

# Dimension reduction and parameter estimation for additive index models\*

LINGYAN RUAN AND MING YUAN†

In this paper, we consider simultaneous model selection and estimation for the additive index model. The additive index model is a class of structured nonparametric models that can be expressed as additive models of a set of unknown linear transformation of the original predictor variables. We introduce a penalized least squares estimator and discuss how it can be efficiently computed in practice. Both theoretical and empirical properties of the estimate are presented to demonstrate its merits. Extensions to more general prediction framework are also discussed.

KEYWORDS AND PHRASES: Additive model, Index model, Model selection, Projection pursuit, Smoothing splines.

## 1. INTRODUCTION

In the additive index model, a response  $y$  is related to a predictor  $\mathbf{x} \in \mathbf{R}^p$  through

$$(1) \quad f(\mathbf{x}) = h_1(\alpha'_1 \mathbf{x}) + h_2(\alpha'_2 \mathbf{x}) + \cdots + h_q(\alpha'_q \mathbf{x}),$$

with some projection indices  $\alpha_1, \alpha_2, \dots, \alpha_q$  and ridge functions  $h_1, h_2, \dots, h_q$ . The goal is to infer  $f$  based on  $n$  independent copies  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  of  $(\mathbf{x}, y)$  without assuming any parametric form for the ridge functions. The additive model has several popular structured nonparametric regression models, most notably additive models or single index models, as special cases. It is also known to be much more flexible than these more specialized examples. In particular, it has been shown that any square integrable function can be approximated to an arbitrary precision by a function of form (1) (Diaconis and Shahshahani, 1984).

To estimate  $f$ , the dimensionality  $q$  is often assumed to be known apriori. Various methods have been proposed to estimate  $f$  when  $q$  is given. See, e.g., Chen (1991) and Chiou and Müller (2004) among others. The assumption that  $q$  is known in advance, however, can be unrealistic. Its choice essentially amounts to a model selection or dimension reduction problem. It is clear that if  $q$  is too small, the additive index model cannot sufficiently capture the relationship between the response and the predictors. On the other hand,

when  $q$  is too big, some of the components in the additive index model will be close to zero, which could cause identifiability problems (see, e.g., Yuan, 2008), and subsequently troubles in parameter estimation. The choice of  $q$ , albeit critical, is notoriously difficult in practice.

In this paper, we propose a penalized least squares method to simultaneously select the dimensionality  $q$  and estimate the regression function  $f$  for the additive index model. Our method is inspired by the COSSO recently proposed by Lin and Zhang (2006) for the purpose of variable selection in additive models. The penalty we employ encourages some of the ridge functions to be exactly zero instead of being estimated as close to zero, which avoids the potential problem of unidentifiability. We prove that the additive index model estimated using the proposed estimator achieves the same rate of convergence as the univariate nonparametric regression, which suggests that models of form (1) have the potential to mitigate the curse of dimensionality.

The rest of this paper is organized as follows. In Section 2, we introduce the proposed estimator for the additive index model and study its asymptotic properties. We show that given that all ridge functions come from a Sobolev space, the additive index model estimated in this fashion achieves the optimal convergence rate. Section 3 discusses an iterative algorithm that can be used to efficiently compute the proposed estimate in practice. Examples, both simulated and real data, are presented in Section 4 to demonstrate the merits of the proposed methodology. We conclude with some discussions in Section 5.

## 2. PENALIZED LEAST SQUARES ESTIMATE

To fix ideas, we shall begin by focusing on the usual mean regression model:

$$(2) \quad y_i = f_0(\mathbf{x}_i) + \epsilon_i, \quad i = 1, 2, \dots, n,$$

where the predictors are properly scaled so that  $\mathbf{x}_i \in [0, 1]^p$  and

$$(3) \quad f_0(\mathbf{x}) = \mu + h_{01}(\alpha'_1 \mathbf{x}) + h_{02}(\alpha'_2 \mathbf{x}) + \cdots + h_{0q}(\alpha'_q \mathbf{x}),$$

where the dimensionality  $q$ , although unknown, is assumed to be upper bounded by a known value  $M$ . Such an upper bound can be easily obtained in practice. For example, it

\*This research was supported in part by NSF Grant DMS-0706724, DMS-0846234, and a grant from the Georgia Cancer Coalition.

†Corresponding author.

can be the final dimensionality obtained by the projection pursuit regression. To avoid ambiguity, we shall also assume that the projection indices  $\alpha_j$ s satisfy  $\|\alpha_j\| = 1$ . To model the ridge functions nonparametrically, we consider them as members of the usual Sobolev space of order  $m$ :

$$(4) \quad \mathcal{S}_m = \{h : h^{(s)} \text{ are absolutely continuous for } s = 0, \dots, m-1, h^{(m)} \in \mathcal{L}_2[0, 1]\}.$$

When equipped with a norm:

$$(5) \quad \|h\|^2 = \sum_{v=0}^{m-1} \left( \int h^{(v)} \right)^2 + \int_0^1 \left( h^{(m)} \right)^2,$$

