

Support Vector Machines, Kernel Logistic Regression and Boosting

Ji Zhu and Trevor Hastie

Department of Statistics
Stanford University
Stanford, CA 94305
`{jzhu,hastie}@stat.stanford.edu`

Abstract. The support vector machine is known for its excellent performance in binary classification, i.e., the response $y \in \{-1, 1\}$, but its appropriate extension to the multi-class case is still an on-going research issue. Another weakness of the SVM is that it only estimates $\text{sign}[p(x) - 1/2]$, while the probability $p(x)$ is often of interest itself, where $p(x) = P(Y = 1|X = x)$ is the conditional probability of a point being in class 1 given $X = x$. We propose a new approach for classification, called the import vector machine, which is built on kernel logistic regression (KLR). We show on some examples that the IVM performs as well as the SVM in binary classification. The IVM can naturally be generalized to the multi-class case. Furthermore, the IVM provides an estimate of the underlying class probabilities. 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 can give the IVM a computational advantage over the SVM, especially when the size of the training data set is large. We illustrate these techniques on some examples, and make connections with boosting, another popular machine-learning method for classification.

Keywords: classification, kernel methods, logistic regression, multi-class 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 $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$, where the output y_i is qualitative and assumes values in a finite set \mathcal{C} . We wish to find a classification rule from the training data, so that when given a new input 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 binary classification, i.e. $y \in \{0, 1\}$, but its appropriate extension to the multi-class case is still an on-going research issue. Another weakness of the SVM is that it only estimates

$\text{sign}[p(x) - 1/2]$, while the probability $p(x)$ is often of interest itself, where $p(x) = P(Y = 1|X = x)$ is the conditional probability of a point being in class 1 given $X = x$. In this paper, 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 binary classification, but also can naturally be generalized to the multi-class case. Furthermore, the IVM provides an estimate of the probability $p(x)$. 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 $O(N^3)$, while the computational cost of the IVM is $O(N^2q^2)$, where q is the number of import points. Since q does not tend to increase as N increases, the IVM can be faster than the SVM, especially for large training data sets. 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 binary classification and compare it with kernel logistic regression (KLR). In section (3), we propose our IVM algorithm. In section (4), we show some simulation results. In section (5), we generalize the IVM to the multi-class case.

2 Support Vector Machines and Kernel Logistic Regression

The standard SVM produces a non-linear 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 , 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 Evgeniou, Pontil, and Poggio (1999), Hastie, Tibshirani, and Friedman (2001) and Wahba, Lin, and Zhang (2000). Fitting an SVM is equivalent to minimizing:

$$\frac{1}{N} \sum_{i=1}^N (1 - y_i f(x_i))_+ + \lambda \|f\|_{\mathcal{H}_K}^2. \quad (1)$$

with $f = b + h$, $h \in \mathcal{H}_K$, $b \in \mathcal{R}$. \mathcal{H}_K is the RKHS generated by the kernel K . The classification rule is given by $\text{sign}[f]$.

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

$$f(x) = b + \sum_{i=1}^N a_i K(x, x_i). \quad (2)$$

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 (1). This seems to be an 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 non-zero a_i 's. The corresponding x_i 's are called support points.

Notice that (1) has the form *loss* + *penalty*. The loss function $(1 - yf)_+$ is plotted in Figure 1, along with several traditional loss functions. As we can see, the negative log-likelihood (NLL) of the binomial distribution has a similar shape to that of the SVM. If we replace $(1 - yf)_+$ in (1) with $\ln(1 + e^{-yf})$, the NLL of the binomial distribution, the problem becomes a KLR problem. We expect that the fitted function performs similarly to the SVM for binary classification.

