

$$p_1(\mathbf{x}) = \frac{e^{f_1(\mathbf{x})}}{\sum_{c=1}^C e^{f_c(\mathbf{x})}}, \quad (5.1)$$

$$p_2(\mathbf{x}) = \frac{e^{f_2(\mathbf{x})}}{\sum_{c=1}^C e^{f_c(\mathbf{x})}}, \quad (5.2)$$

$$\vdots \quad (5.3)$$

$$p_C(\mathbf{x}) = \frac{e^{f_C(\mathbf{x})}}{\sum_{c=1}^C e^{f_c(\mathbf{x})}}, \quad (5.4)$$

where  $f_c(\mathbf{x}) \in \mathcal{H}_K$ ;  $\mathcal{H}_K$  is the RKHS generated by a positive definite kernel  $K(\cdot, \cdot)$ . Notice that  $f_1(\mathbf{x}), \dots, f_C(\mathbf{x})$  are not identifiable in this model, for if we add a common term to each  $f_c(\mathbf{x}), p_1(\mathbf{x}), \dots, p_C(\mathbf{x})$  will not change. To make  $f_c(\mathbf{x})$  identifiable, we consider the symmetric constraint

$$\sum_{c=1}^C f_c(\mathbf{x}) = 0. \quad (5.5)$$

Table 2. Comparison of Classification Performance of SVM and IVM on Ten Benchmark Datasets

<i>Dataset</i>	<i>SVM Error (%)</i>	<i>IVM Error (%)</i>
Banana	10.78 ( $\pm$ .68)	10.34 ( $\pm$ .46)
Breast-cancer	25.58 ( $\pm$ 4.50)	25.92 ( $\pm$ 4.79)
Flare-solar	32.65 ( $\pm$ 1.42)	33.66 ( $\pm$ 1.64)
German	22.88 ( $\pm$ 2.28)	23.53 ( $\pm$ 2.48)
Heart	15.95 ( $\pm$ 3.14)	15.80 ( $\pm$ 3.49)
Image	3.34 (.70)	3.31 ( $\pm$ .80)
Ringnorm	2.03 ( $\pm$ .19)	1.97 ( $\pm$ .29)
Thyroid	4.80 ( $\pm$ 2.98)	5.00 ( $\pm$ 3.02)
Titanic	22.16 ( $\pm$ .60)	22.39 ( $\pm$ 1.03)
Twonorm	2.90 ( $\pm$ .25)	2.45 ( $\pm$ .15)
Waveform	9.98 ( $\pm$ .43)	10.13 ( $\pm$ .47)

Then the multiclass KLR fits a model to minimize the regularized negative log-likelihood

$$H = -\frac{1}{n} \sum_{i=1}^n \ln p_{y_i}(\mathbf{x}_i) + \frac{\lambda}{2} \|\mathbf{f}\|_{\mathcal{H}_K}^2 \quad (5.6)$$

$$= \frac{1}{n} \sum_{i=1}^n \left[ -\mathbf{y}_i^T \mathbf{f}(\mathbf{x}_i) + \ln \left( e^{f_1(\mathbf{x}_i)} + \dots + e^{f_C(\mathbf{x}_i)} \right) \right] + \frac{\lambda}{2} \|\mathbf{f}\|_{\mathcal{H}_K}^2, \quad (5.7)$$

where  $\mathbf{y}_i$  is a binary  $C$ -vector with values all zero except a 1 in position  $c$  if the class is  $c$ , and

$$\mathbf{f}(\mathbf{x}_i) = (f_1(\mathbf{x}_i), \dots, f_C(\mathbf{x}_i))^T, \quad (5.8)$$

$$\|\mathbf{f}\|_{\mathcal{H}_K}^2 = \sum_{c=1}^C \|f_c\|_{\mathcal{H}_K}^2. \quad (5.9)$$

Using the representer theorem (Kimeldorf and Wahba 1971), one can show that  $f_c(\mathbf{x})$ , which minimizes  $H$ , has the form

$$f_c(\mathbf{x}) = \sum_{i=1}^n a_{ic} K(\mathbf{x}_i, \mathbf{x}). \quad (5.10)$$