$\mathcal{S}_m$  forms a reproducing kernel Hilbert space (see, e.g., Wahba, 1990). Let  $\bar{\mathcal{S}}_m$  be the orthogonal complement of constant functions in  $\mathcal{S}_m$  in that  $\mathcal{S}_m = \{1\} \oplus \bar{\mathcal{S}}_m$ . It then suffices to consider  $h_j \in \bar{\mathcal{S}}_m$ , i.e.,

$$(6) \quad f_0 \in \mathcal{F}_M := \{\mu + h_1(\alpha'_1 \mathbf{x}) + h_2(\alpha'_2 \mathbf{x}) + \dots + h_M(\alpha'_M \mathbf{x}) : h_j(\cdot) \in \bar{\mathcal{S}}_m, \|\alpha_j\| = 1\}.$$

For any  $f \in \mathcal{F}_M$ , define

$$(7) \quad J(f) = \inf \{\|h_1\| + \|h_2\| + \dots + \|h_M\|\}$$

where the infimum is taken over all possible  $h_j$ s such that  $f(\mathbf{x}) = \mu + h_1(\alpha'_1 \mathbf{x}) + h_2(\alpha'_2 \mathbf{x}) + \dots + h_p(\alpha'_p \mathbf{x})$  and  $h_j \in \bar{\mathcal{S}}_m$ . It is not hard to check that  $J$  is a pseudo-norm on  $\mathcal{F}_M$ .

We now propose to estimate  $f_0$  by the following penalized least squares estimator:

$$(8) \quad \hat{f}_\lambda = \arg \min_{f \in \mathcal{F}_M} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 + \lambda_n J(f) \right\},$$

where  $\lambda_n > 0$  is a tuning parameter. The following theorem states that our estimate achieves the same convergence rate as the usual additive model.

**Theorem 1.** *If  $f_0$  is not a constant, and  $\lambda_n^{-1} = O_p(n^{2m/(2m+1)})J^{(2m-1)/(2m+1)}(f_0)$ , then*

$$(9) \quad \frac{1}{n} \sum_{i=1}^n \left( \hat{f}_\lambda(\mathbf{x}_i) - f_0(\mathbf{x}_i) \right)^2 = O_p(\lambda_n) J(f_0).$$

If  $f_0$  is a constant, then

$$(10) \quad \frac{1}{n} \sum_{i=1}^n \left( \hat{f}_\lambda(\mathbf{x}_i) - f_0(\mathbf{x}_i) \right)^2 = O_p \left( \max\{n^{-1}, \lambda_n^{-1/(2m-1)} n^{-2m/(2m-1)}\} \right).$$

*Proof.* Recall that  $\hat{f}, f_0 \in \mathcal{F}_q$ . We have  $\hat{f} - f_0 \in \mathcal{F}_{2q} \subseteq \mathcal{F}_{2M}$ . Therefore

$$(11) \quad \frac{\hat{f} - f_0}{J(\hat{f}) + J(f_0)} \in \{f \in \mathcal{F}_{2M} : \|f\| \leq 1\} =: \mathcal{G}.$$

Denote the set on the right-hand side  $\mathcal{G}$ . In the light of Theorem 10.2 of van de Geer (2000), it suffices to show that

$$(12) \quad H_\infty(\delta, \mathcal{G}) \leq C\delta^{-1/m}.$$

Hereafter, we use  $C$  to denote a generic positive constant that does not depend on  $\delta$ .

For any  $f_1, f_2 \in \mathcal{G}$ , there exists  $h_j$ s and  $g_j$ s such that

$$(13) \quad f_1(\mathbf{x}) = h_1(\alpha'_1 \mathbf{x}) + \dots + h_{2M}(\alpha'_{2M} \mathbf{x});$$

$$(14) \quad f_2(\mathbf{x}) = g_1(\beta'_1 \mathbf{x}) + \dots + g_{2M}(\beta'_{2M} \mathbf{x}),$$

where  $h_j, g_j \in \bar{\mathcal{S}}_m$  and

$$\|h_1\|^2 + \dots + \|h_{2M}\|^2 \leq 1;$$

$$\|g_1\|^2 + \dots + \|g_{2M}\|^2 \leq 1.$$

Recall that

$$(15) \quad \|h_j\|^2 = \sum_{s=1}^{m-1} \left\{ h_j^{(s)}(1) - h_j^{(s)}(0) \right\}^2 + \int \left( h_j^{(m)} \right)^2, \\ j = 1, 2, \dots, 2p.$$

