

Hinging Hyperplanes for Regression, Classification, and Function Approximation

By

Leo Breiman

**Technical Report No. 324
September 1991**

**Department of Statistics
University of California
Berkeley, California 94720**

Hinging Hyperplanes for Regression, Classification, and Function Approximation

Leo Breiman

Abstract

A hinge function $y = h(\mathbf{x})$ consists of two hyperplanes continuously joined together at a hinge. In regression (prediction), classification (pattern recognition), and noiseless function approximation, use of sums of hinge functions gives a powerful and efficient alternative to neural networks with compute times several orders of magnitude less than fitting neural networks with a comparable number of parameters. The core of the methodology is a simple and effective method for finding good hinges.

1). Introduction

In an M -dimensional space (x_1, \dots, x_M) , a *hinge function* $y = h(\mathbf{x})$ consists of two hyperplanes continuously joined together. Taking $x_0 \equiv 1$ and using \cdot to denote the inner product of two vectors, if the two hyperplanes are given by

$$y = \beta^+ \cdot \mathbf{x}, \quad y = \beta^- \cdot \mathbf{x},$$

they are joined together on $\{\mathbf{x} : (\beta^+ - \beta^-) \cdot \mathbf{x} = 0\}$ and we refer to $\Delta = \beta^+ - \beta^-$, or any multiple of Δ , as the *hinge* for the function. The explicit form of the hinge function is either $\max(\beta^+ \cdot \mathbf{x}, \beta^- \cdot \mathbf{x})$ or $\min(\beta^+ \cdot \mathbf{x}, \beta^- \cdot \mathbf{x})$.

Most of the recently introduced methods for nonlinear regression, classification, and function approximation use expansions into sums of basis functions. The basis functions used are “data selected” from a large parametric class of primitive functions. For instance, CART (Breiman et al. [1984]) uses an expansion into indicator functions of multidimensional rectangles with sides parallel to the coordinate axes. Neural network methods use sigmoid functions of linear functions as primitives. The MARS method (Friedman [1991]) uses products of univariate linear spline functions as its primitive class. In this work, the hinge functions form the primitive class. There are good reasons, as given below, for this approach.

Let P be any measure with compact support on $E^{(M)}$ and $f(\mathbf{x})$ any sufficiently smooth function. Then we show in section 3, using methods developed by Jones [1991] and Barron [1991] that there is a constant $C(f, P)$ such that for any K , there are hinge functions h_1, \dots, h_K with

$$\|f - \sum_{k=1}^K h_k\|^2 \leq \frac{C}{K}.$$

The property that makes the hinge functions effective is that there is a simple and computationally efficient method for locating hinges. Suppose we are told that $y = h(\mathbf{x})$ is a hinge function with unknown hinge Δ^* ; are given data (y_n, \mathbf{x}_n) , $n = 1, \dots, N$; and want to use this data to locate the hinge. One approach is this: for any specified candidate hinge Δ , do a least squares fit of a hinge function with hinge Δ to the data. Let $RSS(\Delta)$ be the residual sum-of-squares. Now search over Δ -space to find the minimizer of $RSS(\Delta)$. This procedure is compute intensive and global searches are not feasible unless M, N are small.

Here is an alternative: start with an arbitrary hinge $\Delta^{(0)}$. Using least squares, fit the data on the side $\Delta^{(0)} \cdot \mathbf{x}_n \geq 0$ to a hyperplane $y = \beta^+ \cdot \mathbf{x}$, and do a similar fit to the data such that $\Delta^{(0)} \cdot \mathbf{x}_n < 0$ getting $y = \beta^- \cdot \mathbf{x}$. Take the new estimate for the hinge as

$$\Delta^{(1)} = \beta^+ - \beta^-,$$

and repeat the procedure, getting a sequence $\Delta^{(k)}$ of estimates.

In section 2 we give evidence that generally $\Delta^{(k)} \rightarrow \Delta^*$ and that the convergence is rapid. The noisy case, $y_n = h(\mathbf{x}_n) + \epsilon_n$, is also looked at and the accuracy of $\lim \Delta^{(k)}$ as an estimate of Δ^* examined. Simulations support the theoretical accuracy results and show that the hinge finding algorithm is computationally efficient and accurate even for large M , and high noise.

In section 3 we state the theorem regarding the approximation of a smooth function by a sum of hinge functions. Then we look at the implementation that at the K^{th} stage adds a new hinge function by using the hinge finding algorithm on $y = f(\mathbf{x}) - \sum_{k=1}^{K-1} h_k(\mathbf{x})$, and then readjusts the sum $\sum_{k=1}^K h_k(\mathbf{x})$. Simulations in dimensions ranging from $M = 2$ to $M = 16$ show good approximation properties and verify the $1/K$ decrease.

The most important applications of hinge functions is to multivariate regression and classification. Section 4 discusses the use of hinge functions to produce a nonlinear prediction function given noisy data $\{(y_n, \mathbf{x}_n), n = 1, \dots, N\}$. Examples are given to show how sums of hinge functions can be used to construct accurate predictors. In particular, to cope with high dimensional spaces, a variable selection method in hinge finding is introduced. It took 2.8 minutes of cpu time to compute the prediction equation in a highly nonlinear 100-dimensional example with a training set of size 2000.

In section 5 we show how the hinge finding algorithm can be extended so that hinge functions can be usefully employed in classification (pattern recognition) problems. The approach is to formulate the J -class problem in terms of J linked regression equations, and then to locate, at each stage, the hinge that is optimum for a combination of the equations. This leads to a $J \times J$ eigenvalue problem at each iteration.

Examples are given, with both real and simulated data to illustrate the effectiveness of this method.

Our concluding remarks and summary appears in section 6, and the proofs of theorems are in Appendix I. Computational aspects of the hinge procedures are discussed in Appendix II. Appendix III gives some equations for computing cross-validated error measures in regression and classification.

On current computing equipment, the procedures are quite feasible for substantial sample sizes and dimensionality and we give cpu times involved for various examples in the text. The examples were run on a variety of machines. To provide uniform time benchmarks, the 16-dimensional example of section 3 was run on all machines including an IBM RS6000/540. The timings given in the text are scaled to the RS 6000, a 12.5 megaflop machine.

Another important application of hinge finding is in determining splits for the construction of classification and regression trees. This is the subject of a sequel paper.

For recent advances in constructing nonlinear prediction functions, see (in the statistical literature) Friedman and Stuetzle [1981], Breiman et. al. [1984], Breiman and Friedman [1985], Friedman [1991], Hastie and Tibshirani [1990], Breiman [1991], Whaba [1990]. In classification, some recent statistical publications are Breiman et. al [1984], Breiman and Ihaka [1987], Hastie and Tibshirani [1990]. In the engineering and computer science fields, where recent attention has been focussed on neural networks, see Lippmann [1989] for a review and references.

2). Finding the Hinge

2.0). Analytic Results

Assume that $y = h(x)$ is a hinge function with unknown hinge, and that there is a distribution of points $x \in E^{(M+1)}$ governed by a probability measure $P(dx)$. For a fixed vector $\Delta^{(0)}$, denote

$$\begin{aligned} S_+ &= \{x: \Delta^{(0)} \cdot x \geq 0\}, \quad S_- = \{x: \Delta^{(0)} \cdot x < 0\} \\ \Gamma_+(m, m') &= \int_{S_+} x_m x_{m'} dP, \quad \Gamma_-(m, m') = \int_{S_-} x_m x_{m'} dP \\ (xy)_+(m) &= \int_{S_+} x_m h(x) dP, \quad (xy)_-(m) = \int_{S_-} x_m h(x) dP. \end{aligned}$$

The least-squares coefficients of a hyperplane fitted to the y -values in S_+ are $\beta_+ = \Gamma_+^{-1}(xy)_+$, and those in the S_- fit are $\beta_- = \Gamma_-^{-1}(xy)_-$. The new hinge value $\Delta^{(1)}$ is $\beta_+ - \beta_-$ and the process is repeated starting from $\Delta^{(1)}$.

Since $y = h(x)$ is a hinge function, there is a hinge Δ^* such that for $S_+^* = \{x: \Delta^* \cdot x \geq 0\}$ and $S_-^* = \{x: \Delta^* \cdot x < 0\}$, if the fitted coefficients on the +, - sides are β_+^*, β_-^* , then $\Delta^* = \beta_+^* - \beta_-^*$. Thus

$$\begin{aligned} (xy)_+(m) &= \int_{S_+ \cap S_+^*} x_m (\beta_+^* \cdot x) dP + \int_{S_+ \cap S_-^*} x_m (\beta_-^* \cdot x) dP \\ (xy)_-(m) &= \int_{S_- \cap S_+^*} x_m (\beta_+^* \cdot x) dP + \int_{S_- \cap S_-^*} x_m (\beta_-^* \cdot x) dP, \end{aligned}$$

Denote

$$\Gamma_{++}(m, m') = \int_{S_+ \cap S_+^*} x_m x_{m'} dP$$

with the analogous definitions for $\Gamma_{+-}, \Gamma_{-+}, \Gamma_{--}$. Then, some manipulations give

$$\Delta^{(1)} = A\Delta^*$$

where

$$A = I - \Gamma_{++}^{-1} \Gamma_{+-} - \Gamma_{-+}^{-1} \Gamma_{--}. \quad (2.1)$$

Hinge Convergence Theorem (2.1). Denote by s any linear function of x_0, \dots, x_M . Assume that P satisfies

- i) There is a $c < \infty$ such that for every s with $Es^2 \leq 1$, $Es^4 \leq c$.
- ii) Let $p(\alpha) = \sup_{Es^2=1} P(|s| \leq \alpha)$. Then $p(\alpha) = O(\alpha)$ as $\alpha \rightarrow 0$.
- iii) If $S_+ = \{x; s \geq 0\}$ and $P(S_+) = c > 0$, the inf over s of the minimum eigenvalue of Γ_+ is $\lambda^*(c) > 0$.

Then, there is a $\delta > 0$ such that if $\|\Delta^{(0)} - \Delta^*\| \leq \delta$, then

$$\|\Delta^{(k)} - \Delta^*\| \rightarrow 0.$$

and convergence is exponentially fast.

The proof is in Appendix I. Both theoretical and simulation results suggest that convergence occurs even for $\Delta^{(0)}$ distant from Δ^* . For example;

Theorem 2.2. If $x_0 \equiv 0$ and (x_1, \dots, x_M) have a joint normal distribution, then $\Delta^{(k)} \rightarrow \Delta^*$ for any starting $\Delta^{(0)}$ such that $\Delta^{(0)} \cdot \Delta^* \neq 0$.

The proof is given in Appendix I.

2.1 Simulations

The simulations were done as follows: the $(\Delta^*(1), \dots, \Delta^*(M))$ were taken to be M i.i.d. $N(0, 1)$'s, and $\Delta^*(0)$ adjusted so that the proportion of data in $\{\Delta^* \cdot x \geq 0\}$ was uniform on $[.15, .85]$. Initial hinge values $\Delta^{(0)}$ were selected by taking

$(\Delta^{(0)}(1), \dots, \Delta^{(0)}(M))$ to be another M -vector of i.i.d. $N(0, 1)$'s and $\Delta^{(0)}(0)$ taken so that $\#\{\mathbf{x}_n: \Delta^{(0)}\mathbf{x}_n \geq 0\} = \text{int}(N/2)$. Starting from $\Delta^{(0)}$, the hinge finding algorithm was run to convergence giving Δ . The run was counted as a hit if $|\cos(\Delta, \Delta^*)| \geq .99$, otherwise a miss. The whole procedure, including the data generation, was repeated 100 times.