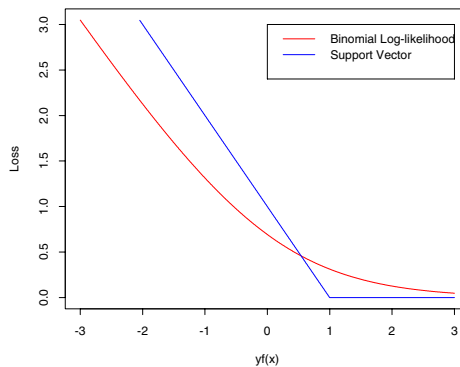


Fig. 1. The binomial log-likelihood and hinge loss function, $y \in \{-1, 1\}$

There are two immediate advantages of making such a replacement: (a) Besides giving a classification rule, the KLR also offers a natural estimate of the probability $p(x) = e^f / (1 + e^f)$, while the SVM only estimates $\text{sign}[p(x) - 1/2]$; (b) The KLR can naturally be generalized to the multi-class case through kernel multi-logit regression, whereas this is not the case for the SVM. However, because the 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) are non-zero.

KLR has been studied by many researchers; see Wahba, Gu, Wang, and Chappell (1995) and references there; see also Green and Silverman (1994) and Hastie and Tibshirani (1990).

The computational cost of the KLR is $O(N^3)$; to save the computational cost, the IVM algorithm will find a sub-model to approximate the full model (2) given by the KLR. The sub-model has the form:

$$f(x) = b + \sum_{x_i \in S} a_i K(x, x_i) \quad (3)$$

where \mathcal{S} is a subset of the training data $\{x_1, x_2, \dots, x_N\}$, and the data in \mathcal{S} are called import points. The advantage of this sub-model is that the computational cost is reduced, especially for large training data sets, while not jeopardizing the performance in classification.

Several other researchers have investigated techniques in selecting the subset \mathcal{S} . Lin, Wahba, Xiang, Gao, Klein, and B. (2001) divide the training data into several clusters, then randomly select a representative from each cluster to make up \mathcal{S} . Smola and Schölkopf (2000) develop a greedy technique to sequentially select q columns of the kernel matrix $[K(x_i, x_j)]_{N \times N}$, such that the span of these q columns approximates the span of $[K(x_i, x_j)]_{N \times N}$ well in the Frobenius norm. Williams and Seeger (2001) propose randomly selecting q points of the training data, then using the Nystrom method to approximate the eigen-decomposition of the kernel matrix $[K(x_i, x_j)]_{N \times N}$, 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 x_i). The IVM algorithm uses both the output y_i and the input x_i to select the subset \mathcal{S} , in such a way that the resulting fit approximates the full model well.

3 Import Vector Machine

Following the tradition of logistic regression, we let $y_i \in \{0, 1\}$ for the rest of this paper. For notational simplicity, the constant term in the fitted function is ignored.

In the KLR, we want to minimize:

$$H = - \sum_{i=1}^N [y_i f(x_i) - \ln(1 + \exp(f(x_i)))] + \frac{\lambda}{2} \|f\|_{\mathcal{H}_K}^2$$

From (2), it can be shown that this is equivalent to the finite dimensional form:

$$H = -\mathbf{y}^T (K_a \mathbf{a}) + \mathbf{1}^T \ln(1 + \exp(K_a \mathbf{a})) + \frac{\lambda}{2} \mathbf{a}^T K_q \mathbf{a} \quad (4)$$

where $\mathbf{a} = (a_1, \dots, a_N)^T$; the regression matrix $K_a = [K(x_i, x_j)]_{N \times N}$; and the regularization matrix $K_q = K_a$.

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:

$$\mathbf{a}^{(k)} = (K_a^T W K_a + \lambda K_q)^{-1} K_a^T W \mathbf{z} \quad (5)$$

where $\mathbf{a}^{(k)}$ is the value of \mathbf{a} in the k th step, $\mathbf{z} = (K_a \mathbf{a}^{(k-1)} + W^{-1}(\mathbf{y} - \mathbf{p}))$. The weight matrix is $W = \text{diag}[p(x_i)(1 - p(x_i))]_{N \times N}$.

As mentioned in section 2, we want to find a subset \mathcal{S} of $\{x_1, x_2, \dots, x_N\}$, such that the sub-model (3) is a good approximation of the full model (2). Since it is impossible to search for every subset \mathcal{S} , we use the following greedy forward strategy:

3.1 Basic Algorithm

(B1) Let $\mathcal{S} = \emptyset$, $\mathcal{R} = \{x_1, x_2, \dots, x_N\}$, $k = 1$.

