## American Society for Quality

# Simultaneous Variable Selection

**Berwin A. TURLACH**

School of Mathematics and Statistics (M019)
The University of Western Australia
Crawley, WA 6009, Australia

**William N. VENABLES**

CSIRO Mathematical and Information Sciences
Cleveland, Qld 4163, Australia

**Stephen J. WRIGHT**

Computer Sciences Department
University of Wisconsin
Madison, WI 53706

We propose a new method for selecting a common subset of explanatory variables where the aim is to model *several* response variables. The idea is a natural extension of the LASSO technique proposed by Tibshirani (1996) and is based on the (joint) residual sum of squares while constraining the parameter estimates to lie within a suitable polyhedral region. The properties of the resulting convex programming problem are analyzed for the special case of an orthonormal design. For the general case, we develop an efficient interior point algorithm. The method is illustrated on a dataset with infrared spectrometry measurements on 14 qualitatively different but correlated responses using 770 wavelengths. The aim is to select a subset of the wavelengths suitable for use as predictors for as many of the responses as possible.

KEY WORDS: Constrained least squares problem; Constrained regression; Convex programming; Infrared spectrometry; Interior-point algorithm; Quadratic programming; Subset selection; Variable selection.

## 1. INTRODUCTION

Many practical (linear) regression problems are ill-conditioned. When the problem contains a large number of highly correlated predictors, the need to select predictors carefully, or otherwise regularize the problem, is well known. Some traditional techniques used for this purpose are direct variable selection (Miller 1990; Burnham and Anderson 1998), ridge regression (see, e.g., Hocking 1996; Draper and Smith 1998), and partial least squares (Wold 1984; Martens and Naes 1989; Brown 1993). The latter is typically used if the number of explanatory variables is large relative to the number of observations. Newer techniques that have been proposed include the *nonnegative garotte* of Breiman (1995) and the *least absolute shrinkage and selection operator* (LASSO) of Tibshirani (1996). Later we discuss the LASSO in more detail and propose a method that extends the LASSO methodology in a natural way to the problem in which several related response variables are observed and the researcher desires, either because of available a priori information or for other reasons, to model these response variables using the same common subset of predictors.

Breiman and Friedman (1997) discussed various applications in which the aim is to model several related response variables using the same set of predictors, and proposed a method that uses the relationship between the response variables to find a "simultaneous" model for each response variable. They showed that such a simultaneous model can outperform an approach in which each response variable is modeled separately. However, their approach uses all available predictors to simultaneously build models for all of the response variables, and these authors did not address the question of variable selection.

Using a Bayesian approach, Brown, Fearn, and Vannucci (1999) and Brown, Vannucci, and Fearn (1998, 2002) addressed the question of variable selection in the setting where one has several related response variables and a (large) set of predictors

to choose from. However, their methods require the use of quite sophisticated Markov chain Monte Carlo (MCMC) algorithms, for which the choice of tuning parameters and the monitoring for convergence do not appear to be trivial. By way of contrast, the method that we propose for variable selection in this setting is based on a regularization approach inspired by the LASSO methodology. Although we are aware that many variable selection procedures that use a regularization approach, including the LASSO, can be explained via a Bayesian framework (see, e.g., Leamer 1978), we do not develop a Bayesian interpretation for our method in this article.

The LASSO technique minimizes the residual sum of squares while bounding the $L_1$-norm of the coefficient vector by a specified value. Suppose that we observe data on a response variable $y_i$ and $p$ explanatory variables $x_{il}$, $i = 1, \ldots, n$ and $l = 1, \ldots, p$, and that the response variable is centered ($\sum_i y_i = 0$) and the explanatory variables are standardized ($\sum_i x_{il} = 0$ and $\sum_i x_{il}^2/n = 1$ for all $l = 1, \ldots, p$). Then the LASSO estimates are given by the solution to the following optimization problem:

$$\underset{b_1,\ldots,b_p}{\text{minimize}} \quad \frac{1}{2} \sum_{i=1}^{n} \left( y_i - \sum_{l=1}^{p} x_{il} b_l \right)^2 \qquad (1a)$$

$$\text{subject to} \quad \sum_{l=1}^{p} |b_l| \leq t. \qquad (1b)$$

Tibshirani (1996) showed that this approach has features in common with both ridge regression and variable selection. As in ridge regression, the solution $\hat{b}_i$ of (1) tends to shrink to 0 as $t$ goes to 0. In contrast, the nonsmooth nature of the $L_1$-norm,

which is nondifferentiable when any components $b_l$ are 0, tends to force some of the solution components $\hat{b}_i$'s to be 0. In this sense, the outcome is similar to variable selection.

In this article we propose a method that extends the LASSO methodology in a natural way to the problem in which several related response variables are observed and the researcher's aim is to find predictors for all of them from a common subset of variables. The dataset that motivated this research and that we use to illustrate our methodology was kindly provided by Dr. Bronwyn Harch of CMIS/CSIRO in Adelaide. In this dataset, experimenters used 24 soil samples to take measurements on 14 quantities (EC, pH, pHCaCl2, CLeco, Org.C, NLeco, extP, Ca, Mg, Na, K, TotalCations, CEC, and CaCO3) at 770 different wavelengths. The aim is to identify those wavelengths (explanatory variables) that are the most informative for detecting and quantifying the presence of a particular quantity, say Org.C.