In the first run, $N = 250$, the dimension $M = 6$ and (x_1, \dots, x_M) were sampled from X_1, \dots, X_M with $X_m = Z_m - Z_{m-1}$, $\{Z_0, \dots, Z_M\}$ independent unit exponentials. To study the affect of dimensionality, the same structure was used in the second run but with $M = 12$. Since this example used long tailed distributions, in the next two runs we again set $X_m = Z_m - Z_{m-1}$ but took the $\{Z_m\}$ to be uniform on $[0, 1]$. The first run was at dimension 6 and the second at dimension 12.

There were no misses in the 400 repetitions. But the algorithm can be made to miss. For instance, at least 13 points are needed to define a hyperplane in 12 dimensions. We set $\Delta^*(0)$ in the 12-dimensional exponential case so that $\#\{\mathbf{x}_n; \Delta^*\mathbf{x}_n \geq 0\} = 20$. Now 20 points to define a corner hinge in 12 dimensions is sparse. In 100 trials there were 13 misses.

2.2 The Effect of Noise

Suppose that the data is of the form $\{y_n, \mathbf{x}_n\}$, $n = 1, \dots, N$ with $y_n = h(\mathbf{x}_n) + \epsilon_n$, $h(\mathbf{x}_n)$ a hinge function and $\{\epsilon_n\}$ noise. Then

$$\Delta^{(k+1)} = \Delta^* - B\Delta^* + \Gamma_+^{-1}(\epsilon x)_+ - \Gamma_-^{-1}(\epsilon x)_-$$

where B is $\Gamma_+^{-1}\Gamma_{+-} + \Gamma_-^{-1}\Gamma_{-+}$. The matrix B is a function $B(\Delta^{(k)}, \Delta^*)$. The proof of the hinge convergence theorem indicates that B is small order of the noise. That is, if the noise is small, and $\Delta^{(k)} \sim \Delta^* + Z$, where Z is a linear function of the noise, then $B \sim o(Z)$. Thus, for small noise, at convergence

$$\Delta \equiv \Delta^* + \Gamma_+^{-1}(\epsilon x)_+ - \Gamma_-^{-1}(\epsilon x)_-$$

From this,

$$E \|\Delta - \Delta^*\|^2 = \frac{\sigma^2}{N} (\text{Tr}(\Gamma_+^{-1}) + \text{Tr}(\Gamma_-^{-1})). \quad (2.3)$$

2.3 Simulations: Noisy Case

The $\{\mathbf{x}_n\}$, $n = 1, \dots, N$ data was generated as in the first example in section 2.1. Then y_n was set equal to $h(\mathbf{x}_n) + \epsilon_n$ with the $\{\epsilon_n\}$ i.i.d. $N(0, \sigma^2)$ noise. The hinge Δ^* was selected as in section 2.1.

The signal/noise ratio (s/n) is here defined by first fitting a least squares hyperplane to $\{h(\mathbf{x}_n)\}$. Then the signal σ_l is defined as the square root of the mean residual

sum-of-squares (MRSS). The noise is defined by its standard deviation σ_N . Then $s/n = \sigma_l/\sigma_N$. Another way to look at this is that the expected MRSS on fitting $h + \epsilon$ by the right hyperplane is $\sigma_l^2 + \sigma_N^2$. The expected MRSS on fitting by the right hinge function is σ_N^2 . The ratio of expected MRSS is

$$\frac{\text{MRSS (hinge)}}{\text{MRSS (linear)}} = \frac{1}{1 + (s/n)^2}.$$

The s/n ratios used were 1, 1/2, 1/3, 1/4. The output of each run consisted of the averages over the repetitions of $|\cos(\Delta, \Delta^*)|$ and of $\|\Delta - \Delta^*\|/D$ where

$$D = \sqrt{\frac{\hat{\sigma}^2}{N} (\text{Tr}(\Gamma_+^{-1}) + \text{Tr}(\Gamma_-^{-1}))},$$

with Γ_+ , Γ_- computing using Δ and

$$\hat{\sigma}^2 = \frac{1}{N - 2(M + 1)} \sum_n (y_n - \hat{f}(x_n))^2$$

where \hat{f} is the estimated hinge function. The discussion leading to (2.3) indicates that the expectation of $\|\Delta - \Delta^*\|/D$ should be close to one in low noise situations. The standard deviation of these two quantities over the repetitions in the run is also given.

Another quantity of interest is $\hat{\sigma}^2 / (\sum_1^N \epsilon_n^2 / N)$. This is the ratio of the MRSS using the estimated hinge function debiased by the factor $1/(N - 2(M + 1))$, to the MRSS using the known ridge function. Denote this ratio by RSSRATIO. We report on four runs of 100 repetitions each using $N = 250$, and the 6 dimensional exponential distribution specified previously.

Table 2.1
Effects of Noise

s/n	av(cos)	Sd(cos)	av($\ \Delta - \Delta^*\ /D$)	Sd($\ \Delta - \Delta^*\ /D$)	RSSRATIO
1	.99	.01	.92	.38	1.00
1/2	.93	.07	1.11	.58	1.00
1/3	.81	.20	1.31	.67	1.00
1/4	.68	.24	1.49	.71	.99

Its surprising how well the algorithm does even with high noise levels. When $s/n = .25$, only a 6% decrease in RSS is gotten by fitting the underlying hinge instead of a hyperplane. Yet the values of the RSSRATIO show that the minimum RSS hinge is consistently being found.

3). Noiseless Function Approximation

On $E^{(M)}$ a sum of hinge functions $\sum_1^K h_k(\mathbf{x})$ is a continuous piecewise linear function. Results of Barron [1991] on sigmoid function and Jones [1991] on sinusoids can be extended to show that sums of hinge functions are effective in function approximation. Let $\tilde{f}(\mathbf{w})$ be the Fourier transform of $f(\mathbf{x})$.

Theorem 3.1. *If the support of P is contained in the sphere of radius R , and if*

$$\int \|\mathbf{w}\|^2 |\tilde{f}(\mathbf{w})| d\mathbf{w} = c < \infty$$

then there are hinge functions h_1, \dots, h_K such that

$$\|f - \sum_1^K h_k\|^2 \leq \frac{(2R)^4 c^2}{K}.$$

The proof of this theorem is an extension of Barron's extension of the Jones result and is given in Appendix I. In the Jones and Barron results, a function ϕ on $E^{(1)}$ is specified, and the approximating sums are of the form $\sum_1^K \phi(\beta_k \cdot \mathbf{x})$. In neural network approximation, the sum which minimizes $\|f(\mathbf{x}) - \sum_1^K \phi(\beta_k \cdot \mathbf{x})\|$ is gotten by using local gradient searches to minimize over the $\{\beta_1, \dots, \beta_K\}$. Not only is this highly compute intensive even for K, M of moderate size, but only local minima are guaranteed.

Barron and Jones show that there is a "greedy" algorithm which at each step enters the next function $\phi(\beta_K \cdot \mathbf{x})$ by minimization over $M + 2$ parameters only, and still achieves the upper bound of the existence proofs.

A similar "greedy" result can be proved for hinge functions but the point may be moot. The hinge finding algorithm makes the optimization over the entire sum computationally efficient. Our simulations show that the decrease in squared norm is inversely linear but that the constant is orders of magnitude less than that given in Theorem 3.1.

3.1 The Approximation Algorithm

The basic algorithm is: given a function specified at the points $\{\mathbf{x}_n\}$, $n = 1, \dots, N$ run the hinge finding algorithm on this data, resulting in a hinge function approximation $h(\mathbf{x})$. Since the function to be fitted is not a hinge function, M starting values of $\Delta^{(0)}$ are used, and the hinge adopted is that with minimum RSS. The m^{th} starting $\Delta^{(0)}$ is given by $\Delta^{(0)}(j) = 0$, $j \neq m$, $\Delta^{(0)}(m) = 1$, and $\Delta^{(0)}(0)$ selected so that the condition $\{\Delta^{(0)} \cdot \mathbf{x} \geq 0\}$ cuts the data in half.

To compute the approximation to $f(\mathbf{x})$ using K hinge functions, at the K^{th} stage find the hinge function approximation $h_K(\mathbf{x})$ to $f(\mathbf{x}) - \sum_1^{K-1} h_k(\mathbf{x})$. Then refit: update h_1

by refitting the difference $f - \sum_{k=1}^K h_k$. Using this updated h_1 , update h_2 by refitting $f - h_1 - \sum_{k=1}^K h_k$. After h_K is refitted, start the cycle again with f_1 . These cycles are continued until there is no further appreciable decrease in RSS. The procedure is made more efficient by using, as the *single* starting $\Delta^{(0)}$ for each refit, the current hinge of the function to be refitted.

Another algorithm was also tested. Each hinge function is the sum of a linear function and a function $h^+(\Delta \cdot x)$ where $h^+(x) = x$, $x \geq 0$, and 0 , $x < 0$. Let $h^+(\Delta_k \cdot x)$, $k=1, \dots, K-1$ be the nonlinear parts of the hinge functions entered at steps $1, 2, \dots, K-1$. Do a linear regression of $f(x)$ on the $M + K$ variables $1, x_1, \dots, x_M, h^+(\Delta_1 \cdot x), \dots, h^+(\Delta_{K-1} \cdot x)$, getting

$$\hat{f}_{K-1}(x) = \sum_{m=0}^M \beta_m x_m + \sum_{k=1}^{K-1} \gamma_k h^+(\Delta_k \cdot x).$$

Use the hinge finding algorithm on $f - \hat{f}_{K-1}$ to find h_K and suppose the nonlinear part of h_K is $h^+(\Delta_K \cdot x)$. Do a linear regression of $f(x)$ on the $M + K + 1$ variables

$$1, x_1, \dots, x_M, h^+(\Delta_1 \cdot x), \dots, h^+(\Delta_K \cdot x)$$

getting new coefficients β_0, \dots, β_M and $\gamma_1, \dots, \gamma_K$. Take \hat{f}_K to be the linear combination using the new coefficients.

Although refitting is computationally fairly efficient, the cpu time to fit, say, 50 hinge functions in 16 dimensions to 1000 data points is considerably larger than the regression type algorithm described in the above paragraph. The trade off in cpu time and accuracy is explored in the simulations reported on in the next section.

3.2 Simulations: Function Approximation

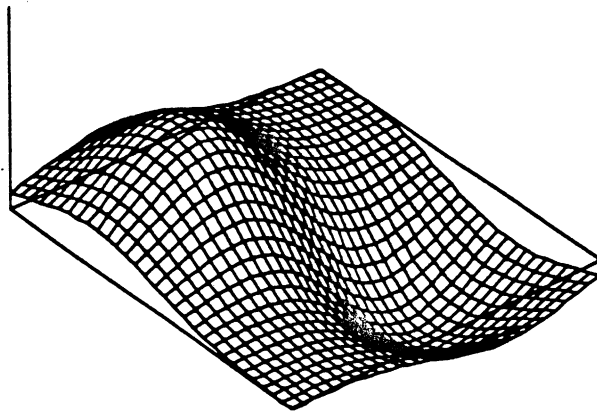
The first example is chosen for visual inspection. There are 1000 $\{x_n\}$ values uniformly distributed on the square $[0, 1]^2$. Let $e = (1, 1)$, then

$$f(x) = e^{-\|x - 3e\|^2} - e^{-\|x - .7e\|^2}.$$

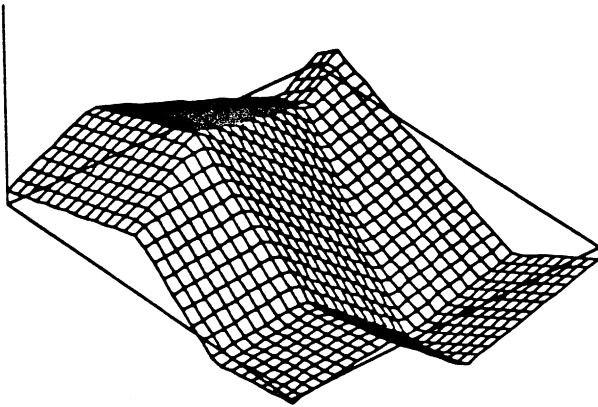
Figure 1a gives the surface plot of f . Hinge functions are fitted and refitted. The surface plots of $\sum_{k=1}^K h_k(x)$, $K = 4, 8, 16, 32$ are given in Figure 1 b,c,d,e.