Following the same argument as that of Lemma A.1 in Lin and Zhang (2006), one can show that  $\|h'_j\|_\infty \equiv \max |h'_j(u)| \leq 1$ ,  $\|h_j\|_\infty \leq 1$ ,  $j = 1, 2, \dots, 2M$  and the same holds true for  $g_j$ s. Denote  $\mathcal{B}_m$  the unit ball in  $\mathcal{S}_m$ . Then  $h_j, g_j \in \mathcal{B}_m$ . It is well known that the  $\delta$  entropy of  $\mathcal{B}_m$  for the supreme norm is bounded by

$$(16) \quad H_\infty(\delta, \mathcal{B}_m) \leq C\delta^{-1/m}.$$

In other words,  $\mathcal{B}_m$  can be covered by  $N = \exp(C\delta^{-1/m})$  balls with radius  $\delta/4M$ . Note that

$$(17) \quad \|f_1 - f_2\|_\infty \leq \sum_{j=1}^{2M} \|g_j(\beta'_j \cdot) - h_j(\alpha'_j \cdot)\|_\infty \\ \leq \sum_{j=1}^{2M} (\|g_j(\beta'_j \cdot) - g_j(\alpha'_j \cdot)\|_\infty \\ + \|g_j(\alpha'_j \cdot) - h_j(\alpha'_j \cdot)\|_\infty) \\ \leq \sum_{j=1}^{2M} (\|g'_j\|_\infty \|\beta_j - \alpha_j\|_1 + \|g_j - h_j\|_\infty) \\ \leq \sum_{j=1}^{2M} (\|\beta_j - \alpha_j\|_1 + \|f_j - h_j\|_\infty)$$

where  $\|\beta_j - \alpha_j\|_1 = \sum_k |\beta_{jk} - \alpha_{jk}|$ . Because a unit sphere in  $\mathbf{R}^p$  can be covered by  $C\delta^{-p}$  balls of radius  $\delta/4p$ ,  $\mathcal{G}$  can be covered by  $CN^{2M}\delta^{-p}$ , which implies the desired entropy condition (12).  $\square$

Note that additive models are also members of  $\mathcal{F}$ . According to Stone (1985),  $n^{-2m/(2m+1)}$  is the optimal rate in

estimating  $f_0$  if it follows an additive model. This suggests that the rate obtained in Theorem 1 is optimal and can not be improved. Allowing projection directions renders the additive index model more flexibility than the additive model. Theorem 1 shows that such flexibility can be gained without loss in terms of estimability asymptotically.

### 3. COMPUTATION

We now discuss the practical aspects in computing the penalized least squares estimator (8). The objective function on the right-hand side of (8) can be minimized in an iterative fashion. Similar to the COSSO, minimizing the objective function on the right-hand side of (8) is equivalent to minimizing

$$(18) \quad \frac{1}{n} \sum_{i=1}^n \{y_i - f(\mathbf{x}_i)\}^2 + \lambda_0 \sum_{j=1}^M \tau_j^{-1} \|h_j\|^2 + \frac{\lambda^2}{4\lambda_0} \sum_{j=1}^M \tau_j$$

for any  $\lambda_0 > 0$  with respect to both  $f$  and  $\tau_j$ s under the constraint that  $\tau_j \geq 0$ ,  $j = 1, \dots, M$ . Hereafter, we abbreviate the subscript of the tuning parameters that signifies their dependence on the sample size. Now we can minimize (18) with respect to the projection directions, ridge functions and  $\tau_j$ s iteratively.

Write  $A = (\alpha_1, \dots, \alpha_M)$  and  $\mathbf{z} = A'\mathbf{x}$ . When  $A$  and  $\tau_j$ s are known, (18) becomes

$$(19) \quad \frac{1}{n} \sum_{i=1}^n \left( y_i - \mu - \sum_{j=1}^p h_j(z_{ij}) \right)^2 + \lambda_0 \sum_{j=1}^M \tau_j^{-1} \|h_j\|^2,$$

the usual smoothing spline estimate of the additive model (Wahba, 1990). The representer Lemma (Wahba, 1990) implies that the minimizer of (19) can be given as

$$(20) \quad h_j(u) = \tau_j \sum_{i=1}^n c_i K(u; z_{ij}), \quad j = 1, \dots, M,$$

where  $K(\cdot, \cdot)$  is the reproducing kernel associated with  $\bar{\mathcal{S}}_m$  and  $z_{ij}$  is the  $j$ th factor of the  $i$ th observation, i.e.,  $z_{ij} = \alpha'_j \mathbf{x}_i$ . Furthermore,

$$(21) \quad \|h_j\|^2 = \tau_j \sum_{i,k=1}^n c_i c_k K(z_{ij}, z_{kj}).$$

Plugging both formulae back to (19), we can get a closed form solution for  $\mu$  and  $\mathbf{c} = (c_1, \dots, c_n)'$ . Readers are referred to Wahba (1990) for details.

