

Variable Selection and Model Building via Likelihood Basis Pursuit

Hao Helen ZHANG, Grace WAHBA, Yi LIN, Meta VOELKER, Michael FERRIS, Ronald KLEIN, and Barbara KLEIN

This article presents a nonparametric penalized likelihood approach for variable selection and model building, called likelihood basis pursuit (LBP). In the setting of a tensor product reproducing kernel Hilbert space, we decompose the log-likelihood into the sum of different functional components such as main effects and interactions, with each component represented by appropriate basis functions. Basis functions are chosen to be compatible with variable selection and model building in the context of a smoothing spline ANOVA model. Basis pursuit is applied to obtain the optimal decomposition in terms of having the smallest l_1 norm on the coefficients. We use the functional L_1 norm to measure the importance of each component and determine the “threshold” value by a sequential Monte Carlo bootstrap test algorithm. As a generalized LASSO-type method, LBP produces shrinkage estimates for the coefficients, which greatly facilitates the variable selection process and provides highly interpretable multivariate functional estimates at the same time. To choose the regularization parameters appearing in the LBP models, generalized approximate cross-validation (GACV) is derived as a tuning criterion. To make GACV widely applicable to large datasets, its randomized version is proposed as well. A technique “slice modeling” is used to solve the optimization problem and makes the computation more efficient. LBP has great potential for a wide range of research and application areas such as medical studies, and in this article we apply it to two large ongoing epidemiologic studies, the Wisconsin Epidemiologic Study of Diabetic Retinopathy (WESDR) and the Beaver Dam Eye Study (BDES).

KEY WORDS: Generalized approximate cross-validation; LASSO; Monte Carlo bootstrap test; Nonparametric variable selection; Slice modeling; Smoothing spline ANOVA.

1. INTRODUCTION

Variable selection, or dimension reduction, is fundamental to multivariate statistical model building. Not only does judicious variable selection improve the model’s predictive ability, it also generally provides a better understanding of the underlying concept that generates the data. Due to recent proliferation of large, high-dimensional databases, variable selection has become the focus of intensive research in several areas such as text processing, environmental sciences, and genomics, particularly gene expression array data involving tens or hundreds of thousands of variables.

Traditional variable selection approaches, such as stepwise selection and best subset selection, are built in linear regression models, and the well-known criteria like Mallows’s C_p , the Akaike information criterion (AIC) and the Bayesian information criterion (BIC) are often used to penalize the number of nonzero parameters (see Linhart and Zucchini 1986 for an introduction). To achieve better prediction and reduce the variances of estimators, many shrinkage estimation approaches have been proposed. Bridge regression was introduced by Frank and Friedman (1993), which is a constrained least squares method subject to an L_p penalty with $p \geq 1$. Two special cases of bridge regression are the LASSO proposed by Tibshirani (1996) when $p = 1$ and the ridge regression when $p = 2$. Due

to the nature of the L_1 penalty, the LASSO tends to shrink small coefficients to 0 and hence gives concise models. It also exhibits the stability of ridge regression estimates. Fu (1998) made a thorough comparison between the bridge model and the LASSO. Knight and Fu (2000) proved some asymptotic results for LASSO-type estimators. In the case of wavelet regression, this L_1 penalty approach is called “basis pursuit.” Chen, Donoho, and Saunders (1998) discussed atomic decomposition by basis pursuit in some detail. (A related development can be found in Bakin 1999.) Gunn and Kandola (2002) proposed a structural modeling approach with sparse kernels. Recently, Fan and Li (2001) suggested a nonconcave penalized likelihood approach with the smoothly clipped absolute deviation (SCAD) penalty function, which resulted in an unbiased, sparse, and continuous estimator. Our motivation of this study is to provide a flexible nonparametric alternative to the parametric approaches for variable selection as well as model building. Yau, Kohn, and Wood (2003) presented a Bayesian method for variable selection in a nonparametric manner.

Smoothing spline analysis of variance (SS-ANOVA) provides a general framework for nonparametric multivariate function estimation and has been studied intensively for Gaussian data. Wahba, Wang, Gu, Klein, and Klein (1995) gave a general setting for applying the SS-ANOVA model to exponential families. Gu (2002) provided a comprehensive review of the SS-ANOVA and some recent progress. In this work we have developed a unified model that appropriately combines the SS-ANOVA model and basis pursuit for variable selection and model building.

The article is organized as follows. Section 2 introduces the notation and illustrates the general structure of the likelihood basis pursuit (LBP) model. We focus on the main-effects model and the two-factor interaction model, then generalize the models to incorporate categorical variables. Section 3 discusses the important issue of adaptively choosing regularization

Hao Helen Zhang is Assistant Professor, Department of Statistics, North Carolina State University, Raleigh, NC 27695 (E-mail: hzhang2@stat.ncsu.edu). Grace Wahba is IJ Schoenberg and Bascom Professor (E-mail: wahba@stat.wisc.edu) and Yi Lin is Associate Professor (E-mail: yilin@stat.wisc.edu), Department of Statistics, and Michael Ferris is Professor (E-mail: ferris@cs.wisc.edu), Department of Computer Sciences and Industrial Engineering, University of Wisconsin-Madison, Madison, WI 53706. Meta Voelker is Senior Analyst, Alphatech Inc., Arlington, VA 22203 (E-mail: meta.voelker@dc.alphatech.com). Ronald Klein is Professor (E-mail: kleinr@epi.ophth.wisc.edu) and Barbara Klein is Professor (E-mail: kleinB@epi.ophth.wisc.edu), Department of Ophthalmology, Medical School, University of Wisconsin-Madison, Madison, WI 53726. This work was supported in part by National Science Foundation grants DMS-00-72292, DMS-01-34987, DMS-04-05913, and CCR-9972372; National Institutes of Health grants EY09946 and EY03083; and AFOSR grant F49620-01-1-0040. The authors thank the editor, the associate editor, and the two referees for their constructive comments and suggestions that have led to significant improvement of this article.

parameters. An extension of generalized approximate cross-validation (GACV) proposed by Xiang and Wahba (1996) is derived as a tuning criterion. Section 4 proposes the measure of importance for the variables and, if desired, their interactions. We develop sequential Monte Carlo bootstrap test algorithm to determine the selection threshold. Section 5 covers the numerical computation issue. Sections 6–8 present several simulation examples and the applications of LBP to two large epidemiologic studies. We carry out a data analysis for the 4-year risk of progression of diabetic retinopathy in the Wisconsin Epidemiologic Study of Diabetic Retinopathy (WESDR) and for the 5-year risk of mortality in the Beaver Dam Eye Study (BDES). The final section contains some concluding remarks. Proofs are relegated to Appendixes A and B.

2. LIKELIHOOD BASIS PURSUIT

2.1 Smoothing Spline ANOVA for Exponential Families

We are interested in estimating the dependence of Y on the covariates $\mathbf{X} = (X^{(1)}, \dots, X^{(d)})$. Typically, \mathbf{X} is in a high-dimensional space $\mathcal{X} = \mathcal{X}^{(1)} \otimes \dots \otimes \mathcal{X}^{(d)}$, where $\mathcal{X}^{(\alpha)}$, $\alpha = 1, \dots, d$, is some measurable space and \otimes denotes the tensor product operation. Conditioning on \mathbf{x} , assume that Y is from an exponential family with density of form $h(y|\mathbf{x}) = \exp\{yf(\mathbf{x}) - b(f(\mathbf{x}))/a(\phi) + c(y, \phi)\}$, where $a > 0$, b , and c are known functions, $f(\mathbf{x})$ is the parameter of interest dependent on \mathbf{x} , and ϕ is either known or a nuisance parameter independent of \mathbf{x} . We denote the observations of $Y_i|\mathbf{x}_i \sim h(y|\mathbf{x}_i)$, $i = 1, \dots, n$, by the vector $\mathbf{y} = (y_1, \dots, y_n)'$. The scaled conditional log-likelihood is

$$\begin{aligned} \mathcal{L}(\mathbf{y}, f) &= \frac{1}{n} \sum_{i=1}^n [-l\{y_i, f(\mathbf{x}_i)\}] \\ &\equiv \frac{1}{n} \sum_{i=1}^n [-y_i f(\mathbf{x}_i) + b\{f(\mathbf{x}_i)\}]. \end{aligned} \quad (1)$$

Although the methodology proposed here is general and valid for any exponential family, we use the Bernoulli case as our working example. In Bernoulli data, Y takes on values $\{0, 1\}$ with the conditional probability $p(\mathbf{x}) \equiv \Pr(Y = 1|\mathbf{X} = \mathbf{x})$. The logit function $f(\mathbf{x}) = \log(\frac{p(\mathbf{x})}{1-p(\mathbf{x})})$, and the log-likelihood $l(y, f) = yf - \log(1 + e^f)$. Many parametric approaches, such as those of Tibshirani (1996), Fu (1998), and Fan and Li (2001), assume $f(\mathbf{x})$ to be a linear function of \mathbf{x} . Instead, we allow f to vary in a high-dimensional function space, which leads to a more flexible estimate for the target function. In this section and Section 2.2, we assume that all of the covariates are continuous. Later in Section 2.3, we take categorical variables into account. Similar to the classical ANOVA, for any function $f(\mathbf{x}) = f(x^{(1)}, \dots, x^{(d)})$ on a product domain \mathcal{X} , we can define its functional ANOVA decomposition as

$$\begin{aligned} f(\mathbf{x}) &= b_0 + \sum_{\alpha=1}^d f_{\alpha}(x^{(\alpha)}) \\ &+ \sum_{\alpha < \beta} f_{\alpha\beta}(x^{(\alpha)}, x^{(\beta)}) \\ &+ \text{all higher-order interactions,} \end{aligned} \quad (2)$$

where b_0 is constant, f_{α} 's are the main effects, and $f_{\alpha\beta}$'s are the two-factor interactions. The identifiability of the terms is ensured by side conditions through averaging operators. We are to estimate each f_{α} in an reproducing kernel Hilbert space (RKHS) $\mathcal{H}^{(\alpha)}$, each $f_{\alpha\beta}$ in the tensor product space $\mathcal{H}^{(\alpha)} \otimes \mathcal{H}^{(\beta)}$, and so on. The full model space is then the d th-order tensor product space $\bigotimes_{\alpha=1}^d \mathcal{H}^{(\alpha)}$. Each functional component in the decomposition (2) falls in the corresponding subspace of $\bigotimes_{\alpha=1}^d \mathcal{H}^{(\alpha)}$.

For any continuous covariate $x^{(\alpha)}$, we scale it onto $[0, 1]$ and choose $\mathcal{H}^{(\alpha)}$ to be the second-order Sobolev space $W_2[0, 1]$, which is defined as $\{g: g(t), g'(t) \text{ are absolutely continuous, } g''(t) \in \mathcal{L}_2[0, 1]\}$. When endowed with a certain inner product, $W_2[0, 1]$ is an RKHS with the reproducing kernel (RK) $1 + K_0(s, t) + K_1(s, t)$. Here $K_0(s, t) = k_1(s)k_1(t)$ and $K_1(s, t) = k_2(s)k_2(t) - k_4(|s - t|)$, with $k_1(t) = t - \frac{1}{2}$, $k_2(t) = \frac{1}{2}(k_1^2(t) - \frac{1}{12})$, and $k_4(t) = \frac{1}{24}(k_1^4(t) - \frac{1}{2}k_1^2(t) + \frac{7}{240})$. Notice that $\mathcal{H}^{(\alpha)} = [1] \oplus \mathcal{H}_{\pi}^{(\alpha)} \oplus \mathcal{H}_s^{(\alpha)}$, with $[1]$ the constant subspace, $\mathcal{H}_{\pi}^{(\alpha)}$ the "parametric" subspace generated by K_0 consisting of linear functions, and $\mathcal{H}_s^{(\alpha)}$ the "nonparametric" subspace generated by K_1 consisting of smooth functions. The reproducing kernel of $\bigotimes_{\alpha=1}^d \mathcal{H}^{(\alpha)}$ is

$$\prod_{\alpha=1}^d (1 + k_1(s^{(\alpha)})k_1(t^{(\alpha)}) + K_1(s^{(\alpha)}, t^{(\alpha)})). \quad (3)$$

Correspondingly, $\bigotimes_{\alpha=1}^d \mathcal{H}^{(\alpha)}$ can be decomposed into the tensor sum of parametric main-effect subspaces, and smooth main-effect subspaces, and two-factor interaction subspaces of three possible forms: parametric \otimes parametric, smooth \otimes parametric, and smooth \otimes smooth, and similarly for three-factor or higher interaction subspaces. The RKs for these subspaces are given by the corresponding terms in the expansion of (3) and, the terms in (2) can be expanded in terms corresponding to those RKs in the expansion of (3). Therefore, our model encompasses the linear model as a special case (see Wahba 1990 for more details). In various situations, the ANOVA decomposition in (2) and in the expansion of (3) is truncated at some point. In this work we consider truncation for the continuous variables no later than after the two-factor interactions. The remaining RKs will be used to construct an overcomplete set of basis functions for the likelihood basis pursuit, via details given later.