We are not advocating the use of hinge functions to fit smooth surfaces in low dimensions. Other methods are available which give smoother and more accurate fits in 2 or 3 dimensions (see, for instance, Breiman [1991], Friedman [1991], Whaba [1990]). The above 2-dimensional example is given only because visual inspection is possible.

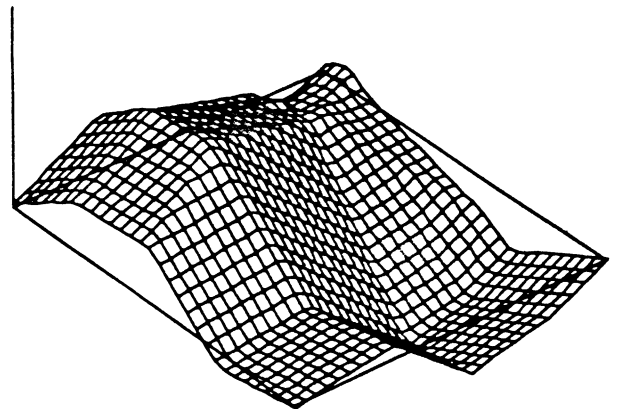
Figure 1. Surface Plots of Hinge Approximations



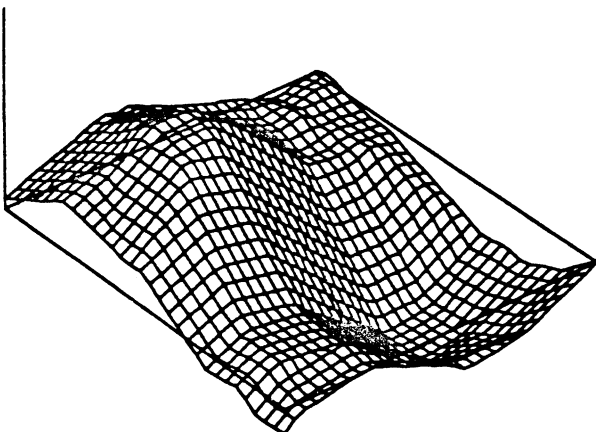
1(a) Original Function



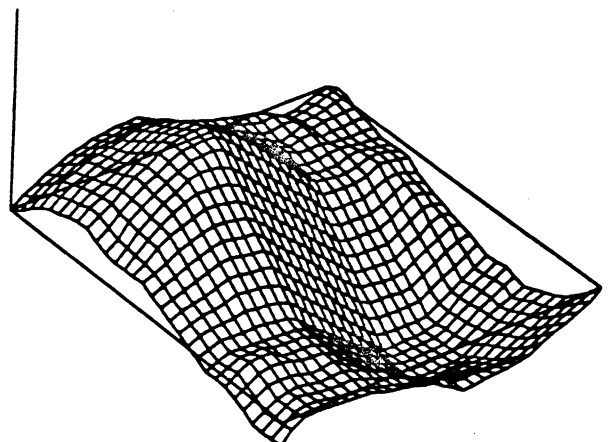
1(b) 4 Hinges



1(c) 8 Hinges



1(d) 16 Hinges



1(e) 32 Hinges

Figure 3

Inverse MRSS Using Refitting

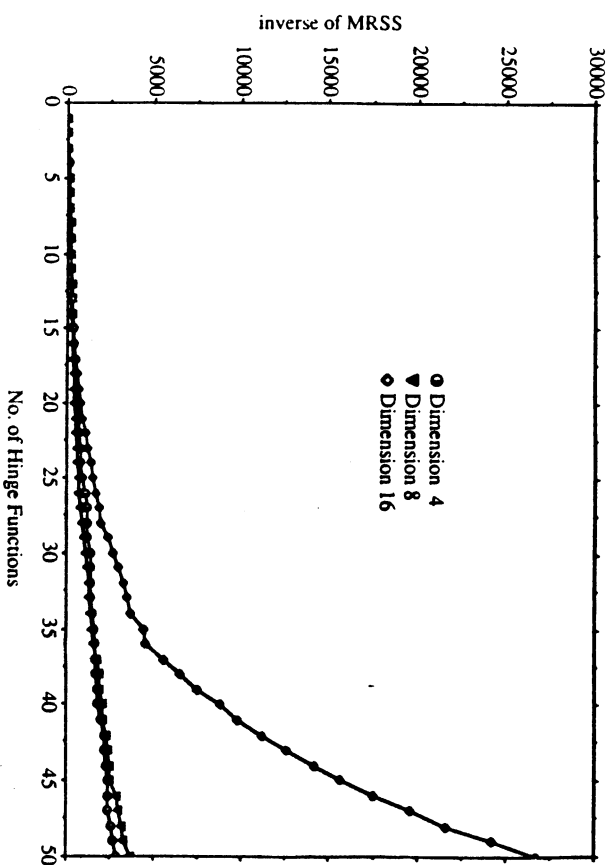


Figure 2

Comparison of Methods in 4 Dimensions

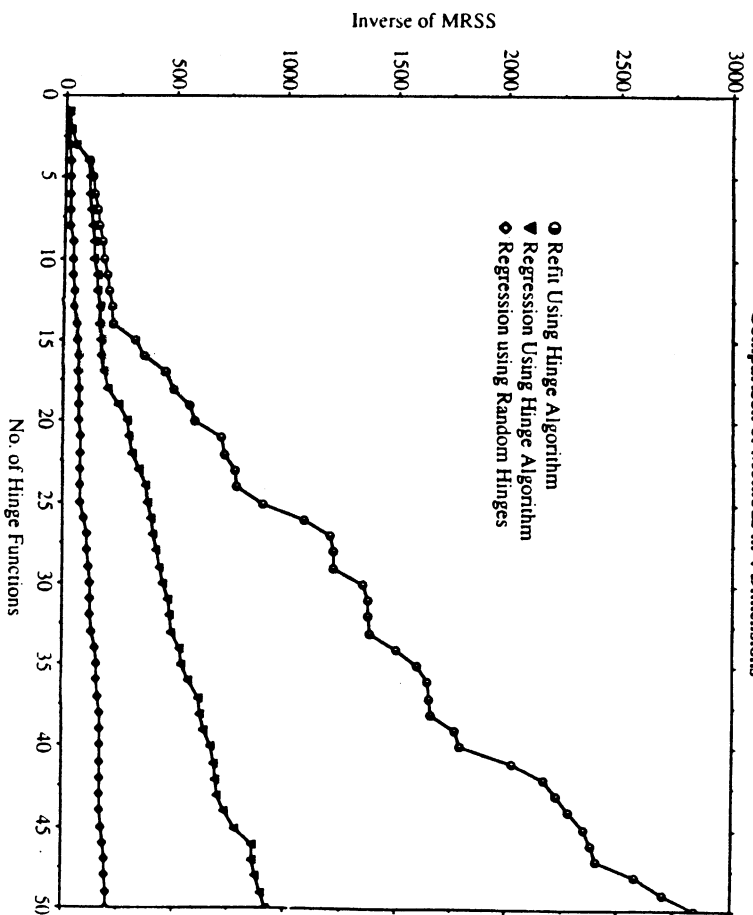
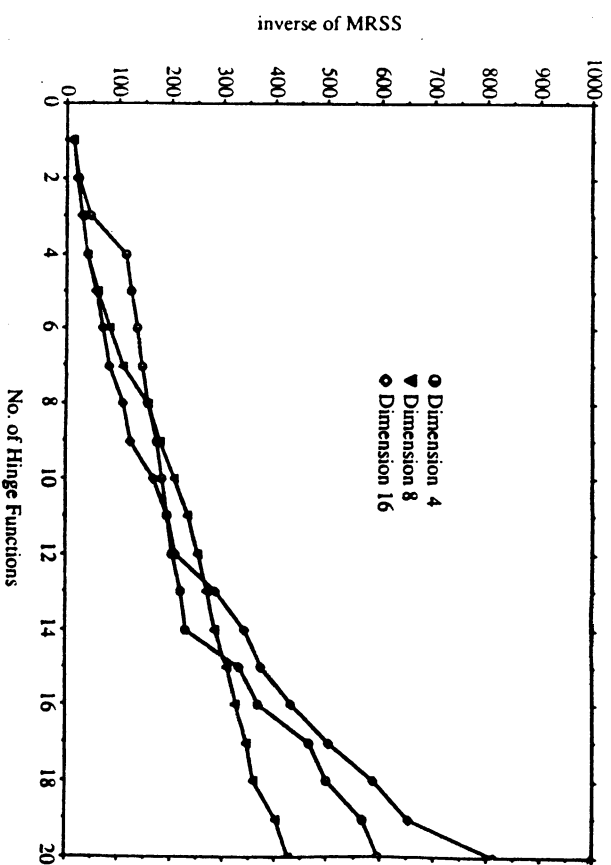


Figure 4

Inverse MRSS Using Refitting



The second series of examples are similar except for dimension. The function is

$$f(\mathbf{x}) = e^{-\|\mathbf{x}\|^2/2}.$$

The points $\{\mathbf{x}_n\}$ are distributed on the sphere $\|\mathbf{x}\| \leq 3$ using a spherically symmetric distribution such that $\|\mathbf{x}\|$ is uniform on $[0, 3]$. For this function

$$M = \int \|\mathbf{w}\|^2 |\tilde{f}(\mathbf{w})| d\mathbf{w},$$

so that the upper bound of theorem 3.1 is $1296 M^2 / K$.

For $M = 4$, the refit and regression algorithm were run up to $K = 50$. Figure 2 gives the plots of $1/\|f - \sum_1^K h_k\|^2$ vs K for both the refit and regression algorithms for $K = 1$ to 50. Fitting the $1/\|f - \sum_1^K h_k\|^2$ data by K/b using least squares gives $b = .022$ (refit) and $b = .061$ (regression). Either one is orders of magnitude smaller than the constant given by Theorem 3.1.

There is a third set of points plotted in Figure 2. The suspicion arises that if one put down a series $\Delta_1, \dots, \Delta_K$ of hinges chosen at random and regressed $f(\mathbf{x})$ on the variables $1, x_1, \dots, x_M, h^+(\Delta_1 \cdot \mathbf{x}), \dots, h^+(\Delta_K \cdot \mathbf{x})$ one might do almost as well as using the hinge finding algorithm. To check this, random hinges $\Delta_1, \dots, \Delta_K$ were generated with $(\Delta_K(1), \dots, \Delta_K(M))$ being i.i.d. $N(0, 1)$ and $\Delta_K(0)$ taken such that the proportion of data satisfying $\{\Delta_K \cdot \mathbf{x}_n \geq 0\}$ is uniform on $[.1, .9]$. The values $\|f - \sum_1^K h_k\|^{-2}$ generated this way form the third graph on Figure 2.