(B2) For each $x_l \in \mathcal{R}$, let

$$f_l(x) = \sum_{x_j \in \mathcal{S} \cup \{x_l\}} a_j K(x, x_j)$$

Find \mathbf{a} to minimize

$$\begin{aligned} H(x_l) &= - \sum_{i=1}^N [y_i f_l(x_i) - \ln(1 + \exp(f_l(x_i)))] + \frac{\lambda}{2} \|f_l(x)\|_{\mathcal{H}_K}^2 \\ &= -\mathbf{y}^T (K_a^l \mathbf{a}^l) + \mathbf{1}^T \ln(1 + \exp(K_a^l \mathbf{a}^l)) + \frac{\lambda}{2} \mathbf{a}^{lT} K_q^l \mathbf{a}^l \end{aligned} \quad (6)$$

where the regression matrix $K_a^l = [K(x_i, x_j)]_{N \times (q+1)}$, $x_i \in \{x_1, x_2, \dots, x_N\}$, $x_j \in \mathcal{S} \cup \{x_l\}$; the regularization matrix $K_q^l = [K(x_j, x_l)]_{(q+1) \times (q+1)}$, $x_j, x_l \in \mathcal{S} \cup \{x_l\}$; $q = |\mathcal{S}|$.

(B3) Let

$$x_{l^*} = \operatorname{argmin}_{x_l \in \mathcal{R}} H(x_l).$$

Let $\mathcal{S} = \mathcal{S} \cup \{x_{l^*}\}$, $\mathcal{R} = \mathcal{R} \setminus \{x_{l^*}\}$, $H_k = H(x_{l^*})$, $k = k + 1$.

(B4) Repeat steps (B2) and (B3) until H_k converges.

We call the points in \mathcal{S} import points.

3.2 Revised Algorithm

The above algorithm is computationally feasible, but in step (B2) we need to use the Newton-Raphson method to find \mathbf{a} iteratively. When the number of import points q 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. To get a good approximation, we take advantage of the fitted result from the current “optimal” \mathcal{S} , i.e., the sub-model when $|\mathcal{S}| = q$, and use it as the initial value. 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 (5) to be computed in $O(Nq)$ time.

Hence, we have the revised step (B2) for the basic algorithm:

(B2*) For each $x_l \in \mathcal{R}$, correspondingly augment K_a with a column, and K_q with a column and a row. Use the updating formula to find \mathbf{a} in (5). Compute (6).

3.3 Stopping Rule for Adding Point to \mathcal{S}

In step (B4) 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 (B4). At each step k , we compare H_k with H_{k-r} , where r is a pre-chosen small integer, for example $r = 1$. If the ratio $\frac{|H_k - H_{k-r}|}{|H_k|}$ is less than some pre-chosen small number α , for example, $\alpha = 0.001$, we stop adding new import points to \mathcal{S} .

3.4 Choosing the Regularization Parameter λ

So far, we have assumed that the regularization parameter λ is fixed. In practice, we also need to choose an "optimal" λ . 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 λ . To reduce the computation, we take advantage of the fact that the regularized NLL converges faster for a larger λ . Thus, instead of running the entire revised algorithm for each λ , we propose the following procedure, which combines both adding import points to \mathcal{S} and choosing the optimal λ :

- (C1) Start with a large regularization parameter λ .
- (C2) Let $\mathcal{S} = \emptyset$, $\mathcal{R} = \{x_1, x_2, \dots, x_N\}$, $k = 1$.
- (C3) Run steps (B2*), (B3) and (B4) of the revised algorithm, until the stopping criterion is satisfied at $\mathcal{S} = \{x_{i1}, \dots, x_{iq_k}\}$. Along the way, also compute the misclassification error on the tuning set.
- (C4) Decrease λ to a smaller value.
- (C5) Repeat steps (C3) and (C4), starting with $\mathcal{S} = \{x_{i1}, \dots, x_{iq_k}\}$.

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

4 Simulation

In this section, we use a simulation to illustrate the IVM method. The data in each class are generated from a mixture of Gaussians (Hastie, Tibshirani, and Friedman 2001). The simulation results are shown in Figure 2.

4.1 Remarks

The support points of the SVM are those which are close to the classification boundary or misclassified and usually have large weights $[p(x)(1 - p(x))]$. 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 only concerned with the classification $\text{sign}[p(x) - 1/2]$, while the IVM also focuses on the unknown probability $p(x)$.