Applying the LASSO methodology to each response variable of this dataset, choosing values for $t$ between 0 and 1, we obtain the results shown in Figure 1. Here the 770 different wavelengths, spaced approximately equally along a spectrum, are labeled X1 to X770 for simplicity. In each panel, the abscissa is the constraint bound $t$ in (1b) and the ordinate is the coefficient value. (Because the predictors are scaled to mean 0 and unit variance, the coefficients are also on comparable scales.) Note that for most of the response variables, only a very few of the coefficients are nonzero for any given value of $t$, but the set of nonzero coefficients depends strongly on $t$. Typically a regressor enters the model (i.e., has a nonzero coefficient) and drops out again to be replaced by a "nearby" regressor. For example, consider the panel for CaCO3. Initially, for small $t$, X116 is selected by the LASSO. As $t$ increases to about .8, X116 is replaced by X117, which in turn is replaced by X119 as $t$ approaches 1. These three regressors are highly correlated. Given this high correlation among regressors, another recently proposed modification of the LASSO, the "elastic net" (Zou and Hastie 2005), might be more appropriate if one wants to model the response variables separately.

A feature in Figure 1 that is not so obvious is that at $t = 1$, the set of variables selected for most response variables includes only a few regressor variables. Moreover, the indices of these regressor variables lie mostly in a few particular regions of the spectrum.

Figure 2 shows a pairwise scatterplot of the 14 response variables. We observe at least some correlation between most of the response variables, in particular among EC, pH, pH-CaCl2, Cleco, Org.C, NLeco, Ca, TotalCations, and CEC. With regard to NLeco, we would also like to point out that one observation—namely, the 10th—appears to be an outlier.

The correlation among the response variables and the results of the separate LASSO analysis suggest that possibly a single set of regressor variables is sufficient to model all (or at least most) of the response variables. One may also expect that by using all response variables *simultaneously* to select a single set of regressor variables, it is possible to avoid overfitting, which is a potentially serious problem if we select a separate set of regressor variables for each response variable.

## 1.1 Extending the LASSO to Multiple Responses

We propose extending the LASSO methodology to achieve simultaneous variable selection. To fix notation, suppose that we have $n$ observations on $k$ response variables $y_{ij}$ and $p$ explanatory variables $x_{il}$ ($i = 1, \ldots, n$, $j = 1, \ldots, k$, and $l = 1, \ldots, p$). We assume not only that the explanatory variables are standardized, as described earlier, but also that the response variables are standardized, that is, $\sum_i y_{ij} = 0$ and $\sum_i y_{ij}^2/n = 1$ for all $j = 1, \ldots, k$. We might interpret the regression parameter $b_{lj}$ as the "explanatory power" that the $l$th regressor variable has on the $j$th response variable. It seems natural to take $b_{l,\max} = \max(|b_{l1}|, \ldots, |b_{lk}|)$ as a measure of the "simultaneous explanatory power" of the $l$th regressor on all $k$ response variables. Following the approach of Tibshirani (1996), we may now impose a constraint on the sum of the $b_{l,\max}$, $l = 1, \ldots, p$, to identify the regressor variables that simultaneously best explain all response variables. Thus we arrive at the following problem:

$$\text{minimize}_{b_{11}, \ldots, b_{pk}} \quad \frac{1}{2} \sum_{j=1}^{k} \sum_{i=1}^{n} \left( y_{ij} - \sum_{l=1}^{p} x_{il} b_{lj} \right)^2, \qquad (2a)$$

$$\text{subject to} \quad \sum_{l=1}^{p} \max(|b_{l1}|, \ldots, |b_{lk}|) \leq t. \qquad (2b)$$

Note that if $k = 1$, then (2) reduces to the LASSO (1). It should also be noted that we propose to use (2) as an exploratory tool to identify a suitable subset of regressor variables. Once this subset is identified, we suggest that its suitability for modeling most (or all) of the response variables is assessed further using standard statistical techniques, and that the selected regressor variables be used in *unconstrained* (linear) models. It is not clear to us that the actual parameter estimates at the solution of (2) have any inherent meaning or use.

Although for the dataset that we use to illustrate our methodology, it was not crucial that the same set of predictors can be used to model all response variables, there may be situations where the ability to identify a set of predictors to model all response variables is crucial. For example, a manufacturer of high-frequency measurement devices produces an instrument that is designed to meet several different specifications (likely correlated) for all carrier frequencies in a given range. (We thank K. Kafadar for bringing this example to our attention.) However, production engineers would not be able to afford to test every single frequency (e.g., 1 MHz, ..., 500 MHz) to verify that the instrument coming off the production line passes all specifications at all frequencies. It would be to their advantage to find a small subset of frequencies that can be used to verify performance at *all* specifications. The engineers would be able to save enormous amounts of time by setting their signal generators to only a few frequencies and testing the instrument response to all specifications as the frequency is changed from one setting to the next. We suggest that our methodology might be helpful in identifying such a subset of frequencies.

The rest of the article is structured as follows. In Section 2 we provide some further motivation for the method that we propose and discuss how it relates to similar work by others.

*Figure 1. Applying the LASSO to Each Variable Separately.*

Figure 2. Pairwise Plot of the 14 Response Variables.

In Section 3 we give an exact characterization of the solutions to (2). We also describe a homotopy algorithm that calculates all solutions (as functions of $t$) in the case where the design matrix is orthonormal and develop an interior point algorithm for the general case. Using the latter algorithm, we reanalyze the infrared spectrometry data in Section 4. We provide further discussion on how this method can be extended in Section 5, and offer some conclusions in Section 6.