2.2 Likelihood Basis Pursuit Models

Basis pursuit is a principle for decomposing a signal into an optimal superposition of dictionary elements, where "optimal" means having the smallest l_1 norm of the coefficients among all such decompositions. Chen et al. (1998) illustrated atomic decomposition by basis pursuit in wavelet regression. In this article we apply basis pursuit to the negative log-likelihood in the context of a dictionary based on the SS-ANOVA decomposition, then select the important components using the multivariate function estimates. Let \mathcal{H} be the model space after truncation. The variational problem for the LBP model is

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n [-l(y_i, f(\mathbf{x}_i))] + J_{\lambda}(f), \quad (4)$$

where $J_\lambda(f)$ denotes the l_1 norm of the coefficients of the basis functions in the representation of f . It is a generalization of the LASSO penalty in the context of nonparametric models. The l_1 penalty often produces coefficients that are exactly 0 and therefore gives sparse solutions. This sparsity helps distinguish important variables from unimportant ones easily and more effectively. (See Tibshirani 1996 and Fan and Li 2001 for comparison of the l_1 penalty with other forms of penalty.) The regularization parameter λ balances the fitness in the likelihood and the penalty part.

For the usual smoothing spline modeling, the penalty J_λ is a quadratic norm or seminorm in an RKHS. Kimeldorf and Wahba (1971) showed that the minimizer f_λ of the traditional smoothing spline model falls in a finite-dimensional space, although the model space is of infinite dimension. For the penalized likelihood approach with a nonquadratic penalty like the l_1 penalty, it is very hard to obtain analytic solutions. In light of the results for the quadratic penalty situation, we propose using a sufficiently large number of basis functions to span the model space and estimate the target function in that space. If all of the data $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ are used to generate the bases, the resulting functional space demands intensive computation, and the application is limited for large-scale problems. Thus we adopt the parsimonious bases approach used by Xiang and Wahba (1998), Ruppert and Carroll (2000), Lin et al. (2000), and Yau et al. (2003). Gu and Kim (2001) showed that the number of basis terms can be much smaller than n without degrading the performance of the estimation. For $N \leq n$, we subsample N points $\{\mathbf{x}_{1*}, \dots, \mathbf{x}_{N*}\}$ from the original data and use them to generate basis functions that then span the model space \mathcal{H}_* . Notice that we are not wasting any data resource here, because all of the data points are used for model fitting, although only a subset of them are selected to generate basis functions.

The issue of choosing N and the subsamples is important. In practice, we generally start with a reasonably large N . It is well known that “reasonably large” is not actually very large (see Lin et al. 2000). In principle, the subspace spanned by the chosen basis terms needs to be rich enough to provide a decent fit to the true curve. In this article we use the simple random sampling technique to choose the subsamples. Alternatively, a cluster algorithm may be used, as done by Xiang and Wahba (1998) and Yau et al. (2003). Its basic idea involves two steps. First, we group the data into N clusters that have maximum separation, by some good algorithm. Then within each cluster we randomly choose one data point as a representative to be included in the base pool. This scheme usually provides well-separated subsamples.

2.2.1 Main-Effects Model. The main-effects model, also known as the additive model, is a sum of d functions of one variable. The function space is the tensor sum of constant and the main effect subspaces of $\bigotimes_{\alpha=1}^d \mathcal{H}^{(\alpha)}$. Define $\mathcal{H}_* = [1] \bigoplus_{\alpha=1}^d \text{span}\{k_1(x^{(\alpha)}), K_1(x^{(\alpha)}, x_{j*}^{(\alpha)}), j = 1, \dots, N\} \equiv [1] \bigoplus_{\alpha=1}^d \mathcal{H}_*^{(\alpha)}$, where $k_1(\cdot)$ and $K_1(\cdot, \cdot)$ are as defined in Section 2.1. Then any component function $f_\alpha \in \mathcal{H}_*^{(\alpha)}$ has the representation $f_\alpha(x^{(\alpha)}) = b_\alpha k_1(x^{(\alpha)}) + \sum_{j=1}^N c_{\alpha,j} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)})$. The LBP estimate $f \in \mathcal{H}_*$ is obtained by minimizing

$$\frac{1}{n} \sum_{i=1}^n [-l(y_i, f(\mathbf{x}_i))] + \lambda_\pi \sum_{\alpha=1}^d |b_\alpha| + \lambda_s \sum_{\alpha=1}^d \sum_{j=1}^N |c_{\alpha,j}|, \quad (5)$$

where $f(\mathbf{x}) = b_0 + \sum_{\alpha=1}^d f_\alpha(x^{(\alpha)})$ and (λ_π, λ_s) are the regularization parameters. Here and in the sequel we have chosen to group terms of similar types (here “parametric” and “smooth”) and to allow distinct λ ’s for different groups. Of course, we could choose to set $\lambda_\pi = \lambda_s$. Note that $\lambda_s = \infty$ yields a parametric model. Alternatively, we could choose different λ ’s for each coefficient if we decide to incorporate some prior knowledge of certain variables or their effects.

2.2.2 Two-Factor Interaction Model. Two-factor interactions arise in many practical problems. (See Hastie and Tibshirani 1990, sec. 9.5.5, or Lin et al. 2000, figs. 9 and 10, for interpretable plots of two-factor interactions with continuous variables.) In the LBP model, the two-factor interaction space consists of the “parametric” part and the “smooth” part. The parametric part is generated by d parametric main-effect terms and $\frac{d(d-1)}{2}$ parametric–parametric interaction terms. The smooth part is the tensor sum of the subspaces generated by smooth main-effect terms, parametric–smooth interaction terms, and smooth–smooth interaction terms. For each pair $\alpha \neq \beta$, the two-factor interaction subspace is $\mathcal{H}_*^{(\alpha\beta)} = \text{span}\{k_1(x^{(\alpha)})k_1(x^{(\beta)}), K_1(x^{(\alpha)}, x_{j*}^{(\alpha)})k_1(x^{(\beta)}), K_1(x^{(\alpha)}, x_{j*}^{(\alpha)})K_1(x^{(\beta)}, x_{j*}^{(\beta)}), j = 1, \dots, N\}$, and the interaction term $f_{\alpha\beta}(x^{(\alpha)}, x^{(\beta)})$ has the representation

$$\begin{aligned} f_{\alpha\beta} &= b_{\alpha\beta} k_1(x^{(\alpha)})k_1(x^{(\beta)}) \\ &+ \sum_{j=1}^N c_{\alpha\beta,j}^{\pi s} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)})k_1(x^{(\beta)}) \\ &+ \sum_{j=1}^N c_{\alpha\beta,j}^{ss} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)})K_1(x^{(\beta)}, x_{j*}^{(\beta)}). \end{aligned}$$

The whole function space is $\mathcal{H}_* \equiv [1] \bigoplus_{\alpha=1}^d \mathcal{H}_*^{(\alpha)} + \bigoplus_{\alpha < \beta} \mathcal{H}_*^{(\alpha\beta)}$. The LBP optimization problem is

$$\begin{aligned} \min_{f \in \mathcal{H}_*} \frac{1}{n} \sum_{i=1}^n [-l(y_i, f(\mathbf{x}_i))] &+ \lambda_\pi \left(\sum_{\alpha=1}^d |b_\alpha| \right) \\ &+ \lambda_{\pi\pi} \left(\sum_{\alpha < \beta} |b_{\alpha\beta}| \right) + \lambda_{\pi s} \left(\sum_{\alpha \neq \beta} \sum_{j=1}^N |c_{\alpha\beta,j}^{\pi s}| \right) \\ &+ \lambda_s \left(\sum_{\alpha=1}^d \sum_{j=1}^N |c_{\alpha,j}| \right) + \lambda_{ss} \left(\sum_{\alpha < \beta} \sum_{j=1}^N |c_{\alpha\beta,j}^{ss}| \right). \quad (6) \end{aligned}$$

Note that different penalties are allowed for the five different types of terms.

2.3 Incorporating Categorical Variables

In real applications, some covariates may be categorical, such as sex, race, and smoking history in many medical studies. In previous sections we proposed the main-effects model (5) and the two-factor interaction model (6) for continuous variables only. Now we generalize these models to incorporate categorical variables, which are denoted by $\mathbf{Z} = (Z^{(1)}, \dots, Z^{(r)})$. For simplicity, we assume that each $Z^{(\gamma)}$ has two categories $\{T, F\}$ for $\gamma = 1, \dots, r$. The following idea is easily extended to the

situation when some variables have more than two categories. Define the mapping Φ_γ by

$$\Phi_\gamma(z^{(\gamma)}) = \begin{cases} \frac{1}{2} & \text{if } z^{(\gamma)} = T \\ -\frac{1}{2} & \text{if } z^{(\gamma)} = F. \end{cases}$$

Generally the mapping is chosen to make the range of categorical variables comparable with that of continuous variables. For any variable with $C > 2$ categories, $C - 1$ contrasts are needed.

The main-effects model that incorporates the categorical variables is as follows:

Minimize

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n [-l(y_i, f(\mathbf{x}_i, \mathbf{z}_i))] + \lambda_\pi \left(\sum_{\alpha=1}^d |b_\alpha| + \sum_{\gamma=1}^r |B_\gamma| \right) \\ + \lambda_s \sum_{\alpha=1}^d \sum_{j=1}^N |c_{\alpha,j}|, \quad (7) \end{aligned}$$

where $f(\mathbf{x}, \mathbf{z}) = b_0 + \sum_{\alpha=1}^d b_\alpha k_1(x^{(\alpha)}) + \sum_{\gamma=1}^r B_\gamma \times \Phi_\gamma(z^{(\gamma)}) + \sum_{\alpha=1}^d \sum_{j=1}^N c_{\alpha,j} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)})$. For each γ , the function Φ_γ can be considered as the main effect of the covariate $Z^{(\gamma)}$. Thus we choose to associate the coefficients $|B|$'s and $|b|$'s with the same parameter λ_π .

Adding two-factor interactions with categorical variables to a model that already includes parametric and smooth terms adds a number of additional terms to the general model. Compared with (6), four new types of terms are involved when we take into account categorical variables: categorical main effects, categorical–categorical interactions, “parametric continuous”–categorical interactions, and “smooth continuous”–categorical interactions. The modified two-factor interaction model is as follows:

Minimize

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n [-l(y_i, f(\mathbf{x}_i, \mathbf{z}_i))] + \lambda_\pi \left(\sum_{\alpha=1}^d |b_\alpha| + \sum_{\gamma=1}^r |B_\gamma| \right) \\ + \lambda_{\pi\pi} \left(\sum_{\alpha < \beta} |b_{\alpha\beta}| + \sum_{\gamma < \theta} |B_{\gamma\theta}| + \sum_{\alpha=1}^d \sum_{\gamma=1}^r |P_{\alpha\gamma}| \right) \\ + \lambda_{\pi s} \left(\sum_{\alpha \neq \beta} \sum_{j=1}^N |c_{\alpha\beta,j}^{\pi s}| + \sum_{\alpha=1}^d \sum_{\gamma=1}^r \sum_{j=1}^N |c_{\alpha\gamma,j}^{\pi s}| \right) \\ + \lambda_s \left(\sum_{\alpha=1}^d \sum_{j=1}^N |c_{\alpha,j}| \right) + \lambda_{ss} \left(\sum_{\alpha < \beta} \sum_{j=1}^N |c_{\alpha\beta,j}^{ss}| \right), \quad (8) \end{aligned}$$

where

$$\begin{aligned} f(\mathbf{x}, \mathbf{z}) = b_0 + \sum_{\alpha=1}^d b_\alpha k_1(x^{(\alpha)}) + \sum_{\gamma=1}^r B_\gamma \Phi_\gamma(z^{(\gamma)}) \\ + \sum_{\alpha < \beta} b_{\alpha\beta} k_1(x^{(\alpha)}) k_1(x^{(\beta)}) \\ + \sum_{\gamma < \theta} B_{\gamma\theta} \Phi_\gamma(z^{(\gamma)}) \Phi_\theta(z^{(\theta)}) \end{aligned}$$

$$\begin{aligned} + \sum_{\alpha=1}^d \sum_{\gamma=1}^r P_{\alpha\gamma} k_1(x^{(\alpha)}) \Phi_\gamma(z^{(\gamma)}) \\ + \sum_{\alpha \neq \beta} \sum_{j=1}^N c_{\alpha\beta,j}^{\pi s} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)}) k_1(x^{(\beta)}) k_1(x_{j*}^{(\beta)}) \\ + \sum_{\alpha=1}^d \sum_{\gamma=1}^r \sum_{j=1}^N c_{\alpha\gamma,j}^{\pi s} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)}) \Phi_\gamma(z^{(\gamma)}) \\ + \sum_{\alpha=1}^d \sum_{j=1}^N c_{\alpha,j} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)}) \\ + \sum_{\alpha < \beta} \sum_{j=1}^N c_{\alpha\beta,j}^{ss} K_1(x^{(\alpha)}, x_{j*}^{(\alpha)}) K_1(x^{(\beta)}, x_{j*}^{(\beta)}). \end{aligned}$$

We assign different regularization parameters for main-effect terms, parametric–parametric interaction terms, parametric–smooth interaction terms, and smooth–smooth interaction terms. In particular, the coefficients $|B_{\gamma\theta}|$'s and $|P_{\alpha\gamma}|$'s are associated with the same parameter $\lambda_{\pi\pi}$, and the coefficients $|c_{\alpha\gamma,j}^{\pi s}|$'s and $|c_{\alpha\beta,j}^{\pi s}|$'s are associated with the same parameter $\lambda_{\pi s}$.