To explore the effect of dimensionality, we ran the refit algorithm using dimensions 4, 8, 16 and going up to $K = 50$. The data was generated as described above with sample size 1000. Figure 3 gives the graphs of $1/\|f - \sum_1^K h_k\|^2$ vs K for $K = 1, \dots, 50$, including the results for dimension four. The MRSS for dimension 16 decreases rapidly when more than 20 hinge functions are fitted. This is probably due to the fact that 1000 data points are sparse in 16 dimensions. Fitting 20 hinge functions involves optimizing over almost 400 parameters, and 50 hinge functions, over almost 900. As the number of parameters approach the number of data points, the error of the *fit at these points* drops rapidly to zero.

A more interesting comparison of the effect of dimensionality is when the number of parameters is constrained to be small compared with the number of data points. Figure 4 graphs the same data as Figure 3 but only up to 20 hinge functions. There is very little effect due to dimensionality.

The run doing fitting and refitting of the 50 hinge functions in 16 dimensions took 11.4 cpu minutes (RS 6000), pretty fast for a highly nonlinear optimization involving almost 900 parameters. The regression type algorithm took 2.3 cpu minutes.

4). Using Hinge Functions in Regression

Hinge functions give a potentially powerful new approach to nonlinear regression. Given data of the form $\{(y_n, \mathbf{x}_n), n = 1, \dots, N\}$ the problem is to use this data to construct a function $\hat{f}(\mathbf{x})$ that will give accurate predictions of future unknown y -values when the corresponding \mathbf{x} -values are known.

If $\hat{f}(\mathbf{x})$ is restricted to be linear, then squared error loss leads to classical least squares regression. If the number of variables is nominal, say ≤ 20 , then the methods proposed by Friedman [1991] and Breiman [1991] give continuous predictors $\hat{f}(\mathbf{x})$ restricted to be the sum of nonlinear functions each depending on only a few variables, say one, two, or at most three. But with some data, predicting y by sums of functions of a small number of the \mathbf{x} -variables may not give accurate results. Consider, for instance, the function $e^{-\|\mathbf{x}\|^2}$ in M -dimensions.

Methods for fitting continuous functions to high interaction, high-dimensional data are rare. One early and remarkable result due to Meisel and Collins [1973] derives a piecewise continuous hyperplane estimate using a method much different from hinging. Friedman and Stuetzle [1981] originated “projection pursuit” regression which uses a sum of estimated smooth functions of linear functions. The tree-structured approach (Breiman, et. al [1984]) fits a discontinuous histogram-like functions.

The approach using hinge functions in prediction is the same as in the noiseless case -- find the best fit to y of $\sum_{k=1}^K h_k(\mathbf{x})$. One important issue is how many hinge functions to fit. If K is too large, overfitting occurs and the fit loses accuracy on future data. Two methods are useful for selecting K . The simplest uses the PE_{GCV} criterion:

Let $MRSS(K)$ be the mean residual sum-of-squares. The PE_{GCV} estimate for the test set prediction error is given by

$$PE_{GCV}(K) = MRSS(K) / (1 - c(K+1)(M+1)/N)^2$$

where N is the sample size, and M the dimension and c a parameter in the $[1,3]$ range. Now take K to minimize $PE_{GCV}(K)$. The most accurate estimate of test set error is given by cross-validation, which requires much more computing. For derivations see Appendix III.

4.1 A High Interaction Example

As an example of a high interaction problem in 10-dimensional space, take

$$g(\mathbf{x}) = \text{logit}(v_1(\mathbf{x})) + \text{logit}(v_2(\mathbf{x})) + \text{logit}(v_3(\mathbf{x})).$$

The $v_1(\mathbf{x})$, $v_2(\mathbf{x})$, $v_3(\mathbf{x})$ are linear in (x_1, \dots, x_{10}) defined by letting

$$l_1 = 10x_1 + 9x_2 + 3x_3 + 7x_4 - 6x_5 - 5x_6 - 9x_7 - 3x_8 - 2x_9 - x_{10}$$

$$l_2 = -x_1 - 2x_2 - 3x_3 - 4x_4 - 5x_5 - 6x_6 + 7x_7 + 8x_8 + 9x_9 + 10x_{10}$$

$$l_3 = -x_1 - 2x_2 - 3x_3 + 4x_4 + 5x_5 + 4x_6 - 3x_7 - 2x_8 - x_9.$$

The $\{x_n\}$ are selected from a uniform distribution on $[0, 1]^{10}$. Normalize the l_1, l_2, l_3 to have upper 97.5% point = 2.0, and then put v_i equal to $2(l_i - 2)$.

Sample size was 400 and the $\{y_n\}$ generated as

$$y_n = a \cdot g(x_n) + \epsilon_n$$

where the $\{\epsilon_n\}$ are unit normal noise. The constant a was defined so that for $f(x_n) = a \cdot g(x_n)$, the standard deviation of the $\{f(x_n)\}$ was 4.0, giving a s/n ratio of 4.0.

A 4000 member test set was generated and used to estimate the prediction error (PE_{TS}) as successive hinge functions were added. The PE_{GCV} for $c = 1.5$ was also computed. The initial value of the sample variance of the $\{y_n\}$ was 17.4. A linear regression fitted to the data resulted in a mean residual-sum-of-squares of 10.3.

The number of hinge functions fitted was increased from 1 up to 7 with the results summarized in Table 4.1

Table 4.1
Summary of Fit by Hinge Functions

# Hinge Functions	MRSS	PE_{GCV}	PE_{TS}
1	5.22	6.20	5.69
2	2.07	2.69	2.49
3	1.09	1.56	1.52
4	1.00	1.58	1.52
5	.96	1.59	1.53
6	.82	1.62	1.54
7	.78	1.75	1.59

The coefficients of the hinges in the 3-hinge fit are given in Table 4.2.

Table 4.2
Coefficients of Hinges in 3-Hinge Fit

Variable	First Hinge	Second Hinge	Third Hinge
1	-2.2	9.8	.4
2	-6.0	10.6	.5
3	-8.1	8.2	1.4
4	10.9	6.8	3.5
5	13.4	-5.7	5.9
6	10.9	-5.2	6.3
7	-9.0	-5.0	-6.6
8	-6.5	-3.1	-7.4
9	-3.1	-2.6	-6.8
10	.7	-2.5	-8.2

The action of these three hinges is clear from inspection. The first hinge is fitting $\text{logit}(v_2)$ on one side and the other two logits on the other side. The second hinge is fitting $\text{logit}(v_1)$ on one side and the other two logits on the other side. Similarly, the third hinge fits $\text{logit}(v_3)$ on one side and the others on the other side.

The minimum PE_{GCV} selection criterion picks the same 3-hinge fit as the minimum PE_{TS} criterion, (although PE_{TS} ties between the 3 and 4-hinge fits). In fact, the PE_{GCV} is a minimum at the 3 hinge fit for every value of the parameter c in $[1,3]$.

The running time for this example is 7.4 cpu seconds.

4.2 Stepwise Forward Selection of Variables and a Higher Dimensional Examples

If the data set has, say, 1000 cases and 100 variables, then fitting 4 hinge functions involves the estimation of 500 parameters. This is only 2 cases per parameter estimated and will probably result in a noisy estimate. If possible, one would want to keep a tighter control on the number of parameters estimated. In addition, with larger dimensionality, fitting hinge functions becomes slower.

To deal with these two issues, a stepwise forward variable selection method for entering hinge functions is used. Here is the idea: start with a search through all M variables to find the single variable hinge that gives minimum RSS. Call the variable used x_1 . Now search among all x_m , $m \geq 2$ to find the lowest RSS hinge based on the pairs of variables (x_1, x_m) . Keep adding variables until a minimum is found in the GCV estimate of prediction error. Then start this process over to find the next hinge: refit and keep adding hinges until the PE_{GCV} hits a minimum.

To illustrate, the following 100 dimensional example was run: the sample size was 2000 and the $\{\mathbf{x}_n\}$ were sampled from the uniform distribution on $\{0,1\}^{100}$. Define

$$g(\mathbf{x}) = \sum_1^K \text{logit}(v_i(\mathbf{x}))$$

where the $\{(v_i(\mathbf{x}))\}$ are linear functions of x_1, \dots, x_{100} . The y_n are given by

$$y_n = a \cdot g(\mathbf{x}_n) + \epsilon_n$$

where the $\{\epsilon_n\}$ are unit normals and a is taken to make the s/n rate 4.0.

The $\{v_1\}$ are defined as follows: For $k = 20, 40, 60, 80$, let

$$c(k, m) = \exp(-.7 |m - k|),$$

define

$$u_k(\mathbf{x}) = \sum_m c(k, m) x_m,$$

and

$$l_1 = 3u_1 - u_2 - u_3 - u_4$$

$$l_2 = -u_1 + 3u_2 - u_3 - u_4$$

$$l_3 = -u_1 - u_2 + 3u_3 - u_4$$

$$l_4 = -u_1 - u_2 - u_3 + 3u_4.$$

Normalize the l_i to have unit variance and take $v_i = l_i - 1$.

The sample variance of the $\{y_n\}$ was 17.1. A linear regression on all 100 variables give a mean residual-sum-of-squares of 16.3, so the regression surface was predominantly nonlinear. The time needed to run this example was 2.8 cpu minutes (RS 6000). Table 4.3 summarizes the results.

Table 4.3
Summary of Results for the 100-dimensional Example

Hinge	# of variables*	MRSS	PE _{GCV} **	PE _{TS}
1	21	11.01	11.78	12.70
2	20	6.71	7.67	8.20
3	25	2.83	3.52	3.57
4	13	2.16	2.82	2.69
5	4	2.13	2.83	2.71
6	5	2.09	2.84	2.73
7	3	2.08	2.86	2.75
8	3	2.06	2.88	2.79
9	3	2.04	2.90	2.83
10	4	2.02	2.90	2.85

*the minimum number of variables in a hinge was set at 3. ** $c = 1.5$

The structure of both of the above examples was devised so that an optimum fit could be gotten using a neural network. It would be interesting to see what accuracies and compute times are produced by running these examples on a neural network program.

5). Using Hinges in Classification

To use hinges in classification, the problem needs to be reformulated into a regression context. Suppose there are J classes numbered $1, \dots, J$ with the probability $P(j|x)$ of being in class j given x assumed known. Then the Bayes optimal classification rule is: classify x into that class j for which $P(j|x)$ is maximum.

Let the random variable Y be one in class j , otherwise zero. Then the function θ of the vector X which minimizes $E(Y - \theta(X))^2$ is $P(j|x)$. With classification type data $\{(j_n, x_n), n = 1, \dots, N\}$ and $j_n \in \{1, \dots, J\}$, these remarks suggest the following approach: define

$$y_{jn} = \begin{cases} 1, & j_n = j \\ 0 & \text{otherwise} \end{cases}$$

Find J functions $\hat{\theta}_j(x)$ such that $RSS_j = \|y_j - \hat{\theta}_j\|^2$ is small. Classify a future object with observed x as class j if

$$\theta_j(x) = \max_i \theta_i(x).$$

The problem is how to estimate the functions $\theta_j(x)$ from the data. If linear θ_j are used, the result is similar to discriminant analysis, and performance can be poor in situations where the $P(j|x)$ are not well approximated by linear functions. The neural network methodology uses $\theta_j(x)$ estimated through the use of sums of sigmoid functions of linear functions.