Table 3. Comparison of Number of Kernel Basis Used by SVM and IVM on Ten Benchmark Datasets

<i>Dataset</i>	<i># of SV</i>	<i># of IV</i>
Banana	90 ( $\pm$ 10)	21 ( $\pm$ 7)
Breast-cancer	115 ( $\pm$ 5)	14 ( $\pm$ 3)
Flare-solar	597 ( $\pm$ 8)	9 ( $\pm$ 1)
German	407 ( $\pm$ 10)	17 ( $\pm$ 2)
Heart	90 ( $\pm$ 4)	12 ( $\pm$ 2)
Image	221 ( $\pm$ 11)	72 ( $\pm$ 18)
Ringnorm	89 ( $\pm$ 5)	72 ( $\pm$ 30)
Thyroid	21 ( $\pm$ 2)	22 ( $\pm$ 3)
Titanic	69 ( $\pm$ 9)	8 ( $\pm$ 2)
Twonorm	70 ( $\pm$ 5)	24 ( $\pm$ 4)
Waveform	151 ( $\pm$ 9)	26 ( $\pm$ 3)

Hence, (5.6) becomes

$$H = \frac{1}{n} \sum_{i=1}^n \left[ -\mathbf{y}_i^T (\mathbf{K}_1(i, \cdot) \mathbf{A})^T + \ln \left( \mathbf{1}^T e^{(\mathbf{K}_1(i, \cdot) \mathbf{A})^T} \right) \right] + \frac{\lambda}{2} \sum_{c=1}^C \mathbf{a}_c^T \mathbf{K}_2 \mathbf{a}_c, \quad (5.11)$$

where  $\mathbf{A} = (\mathbf{a}_1 \dots \mathbf{a}_C)$ ,  $\mathbf{K}_1$  and  $\mathbf{K}_2$  are defined in the same way as in the two-class case; and  $\mathbf{K}_1(i, \cdot)$  is the  $i$ th row of  $\mathbf{K}_1$ . Notice that in this model, the constraint (5.5) is not necessary anymore, for at the minimum of (5.6),  $\sum_{c=1}^C f_c(\mathbf{x}) = 0$  is automatically satisfied.

## 5.1 MULTICLASS KLR AND MULTICLASS SVM

Similar to Theorem 1, a connection between the multiclass KLR and a multiclass SVM also exists.

In going from the two-class SVM to the multiclass classification, many researchers have proposed various procedures.

In practice, the one-versus-rest scheme is often used: given  $C$  classes, the problem is divided into a series of  $C$  one-versus-rest problems, and each one-versus-rest problem is addressed by a different class-specific SVM classifier (e.g., “class 1” vs. “not class 1”); then a new sample takes the class of the classifier with the largest real valued output  $c^* = \arg\max_{c=1, \dots, C} f_c$ , where  $f_c$  is the real valued output of the  $c$ th SVM classifier.

Instead of solving  $C$  problems, Vapnik (1998) and Weston and Watkins (1999) generalized (2.1) by solving one single optimization problem:

$$\max_{f_c} D \quad (5.12)$$

$$\text{subject to } f_{y_i}(\mathbf{x}_i) - f_c(\mathbf{x}_i) \geq D(1 - \xi_{ic}), \quad (5.13)$$

$$i = 1, \dots, n, \quad c = 1, \dots, C, \quad c \neq y_i \quad (5.14)$$

$$\xi_{ic} \geq 0, \quad \sum_i \sum_{c \neq y_i} \xi_{ic} \leq \lambda \quad (5.15)$$

$$\sum_{c=1}^C \|f_c\|_{\mathcal{H}_K}^2 = 1. \quad (5.16)$$

Recently, Lee et al. (2002) pointed out that (5.12)–(5.16) is not always Bayes optimal. They proposed an algorithm that implements the Bayes classification rule and estimates  $\arg\max_c P(Y = c | X = \mathbf{x})$  directly.

Here we propose a theorem that illustrates the connection between the multiclass KLR and one version of the multiclass SVM.