3. GENERALIZED APPROXIMATE CROSS-VALIDATION

Regularization parameter selection has been a very active field of research. Ordinary cross-validation (OCV) (Wahba and Wold 1975), generalized cross-validation (GCV) (Craven and Wahba 1979), and GACV (Xiang and Wahba 1996) are widely used in various contexts of smoothing spline models. We derive the GACV to select λ 's in the LBP models. Here we present the GACV score and its derivation only for the main-effects model; the extension to high-order interaction models is straightforward. With an abuse of notation, we use λ to represent the collective set of tuning parameters; in particular, $\lambda = (\lambda_\pi, \lambda_s)$ for the main-effects model and $\lambda = (\lambda_\pi, \lambda_{\pi\pi}, \lambda_{\pi s}, \lambda_s, \lambda_{ss})$ for the two-factor interaction model.

3.1 Approximate Cross-Validation

Let p be the “true” but unknown probability function and let p_λ be its estimate associated with λ . Similarly, let f and μ denote the true logit and mean functions, and let f_λ and μ_λ denote the corresponding estimates. Kullback–Leibler (KL) distance, also known as the relative entropy, is often used to measure the distance between two probability distributions. For Bernoulli data, we have $KL(p, p_\lambda) = E_{\mathbf{X}}[\frac{1}{2}\{\mu(f - f_\lambda) - (b(f) - b(f_\lambda))\}]$, with $b(f) = \log(1 + e^f)$. Removing the quantity that does not depend on λ from the KL distance expression, we get the comparative KL (CKL) distance $E_{\mathbf{X}}[-\mu f_\lambda + b(f_\lambda)]$. The ordinary leave-one-out cross validation (CV) for CKL is

$$CV(\lambda) = \frac{1}{n} \sum_{i=1}^n [-y_i f_\lambda^{[-i]}(\mathbf{x}_i) + b(f_\lambda(\mathbf{x}_i))], \quad (9)$$

where $f_\lambda^{[-i]}$ is the minimizer of the objective function with the i th data point omitted. CV is commonly used as a roughly unbiased estimate for CKL (see Xiang and Wahba 1996;

Gu 2002). Direct calculation of CV involves computing n leave-one-out estimates, which is expensive and almost infeasible for large-scale problems. Thus we derive a second-order approximate cross-validation (ACV) score to the CV. (See Xiang and Wahba 1996; Lin et al. 2000; Gao, Wahba, Klein, and Klein 2001 for the ACV in traditional smoothing spline models.) We first establish the leave-one-out lemma for the LBP models. The proof of Lemma 1 is given in Appendix A.

Lemma 1 (Leave-One-Out Lemma). Denote the objective function in the LBP model by

$$I_\lambda(f, \mathbf{y}) = \mathcal{L}(f, \mathbf{y}) + J_\lambda(f). \quad (10)$$

Let $f_\lambda^{[-i]}$ be the minimizer of $I_\lambda(f, \mathbf{y})$ with the i th observation omitted and let $\mu_\lambda^{[-i]}(\cdot)$ be the mean function corresponding to $f_\lambda^{[-i]}(\cdot)$. For any $v \in \mathbb{R}$, we define the vector $\mathbf{V} = (y_1, \dots, y_{i-1}, v, y_{i+1}, \dots, y_n)'$. Let $h_\lambda(i, v, \cdot)$ be the minimizer of $I_\lambda(f, \mathbf{V})$; then $h_\lambda(i, \mu_\lambda^{[-i]}(\mathbf{x}_i), \cdot) = f_\lambda^{[-i]}(\cdot)$.

Using Taylor series approximations and Lemma 1, we can derive (in App. B) the ACV score,

$$\begin{aligned} ACV(\lambda) = & \frac{1}{n} \sum_{i=1}^n [-y_i f_\lambda(\mathbf{x}_i) + b(f_\lambda(\mathbf{x}_i))] \\ & + \frac{1}{n} \sum_{i=1}^n h_{ii} \frac{y_i(y_i - \mu_\lambda(\mathbf{x}_i))}{1 - \sigma_{\lambda i}^2 h_{ii}}, \end{aligned} \quad (11)$$

where $\sigma_{\lambda i}^2 \equiv p_\lambda(\mathbf{x}_i)(1 - p_\lambda(\mathbf{x}_i))$ and h_{ii} is the ii th entry of the matrix \mathbf{H} defined in (B.9) (more details have been given in Zhang 2002). Let \mathbf{W} be the diagonal matrix with $\sigma_{\lambda i}^2$ in the ii th position. By replacing h_{ii} with $\frac{1}{n} \sum_{i=1}^n h_{ii} \equiv \frac{1}{n} \text{tr}(\mathbf{H})$ and replacing $1 - \sigma_{\lambda i}^2 h_{ii}$ with $\frac{1}{n} \text{tr}[\mathbf{I} - (\mathbf{W}^{1/2} \mathbf{H} \mathbf{W}^{1/2})]$ in (11), we obtain the GACV score,

$$\begin{aligned} GACV(\lambda) = & \frac{1}{n} \sum_{i=1}^n [-y_i f_\lambda(\mathbf{x}_i) + b(f_\lambda(\mathbf{x}_i))] \\ & + \frac{\text{tr}(\mathbf{H})}{n} \frac{\sum_{i=1}^n y_i(y_i - \mu_\lambda(\mathbf{x}_i))}{\text{tr}[\mathbf{I} - \mathbf{W}^{1/2} \mathbf{H} \mathbf{W}^{1/2}]}. \end{aligned} \quad (12)$$

3.2 Randomized GACV

Direct computation of (12) involves the inversion of a large-scale matrix, whose size depends on sample size n , basis size N , and dimension d . Large values of n , N , or d may make the computation expensive and produce unstable solutions. Thus the randomized GACV (ranGACV) score is proposed as a computable proxy for GACV. Essentially, we use the randomized trace estimates for $\text{tr}(\mathbf{H})$ and $\text{tr}[\mathbf{I} - \frac{1}{2}(\mathbf{W}^{1/2} \mathbf{H} \mathbf{W}^{1/2})]$ based on the following theorem (which has been exploited by numerous authors, e.g., Girard 1998):

If \mathbf{A} is any square matrix and $\boldsymbol{\epsilon}$ is a mean 0 random n -vector with independent components with variance σ_ϵ^2 , then $\frac{1}{\sigma_\epsilon^2} E \boldsymbol{\epsilon}^T \mathbf{A} \boldsymbol{\epsilon} = \text{tr}(\mathbf{A})$.

Let $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)'$ be a mean 0 random n -vector of independent components with variance σ_ϵ^2 . Let $f_\lambda^{\mathbf{y}}$ and $f_\lambda^{\mathbf{y}+\boldsymbol{\epsilon}}$ be the minimizer of (5) using the original data \mathbf{y} and the perturbed

data $\mathbf{y} + \boldsymbol{\epsilon}$. Using the derivation procedure of Lin et al. (2000), we get the ranGACV score for the LBP estimates,

$$\begin{aligned} \text{ranGACV}(\lambda) = & \frac{1}{n} \sum_{i=1}^n [-y_i f_\lambda(\mathbf{x}_i) + b(f_\lambda(\mathbf{x}_i))] \\ & + \frac{\boldsymbol{\epsilon}^T (f_\lambda^{\mathbf{y}+\boldsymbol{\epsilon}} - f_\lambda^{\mathbf{y}})}{n} \frac{\sum_{i=1}^n y_i (y_i - \mu_\lambda(\mathbf{x}_i))}{\boldsymbol{\epsilon}^T \boldsymbol{\epsilon} - \boldsymbol{\epsilon}^T \mathbf{W} (f_\lambda^{\mathbf{y}+\boldsymbol{\epsilon}} - f_\lambda^{\mathbf{y}})}. \end{aligned} \quad (13)$$

In addition, two ideas help reduce the variance of the second term in (13):

1. Choose $\boldsymbol{\epsilon}$ as Bernoulli (.5), taking values in $\{+\sigma_\epsilon, -\sigma_\epsilon\}$. This guarantees that the randomized trace estimate has the minimal variance given a fixed σ_ϵ^2 (see Hutchinson 1989).
2. Generate U independent perturbations $\boldsymbol{\epsilon}^{(u)}$, $u = 1, \dots, U$, and compute U replicated ranGACVs. Then use their average to compute the GACV estimate.

4. SELECTION CRITERIA FOR MAIN-EFFECTS AND TWO-FACTOR INTERACTIONS

4.1 The L_1 Importance Measure

After choosing the optimal $\hat{\lambda}$ by the GACV or ranGACV criteria, the LBP estimate $f_{\hat{\lambda}}$ is obtained by minimizing (5), (6), (7), or (8). How to measure the importance of a particular component in the fitted model is a key question. We propose using the functional L_1 norm as the importance measure. The sparsity in the solutions will help distinguish the significant terms from insignificant ones effectively, and thus improve the performance of our importance measure. In practice, for each functional component, its L_1 norm is empirically calculated as the average of the function values evaluated at all of the data points, for instance, $L_1(f_\alpha) = \frac{1}{n} \sum_{i=1}^n |f_\alpha(x_i^{(\alpha)})|$ and $L_1(f_{\alpha\beta}) = \frac{1}{n} \sum_{i=1}^n |f_{\alpha\beta}(x_i^{(\alpha)}, x_i^{(\beta)})|$. For the categorical variables in model (7), the empirical L_1 norm of the main effect f_γ is $L_1(f_\gamma) = \frac{1}{n} \sum_{i=1}^n |B_\gamma \Phi_\gamma(z_i^{(\gamma)})|$ for $\gamma = 1, \dots, r$. The norms of the interaction terms involved with categorical variables are defined similarly. The rank of the L_1 norm scores is used to rank the relative importance of functional components. For instance, the component with the largest L_1 norm is ranked as the most important one, and any component with a zero or tiny L_1 norm is ranked as an unimportant one. We have also tried using the functional L_2 norm to rank the components; this gave almost identical results in terms of the set of variables selected in numerous simulation studies (not reproduced here).

4.2 Choosing the Threshold

In this section we focus on the main-effects model. Using the chosen parameter $\hat{\lambda}$, we obtain the estimated main-effects components $\hat{f}_1, \dots, \hat{f}_d$ and calculate their L_1 norms $L_1(\hat{f}_1), \dots, L_1(\hat{f}_d)$. We use a sequential procedure to select important terms. Denote the decreasingly ordered norms by $\hat{L}_{(1)}, \dots, \hat{L}_{(d)}$ and the corresponding components by $\hat{f}_{(1)}, \dots, \hat{f}_{(d)}$. A universal threshold value is needed to differentiate the important components from unimportant ones. Call the threshold q . Only variables with their L_1 norms greater than or equal to q are “important.”

Now we develop a sequential Monte Carlo bootstrap test procedure to determine q . Essentially we test the variables' importance one by one in their L_1 norm rank order. If one variable passes the test (hence "important"), then it enters the null model for testing the next variable; otherwise, the procedure stops. After the first η ($0 \leq \eta \leq d - 1$) variables enter the model, it is a one-sided hypothesis-testing problem to decide whether the next component $\hat{f}_{(\eta+1)}$ is important or not. When $\eta = 0$, the null model f is the constant, say $f = \hat{b}_0$, and the hypotheses are $H_0: L_{(1)} = 0$ versus $H_1: L_{(1)} > 0$. When $1 \leq \eta \leq d - 1$, the null model is $f = \hat{b}_0 + \hat{f}_{(1)} + \dots + \hat{f}_{(\eta)}$, and the hypotheses are $H_0: L_{(\eta+1)} = 0$ versus $H_1: L_{(\eta+1)} > 0$. Let the desired one-sided test level be α . If the null distribution of $\hat{L}_{(\eta+1)}$ were known, then we could get the critical value α -percentile and make a decision of rejection or acceptance. In practice, calculating the exact α -percentile is difficult or impossible; however, the Monte Carlo bootstrap test provides a convenient approximation to the full test. Conditional on the original covariates $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, we generate $\{y_1^{*(\eta)}, \dots, y_n^{*(\eta)}\}$ (response 0 or 1) using the null model $f = \hat{b}_0 + \hat{f}_{(1)} + \dots + \hat{f}_{(\eta)}$ as the true logit function. We sample a total of T independent sets of data $(\mathbf{x}_1, y_{1,t}^{*(\eta)}), \dots, (\mathbf{x}_n, y_{n,t}^{*(\eta)})$, $t = 1, \dots, T$, from the null model f , fit the main-effects model for each set, and compute $\hat{L}_t^{*(\eta+1)}$, $t = 1, \dots, T$. If exactly k of the simulated $\hat{L}_t^{*(\eta+1)}$ values exceed $\hat{L}_{(\eta+1)}$ and none equals it, then the Monte Carlo p value is $\frac{k+1}{T+1}$. (See Davison and Hinkley 1997 for an introduction to the Monte Carlo bootstrap test.)