One approach is -- for each j approximate the $\{y_{jn}\}$ by a sum of hinge functions $\sum_1^K h_{jk}(x_n)$. This uses up parameters at an alarming rate. A different approach, employed in classification trees, for instance, is to use a common basis of functions, h_1, \dots, h_K , to approximate each set of $\{y_{jn}\}$, with the approximations differing from class to class through the coefficients of the $\{h_k\}$, i.e. the approximations

$$y_j \sim \sum_k \alpha_{jk} h_k(x) \quad (5.1)$$

are used. Then the problem becomes to find those hinge functions that make

$$RSS = \sum_j \min_{\{\alpha_{jk}\}} \|y_j - \sum_{k=1}^K \alpha_{jk} h_k\|^2$$

as small as possible.

Suppose h_1, \dots, h_{K-1} have been selected and residuals

$$r_{jn} = y_{jn} - \sum_{k=1}^{K-1} \alpha_{jk} h_k(x_n)$$

computed. What is wanted is a hinge function f_K , and coefficients $\{\gamma_j\}$ to minimize

$$\sum_{j,n} (r_{jn} - \gamma_j f_K(x_n))^2. \quad (5.2)$$

Start with a hinge $\Delta^{(0)}$, and denote by (+), that data for which $\Delta^{(0)}x \geq 0$ and (-) for the other data. Consider fitting a hyperplane $\beta^+ \cdot x$ to the (+) data, and a hyperplane β^- to the (-) data and selecting $\{\gamma_j\}$, β^+ , β^- to minimize

$$\sum_{-j} (r_{jn} - \gamma_j \beta^- \cdot x_n)^2 + \sum_{+j} (r_{jn} - \gamma_j \beta^+ \cdot x_n)^2. \quad (5.3)$$

After this minimization has been carried out, the next hinge is $\Delta^{(1)} = \beta^+ - \beta^-$, and the process is repeated until convergence.

Minimization of (5.3) leads to either a $J \times J$ or $M \times M$ eigenvalue problem. We derive the $J \times J$ problem. To begin note that there is an indeterminacy in the scales of γ , β^+ , β^- . This is resolved by taking $\|\gamma\|^2 = 1$. Partial derivatives of (5.2) with respect to β_m^+ , β_m^- leads to the equations

$$S^+ \beta^+ = Z^+ \gamma,$$

$$S^- \beta^- = Z^- \gamma$$

where S^+ is the $M \times M$ matrix, $S_{mn}^+ = \sum_{+} x_{mn} x_{m'n}$, Z^+ is the $M \times J$ matrix $Z_{mj}^+ = \sum_{+} x_{mn} r_{jn}$; similarly for S^- and Z^- . Substituting these expressions for β^+ , β^- into (5.3) gives the expression

$$\sum_{j,n} r_{jn}^2 - \sum_{jj'} \gamma_j \gamma_{j'} H(j, j')$$

where

$$H = (Z^+)^t (S^+)^{-1} Z^+ + (Z^-)^t (S^-)^{-1} Z^-.$$

To minimize the RSS we want to maximize $\gamma^t H \gamma$ under the constraint $\|\gamma\| = 1$. This solution is the eigenvector of H corresponding to its largest eigenvalue. Once γ is known, then $Z^+ \gamma$, $Z^- \gamma$ are computed, $\beta^+ = (S^+)^{-1} Z^+ \gamma$, $\beta^- = (S^-)^{-1} Z^- \gamma$ and $\Delta^{(1)} = \beta^+ - \beta^-$. After h_k is entered, then refitting cycles are carried out refitting both the $\{h_k\}$ and the $\{\alpha_{jk}\}$. The decision on how many hinge functions to use in the fit can be based either on a test set or cross-validation.

5.1 Two Examples:

The first data set used as an example was provided to us by Richard Lippmann. To quote him "The database consists of the first two formants frequencies (F1 and F2) of vowels [Peterson and Barney, "Control Methods Used in a Study of Vowels", THE JOURNAL OF THE ACOUSTICAL SOCIETY OF AMERICA, Vol. 24, No. 2, March 1952, pp 175-184. in Hz. The actual data was gathered by digitizing a figure from Rabiner and Schafer, "Digital Processing of Speech", Prentice-Hall, 1978. The digitized data was arbitrarily divided into 338 training samples and 333 testing samples."

There are ten vowel classes having roughly equal representations in both the training and test set. This two dimensional data was used as a benchmark for various classifiers by Lee and Lippman [1989]. Table 5.1 is an excerpt of their results.

Table 5.1
Classifier Performance

Classifier	Test Set	% Error Training Set	Times (cpu seconds)
Unimodal Gaussian	20.4	22.5	31.4
Back-Prop.	21.0	21.0	2458.4
Hypersphere (RCE)	23.1	17.3	144.9
KNN	17.4	20.4	.5
LVQ	18.0	13.5	425.9
Feature Map	19.5	22.2	113.1
CART Tree	21.9	16.0	8.2

Fitting 22 hinge functions gave a test set error of 18.6% and training set error of 20.4%. The compute time was 41.9 cpu seconds. This translates into 517.5 cpu seconds scaled to the machine used by Lee and Lippman. Although we conceived of the hinge function methodology as primarily useful in high dimensions, it is competitive in this 2-dimensional example.

5.2 Simulated Wave Form Data

The second example consists of simulated data with structure given in pp. 49-55 of the book by Breiman et. al [1984]. It is a 3-class, 21-dimensional problem based on the waveforms $w_1(t)$, $w_2(t)$, $w_3(t)$ graphed in figure 5.

Each class consists of a random convex combination of two of these waveforms sampled at the integers with noise added. More specifically, the measurement vectors are 21 dimensional: $\mathbf{x} = (x_1, \dots, x_{21})$. To generate a class 1 vector \mathbf{x} , independently generate a uniform random number u and 21 random numbers $\epsilon_1, \dots, \epsilon_{21}$ normally distributed with mean zero and variance 1. Then set

$$x_m = uw_1(m) + (1 - u)w_2(m) + \epsilon_m, \quad m = 1, \dots, 21.$$

To generate a class 2 vector, repeat the preceding and set

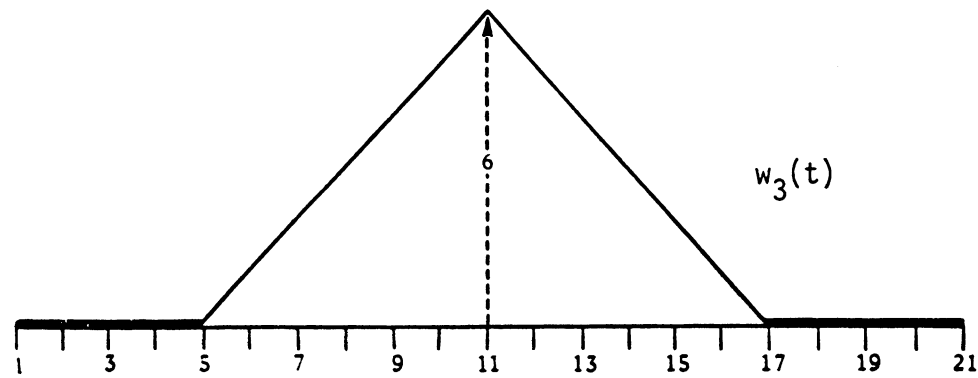
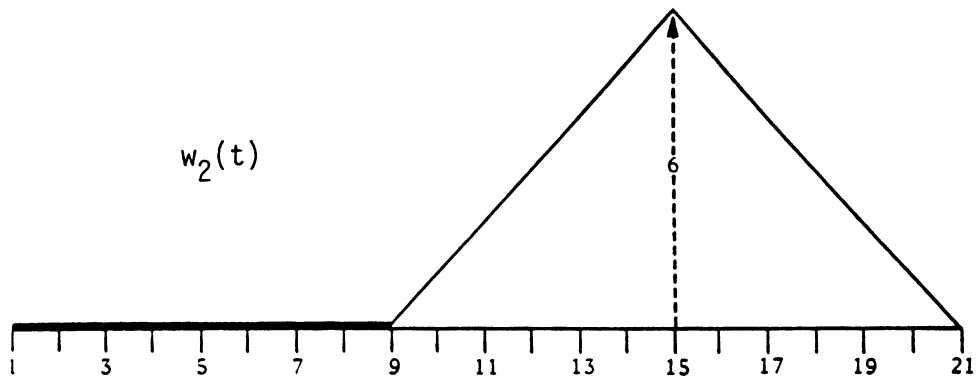
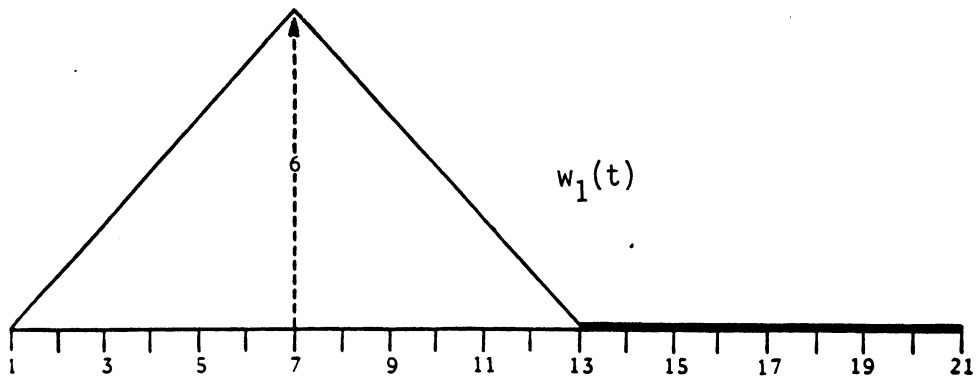
$$x_m = uw_1(m) + (1 - u)w_3(m) + \epsilon_m, \quad m = 1, \dots, 21.$$

Class 3 vectors are generated by

$$x_m = uw_2(m) + (1 - u)w_3(m) + \epsilon_m, \quad m = 1, \dots, 21.$$

Three hundred measurement vectors were generated using prior probabilities of $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$, so there were approximately 100 per class.

Figure 5. Wave Forms



In one set of data generated as above, the CART tree classifier had a test set error rate of 28%. Using linear combinations in the tree construction lowers the test set error rate to 20%. Linear discriminant analysis with stepwise entry of variables gave a test set error rate of 26%.

Ten data sets were generated using the above mechanism, together with ten test sets of size 3000. When each training set hinges were added and refitted. The model selected (i.e. number of hinges fitted) was that with minimum test set error.

For the stochastic structure of this data, the Bayes rule can be computed and applied to the test sets. The averages over the 10 repetitions are:

number of hinges	3.6
training set misclassification rate (%)	9.5
test set misclassification rate (%)	18.1
Bayes rule misclassification rate (%)	13.6
cpu seconds (RS 6000)	39.7

The hinges procedure improves on the tree classifier using linear combination splits. But the error rate is still somewhat above the Bayes rate. We revisit this example in the next section.

5.3 Stepwise Addition of Variables and More Examples

For high dimensional problems a variable selection method is imperative. The method used in classification is an extension of the regression method. Suppose that hinge functions h_1, \dots, h_{K-1} have been selected and residuals $\{r_{jn}\}$ computed. What is wanted is a hinge function h_K and $\{\gamma_j\}$ to minimize (5.2).