Now suppose that  $\mathbf{c}$  and  $\mu$  are known. Minimizing (18) with respect to  $A$  can be done using Newton iteration. We update  $A$  column by column. Consider, for example, updating  $\alpha_j$ . Write

$$(22) \quad r_i = y_i - \left\{ \mu + \sum_{k \neq j} h_k(\alpha'_k \mathbf{x}_i) \right\}$$

with  $\alpha_j$ s from the previous iteration and  $h_j$ s being given by (20). Our goal is to solve

$$(23) \quad \min_{\alpha_j \in \mathbf{R}^p} \frac{1}{n} \sum_{i=1}^n \{r_i - h_j(\alpha'_j \mathbf{x}_i)\}^2, \quad \text{subject to } \|\alpha_j\| = 1.$$

A Lagrange formulation leads to

$$(24) \quad \min_{\alpha_j \in \mathbf{R}^p} \left[ \sum_{i=1}^n \{r_i - h_j(\alpha'_j \mathbf{x}_i)\}^2 + \theta \|\alpha_j\|^2 \right],$$

for some constant  $\theta$ . First order condition yields

$$(25) \quad \sum_{i=1}^n [\{r_i - h_j(\alpha'_j \mathbf{x}_i)\} h'_j(\alpha'_j \mathbf{x}_i) \mathbf{x}_i] = \theta \alpha_j.$$

Multiplying both sides by  $\alpha_j$ ,

$$(26) \quad \theta = \sum_{i=1}^n [\{r_i - h_j(\alpha'_j \mathbf{x})\} h'_j(\alpha'_j \mathbf{x}_i) \alpha'_j \mathbf{x}_i].$$

To this end, we need to be able to compute the gradient and hessian of (24). Note that from (20)

$$(27) \quad h'_j(u) = \tau_j \sum_{i=1}^n c_i \frac{\partial K(u; z_{ij})}{\partial u};$$

$$(28) \quad h''_j(u) = \tau_j \sum_{i=1}^n c_i \frac{\partial^2 K(u; z_{ij})}{\partial u^2}.$$

Both can be easily evaluated. The gradient and hessian of (24) can then be deduced from

$$(29) \quad G(\alpha_j) = -2 \sum_{i=1}^n [\{r_i - h_j(\alpha'_j \mathbf{x})\} h'_j(\alpha'_j \mathbf{x}_i) \mathbf{x}_i] + 2\theta \alpha_j;$$

$$(30) \quad H(\alpha_j) = 2 \sum_{i=1}^n \left[ \{h'_j(\alpha'_j \mathbf{x}_i)\}^2 - \{r_i - h_j(\alpha'_j \mathbf{x})\} h''_j(\alpha'_j \mathbf{x}_i) \right] \mathbf{x}_i \mathbf{x}'_i + 2\theta I.$$

Now a Newton iteration would update the current estimate  $\alpha_j$  by  $\beta_j / \|\beta_j\|$  where

$$(31) \quad \beta_j = \alpha_j - H^{-1}(\alpha_j) G(\alpha_j).$$

Lastly, we update  $\tau_j$ s. This can be done in a similar fashion as the COSSO. Denote  $K_j$  an  $n \times n$  matrix whose  $(i, l)$  entry is  $K(\alpha'_j \mathbf{x}_i, \alpha'_j \mathbf{x}_l)$  and

$$(32) \quad F = (K_1 \mathbf{c}, \dots, K_p \mathbf{c}).$$

Define  $\tilde{\mathbf{y}} = \mathbf{y} - n\lambda_0\mathbf{c}/2 - \mu\mathbf{1}$ . It is not hard to show that the  $\tau_j$ s that minimize (18) are also the solution to

$$(33) \quad \min_{\tau_1, \dots, \tau_p} \{\tilde{\mathbf{y}} - F(\tau_1, \dots, \tau_p)\}' \{\tilde{\mathbf{y}} - F(\tau_1, \dots, \tau_p)\}' + \frac{\lambda^2}{4\lambda_0} \sum_{j=1}^p \tau_j$$

subject to  $\tau_j \geq 0$ .

This is the Lagrange formulation of

$$(34) \quad \min_{\tau_1, \dots, \tau_p} \{\tilde{\mathbf{y}} - F(\tau_1, \dots, \tau_p)\}' \{\tilde{\mathbf{y}} - F(\tau_1, \dots, \tau_p)\}'$$

subject to  $\tau_j \geq 0, \sum_{j=1}^p \tau_j \leq T$

for some  $T > 0$ , which is in the same form as the so-called nonnegative garrote proposed for variable selection in linear regression (Breiman, 1995). This optimization problem is a quadratic program and can be solved efficiently using the standard quadratic program solver.

In summary, (18) can be minimized using the following algorithm.

**Algorithm.**

- (I) Initialize  $A$  with a random  $p \times M$  matrix normalized so that each column has norm one.
- (II) Initialize  $\tau_1 = \dots = \tau_M = 1$ .
- (III) Choose  $\lambda_0$  by the generalized cross validation and compute the initial estimate of  $\mu$  and  $h_j$ s using the smoothing spline algorithm (Wahba, 1990).
- (IV) Compute the first and second derivative of  $h_j$ s using (27) and (28).
- (V) Update  $\alpha_j$ s by a one-step Newton iteration as described before.
- (VI) Given  $A$ , update  $\mu$  and  $h_j$ s using the smoothing spline algorithm.
- (VII) Update  $\tau_j$ s by solving (34).
- (VIII) Go back to step (IV) until a certain convergence criterion is met.