**Theorem 2.** *Suppose the training data are pairwise separable, that is,  $\exists f_c(\mathbf{x})$ , s.t.  $f_{y_i}(\mathbf{x}_i) - f_c(\mathbf{x}_i) > 0, \forall i, \forall c \neq y_i$ . Then as  $\lambda \rightarrow 0$ , the classification boundary given by the multi-class KLR (5.6) will converge to that given by the multiclass SVM (5.12)–(5.16), if the latter is unique.*

The proof of the theorem is very similar to that of Theorem 1, we omit it here. Note that in the case of separable classes, (5.12)–(5.16) is guaranteed to be Bayes optimal.



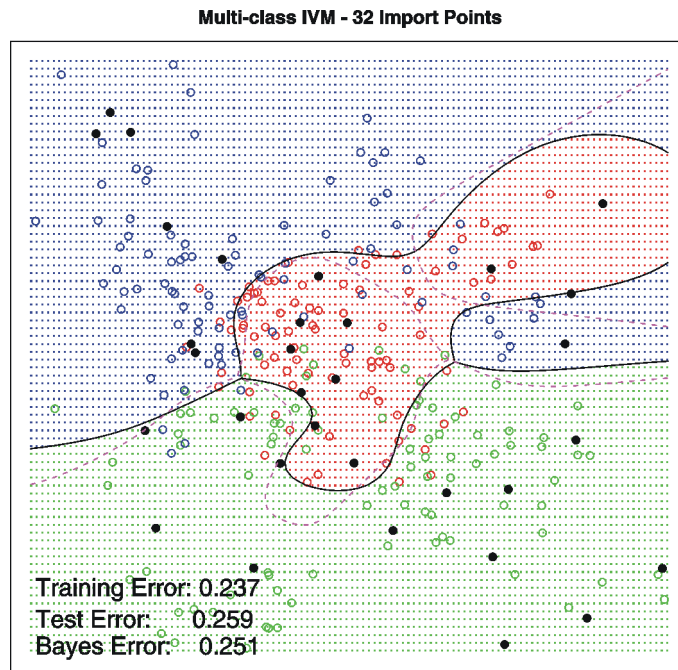


Figure 6. Radial kernel is used.  $C = 3$ ,  $n = 300$ ,  $\lambda = .368$ ,  $|S| = 32$ .

## 5.2 MULTICLASS IVM

The multiclass IVM procedure is similar to the two-class case (we omit the details), and the computational cost is  $O(Cn^2m^2)$ . Figure 6 is a simulation of the multiclass IVM. The data in each class are generated from a mixture of Gaussians (Hastie et al. 2001).

## 6. DISCUSSION

The support vector machine has been very successful for two-class classification and gained a lot of attention in the machine learning society in the past 10 years. Many articles have been published to explain why the support vector machine performs well in two-class classification. Most of this literature concentrates on the concept of margin. Various misclassification error bounds have been derived based on the margin (e.g., Vapnik 1995, 1998; Bartlett and Shawe-Taylor 1999; Shawe-Taylor and Cristianini 1999).

However, our view in this article is a little different from that based on the concept of margin. Several researchers have noted the relationship between the support vector machine and regularized function estimation in RKHS (e.g., Evgeniou et al. 1999; Wahba 1999; Hastie et al. 2001). The regularized function estimation problem contains two parts: a loss function and a regularization term. What is special with the support vector machine is the loss function, that is, the hinge loss. The margin maximizing property of the support vector machine derives from the hinge loss function. Lin (2002) pointed out that the hinge

loss is Bayes consistent, that is, the population minimizer of the loss function agrees with the Bayes rule in terms of classification. We believe this is a big step in explaining the success of the support vector machine, because it implies the support vector machine is trying to implement the Bayes rule. However, this is only part of the story; we believe that the regularization term has also played an important role in the support vector machine's success.