The stepwise procedure first finds the single variable hinge (say on x_1) and coefficients $\{\gamma_j\}$ that minimize (5.2). Then it finds the two variable hinge based on (x_1, x_m) , $m \neq 1$, and coefficients $\{\gamma_j\}$ that minimize (5.2). Variables are added until the PE_{GCV} criterion derived in Appendix III becomes minimum. Then refitting is carried out among all K hinge functions. The number of hinge functions fitted is determined by test set or cross validation.

5.4 The Wave Form Example Revisited

The simulated data of section 5.2 was rerun using the stepwise procedure. The results (averaged over the 10 repetitions) were:

number of hinges	3.7
number of variables per hinge	9.3
training set misclassification rate (%)	10.5
test set misclassification rate (%)	17.3
Bayes rule misclassification rate (%)	13.6
cpu seconds	7.3

Note that the accuracy increases by .8% while the compute time has been cut from 40 cpu seconds to 7 cpu seconds. The decrease of .8% seems insignificant, but another way to look at it is that the all variable procedure is 4.5% above the Bayes rate while the stepwise procedure is 3.7% above.

5.5 High Dimensional Simulated Data

To test the stepwise entry method on a larger problem, we generated 10 class, 61 dimensional data with 1000 cases. In structure, it is an extension of the 3 class wave form data.

The wave form function $w(x) = 10 - \min(10, |x|)$, has a peak at zero of height 10, and is zero outside of $[-10, 10]$. Consider five functions

$$w_k(m) = w(m - 10k - 1), \quad k = 1, \dots, 5$$

where $m = 1, \dots, 61$. These functions are centered at the points 11, 21, 31, 41, 51. Let S_1, \dots, S_{10} consist of all subsets of size three of $\{1, \dots, 5\}$, i.e. $S_1 = \{3, 4, 5\}$ $S_2 = \{2, 4, 5\}, \dots, S_{10} = \{1, 2, 3\}$.

The data is generated as follows: For $n = 1$ to 1000, the class $j \in 1, \dots, 10$ is selected with probability .1 of each choice. Three uniform random numbers u_1, u_2, u_3 are generated. For the lowest index k_1 in S_j , $q(k_1)$ is defined as $u_1 / (u_1 + u_2 + u_3)$; for the second lowest index k_2 in S_j , $q(k_2) = u_2 / (u_1 + u_2 + u_3)$ and for the third index $q(k_3) = u_3 / (u_1 + u_2 + u_3)$. For $k \notin S_j$, $q(k) = 0$. Then, for $m = 1$ to 61

$$x(m, n) = \sum_k q(k)w(m, k) + \epsilon_m,$$

where the $\{\epsilon_m\}$ are independent unit normals.

Two areas were investigated with this data. First was the question of how accurately could classification be done using linear methods only. The second was how accuracy was effected by changing the constant in the criterion which governs the number of variables entered into a hinge.

To explore the accuracy of linear methods, a stepwise entry of variables into a linear classifier was employed. Let $S \subset \{0, \dots, 61\}$ be the subset of variables already entered, with $|S| = K$. Then the $(K + 1)$ st variable entered is that x_k , $k \notin S$ which minimizes

$$\min_{\{c_{jm}, c_{jk}\}} \sum_{j,n} (y_{n,j} - \sum_{m \in S} c_{j,m} x_{mn} - c_{j,k} x_{kn})^2.$$

For each K a 3000 case test set was used to estimate the classification error. The minimum test set misclassification rate was 37.3% occurring when 12 variables were entered.

The second question is this: the PE_{GCV} criterion given in appendix II keeps adding variables to a hinge as long as the PE_{GCV} value is decreasing. Using $c = 1.5$, the decreases are quite small for a long string of entries. In this case, is the program using up degrees of freedom for small gain, and could more accuracy be accomplished by going on to the next hinge? To investigate this, we ran the program using $c = 1.5, 2.0, 3.0, 6.0$. Table 5.1 gives the results for the number of hinges giving the minimum test set error.

Table 5.1
Summary of Runs on Extended Waveform Data

	Value of Constant			
	1.5	2.0	3.0	6.0
Number of Hinges	4	4	4	4
Average No. Variables/Hinge	25.5	16.8	11.5	6.5
Training Set Error (%)	14.0	13.9	15.5	17.3
Test Set Error (%)	19.2	18.6	18.6	19.8
CPU seconds	135.4	84.0	70.1	40.4

The accuracy results are reasonably insensitive to the number of variables per hinge — equivalently, to the value of c between 1.5 and 6.0. The evidence of this example argues for a value of c between 2.0 and 3.0. But in the smaller waveform example increasing c from 1.5 to 2.0 increases the test set error to 17.7%.

The nonlinear hinge procedure improves considerably on the linear error rate. Unfortunately, the Bayes procedure is difficult to calculate, so there is no way of telling how much better can be done.

6). Conclusions

Fitting hinge functions to data can give good predictive results in regression, classification and noiseless function approximation. Using stepwise entry of variables into the hinges, accurate predictors in fairly large problems can be derived in short compute times. These preliminary results are encouraging, but there are important questions that need further work. A short list is:

1. Choice of a good initial starting hinge is important. The current strategy uses M different starting hinges in M dimensions and takes the best result. Is it possible to find a single starting hinge that will provide uniformly good results?

2. The stepwise entry method used is crude. A more refined procedure, for example, is to enter two variables at a time and then delete from the set of entered variables that single variable whose deletion causes the least rise in RSS. Can more sophisticated stepwise entry procedures significantly improve accuracy?

3. In stepwise entry, the variables in the first hinge sometimes appear to be selected almost by happenstance, as the algorithm is hunting around to get some purchase. Is there a way to improve the variable selection in the first hinge?

4. The hinge algorithm is not guaranteed to find the global minimum at $\|f - \sum_1^K h_k\|$ over h_1, \dots, h_K . If $f(\mathbf{x})$ is badly behaved, the hinge algorithm may converge to a poor fit. Is there a way to consistently get near the global minimum?

This methodology will be explored in the binary tree context, and using a related iterative approach, research into fitting by sums and products of ramp functions is being carried out in collaboration with Jerome Friedman. (A ramp function is a continuous function $y = r(\mathbf{x})$ formed by the intersection of a hyperplane with two planes of constant height).

References

- Barron, A.R. (1991). "Universal Approximation Bounds for Superpositions of a Sigmoidal Function". Technical Report #58, Department of Statistics, University of Illinois at Urbana-Champaign.
- Breiman, L. (1991). "The π Method for Estimating Multivariate Functions from Noisy Data" (with discussion). *Technometrics* Vol. 33, No. 2, 125-160.
- Breiman, L. (1990). "The Little Bootstrap and Other Methods for Dimensionality Selection in Regression: X-Fixed Prediction Error". Technical Report No. 169, Statistics Department, University of California, Berkeley, CA (to be published, Journal of the American Statistical Association).
- Breiman, L., and Spector, P. (1990). "Submodel Selection and Evaluation in Regression: the X-Random Case". Technical Report No. 197, Statistics Department, University of California, Berkeley, CA (to be published, International Statistical Review).
- Breiman, L., Friedman, J., Olshen, R. and Stone, C. (1984). "Classification and Regression Trees". Wadsworth Inc.
- Breiman, L. and Friedman, J. (1985). "Estimating Optimal Transformations in Regression and Correlation" (with discussion). *Journal of the American Statistical Association* 80, 580-619.
- Breiman, L. and Ihaka, R. (1987). "Nonlinear Discriminant Analysis via Scaling and ACE". Technical Report #40, Statistics Department, University of California, Berkeley.
- Friedman, J. (1991). "Multivariate Adaptive Regression Splines" (with discussion). *Annals of Statistics*, Vol. 19, No. 1, 1-141.
- Friedman, J. and Stuetzle, W. (1981). "Projection Pursuit Regression". *Journal of the American Statistical Association* 76, 817-823.
- Hastie, T. and Tibshirani, R. (1990). "Generalized Additive Models". Chapman and Hall.
- Jones, L. (1990). "A Simple Lemma on Greedy Approximation in Hilbert Space and Convergence Rates for Projection Pursuit Regression and Neural Network Training". Technical Report No. 16, Department of Mathematics, University of Lowell, Lowell, MA (to appear in Annals of Statistics).

Lee, C. and Lippmann, R. (1989). "Practical Characteristics of Neural Network and Conventional Pattern Classifiers on Artificial and Speech Problems". NIPS 89, Nov. 1989, Denver, CO.

Lippmann, R. (1989). "Pattern Classification Using Neural Networks". *IEEE Communications Magazine* **11**, 47-64.

Meisel, W. and Collins, C. (1973). "Repro-Modeling: An Approach to Efficient Model Utilization and Interpretation". *IEEE Trans. on Systems, Man, and Cybernetics* Vol. 3, No. 4, 349-358.

Wahba, G. (1990). "Spline Functions for Observational Data". SIAM, Philadelphia.

APPENDIX I

Proofs

Proof of Theorem 2.1.

W.l.o.g. we can assume that $\Gamma = I$ and $\|\Delta^*\| = 1$. For any $\Delta^{(0)}$, let $S_+ = \{x; \Delta^{(0)} \cdot x \geq 0\}$, S_- the complement space, and define Γ_+ , Γ_- as previously. Using an orthogonal transformation on x_0, x_1, \dots, x_m we can reduce Γ_+ , Γ_- to diagonal form D_+ , D_- , where $D_+ = I - D_-$. By (2.1)

$$\Delta^* - \Delta^{(1)} = D_+^{-1} \Gamma_+ \Delta^* + D_-^{-1} \Gamma_- \Delta^*.$$

Let $s = \Delta^* \cdot x$, and $\Delta^{(0)} \cdot x = \alpha s + \bar{\alpha} t$, $E(st) = 0$, $E(t^2) = E(s^2) = 1$. Then $(\delta^{(0)})^2 = \|\Delta^{(0)} - \Delta^*\|^2 = (1 - \alpha)^2 + \bar{\alpha}^2$. For $\mu = \Delta^* - \Delta^{(1)}$,

$$\begin{aligned} \mu \cdot (\Delta^* - \Delta^{(1)}) &= \mu D_+^{-1} \Gamma_+ \Delta^* + \mu D_-^{-1} \Gamma_- \Delta^* \\ &= \mu^+ \Gamma_+ \Delta^* + \mu^- \Gamma_- \Delta^* \\ &= \int_{S_{+-}} (\mu^+ \cdot x) s dP + \int_{S_{-+}} (\mu^- \cdot x) s dP. \end{aligned}$$

Take $\mu^+ = \|\mu^+\| e^+$, $\|e^+\| = 1$, and similarly for μ^- and note that $\|\mu^+\| \leq \|\mu\|/\lambda_+$, $\|\mu^-\| \leq \|\mu\|/\lambda_-$ where λ_+ , λ_- are the minimum eigenvalues of Γ_+ , Γ_- . Thus

$$\delta^{(1)} = \|\Delta^* - \Delta^{(1)}\| \leq \frac{1}{\lambda_+} \int_{S_{+-}} |v^+ s| dP + \frac{1}{\lambda_-} \int_{S_{-+}} |v^- s| dP$$

where $v^+ = e^+ \cdot x$, $v^- = e^- \cdot x$. Assume that $\delta^{(0)} < 1$, then $\alpha > 0$. S_{-+} is the set $s \geq 0$, $\alpha s + \bar{\alpha} t < 0$, which combines into $0 \leq s < -(\bar{\alpha}/\alpha)t$. S_{+-} is the set $0 > s \geq -(\bar{\alpha}/\alpha)t$. Letting $\gamma = |\bar{\alpha}/\alpha|$,

$$\delta^{(1)} \leq \frac{\gamma}{\lambda_+} \int_{S_{+-}} |v^+ t| dP + \frac{\gamma}{\lambda_-} \int_{S_{-+}} |v^- t| dP.$$