Thus far, we have assumed that the tuning parameter  $\lambda$ , or equivalently  $T$  of (34) is fixed. In practice, it can be selected using cross-validation. In cross-validation, the full data set  $\mathcal{D}$  is randomly split into  $K$  subsets of about the same size, denoted by  $\mathcal{D}^{(1)}, \dots, \mathcal{D}^{(K)}$ . For each  $k = 1, \dots, K$ , we use the data in  $\mathcal{D} - \mathcal{D}^{(k)}$  to estimate the model parameters and the data in  $\mathcal{D}^{(k)}$  to validate. The squared prediction error could be used as the performance measure in our case. For each candidate value of  $M$ , the  $K$ -fold cross-validated score is defined as

$$(35) \quad CV(T) = \sum_{k=1}^K \sum_{i \in I_k} (y_i - \hat{f}^{(-k)}(\mathbf{x}_i))^2$$

where  $I_k$  is the index set of the data in  $\mathcal{D}^{(k)}$ , and  $\hat{f}^{(-k)}$  is the estimated regression function  $f$  using the training data set

$\mathcal{D} - \mathcal{D}^{(k)}$ . Typically,  $K$  is set to be five or ten and  $CV(T)$  is minimised over a grid of values of  $T$ . Let  $\hat{T}$  be the minimiser of  $CV(T)$ . Our final estimate of  $f$  is based on  $\hat{T}$  and the full data set.

## 4. NUMERICAL EXAMPLES

We now illustrate the finite sample performance of the proposed method through several sets of numerical examples.

### 4.1 Partial single index model

We begin with a partial single index model from Carroll et al. (1997).

$$(36) \quad y = \sin \left[ \frac{\pi(\alpha' \mathbf{x} - A)}{B - A} \right] + \beta z + \epsilon,$$

where  $z$  is a binary variable taking value 1 and 0 respectively for half of the observations,  $\mathbf{x} = (x_1, x_2, x_3)'$  and  $x_i$  are independent observations from a uniform distribution  $U[0, 1]$ ,  $A = \sqrt{3}/2 - 1.645/\sqrt{12}$ ,  $B = \sqrt{3}/2 + 1.645/\sqrt{12}$ , and  $\epsilon \sim N(0, 0.1^2)$ . The parameters used in this example are  $\alpha = (1, 1, 1)'/\sqrt{3}$  and  $\beta = 0.3$ . Following Carroll et al. (1997), the sample size is  $n = 200$ .

The estimate we developed in Section 2 can be easily extended to accommodate the linear component  $\beta z$ . To estimate  $f$  and  $\beta$ , we minimize

$$(37) \quad \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - f(\mathbf{x}_i) - \beta z)^2 + \lambda_n J(f) \right\},$$

with respect to both  $f$  and  $\beta$ . We apply this estimate to one hundred simulated data sets. The left panel of Figure 1 shows a typical simulated data set together with the estimated regression functions for both values of  $Z$ . The overall performance is also summarized in Table 1 and the box plots presented in the right panel of Figure 1. These results compare favorably with those reported by Carroll et al. (1997).

### 4.2 Additive index model

To gain further insight of predictive performance of the additive index model and the proposed penalized least squares estimate, we consider a more complex model. The simulated data consist of observations from the regression model  $y = f_0(\mathbf{x}) + \epsilon$  where

$$f_0(\mathbf{x}) = 5 \sin(2\pi\alpha'_1 \mathbf{x}) + \frac{4 \sin(2\pi\alpha'_2 \mathbf{x})}{2 - \sin(2\pi\alpha'_2 \mathbf{x})} + 6 \times \{ 0.1 \sin(2\pi\alpha'_3 \mathbf{x}) + 0.2 \cos(2\pi\alpha'_3 \mathbf{x}) + 0.3 \sin^2(2\pi\alpha'_3 \mathbf{x}) + 0.4 \cos^3(2\pi\alpha'_3 \mathbf{x}) + 0.5 \sin^3(2\pi\alpha'_3 \mathbf{x}) \}.$$

For each simulated data set, the predictor  $\mathbf{x}$  is generated in two steps: we first sample  $(\mathbf{z}, w)$  from a uniform distribution