Regularization is an essential component in modern data analysis, in particular when the number of basis functions is large, possibly larger than the number of observations, and nonregularized fitting is guaranteed to give badly overfitted models. The enlarged feature space in the support vector machine allows the fitted model to be flexible, and the regularization term controls the complexity of the fitted model; the  $L_2$  nature of the regularization term in the support vector machine allows the fitted model to have a finite representation, even if the fitted model is in an infinite dimensional space. Hence we propose that by replacing the hinge loss of the support vector machine with the negative binomial log-likelihood, which is also Bayes consistent, we should be able to get a fitted model that performs similarly to the support vector machine. The resulting kernel logistic regression is something our statisticians are very familiar with (e.g., Green and Yandell 1985; Hastie and Tibshirani 1990; Wahba et al. 1995). We all understand why it can work well. The same reasoning could be applied to the support vector machine. The import vector machine algorithm is just a way to compress the data and reduce the computational cost.

Kernel logistic regression is not the only model that performs similarly to the support vector machine, replacing the hinge loss with any sensible loss function will give similar result, for example, the exponential loss function of boosting (Freund and Schapire 1997), the squared error loss (e.g., Buhlmann and Yu 2001; Zhang and Oles 2001; Mannor, Meir, and Zhang 2002) and the  $1/yf$  loss for distance weighted discrimination (Marron and Todd 2002). These loss functions are all Bayes consistent. The negative binomial log-likelihood and the exponential loss are also margin-maximizing loss functions; but the squared error loss and the  $1/yf$  loss are not.

To summarize, margin maximization is by nature a nonregularized objective, and solving it in high-dimensional space is likely to lead to overfitting and bad prediction performance. This has been observed in practice by many researchers, in particular Breiman (1999). Our conclusion is that margin maximization is not the only key to the support vector machine's success; the regularization term has played an important role.

## APPENDIX

### A.1 PROOF OF THEOREM 1

For the purpose of simple notation, we omit the constant  $\beta_0$  in the proof. We define

$$G(\beta) \equiv \sum_{i=1}^n \ln \left( 1 + e^{-y_i \mathbf{h}(\mathbf{x}_i)^T \beta} \right).$$

**Lemma A.1.** Consider the optimization problem (2.7)–(2.9), let the solution be denoted by  $\hat{\beta}(s)$ . If the training data are separable, that is,  $\exists \beta$ , s.t.  $y_i \mathbf{h}(\mathbf{x}_i)^T \beta > 0$ ,  $\forall i$ , then  $y_i \mathbf{h}(\mathbf{x}_i)^T \hat{\beta}(s) > 0$ ,  $\forall i$  and  $\|\hat{\beta}(s)\|_2 = s$  for all  $s > s_0$ , where  $s_0$  is a fixed positive number. Hence,  $\left\| \frac{\hat{\beta}(s)}{s} \right\|_2 = 1$ .

**Proof:** Suppose  $\exists i^*$ , s.t.  $y_{i^*} \mathbf{h}(\mathbf{x}_{i^*})^T \beta \leq 0$ , then

$$G(\beta) \geq \ln \left( 1 + e^{-y_{i^*} \mathbf{h}(\mathbf{x}_{i^*})^T \beta} \right) \quad (\text{A.1})$$

$$\geq \ln 2. \quad (\text{A.2})$$

On the other hand, by the separability assumption, we know there exists  $\beta^*$ ,  $\|\beta^*\|_2 = 1$ , s.t.  $y_i \mathbf{h}(\mathbf{x}_i)^T \beta^* > 0$ ,  $\forall i$ . Then for  $s > s_0 = -\ln(2^{1/n} - 1) / \min_i (y_i \mathbf{h}(\mathbf{x}_i)^T \beta^*)$ , we have

$$G(s\beta^*) = \sum_{i=1}^n \ln \left( 1 + e^{-y_i \mathbf{h}(\mathbf{x}_i)^T \beta^* s} \right) \quad (\text{A.3})$$

$$< \sum_{i=1}^n \frac{\ln 2}{n} = \ln 2. \quad (\text{A.4})$$

Since  $G(\hat{\beta}(s)) \leq G(s\beta^*)$ , we have, for  $s > s_0$ ,  $y_i \mathbf{h}(\mathbf{x}_i)^T \hat{\beta}(s) > 0$ ,  $\forall i$ .