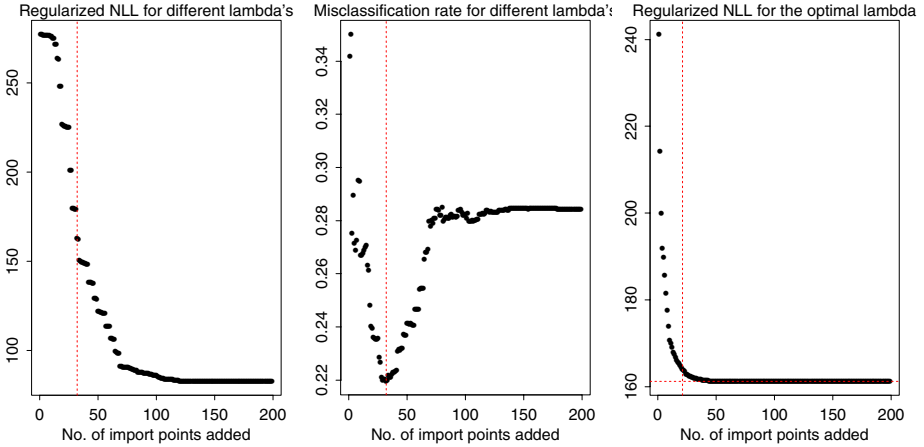


Fig. 2. Radial kernel is used. $N = 200$. The left and middle panels illustrate how to choose the optimal λ . $r = 1$, $\alpha = 0.001$, λ decreases from e^{10} to e^{-10} . The minimum misclassification rate 0.219 is found to correspond to $\lambda = 0.135$. The right panel is for the optimal $\lambda = 0.135$. The stopping criterion is satisfied when $|\mathcal{S}| = 21$.

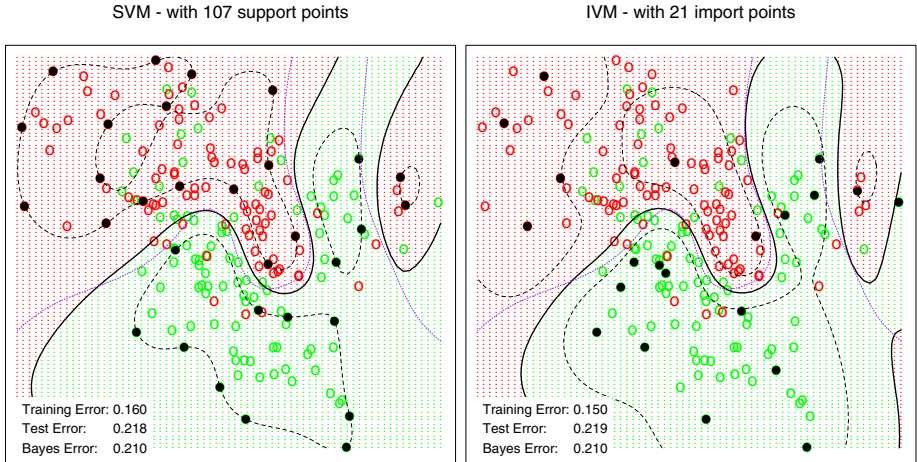


Fig. 3. The solid black lines are the classification boundaries; the dashed purple lines are the Bayes rule boundaries. For the SVM, the dashed black lines are the edges of the margin. For the IVM, the dashed black lines are the $p(x) = 0.25$ and 0.75 lines, and the black points are the import points.

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 $p(x)$. Figure 3 shows a comparison of the SVM and the IVM. The total computational cost of the SVM is $O(N^3)$, while the computational cost of the IVM method is $O(N^2q^2)$, where q is the number

of import points. Since q does not tend to increase as N increases, the computational cost of the IVM can be smaller than that of the SVM, especially for large training data sets.