Sequential Monte Carlo Bootstrap Test Algorithm

- Step 1. Let $\eta = 0$ and $f = \hat{b}_0$. Test $H_0: L_{(1)} = 0$ versus $H_1: L_{(1)} > 0$. Generate T independent sets of data, $(\mathbf{x}_1, y_{1,t}^{*(0)}), \dots, (\mathbf{x}_n, y_{n,t}^{*(0)})$, $t = 1, \dots, T$, from $f = \hat{b}_0$. Fit the LBP main-effects model and compute the Monte Carlo p value p_0 . If $p_0 < \alpha$, then go to step 2; otherwise stop and define q as any number slightly larger than $\hat{L}_{(1)}$.
- Step 2. Let $\eta = \eta + 1$ and $f = \hat{b}_0 + \hat{f}_{(1)} + \dots + \hat{f}_{(\eta)}$. Test $H_0: L_{(\eta+1)} = 0$ versus $H_1: L_{(\eta+1)} > 0$. Generate T independent sets of data $(\mathbf{x}_1, y_{1,t}^{*(\eta)}), \dots, (\mathbf{x}_n, y_{n,t}^{*(\eta)})$ based on f , fit the main-effects model, and compute the Monte Carlo p value, p_η . If $p_\eta < \alpha$ and $\eta < d - 1$, then repeat step 2, and if $p_\eta < \alpha$ and $\eta = d - 1$, then go to step 3; otherwise, stop and define $q = \hat{L}_{(\eta)}$.
- Step 3. Stop the procedure and define $q = \hat{L}_{(d)}$.

The order of entry for sequential testing of the terms being entertained for the model is determined by the magnitude of the component L_1 norms. There are other reasonable ways to determine the order of entry, and the particular strategy used can affect the results, as is the case for any stepwise procedure. For the LBP approach, the relative ranking among the important terms usually does not affect the final component selection solution as long as the important terms are all ranked higher than the unimportant terms. Thus our procedure usually can distinguish between important and unimportant terms, as in most of our examples. When the distinction between important and

unimportant terms is ambiguous with our method, the ambiguous terms can be recognized, and further investigation will be needed on these terms.

5. NUMERICAL COMPUTATION

Because the objective function in either (5) or (6) is not differentiable with respect to the coefficients b 's and c 's, some numerical methods for optimization fail to solve this kind of problem. We can replace the l_1 norms in the objective function by nonnegative variables constrained linearly to be the corresponding absolute values using standard mathematical programming techniques, and then solve a series of programs with nonlinear objective functions and linear constraints. Many optimization methods can solve such problems. We use MINOS (Murtagh and Saunders 1983), because it generally performs well with the linearly constrained models and returns consistent results.

Consider choosing the optimal λ by selecting the grid point for which ranGACV achieves a minimum. For each λ , to find ranGACV , the program (5), (6), (7), or (8) must be solved twice—once with \mathbf{y} (the original problem) and once with $\mathbf{y} + \epsilon$ (the perturbed problem). This often results in hundreds or thousands of individual solves, depending on the range for λ . To obtain solutions in a reasonable amount of time, we use an efficient solving approach known as slice modeling (see Ferris and Voelker 2000, 2001). Slice modeling is an approach to solving a series of mathematical programs with the same structure but different data. Because for LBP we can consider the values (λ, \mathbf{y}) to be individual slices of data (because only these values change between solves), the program can be reduced to an example of nonlinear slice modeling. By applying slice modeling ideas (namely, maintaining the common program structure and "core" data shared between solves and using previous solutions as starting points for late solves), we can improve efficiency and make the grid search feasible. We have developed efficient code for the main-effects and two-way interaction LBP models. The code is easy to use and runs fast. For example, in one simulation example with $n = 1,000$, $d = 10$, and λ fixed, the main-effects model takes less than 2 seconds, and the two-way interaction model takes less than 50 seconds.

6. SIMULATION

6.1 Simulation 1: Main-Effects Model

In this example there are a total of $d = 10$ covariates: $X^{(1)}, \dots, X^{(10)}$. They are taken to be uniformly distributed in $[0, 1]$ independently. The sample size $n = 1,000$. We use the simple random subsampling technique to select $N = 50$ basis functions. The perturbation ϵ is distributed as Bernoulli(.5) taking two values $\{+.25, -.25\}$. Four variables, $X^{(1)}$, $X^{(3)}$, $X^{(6)}$, and $X^{(8)}$, are important, and the others are noise variables. The true conditional logit function is

$$f(\mathbf{x}) = \frac{4}{3}x^{(1)} + \pi \sin(\pi x^{(3)}) + 8(x^{(6)})^5 + \frac{2}{e-1}e^{x^{(8)}} - 5.$$

We fit the main-effects LBP model and search the parameters (λ_π, λ_s) globally. Because the true f is known, both CKL and ranGACV can be used for choosing the λ 's. Figure 1 depicts the values of $CKL(\lambda)$ and $\text{ranGACV}(\lambda)$ as functions of (λ_π, λ_s)

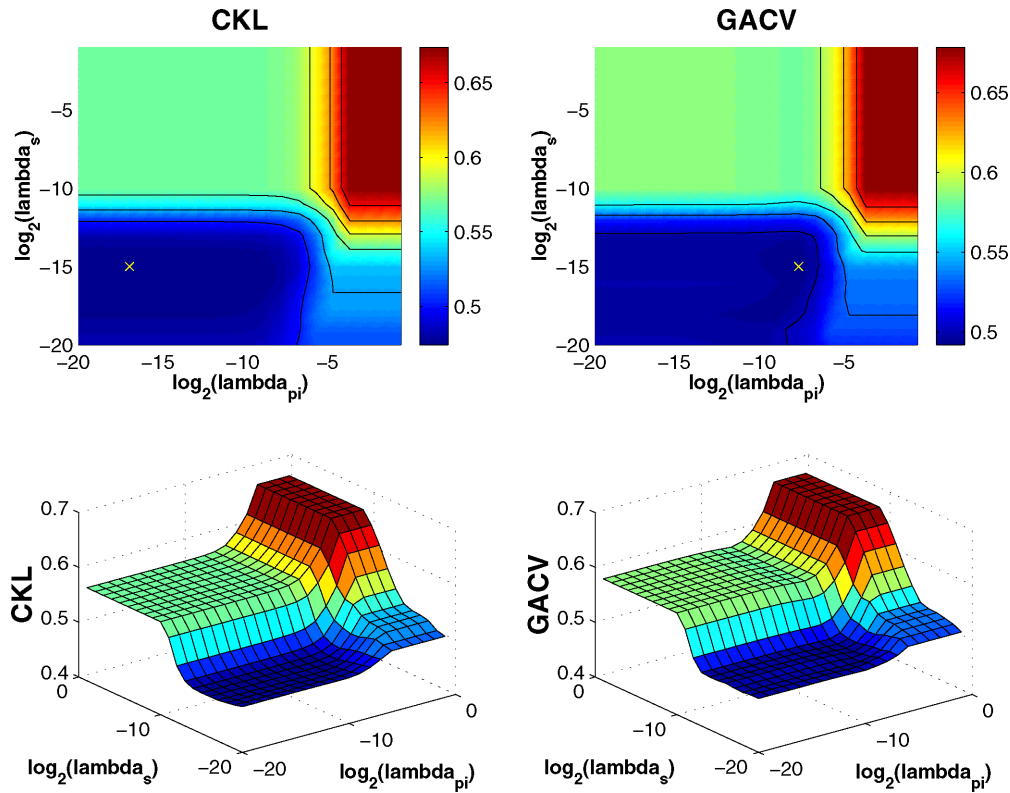


Figure 1. Contours and Three-Dimensional Plots for $CKL(\lambda)$ and $GACV(\lambda)$.

within the region of interest $[2^{-20}, 2^{-1}] \times [2^{-20}, 2^{-1}]$. In the top row are the contours for $CKL(\lambda)$ and $ranGACV(\lambda)$, where the white cross “x” indicates the location of the optimal regularization parameter. Here $\hat{\lambda}_{CKL} = (2^{-17}, 2^{-15})$ and $\hat{\lambda}_{ranGACV} = (2^{-8}, 2^{-15})$. The bottom row shows their three-dimensional plots. In general, $ranGACV(\lambda)$ approximates $CKL(\lambda)$ quite well globally.

Using the optimal parameters we fit the main-effects model and calculate the L_1 norm scores for the individual components $\hat{f}_1, \dots, \hat{f}_{10}$. Figure 2 plots two sets of L_1 norm scores, obtained using $\hat{\lambda}_{CKL}$ and $\hat{\lambda}_{ranGACV}$, in their decreasing order. The dashed line indicates the threshold chosen by the proposed sequential Monte Carlo bootstrap test algorithm. Using this threshold,

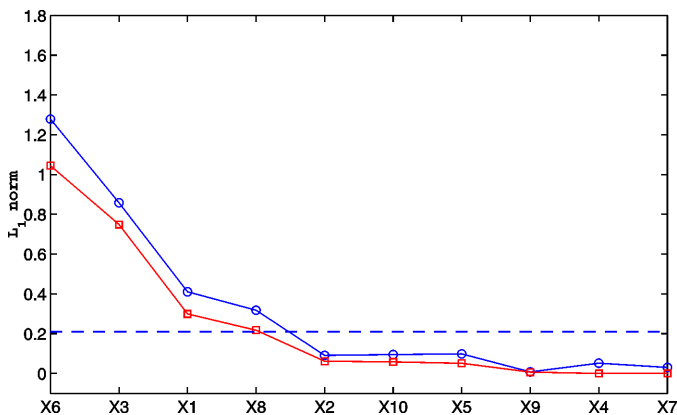


Figure 2. L_1 Norm Scores for the Main-Effects Model (—○— CKL ; —■— $GACV$).

variables $X^{(6)}$, $X^{(3)}$, $X^{(1)}$, and $X^{(8)}$ are selected as “important” variables correctly.

Figure 3 depicts the procedure of the sequential bootstrap tests to determine q . We fit the main-effects model using $\hat{\lambda}_{ranGACV}$ and sequentially test the hypotheses $H_0: L_{(\eta)} = 0$ versus $H_1: L_{(\eta)} > 0$, $\eta = 1, \dots, 10$. In each plot of Figure 3, gray circles denote the L_1 norms for the variables in the null model,

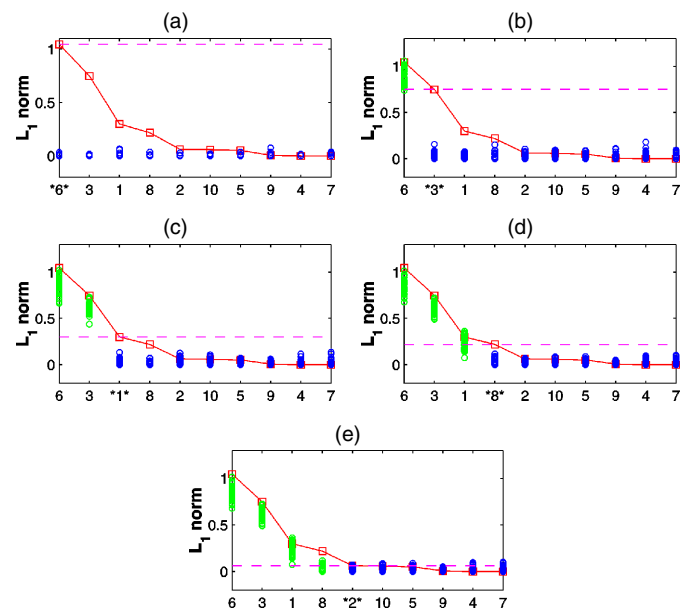


Figure 3. Monte Carlo Bootstrap Tests for Simulation 1. [(a) —■— original, ○ boot 1; (b) —■— original, ○ boot 2; (c) —■— original, ○ boot 3; (d) —■— original, ○ boot 4; (e) —■— original, ○ boot 5.]