## 2. SOME MOTIVATION AND DISCUSSION OF RELATED WORK

The proposed methodology can be viewed as a way to select groups of regression estimates. That is, in a (potentially huge) regression problem with $m$ parameters, $\beta_1, \ldots, \beta_m$, we partition the index set $\sigma = \{1, \ldots, m\}$ into $p$ disjoint sets, $\sigma_l$, such that $\sigma = \bigcup_l \sigma_l$, and seek the "most significant" groups of parameters. We allow $\beta_i$ to be nonzero only if $i$ belongs to one of the selected groups $\sigma_l$.

In Section 3 we show that our problem can be viewed from this perspective. Another problem that fits into this setting is variable selection in generalized additive models (Hastie and Tibshirani 1990), as discussed by Bakin (1999). In this case, each nonparametric function in the generalized additive model is built from a $B$-spline basis, and the corresponding coefficients are collected into a group. By deciding which groups are "significant," Bakin (1999) essentially identified those regressor variables that have a significant influence on the response variable.

How can such a groupwise selection of regression parameters be achieved? By generalizing other methods studied in the statistical literature (Leamer 1978; Frank and Friedman 1993; Tibshirani 1996), one might consider imposing the following constraint onto the parameter estimates:

$$\sum_{l=1}^{p} \left( \sum_{i \in \sigma_l} |\beta_i|^\alpha \right)^{1/\alpha} = \sum_{l=1}^{p} \left\| \boldsymbol{\beta}_{\sigma_l} \right\|_\alpha \le t, \qquad (3)$$

where $t \ge 0$ is some constant, $\boldsymbol{\beta}_{\sigma_l}$ is the vector consisting of those $\beta_i$'s for which $i \in \sigma_l$, and $\| \cdot \|_\alpha$ is the $L^\alpha$-norm. If $\alpha \ge 1$, then the feasible region in (3) is a convex subset of $\mathbb{R}^m$. This property is advantageous from both numerical and theoretical standpoints. It ensures that a solution exists if the parameter estimates are defined as the minimizer of a (strictly) convex function. In fact, for a strictly convex objective function, the solution is unique. If the objective function is not strictly convex then one can ensure that the solution is unique under further regularity conditions, provided that $t$ is small enough (see, e.g., the discussion in Osborne, Presnell, and Turlach 2000b for the special case of the LASSO).

When $t$ is small enough, the solution of the regression problem with constraint (3) lies on the boundary of the feasible set; that is, equality holds in (3). Thus it is clear that imposing a constraint like (3) shrinks the parameter estimates toward 0 as $t$ goes to 0. The size of this shrinkage and the manner in which the parameter estimates are shrunk to 0, however, depends on the particular choice of $\alpha$. For $\alpha = 1$, we have

$$\sum_{l=1}^{p} \left\| \boldsymbol{\beta}_{\sigma_l} \right\|_\alpha = \sum_{l=1}^{p} \sum_{i \in \sigma_l} |\beta_i| = \sum_{l=1}^{m} |\beta_l| = \|\boldsymbol{\beta}\|_1,$$

and the constraint (3) reduces to the $L^1$-norm of the complete vector of parameter estimates, which is the constraint used by Tibshirani (1996) in his LASSO method. Given the behavior of the LASSO method (Tibshirani 1996; Osborne et al. 2000a, b), it is clear that this choice does not achieve the desired "simultaneous" variable selection, so the choice $\alpha = 1$ is not interesting in this context.

Arguably, the most obvious choices for $\alpha > 1$ would be $\alpha = 2$ and $\alpha = \infty$. The former choice was studied by Bakin (1999), whereas we study the latter choice $\alpha = \infty$ in this article. Bakin (1999) noted that the use of $\alpha = 2$ can be interpreted as a hybrid between the LASSO (if $p = m$, i.e., each $\sigma_l$ contains exactly one index) and ridge regression (if $p = 1$ and $\sigma_l = \sigma$ for the entire set of indices). Likewise, the use of $\alpha = \infty$ leads to a hybrid between the LASSO ($p = m$) and interval-restricted least squares ($p = 1$) (Clark and Osborne 1988). We note that the optimization problem with $\alpha = 2$ cannot be handled as effectively with currently available optimization techniques as can the problem with $\alpha = \infty$. The latter leads to a convex quadratic program, for which interior point methods can be devised that exploit its special structure, as we show in this article. The former leads to a second-order cone program (see, e.g., Lobo, Vandenberghe, Boyd, and Lebret 1998). Although software is now available for problems of this type (see, e.g., Sturm 1999), it is less able to take advantage of the structure of the problem and as a consequence will probably be less efficient in practice.

## 3. NUMERICAL ASPECTS OF THE ESTIMATOR

We now return to (2). To avoid some cumbersome notation, we introduce the following matrix notation:

$$
\begin{aligned}
\mathbf{y}_j &= \begin{pmatrix} y_{1j} \\ \vdots \\ y_{nj} \end{pmatrix} \in \mathbb{R}^n, \qquad j = 1, \ldots, k, \\
\mathbf{y} &= \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_k \end{pmatrix} \in \mathbb{R}^{nk}, \\
\mathbf{X} &= \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix}, \qquad \text{and} \\
\tilde{\mathbf{X}} &= \mathbf{I}_k \otimes \mathbf{X} = \begin{pmatrix} \mathbf{X} & & \\ & \ddots & \\ & & \mathbf{X} \end{pmatrix} \in \mathbb{R}^{nk \times pk},
\end{aligned}
\qquad (4)
$$