For  $s > s_0$ , if  $\|\hat{\beta}(s)\|_2 < s$ , we consider to scale up  $\hat{\beta}(s)$  by letting

$$\hat{\beta}'(s) = \frac{\hat{\beta}(s)}{\|\hat{\beta}(s)\|_2} s.$$

Then  $\|\hat{\beta}'(s)\|_2 = s$ , and

$$G(\hat{\beta}'(s)) = \sum_{i=1}^n \ln \left( 1 + e^{-y_i \mathbf{h}(\mathbf{x}_i)^T \hat{\beta}'(s)} \right) \quad (\text{A.5})$$

$$< \sum_{i=1}^n \ln \left( 1 + e^{-y_i \mathbf{h}(\mathbf{x}_i)^T \hat{\beta}(s)} \right) \quad (\text{A.6})$$

$$= G(\hat{\beta}(s)), \quad (\text{A.7})$$

which is a contradiction. Hence  $\|\hat{\beta}(s)\|_2 = s$ .  $\square$

Now we consider two separating candidates  $\beta_1$  and  $\beta_2$ , such that  $\|\beta_1\|_2 = \|\beta_2\|_2 = 1$ . Assume that  $\beta_1$  separates better, that is:

$$d_1 := \min_i y_i \mathbf{h}(\mathbf{x}_i)^T \beta_1 > d_2 := \min_i y_i \mathbf{h}(\mathbf{x}_i)^T \beta_2 > 0.$$

**Lemma A.2.** There exists some  $s_0 = S(d_1, d_2)$  such that  $\forall s > s_0$ ,  $s\beta_1$  incurs smaller loss than  $s\beta_2$ , in other words:

$$G(s\beta_1) < G(s\beta_2).$$

**Proof:** Let

$$s_0 = S(d_1, d_2) = \frac{\ln n + \ln 2}{d_1 - d_2},$$

then  $\forall s > s_0$ , we have

$$\sum_{i=1}^n \ln \left( 1 + e^{-y_i \mathbf{h}(\mathbf{x}_i)^T \beta_1 s} \right) \leq n \ln (1 + e^{-s \cdot d_1}) \quad (\text{A.8})$$

$$\leq n \exp(-s \cdot d_1) \quad (\text{A.9})$$

$$< \frac{1}{2} \exp(-s \cdot d_2) \quad (\text{A.10})$$

$$\leq \ln (1 + e^{-s \cdot d_2}) \quad (\text{A.11})$$

$$\leq \sum_{i=1}^n \ln \left( 1 + e^{-y_i \mathbf{h}(\mathbf{x}_i)^T \beta_2 s} \right). \quad (\text{A.12})$$

The first and the last inequalities imply

$$G(s\beta_1) < G(s\beta_2).$$

□

Given these two lemmas, we can now prove that any convergence point of  $\frac{\hat{\beta}(s)}{s}$  must be a margin maximizing separator. Assume  $\beta^*$  is a convergence point of  $\frac{\hat{\beta}(s)}{s}$ . Denote  $d^* := \min_i y_i \mathbf{h}(\mathbf{x}_i)^T \beta^*$ . Because the training data are separable, clearly  $d^* > 0$  (since otherwise  $G(s\beta^*)$  does not even converge to 0 as  $s \rightarrow \infty$ ).

Now, assume some  $\tilde{\beta}$  with  $\|\tilde{\beta}\|_2 = 1$  has bigger minimal margin  $\tilde{d} > d^*$ . By continuity of the minimal margin in  $\beta$ , there exists some open neighborhood of  $\beta^*$ :

$$N_{\beta^*} = \{\beta : \|\beta - \beta^*\|_2 < \delta\},$$

and an  $\epsilon > 0$ , such that:

$$\min_i y_i \mathbf{h}(\mathbf{x}_i)^T \beta < \tilde{d} - \epsilon, \quad \forall \beta \in N_{\beta^*}.$$

Now by Lemma A.2 we get that there exists some  $s_0 = S(\tilde{d}, \tilde{d} - \epsilon)$  such that  $s\tilde{\beta}$  incurs smaller loss than  $s\beta$  for any  $s > s_0$ ,  $\beta \in N_{\beta^*}$ . Therefore  $\beta^*$  cannot be a convergence point of  $\frac{\hat{\beta}(s)}{s}$ .