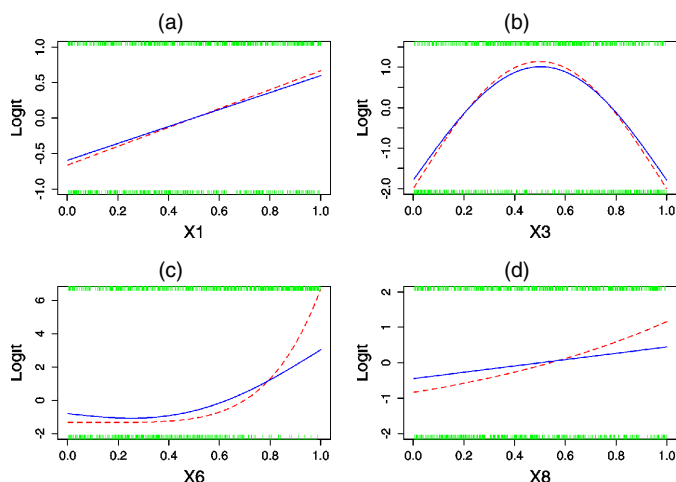


Figure 4. True (---) and Estimated (—) Univariate Logit Components: (a) X_1 ; (b) X_3 ; (c) X_6 ; (d) X_8 .

and black circles denote the L_1 norms for those not in the null model. (In a color plot, gray circles are shown in green and black circles are shown in blue). Along the horizontal axis, the variable being tested for importance is bracketed by a pair of *. Our experiment shows that the null hypotheses of the first four tests are all rejected at level $\alpha = .05$ based on their Monte Carlo p value $1/51 \doteq .02$. However, the null hypothesis for the fifth component f_2 is accepted with the p value $10/51 \doteq .20$. Thus f_6 , f_3 , f_1 , and f_8 are selected as “important” components and $q = L_{(4)} = .21$.

In addition to selecting important variables, LBP also produces functional estimates for the individual components in the model. Figure 4 plots the true main effects f_1 , f_3 , f_6 , and f_8 and their estimates fitted using $\hat{\lambda}_{\text{ranGACV}}$. In each panel, the dashed line denoted the true curve and the solid line denotes the corresponding estimate. In general, the fitted main-effects model provides a reasonably good estimate for each important component. Altogether we generated 20 datasets and fitted the main-effects model for each dataset with regularization parameters tuned separately. Throughout all of these 20 runs, variables $X^{(1)}$, $X^{(3)}$, $X^{(6)}$, and $X^{(8)}$ are always the four top-ranked variables. The results and figures shown here are based on the first dataset.

6.2 Simulation 2: Two-Factor Interaction Model

There are $d = 4$ continuous covariates, independently and uniformly distributed in $[0, 1]$. The true model is a two-factor interaction model, and the important effects are $X^{(1)}$, $X^{(2)}$, and their interaction. The true logit function f is

$$f(\mathbf{x}) = 4x^{(1)} + \pi \sin(\pi x^{(1)}) + 6x^{(2)} - 8(x^{(2)})^3 + \cos(2\pi(x^{(1)} - x^{(2)})) - 4.$$

We choose $n = 1,000$ and $N = 50$, and use the same perturbation ϵ as in the previous example. There are five tuning parameters (λ_π , $\lambda_{\pi\pi}$, λ_s , $\lambda_{\pi s}$, and λ_{ss}) in the two-factor interaction model. In practice, extra constraints may be added on the parameters for different needs. Here we penalize all of the two-factor interaction terms equally by setting $\lambda_{\pi\pi} = \lambda_{\pi s} = \lambda_{ss}$. The optimal parameters are $\hat{\lambda}_{\text{CKL}} = (2^{-10}, 2^{-10}, 2^{-15},$

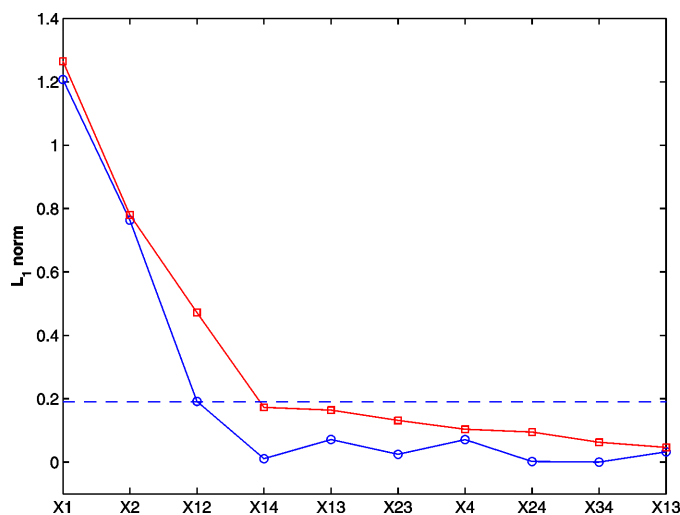


Figure 5. L_1 Norm Scores for the Two-Factor Interaction Model (—○ CKL; —■ GACV).

$2^{-10}, 2^{-10})$ and $\hat{\lambda}_{\text{ranGACV}} = (2^{-20}, 2^{-20}, 2^{-18}, 2^{-20}, 2^{-20})$. Figure 5 plots the ranked L_1 norm scores. The dashed line represents the threshold q chosen by the Monte Carlo bootstrap test procedure. The LBP two-factor interaction model, using either $\hat{\lambda}_{\text{CKL}}$ or $\hat{\lambda}_{\text{GACV}}$, selects all of the important effects $X^{(1)}$, $X^{(2)}$, and their interaction effect correctly.

There is a strong interaction effect between variables $X^{(1)}$ and $X^{(2)}$, which is shown clearly by the cross-sectional plots in Figure 6. The solid lines are the cross-sections of the true logit function $f(x^{(1)}, x^{(2)})$ at distinct values $x^{(1)} = .2, .5, .8$, and the dashed lines are their corresponding estimates given by the LBP model. The parameters are tuned by the ranGACV criterion.

6.3 Simulation 3: Main-Effects Model Incorporating Categorical Variables

In this example, among the 12 covariates, $X^{(1)}, \dots, X^{(10)}$ are continuous and $Z^{(1)}$ and $Z^{(2)}$ are categorical. The continuous

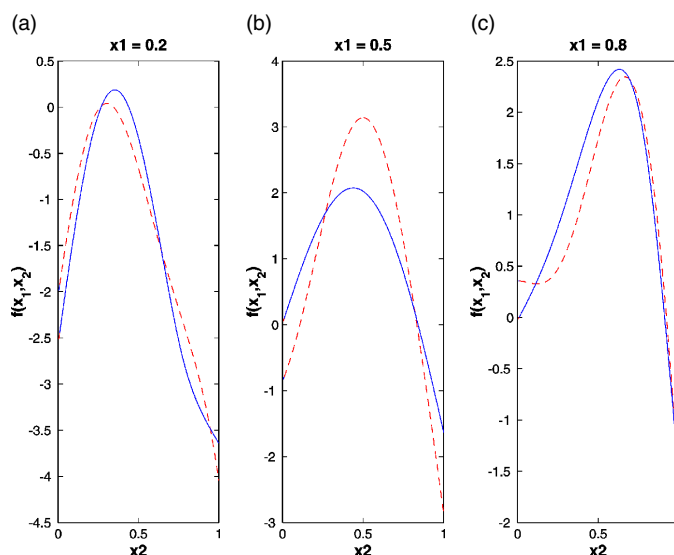


Figure 6. Cross-Sectional Plots for the Two-Factor Interaction Model: (a) $x = .2$; (b) $x = .5$; (c) $x = .8$ (--- true; — estimate).

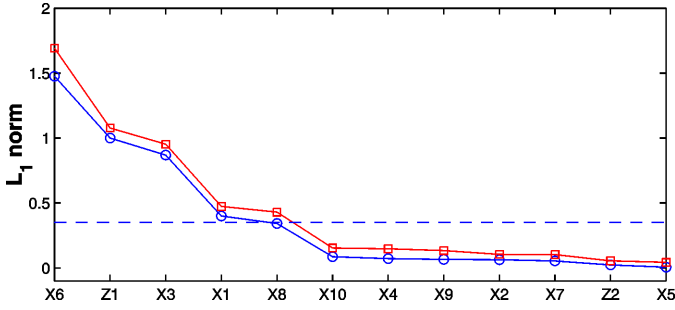


Figure 7. L_1 Norm Scores for the Main-Effects Model Incorporating Categorical Variables (—○— CKL; —■— GACV).

variables are uniformly distributed in $[0, 1]$, and the categorical variables are Bernoulli(.5) with values $\{0, 1\}$. The true logit function is

$$f(\mathbf{x}) = \frac{4}{3}x^{(1)} + \pi \sin(\pi x^{(3)}) + 8(x^{(6)})^5 + \frac{2}{e-1}e^{x^{(8)}} + 4z^{(1)} - 7.$$

The important main effects are $X^{(1)}$, $X^{(3)}$, $X^{(6)}$, $X^{(8)}$, and $Z^{(1)}$. The sample size is $n = 1,000$, and the basis size is $N = 50$. We use the same perturbation ϵ as in the previous examples. The main-effects model incorporating categorical variables in (7) is fitted. Figure 7 plots the ranked L_1 norm scores for all of the covariates. The LBP main-effects models using $\hat{\lambda}_{CKL}$ and $\hat{\lambda}_{GACV}$ both select the important continuous and categorical variables correctly.

7. WISCONSIN EPIDEMIOLOGIC STUDY OF DIABETIC RETINOPATHY

The Wisconsin Epidemiologic Study of Diabetic Retinopathy (WESDR) is an ongoing epidemiologic study of a cohort of patients receiving their medical care southern Wisconsin. All younger-onset diabetic persons (defined as under age 30 of age at diagnosis and taking insulin) and a probability sample of older-onset persons receiving primary medical care in an 11-county area of southwestern Wisconsin in 1979–1980 were invited to participate. Among 1,210 identified younger onset patients, 996 agreed to participate in the baseline examination in 1980–1982; of those, 891 participated in the first follow-up examination. Details about the study have been given by Klein, Klein, Moss, Davis, and DeMets (1984a,b, 1989), Klein, Klein, Moss, and Cruickshanks (1998), and others. A large number of medical, demographic, ocular, and other covariates were recorded in each examination. A multilevel retinopathy score was assigned to each eye based on its stereoscopic color fundus photographs. This dataset has been extensively analyzed using a variety of statistical methods (see, e.g., Craig, Fryback, Klein, and Klein 1999; Kim 1995).

In this section we examine the relation of a large number of possible risk factors at baseline to the 4-year progression of diabetic retinopathy. Each person's retinopathy score was defined as the score for the worse eye, and 4-year progression of retinopathy was considered to occur if the retinopathy score degraded two levels from baseline. Wahba et al. (1995) examined risk factors for progression of diabetic retinopathy on a subset

of the younger-onset group, whose members had no or nonproliferative retinopathy at baseline; the dataset comprised 669 persons. A model of the risk of progression of diabetic retinopathy in this population was built using a SS-ANOVA model (which has a quadratic penalty functional), using the predictor variables glycosylated hemoglobin (*gly*), duration of diabetes (*dur*), and body mass index (*bmi*); these variables are described further in Appendix B. That study began with these variables and two other (nonindependent) variables, age at baseline and age at diagnosis; these latter two were eliminated at the start. Although it was not discussed by Wahba et al. (1995), we report here that that study began with a large number (perhaps about 20) of potential risk factors, which were reduced to *gly*, *dur*, and *bmi* as likely the most important after many extended and laborious parametric and nonparametric regression analyses of small groups of variables at a time and by linear logistic regression by the authors and others. At that time it was recognized that a (smoothly) nonparametric model selection method that could rapidly flag important variables in a dataset with many candidate variables was very desirable. For the purposes of the present study, we make the reasonable assumption that *gly*, *dur*, and *bmi* are the “truth” (i.e., the most important risk factors in the analyzed population), and thus we are presented with a unique opportunity to examine the behavior of the LBP method in a real dataset where, arguably, the truth is known, by giving it many variables in this dataset and comparing the results to those of Wahba et al. (1995). Minor corrections and updatings of that dataset have been made, (but are not believed to affect the conclusions), and we have 648 persons in the updated dataset used here.

We performed some preliminary winnowing of the many potential prediction variables available, reducing the set for examination to 14 potential risk factors. The continuous covariates are *dur*, *gly*, *bmi*, *sys*, *ret*, *pulse*, *ins*, *sch*, *iop*, and the categorical covariates are *smk*, *sex*, *asp*, *famdb*, *mar*; the full names for these are given in Appendix C. Because the true f is not known for real data, only ranGACV is available for tuning λ . Figure 8 plots the L_1 norm scores of the individual functional components in the fitted LBP main-effects model. The dashed line indicates the threshold $q = .39$, which is chosen by the sequential bootstrap tests.

We note that the LBP picks out three most important variables, *gly*, *dur*, and *bmi*, specified by Wahba et al. (1995). The LBP also chose *sch* (highest year of school/college completed). This variable frequently shows up in demographic studies, when one looks for it, because it is likely a proxy for other variables related to disease, such as lifestyle and quality of medical care. It did show up in preliminary studies of Wahba et al.

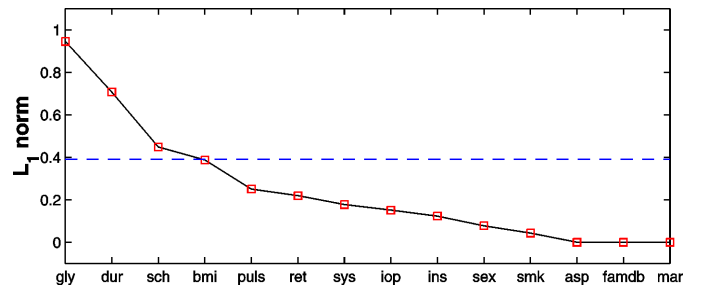


Figure 8. L_1 Norm Scores for the WESDR Main-Effects Model.