where $\otimes$ denotes the Kronecker product. We arrange the regression parameters $b_{lj}$ in a matrix to implicitly define vectors $\mathbf{b}_1, \ldots, \mathbf{b}_k \in \mathbb{R}^p$ and $\mathbf{b}_{(1)}, \ldots, \mathbf{b}_{(p)} \in \mathbb{R}^k$,

$$
\begin{pmatrix} b_{11} & \cdots & b_{1k} \\ \vdots & \ddots & \vdots \\ b_{p1} & \cdots & b_{pk} \end{pmatrix} = \begin{pmatrix} | & & | \\ \mathbf{b}_1 & \cdots & \mathbf{b}_k \\ | & & | \end{pmatrix} = \begin{pmatrix} - \mathbf{b}_{(1)}^\top - \\ \vdots \\ - \mathbf{b}_{(p)}^\top - \end{pmatrix},
$$

and then define the vector $\mathbf{b}$ by

$$\mathbf{b} = \begin{pmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_k \end{pmatrix} \in \mathbb{R}^{pk}.$$

Using this notation, we write (2) as

$$\underset{\mathbf{b} \in \mathbb{R}^{pk}}{\text{minimize}} \quad f(\mathbf{b}) = \frac{1}{2} (\mathbf{y} - \tilde{\mathbf{X}} \mathbf{b})^\top (\mathbf{y} - \tilde{\mathbf{X}} \mathbf{b}) \qquad (5a)$$

$$\text{subject to} \quad g(\mathbf{b}) = t - \sum_{l=1}^{p} \left\| \mathbf{b}_{(l)} \right\|_\infty \ge 0. \qquad (5b)$$

We also define the vector of residuals as

$$r(\underline{b}) = \underline{y} - \tilde{X}\underline{b}.$$

*Remark 1.* A referee pointed out that occasionally a covariance matrix for errors across responses typically is incorporated in multivariate regression problems. This amounts to changing the objective function in (5a) to $\frac{1}{2}(\underline{y} - \tilde{X}\underline{b})^\top W(\underline{y} - \tilde{X}\underline{b})$, where $W$ is a suitable (symmetric positive semidefinite) weight matrix. We can accommodate this generalization easily by premultiplying $\underline{y}$ and $\tilde{X}$ by $W^{1/2}$. Because it is simple to incorporate a weight matrix $W$ in our formulation, our method can be extended to other models in which the objective function is not the residual sum of squares, but rather is obtained from an iteratively (re)weighted least squares procedure, for example, generalized linear models (McCullagh and Nelder 1989).

*Remark 2.* As stated earlier, we generally assume that the explanatory variables and the response variables are centered and standardized to have (sample) mean 0 and (sample) variance 1. Given the previous remark about the possible incorporation of weights, such standardization may seem questionable. Obviously, just as for the LASSO and other statistical techniques that are in wide use, our proposed method is not invariant under rescaling of the variables (be they explanatory variables or response variables). Thus, we suggest that by default, the variables should be centered and standardized unless the researcher has good reasons not to do so.

## 3.1   Characterization of Solutions

We now use results from convex analysis (Rockafellar 1970; Osborne 1985; Clarke 1990) to characterize solutions of (5). Introducing a Lagrange multiplier $\lambda$ for the constraint (5b), we write the Lagrangian for (5) as

$$\mathcal{L}(\underline{b}, \lambda) = f(\underline{b}) - \lambda g(\underline{b}), \qquad (6)$$

where $\lambda \geq 0$. If we fix $\lambda \geq 0$, then $\mathcal{L}(\underline{b}, \lambda)$ is a convex function in $\underline{b}$ and $\underline{b}$ minimizes $\mathcal{L}(\underline{b}, \lambda)$ if and only if the $pk$-dimensional null-vector $\underline{0}$ is an element of the subdifferential $\partial_{\underline{b}} \mathcal{L}(\underline{b}, \lambda)$ (Osborne 1985, p. 23). From (6), we have

$$\partial_{\underline{b}}\mathcal{L}(\underline{b}, \lambda) = -\tilde{X}^\top \underline{r} + \lambda \underline{v},$$

where $\underline{r} = r(\underline{b}) = \underline{y} - \tilde{X}\underline{b}$ denotes the residual vector and $\underline{v} = (v_1, \ldots, v_{pk})^\top$ has the following form:

- If $\|\underline{b}_{(l)}\|_\infty > 0$, then $\underline{v}_{(l)} = (v_{l1}, \ldots, v_{lk})^\top$, where $\sum_{j=1}^k |v_{lj}| = 1$ and, for $j = 1, \ldots, k$, we have $v_{lj} \geq 0$ if $b_{lj} = \|\underline{b}_{(l)}\|_\infty$, $v_{lj} \leq 0$ if $b_{lj} = -\|\underline{b}_{(l)}\|_\infty$ and $v_{lj} = 0$ if $|b_{lj}| \neq \|\underline{b}_{(l)}\|_\infty$.
- If $\|\underline{b}_{(l)}\|_\infty = 0$, then $\underline{v}_{(l)} = (v_{l1}, \ldots, v_{lk})^\top$, where $\sum_{j=1}^k |v_{lj}| \leq 1$.