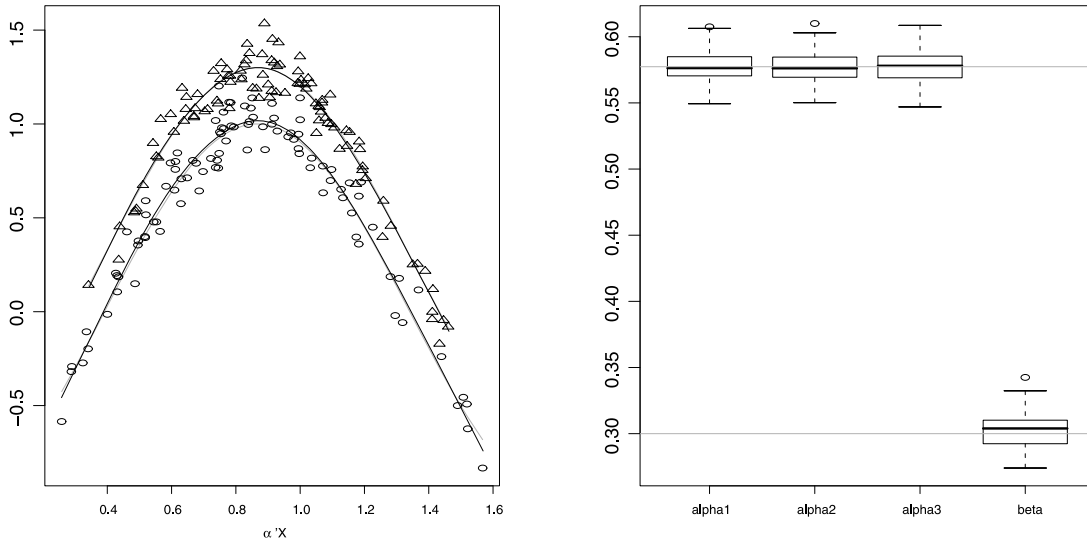


Figure 1. Simulation Example: left panel gives a typical simulated data set with size 200. The open circles correspond to observations with  $z = 0$  and triangles correspond to observations with  $z = 1$ . The solid black lines are the estimated functions for  $z = 0$  and  $1$  respectively. The grey lines, almost coincide with the estimated functions represent the true regression function. The right panel summarizes the parameter estimate for 100 repetitions. The grey horizontal lines correspond to the true values for  $\alpha$  and  $\beta$  respectively.

Table 1. Mean values of the parameter estimation for the simulated example. The numbers in the parentheses are the standard errors

$\alpha_1$	$\alpha_2$	$\alpha_3$	$\beta$
0.575 (0.003)	0.576 (0.001)	0.579 (0.002)	0.302 (0.001)

over  $[0, 1]^{11}$  and then compute  $x_j = (z_j + tu)/(1+t)$  for some  $t$ . We consider  $t = 0$  and  $t = 1$ , corresponding to uncorrelated and correlated predictors respectively. The projection directions are also generated in two steps: a  $10 \times 10$  random matrix was first generated and then we normalize each column so that it has norm one. We discard all the random matrices whose smallest eigenvalue is smaller than 0.1 to ensure that the projection directions are sufficiently linearly independent. The regression noise  $\epsilon \sim N(0, (5/3)^2)$  to yield a signal to noise ratio 3:1. For each simulated data set, we generate 500 observations and estimate  $f$  using the proposed penalized least squares estimate. Note that we assumed that  $q$  is unknown in our estimation procedure. To measure the performance of an estimate of  $f_0$ , we use the mean squared error defined as

$$(38) \quad MSE(\hat{f}) = E_{\mathbf{x}} \left\{ \hat{f}(\mathbf{x}) - f_0(\mathbf{x}) \right\}^2$$

where the expectation is taken with respect to  $\mathbf{x}$ . To evaluate the MSE, we generate an additional 10,000  $\mathbf{x}$ 's in the

same fashion as before. We then estimate the MSE using its sample version:

$$(39) \quad \widehat{MSE}(\hat{f}) = \frac{1}{10000} \sum_{i=1}^{10000} \left\{ \hat{f}(\mathbf{x}_i^*) - f_0(\mathbf{x}_i^*) \right\}^2$$

where  $\mathbf{x}_i^*$ 's are the additionally generated  $\mathbf{x}$ s. For comparison, we also included the MARS (Friedman, 1991) and the projection pursuit regression of Friedman and Stuetzle (1981). The MARS is a nonparametric regression technique with built-in variable selection features. For the projection pursuit regression, we use five fold cross validation to determine the number of components. The left panel of Figure 2 shows the boxplot of the MSE for the three methods when  $t = 0$  summarized over 200 simulated data sets. The right panel corresponds to  $t = 1$ . From Figure 2, we observe that the proposed penalized least squares estimate enjoys superior performance in estimating the regression function  $f$ .

### 4.3 Generalized additive index model

Our last example illustrates how the the proposed method can be applied to more general regression settings. Although we have focused on mean regression, the additive index model can also be extended to more general regression settings such as the generalized linear model. In the generalized regression, the conditional mean  $\mu(\mathbf{x}) = E(y|\mathbf{x})$  is related to a canonical function  $\eta(\mathbf{x})$  via the so-called link function  $\eta = g(\mu)$  which is known. For example the usual Gaussian regression amounts to the  $g(\mu) = \mu$ . A popular example