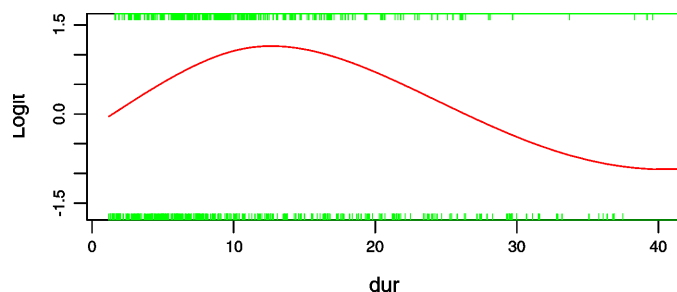


Figure 9. Estimated Logit Component for *dur*.

(1995) (not reported there) but was not included, because it was not considered a direct cause of disease itself. The sequential Monte Carlo bootstrap tests for *gly*, *dur*, *sch*, and *bmi* all have p value $1/51 \doteq .02$; thus these four covariates are selected as important risk factors at the significance level $\alpha = .05$.

Figure 9 plots the estimated logit for *dur*. The risk of progression of diabetic retinopathy increases up to a duration of about 15 years, then decreases thereafter, which generally agrees with the analysis of Wahba et al. (1995). The linear logistic regression model (using the function *glm* in the R package) shows that *dur* is not significant at level $\alpha = .05$. The curve in Figure 9 exhibits a hilly shape, which means that a quadratic function fits the curve better than a linear function. We refit the linear logistic model by intentionally including dur^2 ; the term dur^2 is significant with a p value of .02. This fact confirms the discovery of the LBP, and shows that LBP can be a valid screening tool to aid in determining the appropriate functional form for the individual covariate.

When fitting the two-factor interaction model in (6) with the constraints $\lambda_{\pi\pi} = \lambda_{\pi s} = \lambda_{ss}$, we did not find the *dur*–*bmi* interaction of Wahba et al. (1995) here. We note that the interaction terms tend to be washed out if there are only a few interactions. However, further exploratory analysis may be done by rearranging the constraints and/or varying the tuning parameters subjectively. Note that the solution to the optimization problem is very sparse. In this example, we observed that approximately 90% of the coefficients are 0's in the solution.

8. BEAVER DAM EYE STUDY

The Beaver Dam Eye Study (BDES) is an ongoing population-based study of age-related ocular disorders. It aims to collect information related to the prevalence, incidence, and severity of age-related cataract, macular degeneration, and diabetic retinopathy. Between 1987 and 1988, 5,924 eligible people (age 43–84 years) were identified in Beaver Dam, Wisconsin; and of those, 4,926 (83.1%) participated in the baseline exam. 5- and 10-year follow-up data have been collected, and results are being reported. Many variables of various kinds have been collected, including mortality between baseline and the follow-ups. A detailed description of the study has been given by Klein, Klein, Linton, and DeMets (1991); recent reports include that of Klein, Klein, Lee, Cruickshanks, and Chappell (2001).

We are interested in the relation between the 5-year mortality occurrence for the nondiabetic study participants and possible risk factors at baseline. We focus on the nondiabetic participants, because the pattern of risk factors for people with

diabetes and the rest of the population differs. We consider 10 continuous and 8 categorical covariates; the detailed information for which is given in Appendix D. The abbreviated names of continuous covariates are *pky*, *sch*, *inc*, *bmi*, *glu*, *cal*, *chl*, *hgb*, *sys*, and *age*, and those of the categorical covariates are *cv*, *sex*, *hair*, *hist*, *nout*, *mar*, *sum*, and *vtm*. We deliberately take into account some “noisy” variables in the analysis, including *hair*, *nout*, and *sum*, which are not directly related to mortality in general. We include these to show the performance of the LBP approach, and they are not expected to be picked out eventually by the model. Y is assigned a value of 1 if a person participated in the baseline examination and died before the start of the first 5-year follow-up and 0 otherwise. There are 4,422 nondiabetic study participants in the baseline examination, 395 of whom have missing data in the covariates. For the purpose of this study, we assume that the missing data are missing at random; thus these 335 subjects are not included in our analysis. This assumption is not necessarily valid, because age, blood pressure, body mass index, cholesterol, sex, smoking, and hemoglobin may well affect the missingness, but a further examination of the missingness is beyond the scope of the present study. In addition, we exclude another 10 participants who have either outlier values $pky > 158$ or very abnormal records $bmi > 58$ or $hgb < 6$. Thus we report an analysis of the remaining 4,017 nondiabetic participants from the baseline population.

We fit the main-effects model incorporating categorical variables in (7). The sequential Monte Carlo bootstrap tests for six covariates, *age*, *hgb*, *pky*, *sex*, *sys*, and *cv*, all have Monte Carlo p values $1/51 \doteq .02$, whereas the test for *glu* is not significant with a p value $9/51 = .18$. The threshold is chosen as $q = L_{(6)} = .25$. Figure 10 plots the L_1 norm scores for all of the potential risk factors. Using the threshold (dashed line) .25, chosen by the sequential bootstrap test procedure, the LBP model identifies six important risk factors, *age*, *hgb*, *pky*, *sex*, *sys*, and *cv*, for the 5-year mortality.

Compared with the LBP model, the linear logistic model with stepwise selection using the AIC, implemented by the function *glm* in R, misses the variable *sys* but selects three more variables, *inc*, *bmi*, and *sum*. Figure 11 depicts the estimated univariate logit components for the important continuous variables selected by the LBP model. All of the curves can be approximated reasonably well by linear models except *sys*, whose functional form exhibits a quadratic shape. This explains why *sys* is not selected by the linear logistic model. When we refit the logistic regression model by including sys^2 in the model, the stepwise selection picked out both *sys* and sys^2 .

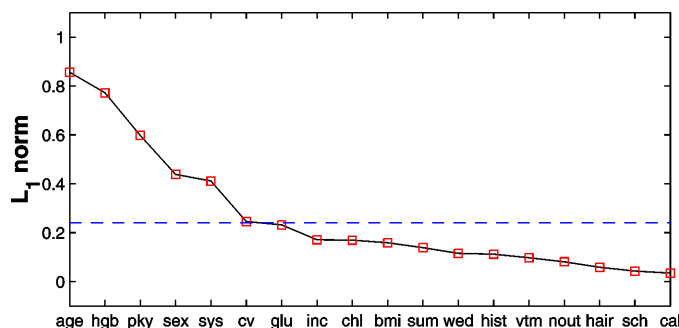


Figure 10. L_1 Norm Scores for the BDES Main-Effects Model.

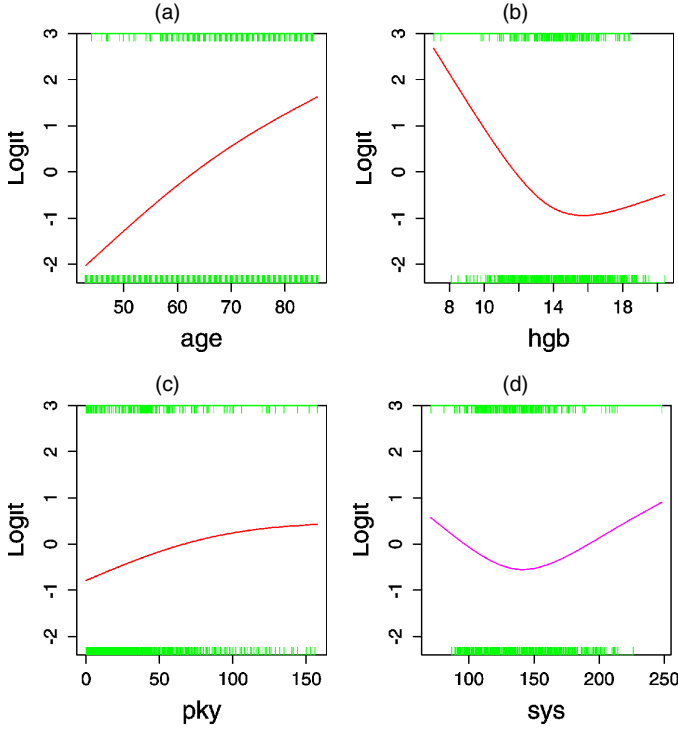


Figure 11. Estimated Univariate Logit Components for Important Variables: (a) age; (b) hgb; (c) pky; (d) sys.

9. DISCUSSION

We have proposed the LBP approach for variable selection in high-dimensional nonparametric model building. In the spirit of LASSO, LBP produces shrinkage functional estimates by imposing the l_1 penalty on the coefficients of the basis functions. Using the proposed measure of importance for the functional components, LBP effectively selects important variables, and the results are highly interpretable. LBP can handle continuous variables and categorical variables simultaneously. Although in this article our continuous variables have all been on subsets of the real line, it is clear that other continuous domains are possible. We have used LBP in the context of the Bernoulli distribution, but it can be extended to other exponential distributions as well (of course, to Gaussian data). We expect that larger numbers of variables than that considered here may be handled, and we expect that there will be many other scientific applications of the method. We plan to provide freeware for public use.

We believe that this method is a useful addition to the toolbox of the data analyst. It provides a way to examine the possible effects of a large number of variables in a nonparametric manner, complementary to standard parametric models in its ability to find nonparametric terms that may be missed by parametric methods. It has an advantage over quadratically penalized likelihood methods when the aim is to examine a large number of variables or terms simultaneously, inasmuch as the l_1 penalties result in sparse solutions. In practice, it can be an efficient tool for examining complex datasets to identify and prioritize variables (and, possibly, interactions) for further study. Based only on the variables or interactions identified by the LBP, one can build more traditional parametric or penalized likelihood models, for which confidence intervals and theoretical properties are known.

APPENDIX A: PROOF OF LEMMA 1

For $i = 1, \dots, n$, we have $-l(\mu_\lambda^{[-i]}(\mathbf{x}_i), \tau) = -\mu_\lambda^{[-i]}(\mathbf{x}_i)\tau + b(\tau)$. Let $f_\lambda^{[-i]}$ be the minimizer of the objective function

$$\frac{1}{n} \sum_{j \neq i} [-l(y_j, f(\mathbf{x}_j))] + J_\lambda(f). \quad (\text{A.1})$$

Because

$$\frac{\partial(-l(\mu_\lambda^{[-i]}(\mathbf{x}_i), \tau))}{\partial \tau} = -\mu_\lambda^{[-i]}(\mathbf{x}_i) + \dot{b}(\tau)$$

and

$$\frac{\partial^2(-l(\mu_\lambda^{[-i]}(\mathbf{x}_i), \tau))}{\partial^2 \tau} = \ddot{b}(\tau) > 0,$$

we see that $-l(\mu_\lambda^{[-i]}(\mathbf{x}_i), \tau)$ achieves its unique minimum at $\hat{\tau}$ that satisfies $\dot{b}(\hat{\tau}) = \mu_\lambda^{[-i]}(\mathbf{x}_i)$. So $\hat{\tau} = f_\lambda^{[-i]}(\mathbf{x}_i)$. Then for any f , we have

$$-l(\mu_\lambda^{[-i]}(\mathbf{x}_i), f_\lambda^{[-i]}(\mathbf{x}_i)) \leq -l(\mu_\lambda^{[-i]}(\mathbf{x}_i), f(\mathbf{x}_i)). \quad (\text{A.2})$$

Define $\mathbf{y}^{-i} = (y_1, \dots, y_{i-1}, \mu_\lambda^{[-i]}(\mathbf{x}_i), y_{i+1}, \dots, y_n)'$. For any f ,

$$\begin{aligned} I_\lambda(f, \mathbf{y}^{-i}) &= \frac{1}{n} \left\{ -l(\mu_\lambda^{[-i]}(\mathbf{x}_i), f(\mathbf{x}_i)) + \sum_{j \neq i} [-l(y_j, f(\mathbf{x}_j))] \right\} + J_\lambda(f) \\ &\geq \frac{1}{n} \left\{ -l(\mu_\lambda^{[-i]}(\mathbf{x}_i), f_\lambda^{[-i]}(\mathbf{x}_i)) + \sum_{j \neq i} [-l(y_j, f(\mathbf{x}_j))] \right\} + J_\lambda(f) \\ &\geq \frac{1}{n} \left\{ -l(\mu_\lambda^{[-i]}(\mathbf{x}_i), f_\lambda^{[-i]}(\mathbf{x}_i)) + \sum_{j \neq i} [-l(y_j, f_\lambda^{[-i]}(\mathbf{x}_j))] \right\} \\ &\quad + J_\lambda(f_\lambda^{[-i]}). \end{aligned}$$

The first inequality comes from (A.2). The second inequality is due to the fact that $f_\lambda^{[-i]}(\cdot)$ is the minimizer of (A.1). Thus we have that $h_\lambda(i, \mu_\lambda^{[-i]}(\mathbf{x}_i), \cdot) = f_\lambda^{[-i]}(\cdot)$.