Thus if $\underline{b}$ minimizes $\mathcal{L}(\underline{b}, \lambda)$ for a given value of $\lambda$, then

$$\underline{0} = -\tilde{X}^\top \underline{r} + \lambda \underline{v}, \qquad (7)$$

for some $\underline{v}$ of the form described earlier, and $\bar{\underline{r}} = r(\bar{\underline{b}}) = \underline{y} - \tilde{X}\bar{\underline{b}}$. The properties of $\bar{\underline{v}}$ imply that $\bar{\underline{v}}^\top \bar{\underline{b}} = \sum_{l=1}^p \|\bar{\underline{b}}_{(l)}\|_\infty$, so it follows from (7) that

$$\lambda = \bar{\underline{r}}^\top \tilde{X}\bar{\underline{b}} \Big/ \sum_{l=1}^p \|\bar{\underline{b}}_{(l)}\|_\infty. \qquad (8)$$

For $\bar{\underline{b}}$ to be a solution of (5), we require not only that (7) holds for some vector $\bar{\underline{v}}$ satisfying the foregoing properties, but also that $\bar{\underline{b}}$ satisfies the constraint (5b), that $\lambda$ satisfying (8) has $\lambda \geq 0$, and that the following complementarity condition holds:

$$\lambda g(\bar{\underline{b}}) = \lambda \left( t - \sum_{l=1}^p \|\bar{\underline{b}}_{(l)}\|_\infty \right) = 0.$$

In the case of $\lambda = 0$, we have from (7) that $\tilde{X}^\top \bar{\underline{r}} = 0$, indicating that $\bar{\underline{b}}$ is the unconstrained least squares minimizer of (5a).

Although equation (7) gives a characterization of the solution for (5), the highly nonlinear way in which $\underline{v}$ depends on $\underline{b}$ makes it impossible to calculate the solution directly from the characterization just described. Some sort of iterative algorithm is needed. For general $X$, an interior-point algorithm is developed in Section 3.3. The next section discusses the special case in which $X$ is an orthonormal matrix.

## 3.2   The Orthonormal Design Case

In this section we assume that $X$ is orthonormal and that $n > p$, so that $X^\top X = I_p$. For the LASSO, we can find explicit formulas for the LASSO estimate based on the unconstrained least squares estimate. Unfortunately, similar formulas do not seem to be available for the current problem. We can, however, develop a homotopy method, in which the constraint bound $t$ becomes the homotopy parameter, and we can examine the behavior of the solution to (5) as $t$ varies. A similar analysis was given by Osborne (1992) for the case of quantile regression, and by Osborne et al. (2000a) and Efron, Hastie, Johnstone, and Tibshirani (2004) for the LASSO. Specifically, the analysis shows that the solution $\underline{b}$ of (2) is piecewise linear as a function of $t$ and gives further insight into how our method selects variables simultaneously.

We start by noting that, because $\tilde{X}^\top \tilde{X} = I_{pk}$, the unconstrained minimizer of (5a) is

$$\underline{b}^0 = (b_1^0, \ldots, b_{pk}^0)^\top = \tilde{X}^\top \underline{y}.$$

Assuming that the $b_{lj}$'s are, for some nonnegative quantities $\rho_l$, $l = 1, \ldots, p$, of the form

$$b_{lj} = \text{sign}(b_{lj}^0) \times \min(|b_{lj}^0|, \rho_l), \qquad l = 1, \ldots, p, \ j = 1, \ldots, k, \qquad (9)$$

we now show, by specifying the dependence of these $\rho_l$'s on $t$, that (9) indeed yields the solution of (5). Observe that as long as $\rho_l \leq \|\underline{b}_{(l)}^0\|_\infty$, for $l = 1, \ldots, p$, we have $\rho_l = \|\underline{b}_{(l)}\|_\infty$. By using $\tilde{X}^\top \tilde{X} = I_{pk}$ again, we rewrite $f(\underline{b})$ as

$$f(\underline{b}) = \frac{1}{2}\{(\underline{y} - \tilde{X}\underline{b}^0)^\top(\underline{y} - \tilde{X}\underline{b}^0) + (\underline{b}^0 - \underline{b})^\top(\underline{b}^0 - \underline{b})\}. \qquad (10)$$

From (9), we have that $|b_{lj}^0 - b_{lj}| = (|b_{lj}^0| - \rho_l)_+$, where $(x)_+ = \max(0, x)$. By using this observation in conjunction with (10), we reformulate (5) as

$$\underset{\rho_1,\dots,\rho_l}{\text{minimize}} \quad \frac{1}{2} \sum_{j=1}^{k} \sum_{l=1}^{p} (|b_{lj}^0| - \rho_l)_+^2 \tag{11a}$$

$$\text{subject to} \quad \sum_{l=1}^{p} \rho_l = t. \tag{11b}$$

We now define $\sigma \subseteq \{1, \dots, p\}$ such that if $l \notin \sigma$, then $\rho_l = 0$; that is, $\sigma$ is the set of indices $l$ for which $\rho_l$ may be different from 0. Furthermore, for each $l = 1, \dots, p$, let $\sigma_l \subseteq \{1, \dots, k\}$ be the set of indices $j$ such that $|b_{lj}^0| > \rho_l$ if and only if $j \in \sigma_l$. We then rewrite (11) as

$$\underset{\rho_1,\dots,\rho_l}{\text{minimize}} \quad \frac{1}{2} \sum_{l \in \sigma} \sum_{j \in \sigma_l} (|b_{lj}^0| - \rho_l)^2,$$

$$\text{subject to} \quad \sum_{l \in \sigma} \rho_l = t.$$

Of course, we also require that $\rho_l \geq 0$ for $l \in \sigma$. By introducing a Lagrange multiplier $\mu$ for the constraint in this problem, we obtain from the optimality conditions that the solution must satisfy the relations

$$\mu = \sum_{j \in \sigma_l} (|b_{lj}^0| - \rho_l) = \left( \sum_{j \in \sigma_l} |b_{lj}^0| \right) - n_l \rho_l,$$

$$l \in \sigma \text{ and } \rho_l > 0, \tag{12a}$$

and

$$t = \sum_{l \in \sigma} \rho_l, \tag{12b}$$

where $n_l = |\sigma_l|$ denotes the number of elements in the set $\sigma_l$, $l = 1, \dots, p$.