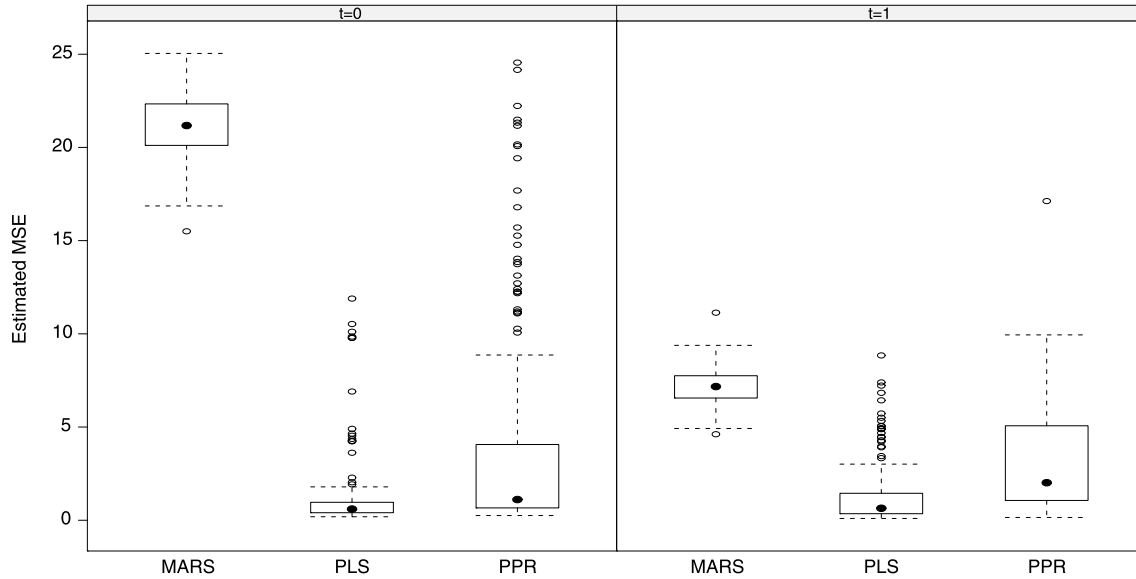


Figure 2. Simulation Example: boxplot of estimated MSE for three methods: MARS, the proposed penalized least squares estimate (PLS) and the original projection pursuit regression (PPR).

Table 2. Estimated projection directions for the Pima Indian Diabetes data

	term 1	term 2	term 3
pregnant	0.01	0.12	-0.15
glucose	-0.24	0.36	-0.05
pressure	0.45	0.26	-0.03
triceps	-0.16	-0.14	0.47
insulin	0.25	0.33	-0.30
mass	0.36	0.78	-0.07
pedigree	-0.32	0.21	-0.09
age	-0.65	-0.09	0.80

of the generalized regression is the logistic regression for binary responses where  $\eta(\mathbf{x})$  is the conditional log-odds ratio

$$(40) \quad \eta(\mathbf{x}) = \log \frac{\Pr(y = 1|\mathbf{x})}{\Pr(y = 0|\mathbf{x})}.$$

The idea of the additive index model can be generalized to these situations by assuming

$$(41) \quad \eta(\mathbf{x}) = \mu + h_1(\alpha'_1 \mathbf{x}) + \dots + h_q(\alpha'_q \mathbf{x}).$$

The generalized additive index model can be estimated in a similar manner as the usual additive index model. A deviance function such as the negative log-likelihood is used in place of the least squares of (8):

$$(42) \quad \mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n \ell[\eta(\mathbf{x}_i); y_i] + \lambda J(f).$$

For example, in the case of logistic regression, we estimate the generalized additive index model by minimizing

$$(43) \quad \mathcal{L}(f) = \frac{1}{n} \sum_{i=1}^n \left[ y_i \eta(\mathbf{x}_i) - \log \left( 1 + e^{\eta(\mathbf{x}_i)} \right) \right] + \lambda J(f),$$

where  $\eta(\cdot)$  is given by (41). The readers are referred to McCullagh and Nelder (1989) for more general discussion regarding the choice of the deviance. Similar to the generalized linear model  $\mathcal{L}$  can be minimized by the iteratively re-weighted least squares. In each iteration, the deviance function is approximated by a weighted least squares and can therefore be optimized using the algorithm given in Section 2. According to our experience, it suffices to run our algorithm only for one iteration at each step in the iteratively re-weighted least squares and the algorithm usually converges in a few iterations.

To illustrate the idea of the generalized additive index model, we consider here a real data example. The Pima Indians Diabetes data have 768 observations on nine variables. The purpose is to predict whether or not a particular subject has diabetes using eight remaining variables. It was often used as a benchmark data set for classification. An application of the additive index model and the penalized least squares estimate yields an additive model with three linear factors. The estimated projection directions are given in Table 2 and their corresponding ridge functions are given in Figure 3. The classification error of the additive index model is estimated as 18.3% using ten fold cross validation, which compares favorably with other popular classification tools (Newman et al., 1998).