We conclude that any convergence point of the sequence  $\frac{\hat{\beta}(s)}{s}$  must be a margin maximizing separator. If the margin maximizing separator is unique then it is the only possible convergence point, and therefore:

$$\lim_{s \rightarrow \infty} \frac{\hat{\beta}(s)}{s} = \arg \max_{\|\beta\|_2=1} \min_i y_i \mathbf{h}(\mathbf{x}_i)^T \beta.$$

In the case that the margin maximizing separating hyperplane is not unique, this conclusion can easily be generalized to characterize a unique solution by defining tie-breakers: if the minimal margin is the same, then the second minimal margin determines which model separates better, and so on. Only in the case that the whole order statistics of the margins is common to many solutions can there really be more than one convergence point for  $\frac{\hat{\beta}(s)}{s}$ .

## ACKNOWLEDGMENTS

We thank Saharon Rosset, Dylan Small, John Storey, Rob Tibshirani, and Jingming Yan for their helpful comments. We are also grateful for the three reviewers and one associate editor for their comments that helped improve the article. Ji Zhu was partially supported by the Stanford Graduate Fellowship. Trevor Hastie is partially supported by grant DMS-9803645 from the National Science Foundation, and grant RO1-CA-72028-01 from the National Institutes of Health.

[Received February 2002. Revised April 2004.]

## REFERENCES

- Bartlett, P., and Shawe-Taylor, J. (1999), "Generalization Performance of Support Vector Machines and Other Pattern Classifiers," in *Advances in Kernel Methods—Support Vector Learning*, eds. B. Schölkopf, C. Burges, and A. Smola, Cambridge, MA: MIT Press, 43–54.
- Bredensteiner, E., and Bennett, K. (1999), "Multicategory Classification by Support Vector Machines," *Computational Optimization and Applications*, 12, 35–46.
- Breiman, L. (1999), "Prediction Games and Arcing Algorithms," *Neural Computation*, 7, 1493–1517.
- Buhlmann, P., and Yu, B. (2001), "Boosting With the  $L_2$  Loss: Regression and Classification," *Journal of American Statistical Association*, 98, 324–340.
- Burges, C. (1998), "A Tutorial on Support Vector Machines for Pattern Recognition," in *Data Mining and Knowledge Discovery*, 2, Boston: Kluwer.
- Evgeniou, T., Pontil, M., and Poggio, T. (1999), "Regularization Networks and Support Vector Machines," in *Advances in Large Margin Classifiers*, eds. A. J. Smola, P. Bartlett, B. Schölkopf, and C. Schuurmans, Cambridge, MA: MIT Press, 171–204.
- Freund, Y., and Schapire, R. (1997), "A Decision-Theoretic Generalization of On-line Learning and an Application to Boosting," *Journal of Computer and System Sciences*, 55, 119–139.
- Green, P., and Yandell, B. (1985), "Semi-parametric Generalized Linear Models," *Proceedings 2nd International GLIM Conference*, Lancaster, Lecture notes in Statistics No. 32, New York: Springer, pp. 44–55.
- Hastie, T., and Tibshirani, R. (1990), *Generalized Additive Models*, New York: Chapman and Hall.
- Hastie, T., Tibshirani, R., and Friedman, J. (2001), *The Elements of Statistical Learning*. New York: Springer.
- Kaufman, L. (1999), "Solving the Quadratic Programming Problem Arising in Support Vector Classification," in *Advances in Kernel Methods—Support Vector Learning*, eds. B. Schölkopf, C. Burges, and A. Smola, Cambridge, MA: MIT Press, 147–168.
- Keerthi, S., Duan, K., Shevade, S., and Poo, A. (2002), "A Fast Dual Algorithm for Kernel Logistic Regression," *International Conference on Machine Learning*, 19.
- Kimeldorf, G., and Wahba, G. (1971), "Some Results on Tchebycheffian Spline Functions," *Journal of Mathematical Analysis and Applications*, 33, 82–95.
- Lee, Y., Lin, Y., and Wahba, G. (2004), "Multicategory Support Vector Machines, Theory, and Application to the Classification of Microarray Data and Satellite Radiance Data," *Journal of American Statistical Association*, 99, 67–81.
- Lin, Y. (2002), "Support Vector Machines and the Bayes Rule in Classification," *Data Mining and Knowledge Discovery*, 6, 259–275.
- Lin, X., Wahba, G., Xiang, D., Gao, F., Klein, R., and Klein B. (2000), "Smoothing Spline ANOVA Models for Large Data Sets with Bernoulli Observations and the Randomized GACV," *The Annals of Statistics*, 28, 1570–1600.
- Luo, Z., and Wahba, G. (1997), "Hybrid Adaptive Splines," *Journal of American Statistical Association*, 92, 107–116.