APPENDIX B: DERIVATION OF ACV

Here we derive the ACV in (11) for the main-effects model; the derivation for two- or higher-order interaction models is similar. Write $\mu(\mathbf{x}) = E(Y|\mathbf{X} = \mathbf{x})$ and $\sigma^2(\mathbf{x}) = \text{var}(Y|\mathbf{X} = \mathbf{x})$. For the functional estimate f , we define $f_i = f(\mathbf{x}_i)$ and $\mu_i = \mu(\mathbf{x}_i)$ for $i = 1, \dots, n$. To emphasize that an estimate is associated with the parameter λ , we also use $f_{\lambda i} = f_\lambda(\mathbf{x}_i)$ and $\mu_{\lambda i} = \mu_\lambda(\mathbf{x}_i)$. Now define

$$OBS(\lambda) = \frac{1}{n} \sum_{i=1}^n [-y_i f_{\lambda i} + b(f_{\lambda i})].$$

This is the observed CKL, which, because y_i and $f_{\lambda i}$ are usually positively correlated, tends to underestimate the CKL. To correct this, we use the leave-one-out CV in (9),

$$\begin{aligned} CV(\lambda) &= \frac{1}{n} \sum_{i=1}^n [-y_i f_{\lambda i}^{[-i]} + b(f_{\lambda i})] \\ &= OBS(\lambda) + \frac{1}{n} \sum_{i=1}^n y_i (f_{\lambda i} - f_{\lambda i}^{[-i]}) \\ &= OBS(\lambda) + \frac{1}{n} \sum_{i=1}^n y_i \frac{f_{\lambda i} - f_{\lambda i}^{[-i]}}{y_i - \mu_{\lambda i}^{[-i]}} \\ &\quad \times \frac{y_i - \mu_{\lambda i}}{1 - (\mu_{\lambda i} - \mu_{\lambda i}^{[-i]})/(y_i - \mu_{\lambda i}^{[-i]})}. \end{aligned} \quad (\text{B.1})$$

For exponential families, we have $\mu_{\lambda i} = \dot{b}(f_{\lambda i})$ and $\sigma_{\lambda i}^2 = \ddot{b}(f_{\lambda i})$. If we approximate the finite difference by the functional derivative, then we get

$$\begin{aligned} \frac{\mu_{\lambda i} - \mu_{\lambda i}^{[-i]}}{y_i - \mu_{\lambda i}^{[-i]}} &= \frac{\dot{b}(f_{\lambda i}) - \dot{b}(f_{\lambda i}^{[-i]})}{y_i - \mu_{\lambda i}^{[-i]}} \\ &\approx \ddot{b}(f_{\lambda i}) \frac{f_{\lambda i} - f_{\lambda i}^{[-i]}}{y_i - \mu_{\lambda i}^{[-i]}} \\ &= \sigma_{\lambda i}^2 \frac{f_{\lambda i} - f_{\lambda i}^{[-i]}}{y_i - \mu_{\lambda i}^{[-i]}}. \end{aligned} \quad (\text{B.2})$$

Denote

$$G_i \equiv \frac{f_{\lambda i} - f_{\lambda i}^{[-i]}}{y_i - \mu_{\lambda i}^{[-i]}}; \quad (\text{B.3})$$

then from (B.2) and (B.1), we get

$$CV(\lambda) \approx OBS(\lambda) + \frac{1}{n} \sum_{i=1}^n G_i \frac{y_i (y_i - \mu_{\lambda i})}{1 - G_i \sigma_{\lambda i}^2}. \quad (\text{B.4})$$

Now we proceed to approximate G_i . Consider the main-effects model with

$$f(\mathbf{x}) = b_0 + \sum_{\alpha=1}^d b_{\alpha} k_1(x^{(\alpha)}) + \sum_{\alpha=1}^d \sum_{j=1}^N c_{\alpha,j} K_1(x^{(\alpha)}, x_j^{(\alpha)}).$$

Let $m = Nd$. Define the vectors $\mathbf{b} = (b_0, b_1, \dots, b_d)'$ and $\mathbf{c} = (c_{1,1}, \dots, c_{d,N})' = (c_1, \dots, c_m)'$. For $\alpha = 1, \dots, d$, define the $n \times N$ matrix $\mathbf{R}_{\alpha} = (K_1(x_i^{(\alpha)}, x_j^{(\alpha)}))$. Let $\mathbf{R} = [\mathbf{R}_1, \dots, \mathbf{R}_d]$ and

$$\mathbf{T} = \begin{bmatrix} 1 & k_1(x_1^{(1)}) & \dots & k_1(x_1^{(d)}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & k_1(x_n^{(1)}) & \dots & k_1(x_n^{(d)}) \end{bmatrix}.$$

Define $\mathbf{f} = (f_1, \dots, f_n)'$, then $\mathbf{f} = \mathbf{T}\mathbf{b} + \mathbf{R}\mathbf{c}$. The objective function in (10) can be expressed as

$$\begin{aligned} I_{\lambda}(\mathbf{b}, \mathbf{c}, \mathbf{y}) &= \frac{1}{n} \sum_{i=1}^n \left[-y_i \left(\sum_{\alpha=1}^{d+1} T_{i\alpha} b_{\alpha} + \sum_{j=1}^m R_{ij} c_j \right) \right. \\ &\quad \left. + b \left(\sum_{\alpha=1}^{d+1} T_{i\alpha} b_{\alpha} + \sum_{j=1}^m R_{ij} c_j \right) \right] \\ &\quad + \lambda_d \sum_{\alpha=2}^{d+1} |b_{\alpha}| + \lambda_s \sum_{j=1}^m |c_j|. \end{aligned} \quad (\text{B.5})$$

With the observed response \mathbf{y} , we denote the minimizer of (B.5) by $(\hat{\mathbf{b}}, \hat{\mathbf{c}})$. When a small perturbation $\boldsymbol{\varepsilon}$ is imposed on the response, we denote the minimizer of $I_{\lambda}(\mathbf{b}, \mathbf{c}, \mathbf{y} + \boldsymbol{\varepsilon})$ by $(\tilde{\mathbf{b}}, \tilde{\mathbf{c}})$. Then we have $\mathbf{f}_{\lambda}^{\mathbf{y}} = \mathbf{T}\hat{\mathbf{b}} + \mathbf{R}\hat{\mathbf{c}}$ and $\mathbf{f}_{\lambda}^{\mathbf{y}+\boldsymbol{\varepsilon}} = \mathbf{T}\tilde{\mathbf{b}} + \mathbf{R}\tilde{\mathbf{c}}$. The l_1 penalty in (B.5) tends to produce sparse solutions; that is, many components of $\hat{\mathbf{b}}$ and $\hat{\mathbf{c}}$ are exactly 0 in the solution. For simplicity of explanation, we assume that the first s components of $\hat{\mathbf{b}}$ and the first r components of $\hat{\mathbf{c}}$ are nonzero, that is,

$$\hat{\mathbf{b}} = \underbrace{(\hat{b}_1, \dots, \hat{b}_s)'}_{\neq \mathbf{0}} \underbrace{(\hat{b}_{s+1}, \dots, \hat{b}_{d+1})'}_{=\mathbf{0}}$$

and

$$\hat{\mathbf{c}} = \underbrace{(\hat{c}_1, \dots, \hat{c}_r)'}_{\neq \mathbf{0}} \underbrace{(\hat{c}_{r+1}, \dots, \hat{c}_m)'}_{=\mathbf{0}}.$$

The general case is similar but notationally more complicated.

The 0's in the solutions are robust against a small perturbation in data. That is, when the magnitude of perturbation $\boldsymbol{\varepsilon}$ is small enough, the 0 elements will stay at 0. This can be seen by transforming the optimization problem (B.5) into a nonlinear programming problem with a differentiable convex objective function and linear constraints, and considering the Karush–Kuhn–Tucker (KKT) conditions (see, e.g., Mangasarian 1969). Thus it is reasonable to assume that the solution of (B.5) with the perturbed data $\mathbf{y} + \boldsymbol{\varepsilon}$ has the form

$$\tilde{\mathbf{b}} = \underbrace{(\tilde{b}_1, \dots, \tilde{b}_s)'}_{\neq \mathbf{0}} \underbrace{(\tilde{b}_{s+1}, \dots, \tilde{b}_{d+1})'}_{=\mathbf{0}}$$

and

$$\tilde{\mathbf{c}} = \underbrace{(\tilde{c}_1, \dots, \tilde{c}_r)'}_{\neq \mathbf{0}} \underbrace{(\tilde{c}_{r+1}, \dots, \tilde{c}_m)'}_{=\mathbf{0}}.$$

For convenience, we denote the first s components of \mathbf{b} by \mathbf{b}^* and the first r components of \mathbf{c} by \mathbf{c}^* . Correspondingly, let \mathbf{T}^* be the submatrix containing the first s columns of \mathbf{T} and let \mathbf{R}^* be the submatrix consisting of the first r columns of \mathbf{R} . Then

$$\mathbf{f}_{\lambda}^{\mathbf{y}+\boldsymbol{\varepsilon}} - \mathbf{f}_{\lambda}^{\mathbf{y}} = \mathbf{T}(\tilde{\mathbf{b}} - \hat{\mathbf{b}}) + \mathbf{R}(\tilde{\mathbf{c}} - \hat{\mathbf{c}}) = [\mathbf{T}^* \quad \mathbf{R}^*] \begin{bmatrix} \tilde{\mathbf{b}}^* - \hat{\mathbf{b}}^* \\ \tilde{\mathbf{c}}^* - \hat{\mathbf{c}}^* \end{bmatrix}. \quad (\text{B.6})$$

Now

$$\begin{aligned} \left[\frac{\partial I_{\lambda}}{\partial \mathbf{b}^*}, \frac{\partial I_{\lambda}}{\partial \mathbf{c}^*} \right]'_{(\tilde{\mathbf{b}}, \tilde{\mathbf{c}}, \mathbf{y}+\boldsymbol{\varepsilon})} &= \mathbf{0} \quad \text{and} \\ \left[\frac{\partial I_{\lambda}}{\partial \mathbf{b}^*}, \frac{\partial I_{\lambda}}{\partial \mathbf{c}^*} \right]'_{(\hat{\mathbf{b}}, \hat{\mathbf{c}}, \mathbf{y})} &= \mathbf{0}. \end{aligned} \quad (\text{B.7})$$

The first-order Taylor approximation of $\left[\frac{\partial I_{\lambda}}{\partial \mathbf{b}^*}, \frac{\partial I_{\lambda}}{\partial \mathbf{c}^*} \right]'_{(\tilde{\mathbf{b}}, \tilde{\mathbf{c}}, \mathbf{y}+\boldsymbol{\varepsilon})}$ at $(\hat{\mathbf{b}}, \hat{\mathbf{c}}, \mathbf{y})$ gives

$$\begin{aligned} &\left[\frac{\partial I_{\lambda}}{\partial \mathbf{b}^*}, \frac{\partial I_{\lambda}}{\partial \mathbf{c}^*} \right]'_{(\tilde{\mathbf{b}}, \tilde{\mathbf{c}}, \mathbf{y}+\boldsymbol{\varepsilon})} \\ &\approx \left[\frac{\partial I_{\lambda}}{\partial \mathbf{b}^*}, \frac{\partial I_{\lambda}}{\partial \mathbf{c}^*} \right]'_{(\hat{\mathbf{b}}, \hat{\mathbf{c}}, \mathbf{y})} + \begin{bmatrix} \frac{\partial^2 I_{\lambda}}{\partial \mathbf{b}^* \partial \mathbf{b}^{*'}} & \frac{\partial^2 I_{\lambda}}{\partial \mathbf{b}^* \partial \mathbf{c}^{*'}} \\ \frac{\partial^2 I_{\lambda}}{\partial \mathbf{c}^* \partial \mathbf{b}^{*'}} & \frac{\partial^2 I_{\lambda}}{\partial \mathbf{c}^* \partial \mathbf{c}^{*'}} \end{bmatrix}_{(\hat{\mathbf{b}}, \hat{\mathbf{c}}, \mathbf{y})} \begin{bmatrix} \tilde{\mathbf{b}}^* - \hat{\mathbf{b}}^* \\ \tilde{\mathbf{c}}^* - \hat{\mathbf{c}}^* \end{bmatrix} \\ &\quad + \left[\frac{\partial^2 I_{\lambda}}{\partial \mathbf{b}^* \partial \mathbf{y}'} \right]_{(\hat{\mathbf{b}}, \hat{\mathbf{c}}, \mathbf{y})} (\mathbf{y} + \boldsymbol{\varepsilon} - \mathbf{y}), \end{aligned} \quad (\text{B.8})$$

when the magnitude of $\boldsymbol{\varepsilon}$ is small. Define

$$\mathbf{U} \equiv \begin{bmatrix} \frac{\partial^2 I_{\lambda}}{\partial \mathbf{b}^* \partial \mathbf{b}^{*'}} & \frac{\partial^2 I_{\lambda}}{\partial \mathbf{b}^* \partial \mathbf{c}^{*'}} \\ \frac{\partial^2 I_{\lambda}}{\partial \mathbf{c}^* \partial \mathbf{b}^{*'}} & \frac{\partial^2 I_{\lambda}}{\partial \mathbf{c}^* \partial \mathbf{c}^{*'}} \end{bmatrix}_{(\hat{\mathbf{b}}, \hat{\mathbf{c}}, \mathbf{y})}$$

and

$$\mathbf{V} \equiv - \left[\frac{\partial^2 I_{\lambda}}{\partial \mathbf{b}^* \partial \mathbf{y}'} \right]_{(\hat{\mathbf{b}}, \hat{\mathbf{c}}, \mathbf{y})}.$$