We now use the Karush–Kuhn–Tucker (KKT) conditions (12) to show that the $\rho_l$'s [and hence, by (9), the $b_{lj}$'s] are piecewise linear functions of $t$. Specifically, we show that there is a sequence of "knots," $0 = t_0 < t_1 < t_2 < \dots < t_m$, such that the $\rho_l$'s are linear in $t$ for $t \in [t_{i-1}, t_i]$ for $i = 1, 2, \dots, m$. At each of the $t_i$'s either $\sigma$ changes by adding one or more indices or one or more $\sigma_l$'s change by dropping one or more indices, or both these events happen.

For $t_0 = 0$, we set $\sigma = \{l : \|\mathbf{b}_{(l)}^0\|_1 = \max_{l=1,\dots,p} \|\mathbf{b}_{(l)}^0\|_1\}$ and $\rho_l(t_0) = 0$ for $l = 1, \dots, p$. Clearly, this is the optimal solution for $t = 0$ and fulfills the KKT conditions (12). An iterative homotopy algorithm proceeds as follows: Assume that we are at $t_i$; then define, for each $l \in \sigma$,

$$\rho_l(t) = \rho_l(t_i) + \frac{1}{\sum_{l \in \sigma} 1/n_l} \frac{1}{n_l} (t - t_i).$$

Note that each $\rho_l$ is a linear function of $t$ and that if $\sum_{l \in \sigma} \rho_l(t_i) = t_i$, then $\sum_{l \in \sigma} \rho_l(t) = t$ for all $t > t_i$, provided neither the set $\sigma$ nor any of the sets $\sigma_l$ changes. Hence these $\rho_l$'s fulfill the KKT condition (12b) for $t$ between $t_i$ and $t_{i+1}$.

Furthermore, for each $l \in \sigma$, we define

$$\mu_l(t) = \left( \sum_{j \in \sigma_l} |b_{lj}^0| \right) - n_l \rho_l(t).$$

Because of the definition of the $\rho_l(t)$'s, this definition ensures that the optimality condition (12a) holds for $t > t_i$ whenever it holds at $t_i$, as long as none of the sets $\sigma_l$ changes. That is, $\mu_{l_1}(t) = \mu_{l_2}(t) = \mu(t)$ for any $l_1, l_2 \in \sigma$.

We conclude that the $\rho_l(t)$'s defined earlier are the solutions to (11) for all $t$'s between $t_i$ and $t_{i+1}$. It remains to determine the next knot $t_{i+1}$. To find $t_{i+1}$, we calculate $\tau_l^*$ such that $\rho_l(\tau_l^*) = \min_{j \in \sigma_l}(|b_{lj}^0|)$, for each $l \in \sigma$. Let $\tau^* = \min_{l \in \sigma} \tau_l^*$ be the constraint bound $t$ at which one or more of the $\sigma_l$'s change by dropping one or more indices.

Furthermore, for some $l_0 \in \sigma$, let $\mu(t) = \mu_{l_0}(t)$. [As noted earlier, the $\mu_l(t)$, $l \in \sigma$, are all identical.] Then, provided that not all variables have yet entered the model (i.e., $\sigma \neq \{1, \dots, p\}$), we calculate $\tau^\dagger$ such that $\mu(\tau^\dagger) = \max_{l \notin \sigma} \|\mathbf{b}_{(l)}^0\|_1$. In other words, $\tau^\dagger$ is the constraint bound $t$ at which $\sigma$ changes by adding one or more indices. In the alternative case of $\sigma = \{1, \dots, p\}$, we calculate $\tau^\dagger$ such that $\mu(\tau^\dagger) = 0$; that is, $\tau^\dagger$ is the constraint bound at which we reach the unconstrained solution.

We then set $t_{i+1} = \min(\tau^*, \tau^\dagger) > t_i$. If $t_{i+1} = \tau^\dagger$ and $\mu(t_{i+1}) \neq 0$, then we update $\sigma$ to $\sigma = \sigma \cup \{l : \|\mathbf{b}_{(l)}^0\|_1 = \mu(t_{i+1})\}$. If $\mu(t_{i+1}) = 0$, then we have reached the unconstrained solution and the algorithm stops; otherwise, it continues as described above.

We conclude from this analysis that for $\mathbf{X}$ orthonormal, the solution vector $\mathbf{b}$ of (2) is a (continuous) piecewise linear function of the constraint parameter $t$. This property of the solution vector $\mathbf{b}$ also holds for general $\mathbf{X}$, if $\mathbf{X}$ has full column rank. (The proof of the more general result is omitted but is available from the authors on request.) We believe that extensions of the continuity and piecewise linearity results hold for general $\mathbf{X}$, as in the LASSO (Osborne et al. 2000a; Efron et al. 2004). However, because the solution for each $t$ is not necessarily unique (except under additional assumptions on $\mathbf{X}$), the analysis of this case is somewhat more difficult.