- Mannor, S., Meir, R., and Zhang, T. (2002), "The Consistency of Greedy Algorithms for Classification," in *Proceedings of the European Conference on Computational Learning Theory*, pp. 319–333.
- Marron, J., and Todd, M. (2002), "Distance Weighted Discrimination," Technical Report No. 1339, School of Operations Research and Industrial Engineering, Cornell University.
- Platt, J. (1999), "Fast Training of Support Vector Machines using Sequential Minimal Optimization," in *Advances in Kernel Methods—Support Vector Learning*, eds. B. Schölkopf, C. Burges, and A. Smola, Cambridge, MA: MIT Press.
- Rätsch, G., Onoda, T., and Müller, K. (2001), "Soft Margins for Adaboost," *Machine Learning*, 42, 287–320.
- Rosset, S., Zhu, J., and Hastie, T. (2004), "Margin Maximizing Loss Functions," *Neural Information Processing Systems*, 16.
- Shawe-Taylor, J., and Cristianini, N. (1999), "Margin Distribution Bounds on Generalization, in *Proceedings of the European Conference on Computational Learning Theory*, 263–273.
- Smola, A., and Schölkopf, B. (2000), "Sparse Greedy Matrix Approximation for Machine Learning," in *Proceedings of the Seventeenth International Conference on Machine Learning*, San Francisco, CA: Morgan Kaufmann, 911–918.
- Vapnik, V. (1995), *The Nature of Statistical Learning Theory*, Berlin: Springer Verlag.
- (1998), *Statistical Learning Theory*, New York: Wiley.
- Wahba, G. (1999), "Support Vector Machines, Reproducing Kernel Hilbert Spaces and the Randomized GACV," in *Advances in Kernel Methods—Support Vector Learning*, eds. B. Schölkopf, C. Burges, and A. Smola, Cambridge, MA: MIT Press, 69–88.
- Wahba, G., Gu, C., Wang, Y., and Chappell, R. (1995), "Soft Classification, a.k.a. Risk Estimation, via Penalized Log-Likelihood and Smoothing Spline Analysis of Variance," in *The Mathematics of Generalization*, ed. D. H. Wolpert, Santa Fe Institute Studies in the Sciences of Complexity, Reading, MA: Addison-Wesley, 329–360.
- Weston, J., and Watkins, C. (1999), "Support Vector Machines for Multiclass Pattern Recognition," in *Proceedings of the Seventh European Symposium on Artificial Neural Networks*.
- Williams, C., and Seeger, M. (2001), "Using the Nystrom Method to Speed up Kernel Machines," in *Advances in Neural Information Processing Systems 13*, eds. T. K. Leen, T. G. Diettrich, and V. Tresp, Cambridge, MA: MIT Press, pp. 682–688.
- Zhang, T., and Oles, F. (2001), "Text Categorization Based on Regularized Linear Classification Methods," *Information Retrieval*, 4, 5–31.
- Zhu, J., and Hastie, T. (2004), "Classification of Gene Microarrays by Penalized Logistic Regression," *Biostatistics*, 25, 427–444.