5 Multi-class Case

In this section, we briefly describe a generalization of the IVM to multi-class classification. Suppose there are $M + 1$ classes. We can write the response as an M -vector \mathbf{y} , with each component being either 0 or 1, indicating which class the observation is in. Therefore $y_k = 1, y_j = 0, j \neq k, j \leq M$ indicates the response is in the k th class, and $y_j = 0, j \leq M$ indicates the response is in the $M + 1$ th class. Using the $M + 1$ th class as the basis, the multi-logit can be written as $f_1 = \ln(p_1/p_{M+1}), \dots, f_M = \ln(p_M/p_{M+1}), f_{M+1} = 0$. Hence the Bayes classification rule is given by:

$$c = \operatorname{argmax}_{k \in \{1, 2, \dots, M+1\}} f_k$$

We use i to index the observations, j to index the classes, i.e. $i = 1, \dots, N, j = 1, \dots, M$. Then the regularized negative log-likelihood is

$$H = - \sum_{i=1}^N [\mathbf{y}_i^T \mathbf{f}(x_i) - \ln(1 + e^{f_1(x_i)} + \dots + e^{f_M(x_i)})] + \frac{\lambda}{2} \|\mathbf{f}\|_{\mathcal{H}_K}^2 \quad (7)$$

where $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iM})^T, \mathbf{f}(x_i) = (f_1(x_i), f_2(x_i), \dots, f_M(x_i))^T$, and

$$\|\mathbf{f}\|_{\mathcal{H}_K}^2 = \sum_{j=1}^M \|f_j\|_{\mathcal{H}_K}^2$$

Using the representer theorem (Kimeldorf and Wahba 1971), the j th element of $\mathbf{f}(x), f_j(x)$, which minimizes H has the form

$$f_j(x) = \sum_{i=1}^N a_{ij} K(x, x_i). \quad (8)$$

Hence, (7) becomes

$$H = - \sum_{i=1}^N [\mathbf{y}_i^T (K_a(i, \cdot) A)^T - \ln(1 + \mathbf{1}^T e^{(K_a(i, \cdot) A)^T})] + \frac{\lambda}{2} \sum_{j=1}^M \mathbf{a}_j^T K_q \mathbf{a}_j \quad (9)$$

where $A = (\mathbf{a}_1 \dots \mathbf{a}_M) = (a_{ij})$, K_a and K_q are defined in the same way as in the binary case; and $K_a(i, \cdot)$ is the i th row of K_a .

The multi-class IVM procedure is similar to the binary case, and the computational cost is $O(MN^2q^2)$. Figure 4 is a simulation of the multi-class IVM. The data in each class are generated from a mixture of Gaussians.

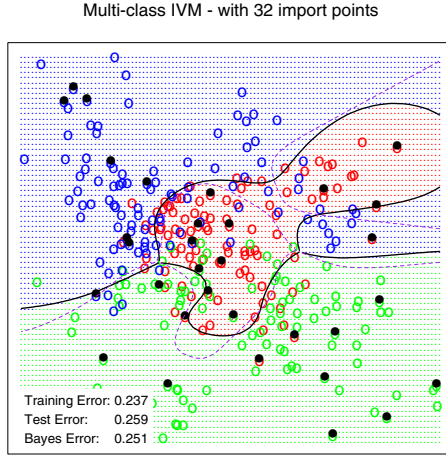


Fig. 4. Radial kernel is used. $M + 1 = 3$, $N = 300$, $\lambda = 0.368$, $|\mathcal{S}| = 32$.

6 Discussion

Although the intuitive motivation of the SVM is via separating hyperplanes, this intuition gets a murky when the classes overlap. In this case it is perhaps more intuitive to pose the problem as that of regularized function estimation with a loss function particularly suited to classification. Furthermore, the “magic” of the kernel finds its proper home when this function estimation takes place in reproducing kernel Hilbert spaces. We have argued in this paper that the binomial loss function offers several advantages over the hinge loss, in particular that it estimates class probabilities, and generalizes naturally to problems with more than two classes. Although KLR lacks the “support vector” property of SVMs, we propose the IVM, a simple and attractive compromise with performance similar to that of the SVM. The computational cost of the IVM is $O(N^2q^2)$ for the binary case and $O(MN^2q^2)$ for the multi-class case, where q is the number of import points.