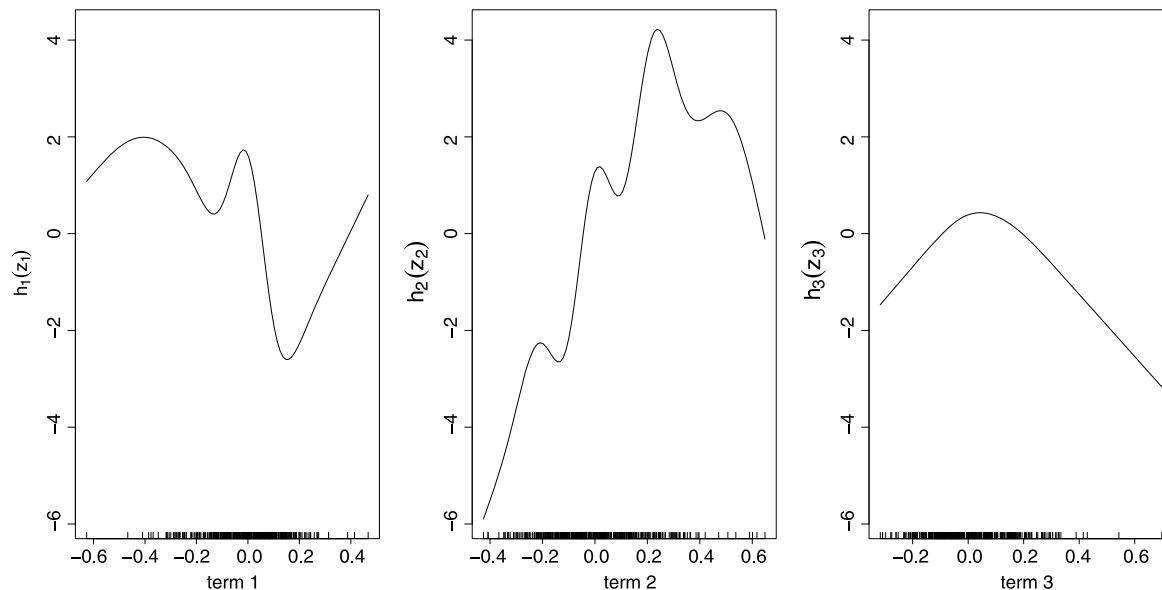


Figure 3. Estimated ridge functions for the Pima Indian Diabetes data.

## 5. CONCLUSION

In this paper, we propose a penalized least squares estimator to simultaneously determine the number of factors and estimate the regression function in the additive index model. The method becomes the recently proposed COSSO estimate if the factors are known in advance and the additive index model reduces to the additive model. We show the regression function estimate can achieve the optimal convergence rate. Thanks to the efficient algorithm for computing the smoothing spline estimate and their derivatives, the proposed estimate can be computed using an iterative algorithm.

Numerical experiments have been conducted to demonstrate the flexibility of the additive index model and the efficiency of the proposed estimate. We also used a real data example to show that the idea of the additive index model can be extended beyond the usual mean regression which is the primary focus of this paper.

Received 17 October 2010

## REFERENCES

- BREIMAN, L. (1995), Better subset regression using the nonnegative garrote, *Technometrics*, **37**, 373–384. [MR1365720](#)
- CARROLL, P., FAN, J., GIJBELS, I. and WAND, M. (1997), Generalized partial linear single-index models, *Journal of the American Statistical Association*, **92(438)**, 477–489. [MR1467842](#)
- CHEN, H. (1991), Estimation of a projection-pursuit type regression model, *Annals of Statistics*, **19(1)**, 142–157. [MR1091843](#)
- CHIOU, J. and MÜLLER, H. (2004), Quasi-likelihood regression with multiple indices and smooth link and variance functions, *Scandinavian J. Statistics*, **31**, 367–386. [MR2087831](#)
- DIACONIS, P. and SHAHSHAHANI, M. (1984), On nonlinear functions of linear combinations, *SIAM Journal of Scientific Computing*, **5(1)**, 175–191. [MR0731890](#)

- FRIEDMAN, J. (1991), Multivariate adaptive regression splines, *Annals of Statistics*, **19**, 1–67. [MR1091842](#)
- FRIEDMAN, J. and STUETZLE, W. (1981), Projection pursuit regression, *Journal of the American Statistical Association*, **76(376)**, 817–823. [MR0650892](#)
- LIN, Y. and ZHANG, H. (2006), Component selection and smoothing in smoothing spline analysis of variance models. *Annals of Statistics*, **34**, 2272–2297. [MR2291500](#)
- MCCULLAGH, P. and NELDER, J. (1989), *Generalized Linear Models*, London: Chapman and Hall. [MR0727836](#)
- NEWMAN, D., HETTICH, S., BLAKE, C. and MERZ, C. (1998), UCI repository of machine learning dataset [http://www.ics.uci.edu/~mllearn/MLRepository.html], Irvine, CA: University of California, Department of Information and Computer Science.
- STONE, C. (1985), Additive regression and other nonparametric models, *Annals of Statistics*, **13**, 689–705. [MR0790566](#)
- VAN DE GEER, S. (2000), *Empirical Processes in M-Estimation*, Cambridge: Cambridge University Press.
- WAHBA, G. (1990), *Spline Models for Observational Data*, Philadelphia: SIAM. [MR1045442](#)
- YUAN, M. (2008), On the identifiability of additive index model, Manuscript.

Lingyan Ruan  
H. Milton Stewart School  
of Industrial and Systems Engineering  
Georgia Institute of Technology  
USA

Ming Yuan  
H. Milton Stewart School  
of Industrial and Systems Engineering  
Georgia Institute of Technology  
USA  
E-mail address: [myuan@isye.gatech.edu](mailto:myuan@isye.gatech.edu)