This analysis also gives some insight into how our method selects variables. In the case of an orthonormal design, it essentially orders the variables such that

$$\left\| \mathbf{b}_{(l_1)}^0 \right\|_1 \geq \left\| \mathbf{b}_{(l_2)}^0 \right\|_1 \geq \left\| \mathbf{b}_{(l_3)}^0 \right\|_1 \geq \dots \geq \left\| \mathbf{b}_{(l_{p-1})}^0 \right\|_1 \geq \left\| \mathbf{b}_{(l_p)}^0 \right\|_1,$$

and then selects the variables $x_{il_1}, x_{il_2}, \dots, x_{il_m}$, where $m$ depends on $t$, using this ordering. Note that the unconstrained coefficient estimates $\mathbf{b}_{(l)}^0$ are sorted according to their $L^1$-norms. This shows that the constraint that we propose achieves its "simultaneous" variable selection by measuring the over all contributions of an explanatory variable by summing its (absolute) contribution over all of the $k$ regressions. The variable that is best with respect to this measure is selected first, followed by the variable that is second best with respect to this measure, and so on.

### 3.3 The General Case

In this section we develop an interior point algorithm for solving (2) for general $\tilde{\mathbf{X}}$. First, to express this problem as a convex quadratic program, we define

$$\mathbf{Q} = \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}, \qquad \mathbf{c} = -\tilde{\mathbf{X}}^\top \mathbf{y} \in \mathbb{R}^{pk}, \qquad d = \frac{1}{2} \mathbf{y}^\top \mathbf{y},$$

and use $\underline{u}_l$ to denote an $l$-dimensional vector with all entries equal to 1. By introducing an auxiliary vector $\underline{z} \in \mathbb{R}^p$, we now write (5) as

$$\underset{\underline{b}}{\text{minimize}} \quad \tfrac{1}{2}\underline{b}^\top Q\underline{b} + \underline{c}^\top \underline{b} + d \tag{13a}$$

$$\text{subject to} \quad \underline{u}_k \otimes \underline{z} - \underline{b} \geq \underline{0}, \tag{13b}$$

$$\underline{u}_k \otimes \underline{z} + \underline{b} \geq \underline{0}, \quad \text{and} \tag{13c}$$

$$t - \underline{u}_p^\top \underline{z} \geq 0. \tag{13d}$$

It is well known that convex quadratic programming problems can be solved efficiently using primal–dual infeasible interior-point algorithms (Roos, Terlaky, and Vial 1997; Wright 1997; Ye 1997). We now present a brief derivation of the interior-point approach, as applied to our specific problem (13).

Using the Lagrange multipliers $\underline{\lambda}_l$, $\underline{\lambda}_u \in \mathbb{R}^{kp}$ and $\tau \in \mathbb{R}$, the Lagrangian for (13) is

$$\mathcal{L}(\underline{b}, \underline{z}, \underline{\lambda}_u, \underline{\lambda}_l, \tau) = \tfrac{1}{2}\underline{b}^\top Q\underline{b} + \underline{c}^\top \underline{b} + d - \underline{\lambda}_u^\top (\underline{u}_k \otimes \underline{z} - \underline{b})$$

$$- \underline{\lambda}_l^\top (\underline{u}_k \otimes \underline{z} + \underline{b}) - \tau(t - \underline{u}_p^\top \underline{z}).$$

The optimality conditions for $\underline{b}$ to solve (13) are

$$Q\underline{b} + \underline{c} + \underline{\lambda}_u - \underline{\lambda}_l = \underline{0},$$

$$-(\underline{u}_k^\top \otimes I_p)\underline{\lambda}_u - (\underline{u}_k^\top \otimes I_p)\underline{\lambda}_l + \tau\underline{u}_p = \underline{0},$$

$$\underline{u}_k \otimes \underline{z} - \underline{b} \geq \underline{0},$$

$$\underline{u}_k \otimes \underline{z} + \underline{b} \geq \underline{0},$$

$$t - \underline{u}_p^\top \underline{z} \geq 0,$$

$$\underline{\lambda}_u \geq \underline{0}, \qquad \underline{\lambda}_l \geq \underline{0}, \qquad \tau \geq 0,$$

and

$$\underline{\lambda}_u^\top (\underline{u}_k \otimes \underline{z} - \underline{b}) = \underline{0},$$

$$\underline{\lambda}_l^\top (\underline{u}_k \otimes \underline{z} + \underline{b}) = \underline{0},$$

$$\tau(t - \underline{u}_p^\top \underline{z}) = 0.$$

Because the problem is a convex quadratic program, these conditions are sufficient as well as necessary. By introducing slack variables $\underline{s}_u$ and $\underline{s}_l$ in $\mathbb{R}^{kp}$ and $\zeta \in \mathbb{R}$, we restate these conditions in a form that is more convenient for development of the interior point approach:

$$Q\underline{b} + \underline{c} + \underline{\lambda}_u - \underline{\lambda}_l = \underline{0}, \tag{14a}$$

$$-(\underline{u}_k^\top \otimes I_p)\underline{\lambda}_u - (\underline{u}_k^\top \otimes I_p)\underline{\lambda}_l + \tau\underline{u}_p = \underline{0}, \tag{14b}$$

$$\underline{u}_k \otimes \underline{z} - \underline{b} - \underline{s}_u = \underline{0}, \tag{14c}$$

$$\underline{u}_k \otimes \underline{z} + \underline{b} - \underline{s}_l = \underline{0}, \tag{14d}$$

$$t - \underline{u}_p^\top \underline{z} - \zeta = 0, \tag{14e}$$

$$\Lambda_u S_u \underline{u}_{pk} = \underline{0}, \qquad \Lambda_l S_l \underline{u}_{pk} = \underline{0}, \qquad \tau\zeta = 0, \tag{14f}$$

and

$$\underline{\lambda}_u \geq \underline{0}, \qquad \underline{\lambda}_l \geq \underline{0}, \qquad \tau \geq 0,$$
$$\tag{14g}$$
$$\underline{s}_u \geq \underline{0}, \qquad \underline{s}_l \geq \underline{0}, \qquad \zeta \geq 0.$$