Define

$$\theta(\gamma) = \sup_{v^+, v^-, t} \max \left[\int_{S_{+-}} |v^+ t| dP, \int_{S_{-+}} |v^- t| dP \right],$$

then $\theta(\gamma) \leq 1$, and $\theta(\gamma) = O(\gamma^{2/5})$ as $\gamma \rightarrow 0$. To verify these statements, note that

$$\int_T |vt| dP \leq \sqrt{E v^2 E t^2} = 1,$$

and

$$\int_T |vt| dP \leq (P(T))^{1/2} (E v^4 E t^4)^{1/4}.$$

Now $\max(P(S_{+-}), P(S_{-+})) \leq P(|s| \leq \gamma |t|)$, and the latter is dominated by

$p(\gamma b) + c/b^4$. Taking $b = \gamma^{-1/5}$ makes this bound equal to

$$q(\gamma) = p(\gamma^{4/5}) + c\gamma^{4/5}.$$

Now

$$\begin{aligned} P(\Delta^{(0)} \cdot x \geq 0) &= P(\Delta^* \cdot x \geq (\Delta^* - \Delta^{(0)}) \cdot x) \\ &\geq P(\Delta^* \cdot x \geq 0) - P(|\Delta^* \cdot x| \leq |(\Delta^* - \Delta^{(0)}) \cdot x|). \end{aligned}$$

The last term is bounded by

$$P(|\Delta^* \cdot x| \leq \alpha) + \frac{c}{\alpha^4} \delta^{(0)4}.$$

With $\alpha = (\delta^{(0)})^{4/5}$,

$$P(\Delta^{(0)} \cdot x \geq 0) \geq P(\Delta^* \cdot x \geq 0) - q(\delta^{(0)}).$$

Let $c^* = P(\Delta^* \cdot x > 0)$, then

$$\delta^{(1)} \leq \gamma \theta(\gamma) \left[\frac{1}{\lambda(c^* - q(\delta^{(0)}))} + \frac{1}{\lambda(1 - c^* - q(\delta^{(0)}))} \right].$$

Noting that $\gamma \leq \delta^{(0)} / \sqrt{1 - (\delta^{(0)})^2}$ gives the result that

$$\delta^{(1)} \leq \delta^{(0)} \phi(\delta^{(0)})$$

where $\phi(\delta^{(0)}) = O((\delta^{(0)})^{2/5})$.

Proof of Theorem 2.2.

W.l.o.g. we can take x_1, \dots, x_M i.i.d. $N(0, 1)$ such that $(\Delta^* \cdot x) = x_1$. For any $\Delta^{(0)}$, let $s = (\Delta^{(0)} \cdot x)$. Then

$$\Gamma_+ = \Gamma_- = I/2$$

and

$$\Delta^{(1)}(m) = \Delta^*(m) - 2\{E(x_m x_1; x_1 \geq 0, s < 0) + E(x_m x_1; x_1 < 0, s > 0)\}.$$

Let $\|\Delta^{(0)}\|_1 = \sqrt{\|\Delta^{(0)}\|^2 - \Delta^{(0)}(1)^2}$, and $t^{(0)} = \Delta^{(0)}(1) / \|\Delta^{(0)}\|_1$. By simple manipulations

$$\begin{aligned} \Delta^{(1)}(1) &= \theta(t^{(0)}) \\ \Delta^{(1)}(m) &= \frac{\Delta^{(0)}(m)}{\|\Delta^{(0)}\|_1} \chi(t^{(0)}), \quad m \geq 2 \end{aligned}$$

where

$$\theta(t) = 1 - 4 \int_0^\infty x^2 \Phi(-tx) f(x) dx$$

$$\chi(t) = \frac{2}{\pi} \cdot \frac{1}{1+t^2},$$

$f(x)$ is the standard normal density and $\Phi(x) = \int_{-\infty}^x f(y) dy$. Therefore

$$t^{(1)} = \phi(t^{(0)})$$

where

$$\phi(t) = \frac{\pi}{2} (1+t^2) \theta(t).$$

Now $\theta'(t) = (4/\pi)(1+t^2)^{-2}$, so for $t > 0$

$$\theta(t) > \frac{4}{\pi t} \int_0^t \frac{\xi d\xi}{(1+\xi^2)^2} = \frac{2t}{\pi(1+t^2)}.$$

In consequence, $\phi(t) > t$, all $t > 0$, and $t^{(k)} \rightarrow \infty$ unless $t^{(0)} = 0$, so $\Delta^{(k)}(1) \rightarrow 1$ and $\Delta^{(k)}(m) \rightarrow 0$, $m \geq 1$.

Proof of Theorem 3.1.

The key is a lemma which we reproduce from Barron [1991].

Lemma. *If f is in the closure of the convex hull of a set G in a Hilbert space, with $\|g\| \leq b$, all $g \in G$, then for every $K \geq 1$ and every $c' > b^2 - \|f\|^2$ there is an f_K in the convex hull of K points in G such that $\|f - f_K\|^2 \leq c'/K$.*

Denote by $h^+(x)$, the x^+ function (i.e., $h^+(x) = x$, $x \geq 0$, else $= 0$). Every hinge function is of the form

$$\beta_0 + \beta \cdot x + h^+(\Delta_0 + \Delta \cdot x).$$

The idea of the proof is to show that

$$f(x) - x \cdot \nabla f(0) - f(0)$$

is in the closure of the convex hull of a set of functions $\{h^+(\Delta_0 + \Delta \cdot x)\}$ with such that each function in the set has norm bounded by $c(2R)^2$.

Begin by noting that

$$\begin{aligned} g(x) &= f(x) - x \cdot \nabla f(0) - f(0) = R \int (e^{iw \cdot x} - iw \cdot x - 1) \tilde{f}(w) dw \\ &= R \int (e^{iw \cdot x} - iw \cdot x - 1) e^{i\theta(w)} |\tilde{f}(w)| dw \\ &= \int [\cos(w \cdot x + \theta(w)) + w \cdot x \sin \theta(w) - \cos \theta(w)] |\tilde{f}(w)| dw. \end{aligned}$$

Write $Q(dw) = \|w\|^2 |\tilde{f}(w)|/c$, so

$$g(x) = c \int [\cos(w \cdot x + \theta(w)) + w \cdot x \sin \theta(w) - \cos \theta(w)] \frac{Q(dw)}{\|w\|^2}$$

and $g(\mathbf{x})$ is in the closure of the convex hull of the functions

$$s_w(\mathbf{x}) = \frac{c [\cos(w \cdot \mathbf{x} + \theta(w)) + w \cdot \mathbf{x} \sin \theta - \cos \theta]}{\|w\|^2}.$$

Let $w = \|w\| \alpha$, $b = \theta(w)$, $y = \alpha \cdot \mathbf{x}$, and

$$\phi(y) = \frac{c [\cos(\|w\| y + b) - \|w\| y \sin b - \cos b]}{\|w\|^2},$$

so $s_w(\mathbf{x}) = \phi(y)$, $\|\alpha\| = 1$, $y \in [-R, R]$, $\phi(0) = 0$, and $\phi''(y)$ is continuous on $[-R, R]$. Consider fitting $\phi(y)$ at the points j/J , $j = 0, \dots, \text{int}(RJ)$ with a function of the form

$$s_J(y) = \sum_0^{\text{Int}(RJ)} \alpha_i h^+(y - i/J).$$

Then $\sum_1^{j-1} \alpha_i = J [\phi(j/J) - \phi(j-1/J)]$, and α_i is J times the second difference of the sequence $\phi(j/J)$ evaluated at $j = i$. Thus, as $J \rightarrow \infty$,

$$\sum_i |\alpha_i| \rightarrow \int_0^R |\phi''(y)| dy.$$

A similar argument holds for approximating $\phi(y)$ on $[-R, 0]$, with sums of functions of the form $h^+(-y - i/J)$. Thus, with

$$\int_{-R}^{+R} |\phi''(y)| dy \leq 2cR$$

$\phi(y)$ is in the closure of the convex hull of the set of functions $\{h^+(\gamma(y - \tau))\}$ with $|\gamma| \leq 2cR$ and $\tau \in [-R, R]$. This implies that $g(\mathbf{x})$ is in the closure of the convex hull of the set of functions $\{h^+(\gamma(\alpha \cdot \mathbf{x} - \tau))\}$ with $\|\alpha\| = 1$. Since

$$\int [h^+(\gamma(\alpha \cdot \mathbf{x} - \tau))]^2 dP \leq \gamma^2 \int (\alpha \cdot \mathbf{x} - \tau)^2 dP \leq (2R)^2 \gamma^2,$$

using the lemma completes the proof.

APPENDIX II

Computational Aspects of Hinge Fitting

There are devices used in the hinge programs both to improve accuracy and improve computing speed. We mention these to assist other researchers in building their own hinge programs.

The basic building block is the hinge finding algorithm. In this, two devices are noted: i) if, in the iterations, the hinge gets too close to the edge of the data, then it is pushed back into the data. More precisely, if the number of points on one side of the hinge (say $\Delta x \leq 0$) falls below a threshold value NH then $\Delta(0)$ is readjusted so that $\#\{x_n; \Delta \cdot x_n \leq 0\} = NH$.

The major computational burden in the algorithm is usually in the recomputation of $S_+ = (X^t X)_+$, S_- after a new hinge has been selected. To reduce this, the sign of the new hinge is taken so that the number of points in S_+ is $\leq \text{int}(N/2)$, S_+ is updated only for those points that have changed sides, and S_- set equal to $S - S_+$.

In the refitting algorithm, the last hinge Δ_k used in a hinge function h_k is stored. Any refitting of h_k starts from Δ_k and the new hinge is then stored in place of Δ_k . In the refitting algorithm, the best performance has been gotten by using one iteration of the hinge finding algorithm per hinge function per cycle. Often, there are many cycles until the reduction in RSS levels off.

In stepwise variable selection the problem was to construct a computationally feasible procedure when dimensionality is high — say over 50 variables. The device used is this: suppose variables x_1, \dots, x_L have been entered into the current hinge. What is stored is an indicator vector $\text{id}(n)$ which tells which side of the current hinge each x_n is on, the matrices S_+^{-1} , S_-^{-1} , and the values of $x_l \cdot x_m$ for all $l \leq L$ and $m > L$. For a variable x_m , $m > L$, the indicator vector id is used to compute $(x_l \cdot x_m)_+$ for all $l \leq L$, $(x_m \cdot x_m)_+$, $(x_m \cdot y)_+$ and the corresponding $()_-$ values. Then a fast update formula uses these values to produce the inverses \tilde{S}_+^{-1} , \tilde{S}_-^{-1} , where the tilde indicates that the m^{th} variable values are included in the \tilde{S}_+ , \tilde{S}_- matrices.

Then hyperplanes based on (x_1, \dots, x_L, x_m) are fitted on either side of the current hinge, the new hinge function computed along with the new residual-sum-of squares RSS_m . The variable added is that with minimum RSS_m . In actuality, computing RSS_m requires only a partial update of S_+^{-1} , S_-^{-1} and after m is selected, then the full update to \tilde{S}_+^{-1} , \tilde{S}_-^{-1} is done. The final step is to run the hinge finding algorithm to convergence using (x_1, \dots, x_L, x_m) and store as starting values for the next round the final \tilde{S}_+^{-1} , \tilde{S}_-^{-1} and $\tilde{\text{id}}$.