The loss function representation of the SVM and KLR encourages a comparison with *boosting* (Freund and Schapire 1996; Hastie, Tibshirani, and Friedman 2001). Figure 5 is similar to Figure 1, and includes the exponential loss function that drives the boosting procedures. Boosting has been shown to fit a logistic regression model by a form of non-parametric gradient descent (Friedman, Hastie, and Tibshirani 2000) using this exponential loss function. It is also motivated as a means for generating a classifier that creates a wide margin between the classes (Schapire and Freund 1997). The comparisons in Figure 5 make it clear that all three methods are similar in this regard, although the exponential left tail in boosting suggests a non-robustness to clumps of observations far from their parent class.

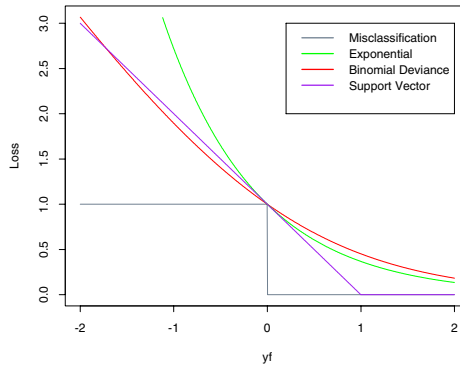


Fig. 5. The two loss functions from Figure 1, along with the exponential loss function implicit in boosting.

Acknowledgments We thank Dylan Small, John Storey, Rob Tibshirani, and Jingming Yan for their helpful comments. Ji Zhu is partially supported by the Stanford Graduate Fellowship. Trevor Hastie is partially supported by grant DMS-9803645 from the National Science Foundation, and grant ROI-CA-72028-01 from the National Institutes of Health.

References

- Evgeniou, T., M. Pontil, and T. Poggio (1999). Regularization networks and support vector machines. *Advances in Computational Mathematics (to appear)*.
- Freund, Y. and R. E. Schapire (1996). Experiments with a new boosting algorithm. In *Machine Learning: Proceedings of the Thirteenth International Conference*, pp. 148–156. Morgan Kaufman, San Francisco.
- Friedman, J., T. Hastie, and R. Tibshirani (2000). Additive logistic regression: a statistical view of boosting (with discussion). *Annals of Statistics* 28, 337–307.
- Green, P. and B. Silverman (1994). *Nonparametric Regression and Generalized Linear Models*. Chapman and Hall.
- Hastie, T. and R. Tibshirani (1990). *Generalized Additive Models*. Chapman and Hall.
- Hastie, T., R. Tibshirani, and J. Friedman (2001). *The Elements of Statistical Learning; Data mining, Inference and Prediction*. New York: Springer Verlag.
- Kimeldorf, G. and G. Wahba (1971). Some results on tchebycheffian spline functions. *J. Math. Anal. Applic.* 33, 82–95.
- Lin, X., G. Wahba, D. Xiang, F. Gao, R. Klein, and K. B. (2001). Smoothing spline anova models for large data sets with bernoulli observations

- and the randomized gacv. Technical Report 998, Department of Statistics, University of Wisconsin.
- Schapire, R. E. and Y. Freund (1997). Boosting the margin: A new explanation for the effectiveness of voting methods. *Proceedings of the Fourteenth International Conference on Machine Learning*.
- Smola, A. and B. Schölkopf (2000). Sparse greedy matrix approximation for machine learning. In *Proceedings of the Seventeenth International Conference on Machine Learning*. Morgan Kaufmann.
- Wahba, G., C. Gu, Y. Wang, and R. Chappell (1995). *The Mathematics of Generalization.*, Chapter Soft Classification, a.k.a. Risk Estimation, via Penalized Log Likelihood and Smoothing Spline Analysis of Variance. Santa Fe Institute Studies in the Sciences of Complexity. Addison-Wesley.
- Wahba, G., Y. Lin, and H. Zhang (2000). Gacv for support vector machines. In A. Smola, P. Bartlett, B. Schölkopf, and D. Schuurmans (Eds.), *Advances in Large Margin Classifiers*, Cambridge, MA, pp. 297–311. MIT Press.
- Williams, C. and M. Seeger (2001). Using the nystrom method to speed up kernel machines. In T. K. Leen, T. G. Diettrich, and V. Tresp (Eds.), *Advances in Neural Information Processing Systems*, Volume 13. MIT Press.