From (B.7) and (B.8), we have

$$\mathbf{U} \begin{bmatrix} \tilde{\mathbf{b}}^* - \hat{\mathbf{b}}^* \\ \tilde{\mathbf{c}}^* - \hat{\mathbf{c}}^* \end{bmatrix} \approx \mathbf{V}\boldsymbol{\varepsilon}.$$

Then (B.6) gives us $\mathbf{f}_{\lambda}^{\mathbf{y}+\boldsymbol{\varepsilon}} - \mathbf{f}_{\lambda}^{\mathbf{y}} \approx \mathbf{H}\boldsymbol{\varepsilon}$, where

$$\mathbf{H} \equiv [\mathbf{T}^* \quad \mathbf{R}^*] \mathbf{U}^{-1} \mathbf{V}. \quad (\text{B.9})$$

Now consider a special form of perturbation $\boldsymbol{\varepsilon}_0 = (0, \dots, \mu_{\lambda i}^{[-i]} - y_i, \dots, 0)'$; then $\mathbf{f}_{\lambda}^{\mathbf{y}+\boldsymbol{\varepsilon}_0} - \mathbf{f}_{\lambda}^{\mathbf{y}} \approx \varepsilon_{0i} H_{.i}$, where $\varepsilon_{0i} = \mu_{\lambda i}^{[-i]} - y_i$. Lemma 1

in Section 3 shows that $f_{\lambda}^{[-i]}$ is the minimizer of $I(f, \mathbf{y} + \mathbf{e}_0)$. Therefore G_i in (B.3) is

$$G_i = \frac{f_{\lambda i} - f_{\lambda i}^{[-i]}}{y_i - \mu_{\lambda i}^{[-i]}} = \frac{f_{\lambda i}^{[-i]} - f_{\lambda i}}{\varepsilon_{0i}} \approx h_{ii}.$$

From (B.4), an ACV score is then

$$ACV(\lambda) = \frac{1}{n} \sum_{i=1}^n [-y_i f_{\lambda i} + b(f_{\lambda i})] + \frac{1}{n} \sum_{i=1}^n h_{ii} \frac{y_i (y_i - \mu_{\lambda i})}{1 - \sigma_{\lambda i}^2 h_{ii}}. \quad (\text{B.10})$$

APPENDIX C: WISCONSIN EPIDEMIOLOGIC STUDY OF DIABETIC RETINOPATHY

• Continuous covariates:

- X_1 : (*dur*) duration of diabetes at the time of baseline examination, years
- X_2 : (*gly*) glycosylated hemoglobin, a measure of hyperglycemia, %
- X_3 : (*bmi*) body mass index, kg/m²
- X_4 : (*sys*) systolic blood pressure, mm Hg
- X_5 : (*ret*) retinopathy level
- X_6 : (*pulse*) pulse rate, count for 30 seconds
- X_7 : (*ins*) insulin dose, kg/day
- X_8 : (*sch*) years of school completed
- X_9 : (*iop*) intraocular pressure, mm Hg

• Categorical covariates:

- Z_1 : (*smk*) smoking status (0 = no, 1 = any)
- Z_2 : (*sex*) gender (0 = female, 1 = male)
- Z_3 : (*asp*) use of at least one aspirin for (0 = no, 1 = yes)
at least three months while diabetic
- Z_4 : (*famdb*) family history of diabetes (0 = none, 1 = yes)
- Z_5 : (*mar*) marital status (0 = no, 1 = yes/ever).

APPENDIX D: BEAVER DAM EYE STUDY

• Continuous covariates:

- X_1 : (*pkv*) pack years smoked, (packs per day/20)*years smoked
- X_2 : (*sch*) highest year of school/college completed, years
- X_3 : (*inc*) total household personal income, thousands/month
- X_4 : (*bmi*) body mass index, kg/m²
- X_5 : (*glu*) glucose (serum), mg/dL
- X_6 : (*cal*) calcium (serum), mg/dL
- X_7 : (*chl*) cholesterol (serum), mg/dL
- X_8 : (*hgb*) hemoglobin (blood), g/dL
- X_9 : (*sys*) systolic blood pressure, mm Hg
- X_{10} : (*age*) age at examination, years.

• Categorical covariates:

- Z_1 : (*cv*) history of cardiovascular disease (0 = no, 1 = yes)
- Z_2 : (*sex*) gender (0 = female, 1 = male)
- Z_3 : (*hair*) hair color (0 = blond/red, 1 = brown/black)
- Z_4 : (*hist*) history of heavy drinking (0 = never, 1 = past/currently)
- Z_5 : (*nout*) winter leisure time (0 = indoors, 1 = outdoors)
- Z_6 : (*mar*) marital status (0 = no, 1 = yes/ever)
- Z_7 : (*sum*) part of day spent outdoors in summer (0 = < 1/4 day, 1 = > 1/4 day)
- Z_8 : (*vtm*) vitamin use (0 = no, 1 = yes).

[Received July 2002. Revised February 2004.]

REFERENCES

- Bakin, S. (1999), "Adaptive Regression and Model Selection in Data Mining Problems," unpublished doctoral thesis, Australia National University.
- Chen, S., Donoho, D., and Saunders, M. (1998), "Atomic Decomposition by Basis Pursuit," *SIAM Journal of the Scientific Computing*, 20, 33–61.
- Craig, B. A., Fryback, D. G., Klein, R., and Klein, B. (1999), "A Bayesian Approach to Modeling the Natural History of a Chronic Condition From Observations With Intervention," *Statistics in Medicine*, 18, 1355–1371.
- Craven, P., and Wahba, G. (1979), "Smoothing Noise Data With Spline Functions: Estimating the Correct Degree of Smoothing by the Method of Generalized Cross-Validation," *Numerische Mathematik*, 31, 377–403.
- Davison, A. C., and Hinkley, D. V. (1997), *Bootstrap Methods and Their Application*, Cambridge, U.K.: Cambridge University Press.
- Fan, J., and Li, R. Z. (2001), "Variable Selection via Penalized Likelihood," *Journal of the American Statistical Association*, 96, 1348–1360.
- Ferris, M. C., and Voelker, M. M. (2000), "Slice Models in General Purpose Modeling Systems," *Optimization Methods and Software*, 17, 1009–1032.
- (2001), "Slice Models in GAMS," in *Operations Research Proceedings*, eds. P. Chamon, R. Leisten, A. Martin, J. Minnemann, and H. Stadler, New York: Springer-Verlag, pp. 239–246.
- Frank, I. E., and Friedman, J. H. (1993), "A Statistical View of Some Chemometrics Regression Tools," *Technometrics*, 35, 109–148.
- Fu, W. J. (1998), "Penalized Regression: The Bridge versus the LASSO," *Journal of Computational and Graphical Statistics*, 7, 397–416.
- Gao, F., Wahba, G., Klein, R., and Klein, B. (2001), "Smoothing Spline ANOVA for Multivariate Bernoulli Observations, With Application to Ophthalmology Data," *Journal of the American Statistical Association*, 96, 127–160.
- Girard, D. (1998), "Asymptotic Comparison of (Partial) Cross-Validation, GCV and Randomized GCV in Nonparametric Regression," *The Annals of Statistics*, 26, 315–334.
- Gu, C. (2002), *Smoothing Spline ANOVA Models*, New York: Springer-Verlag.
- Gu, C., and Kim, Y. J. (2002), "Penalized Likelihood Regression: General Formulation and Efficient Approximation," *Canadian Journal of Statistics*, 30, 619–628.
- Gunn, S. R., and Kandola, J. S. (2002), "Structural Modelling With Sparse Kernels," *Machine Learning*, 48, 115–136.
- Hastie, T. J., and Tibshirani, R. J. (1990), *Generalized Additive Models*, New York: Chapman & Hall.
- Hutchinson, M. (1989), "A Stochastic Estimator for the Trace of the Influence Matrix for Laplacian Smoothing Splines," *Communication in Statistics, Simulation*, 18, 1059–1076.
- Kim, K. (1995), "A Bivariate Cumulative Probit Regression Model for Ordered Categorical Data," *Statistics in Medicine*, 14, 1341–1352.
- Kimeldorf, G., and Wahba, G. (1971), "Some Results on Tchebycheffian Spline Functions," *Journal of Mathematical Analysis and Applications*, 33, 82–95.
- Klein, R., Klein, B., Lee, K., Cruickshanks, K., and Chappell, R. (2001), "Changes in Visual Acuity in a Population Over a 10-Year Period. Beaver Dam Eye Study," *Ophthalmology*, 108, 1757–1766.
- Klein, R., Klein, B., Linton, K., and DeMets, D. L. (1991), "The Beaver Dam Eye Study: Visual Acuity," *Ophthalmology*, 98, 1310–1315.
- Klein, R., Klein, B., Moss, S., and Cruickshanks, K. (1998), "The Wisconsin Epidemiologic Study of Diabetic Retinopathy XVII. The 14-Year Incidence and Progression of Diabetic Retinopathy and Associated Risk Factors in Type 1 Diabetes," *Ophthalmology*, 105, 1801–1815.
- Klein, R., Klein, B., Moss, S. E., Davis, M. D., and DeMets, D. L. (1984a), "The Wisconsin Epidemiologic Study of Diabetic Retinopathy II. Prevalence and Risk of Diabetes When Age at Diagnosis is Less Than 30 Years," *Archives of Ophthalmology*, 102, 520–526.
- (1984b), "The Wisconsin Epidemiologic Study of Diabetic Retinopathy. III. Prevalence and Risk of Diabetes When Age at Diagnosis is 30 or More Years," *Archives of Ophthalmology*, 102, 527–532.
- (1989), "The Wisconsin Epidemiologic Study of Diabetic Retinopathy IX. Four-Year Incidence and Progression of Diabetic Retinopathy When Age at Diagnosis Is Less Than 30 Years," *Archives of Ophthalmology*, 107, 237–243.
- Knight, K., and Fu, W. J. (2000), "Asymptotics for Lasso-Type Estimators," *The Annals of Statistics*, 28, 1356–1378.
- Lin, X., Wahba, G., Xiang, D., Gao, F., Klein, R., and Klein, B. (2000), "Smoothing Spline ANOVA Models for Large Data Sets With Bernoulli Observations and the Randomized GACV," *The Annals of Statistics*, 28, 1570–1600.
- Linhart, H., and Zucchini, W. (1986), *Model Selection*, New York: Wiley.
- Mangasarian, O. (1969), *Nonlinear Programming*, New York: McGraw-Hill.
- Murtagh, B. A., and Saunders, M. A. (1983), "MINOS 5.5 User's Guide," Technical Report SOL 83-20R, Stanford University, OR Dept.

- Ruppert, D., and Carroll, R. J. (2000), "Spatially-Adaptive Penalties for Spline Fitting," *Australian and New Zealand Journal of Statistics*, 45, 205–223.
- Tibshirani, R. J. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society, Ser. B*, 58, 267–288.
- Wahba, G. (1990), *Spline Models for Observational Data*, Vol. 59, Philadelphia: SIAM.
- Wahba, G., Wang, Y., Gu, C., Klein, R., and Klein, B. (1995), "Smoothing Spline ANOVA for Exponential Families, With Application to the Wisconsin Epidemiological Study of Diabetic Retinopathy," *The Annals of Statistics*, 23, 1865–1895.
- Wahba, G., and Wold, S. (1975), "A Completely Automatic French Curve," *Communications in Statistics*, 4, 1–17.
- Xiang, D., and Wahba, G. (1996), "A Generalized Approximate Cross-Validation for Smoothing Splines With Non-Gaussian Data," *Statistica Sinica*, 6, 675–692.
- (1998), "Approximate Smoothing Spline Methods for Large Data Sets in the Binary Case," in *Proceedings of the American Statistical Association Joint Statistical Meetings, Biometrics Section*, pp. 94–98.
- Yau, P., Kohn, R., and Wood, S. (2003), "Bayesian Variable Selection and Model Averaging in High-Dimensional Multinomial Nonparametric Regression," *Journal of Computational and Graphical Statistics*, 12, 23–54.
- Zhang, H. H. (2002), "Nonparametric Variable Selection and Model Building via Likelihood Basis Pursuit," Technical Report 1066, University of Wisconsin-Madison, Dept. of Statistics.