The stepwise forward procedure in classification adds an extra step. After the update to \tilde{S}_+ , \tilde{S}_- and $\tilde{Z}_+ = (x\tilde{y})_+$, $\tilde{Z}_- = (x\tilde{y})_-$, form the updated \tilde{H} matrix (sec 5.2), and

solve the $J \times J$ eigenvalue problem to find the new hinge and RSS value. In the L variable case, either a $J \times J$ or $L \times L$ eigenvalue provides the new hinge. If J is large, it will be faster to solve the sequence of $L \times L$ problems.

APPENDIX III

Cross-Validation Measures

An important question in the use of hinges is: how many hinges to use? If step-wise entry of variables is used, another question is how many variables to enter into the current hinge. If large test sets are available, they can be used to answer these questions. If not, then we use cross-validation measures as our basic tool.

Number of Hinges Used in Regression

One approach is to use 10-fold cross-validation (see Breiman et. al [1984]). However, this increases the computations by a factor of ten. Another approach is to use cross-validation measures based on the leave-one-out method. That is, the case \mathbf{x}_n is left out of the fitting, coefficients derived and then \mathbf{x}_n used as a single case test set. This is repeated for each case in the data and then the mean taken of the sums-of-squares of test set errors. The resulting measure is called the cross-validated prediction error, or PE_{CV} . In ordinary regression, its not difficult to show that

$$PE_{CV} = \frac{1}{N} \sum_1^N (r_n / 1 - h_n)^2$$

where r_n is the usual residual and h_n is the n^{th} diagonal element of the matrix $H_{nn} = \mathbf{x}_n^t S^{-1} \mathbf{x}_n$, where S^{-1} is the inverse of the matrix $\sum_n \mathbf{x}_{nn} \mathbf{x}_{m'n}$. The trace of H is the number of variables M in the equation. If h_n is replaced by its average $\sum_n h_n / N$ in the PE_{CV} expression, the resulting approximation is called the generalized cross-validation estimation of the prediction error and is given by

$$PE_{GCV} = MRSS / (1 - M/N)^2.$$

In fitting hinge functions, recall that the fit is the sum of a linear function and a function of the type $h^+(\Delta \cdot \mathbf{x})$. That is, RSS is the minimum of

$$\sum_n (y_n - \beta \cdot \mathbf{x}_n - \sum_{k=1}^K h^+(\Delta_k \cdot \mathbf{x}_n))^2. \quad (A.3.1)$$

Denote by β , Δ_k the minimizers of (A.3.1), and $\{r_n\}$ the residuals. In the neighborhood of β , Δ_k , say $\beta' = \beta + \epsilon$, $\Delta_k' = \Delta_k + \delta_k$, (A.3.1) is approximated by

$$\sum_n (r_n - \epsilon \cdot \mathbf{x}_n - \sum_{k=1}^K \delta_k \cdot \mathbf{x}_n I(\Delta_k \cdot \mathbf{x}_n \geq 0))^2 \quad (A.3.2)$$

where $I(\cdot)$ is the indicator function.

Define the $(M + 1)(K + 1)$ dimensional variable ξ by

$$\xi = \{x, xI(\Delta_1 \cdot x \geq 0), \dots, xI(\Delta_K \cdot x \geq 0)\}.$$

Then letting $\alpha = (\epsilon, \delta_1, \dots, \delta_K)$, (A.3.2) has the form

$$\sum_n (r_n - \alpha \cdot \xi_n)^2. \quad (\text{A.3.3})$$

If the case numbered n' is deleted and the deleted sum-of-squares in (A.3.1) minimized, the resulting β' , Δ_k' will be close to β , Δ_k . Therefore, the deletion and cross-validation measures for the nonlinear (A.3.1) can be approximated by the same measures for the linearized (A.3.3). Thus, to 1st order, the PE_{GCV} for the hinge fitting is given by

$$PE_{GCV} = MRSS / (1 - N_p / N)^2 \quad (\text{A.3.4})$$

where $N_p = (K + 1)(M + 1)$.

Because of the replacement of h_n by \bar{h} , the expression for PE_{GCV} given by (A.3.4) is usually biased low. The simplest way to correct this bias is to put

$$PE_{GCV} = MRSS (1 - cN_p / N)^2$$

where $c \in [1, 3]$. We have used $c = 1.5$ throughout with fairly accurate results, although the choice of the number of hinges selected is fairly insensitive to the choice of c .

Number of Variables Used in a Hinge (Regression)

Suppose that some hinges have already been fitted, the current residuals are $\{r_n\}$ and stepwise entry of variables into the next hinge function is being carried out to fit the $\{r_n\}$. Consider fitting with a single hinge function of M_0 variables. Then from the previous derivation leading to (A.3.4), we get

$$PE_{GCV} = MRSS / (1 - N_p / N)^2$$

where $N_p = 2(M_0 + 1)$.

In the stepwise selection of variables context, there is an additional reason why this expression is biased low. In brief, the cross-validation measure used above does not cross-validate the variable selection process (see Breiman [1990], Breiman and Spector [1990]). Again, the estimate

$$PE_{GCV} = MRSS / (1 - cN_p / N)^2$$

is used ($c = 1.5$ throughout), and variables are added until this expression starts to increase.

Number of Hinges Used in Classification

An expression can be developed for the leave-one-out miss-classification rate, but it is complex and difficult to evaluate. For K hinges and M variables it requires the inversion of a matrix of order $K(M + 1) \times K(M + 1)$. For 99 variables with 5 hinges, this is a 500×500 matrix. If forward stepwise entry of variables is used, the matrix is sparse, and efficient methods may be found to invert it. In large problems a test set provides a rapid method for determining how many hinges to use. If the sample size is not large enough for a test set, we recommend 10-fold cross-validation.

Number of Variables Used in a Hinge (Classification)

Unless the test set is large, using it to both determine the number of hinges and the number of variables in each hinge may be overfitting the test set. The procedure outlined below has proven a satisfactory alternative in the examples we have run. It assumes that the "true misclassification rate" is decreasing as long as the "true regression error" is decreasing.

Suppose there are currently M_0 variables entered into the K^{th} hinge, and denote the residuals after the $(K - 1)^{\text{st}}$ hinge was fitted by $\{r_{jn}\}$. The PE_{GCV} when M_0 variables are used will be computed and variables added until the PE_{GCV} stops decreasing.

Let the coefficients $\{c_j\}$ and hinge function $h(\mathbf{x}) = \beta \cdot \mathbf{x} + h^+(\Delta \cdot \mathbf{x})$ based on M_0 variables minimize

$$\sum_{j,n} (r_{jn} - c_j h(\mathbf{x}_n))^2. \quad \text{A.3.5}$$

Denote $z_{jn} = r_{jn} - c_j h(\mathbf{x}_n)$. Then $\mathbf{z}_j \perp \mathbf{h}$, $j = 1, \dots, J$; and if $\boldsymbol{\mu} = \sum_j c_j \mathbf{z}_j$, then $\boldsymbol{\mu} \perp \mathbf{x}_m$, $m = 0, \dots, M_0$. Further, $\boldsymbol{\mu} \perp \mathbf{t}_m$, $m = 0, \dots, M_0$ where $\mathbf{t}_m = \mathbf{x}_m I(\Delta \cdot \mathbf{x} > 0)$. Let $\beta' = \beta + \epsilon$, $\Delta' = \Delta + \delta$, and $c_j' = c_j + \sigma_j$. Then the linearized version of (A.3.5) is

$$\sum_{j,n} (z_{jn} - \sigma_j h(\mathbf{x}_n) - c_j (\epsilon \cdot \mathbf{x}_n + \delta \cdot \mathbf{t}_n))^2. \quad \text{(A.3.6)}$$

Define the $2(M_0 + 1)$ dimensional variable $\xi_n = (\mathbf{x}_n, \mathbf{t}_n)$, and $\mathbf{s} = (\epsilon, \delta)$ so that A.3.6 becomes

$$\sum_{j,n} (z_{jn} - \sigma_j h(\mathbf{x}_n) - c_j \mathbf{s} \cdot \xi_n)^2. \quad \text{A.3.7}$$

If the n' case is deleted, then the RSS is

$$\sum_{j,n} (z_{jn} - \sigma_j h(\mathbf{x}_n) - c_j \mathbf{s} \cdot \xi_n)^2 - \sum_j (z_{jn'} - \sigma_j h(\mathbf{x}_{n'}) - c_j \mathbf{s} \cdot \xi_{n'})^2. \quad \text{(A.3.8)}$$

Because of the orthogonality relations, the 1st term is

$$\sum_{j,n} z_{jn}^2 + \sum_{j,n} (\sigma_j h(\mathbf{x}_n) + c_j \mathbf{s} \cdot \xi_n)^2.$$

If the 2nd term in A.3.8 is expanded only to 1st order terms in σ_j , s , and we assume $\sum c_j^2 = 1$, then partial derivatives with respect to σ_j and s give the equations

$$(A) \quad \sigma_j \|h\|^2 + c_j s (\xi, h) = -z_{jn'} h_{n'}$$

$$(B) \quad (\sum_j \sigma_j c_j) (\xi, h) + Ss = -u_{n'} \xi_{n'}$$

where $u_{n'} = \sum_j c_j z_{jn'}$, and $S_{mm'} = (\xi_m, \xi_{m'})$.

Equation (A) gives σ_j in terms of s . The solution for s (to first order) is as follows: define a matrix A by

$$A_{mm'} = S_{mm'} - \frac{(\xi_m, h) (\xi_{m'}, h)}{\|h\|^2}$$

and a vector v_n by

$$v_n = (\xi_n - \frac{(\xi, h) h_n}{\|h\|^2}).$$

Then

$$s = -u_{n'} A^{-1} v_{n'},$$

and

$$PE_{CV} \equiv \sum_{j,n} (z_{jn} - \sigma_{jn} h(x_n) - c_j s_n \cdot \xi_n)^2 \quad A.3.9$$

where the subscripts n on σ_{jn} , s_n indicate that these are the values of σ_j , s gotten by minimizing the RSS with the n^{th} case deleted. Let $RSS = \sum_j (z_j, z_j)$. To terms of order RSS/N , A.3.9 equals

$$RSS + \sum_n u_n^2 v_n^t A^{-1} v_n.$$

To approximate the 2nd term, we replace u_n^2 by \bar{u}^2 . Noting that $\sum_n v_n^t v_n = A$, the approximation equals $2(M_0 + 1) \bar{u}^2$. Therefore

$$PE_{GCV} \equiv MRSS + 2(M_0 + 1) \bar{u}^2 / N.$$

Analysis of 2nd order terms shows that a slightly more accurate value is

$$PE_{GCV} \equiv MRSS - \bar{u}^2 + \bar{u}^2 / (1 - N_p / N)^2$$

with $N_p = 2(M_0 + 1)$. As previously, we correct this for downward bias by using

$$PE_{GCV} = MRSS - \bar{u}^2 + \bar{u}^2 / (1 - cN_p / N)^2.$$

In two examples, the evidence indicates that the more selection is going on (i.e. the higher the dimensionality) the larger c should be. In 21 dimensions, $c = 1.5$ was better than $c = 2$. In 61 dimensions, the best value of c is between 2.0 and 3.0. But in both

cases, the test set error changes only slightly over a wide range of c .