(Here we use a standard notational convention from the interior-point literature, namely that if a lowercase and an uppercase letter are used at the same time, then the lowercase letter indicates a vector and the uppercase letter indicates a diagonal matrix whose diagonal elements are the elements of the corresponding vector.) Primal–dual interior-point methods view (14) as a constrained system of nonlinear equations. They seek a root of the function $\mathbf{F}$ defined by the equality conditions in (14), that is,

$$\mathbf{F}(\underline{b}, \underline{z}, \underline{\lambda}_u, \underline{\lambda}_l, \tau, \underline{s}_u, \underline{s}_l, \zeta)$$

$$= \begin{pmatrix} Q\underline{b} + \underline{c} + \underline{\lambda}_u - \underline{\lambda}_l \\ -(\underline{u}_k^\top \otimes I_p)\underline{\lambda}_u - (\underline{u}_k^\top \otimes I_p)\underline{\lambda}_l + \tau\underline{u}_p \\ \underline{u}_k \otimes \underline{z} - \underline{b} - \underline{s}_u \\ \underline{u}_k \otimes \underline{z} + \underline{b} - \underline{s}_l \\ t - \underline{u}_p^\top \underline{z} - \zeta \\ \Lambda_u S_u \underline{u}_{pk} \\ \Lambda_l S_l \underline{u}_{pk} \\ \tau\zeta \end{pmatrix}$$

$$= \mathbf{0}, \tag{15}$$

over the set defined by the inequalities listed in (14g). An important concept in primal–dual methods is the *central path*, which is defined as the solution of the following perturbed variant of (15), for some parameter $\mu > 0$:

$$\mathbf{F}(\underline{b}, \underline{z}, \underline{\lambda}_u, \underline{\lambda}_l, \tau, \underline{s}_u, \underline{s}_l, \zeta) = \begin{pmatrix} \underline{0} \\ \underline{0} \\ \underline{0} \\ \underline{0} \\ 0 \\ \mu\underline{u}_{pk} \\ \mu\underline{u}_{pk} \\ \mu \end{pmatrix}, \tag{16}$$

over the strict interior of the feasible region defined by (14g), that is,

$$\underline{\lambda}_u > \underline{0}, \qquad \underline{\lambda}_l > \underline{0}, \qquad \tau > 0,$$
$$\tag{17}$$
$$\underline{s}_u > \underline{0}, \qquad \underline{s}_l > \underline{0}, \qquad \zeta > 0.$$

Maintenance of strict positivity of these variables at each iteration [versus the nonnegative constraints listed in (14g)] is the origin of the term "interior point."

Interior-point methods of the path-following type (such as the one that we use here) find the solution of (15), subject to the constraints in (14g), by following the central path (16), (17) as $\mu$ decreases to 0. Rather than calculate the central path point exactly for each value of $\mu$, path-following methods take a single Newton-like step toward a point on the central path that is, in a sense, closer to the solution than the current iterate. We define the central path point corresponding to the current iterate by defining $\mu$ as

$$\mu = \frac{\underline{\lambda}_l^\top \underline{s}_l + \underline{\lambda}_u^\top \underline{s}_u + \tau\zeta}{2pk + 1}. \tag{18}$$

(Note that this $\mu$ is the average value of the pairwise products $\lambda_{u,i}s_{u,i}$ and $\lambda_{l,i}s_{l,i}$, $i = 1, 2, \ldots, pk$, and $\tau\zeta$.) We then choose a *centering parameter* $\sigma \in (0, 1)$, and apply a modified Newton

step toward the central path point defined by (16)–(17) in which $\mu$ is replaced by $\sigma\mu$. The modification (due to Mehrotra 1992 and detailed later) enhances the approximation of (15) on which the Newton step is based, making it approach a second-order approximation rather than the usual first-order (linear) approximation. By requiring $\sigma$ to be strictly less than 1, we ensure that the step aims at a point further along the central path than the point corresponding to the current iterate. A heuristic for choosing $\sigma$ was also described by Mehrotra (1992) and is detailed later. Practical variants of this algorithm may contain other important features, such as techniques for determining the distance to move along the calculated step, a method for calculating the starting point, and possibly third- and higher-order modifications to the search direction.

We now focus on the system of equations obtained from the modified Newton step for (16). By defining the residuals at the current point from (15), with the appropriate adjustments for central-path perturbation (the terms involving $\sigma$) and for higher-order enhancement (the terms $\hat{\mathbf{r}}_{hi}$, $\hat{\mathbf{r}}_{lo}$, and $\hat{r}_{\tau\zeta}$), we obtain

$$
\begin{pmatrix} \mathbf{r}_b \\ \mathbf{r}_z \\ \mathbf{r}_u \\ \mathbf{r}_l \\ r_t \\ \mathbf{r}_{hi} \\ \mathbf{r}_{lo} \\ r_{\tau\zeta} \end{pmatrix} := \begin{pmatrix} \mathbf{Q}\mathbf{b} + \mathbf{c} + \boldsymbol{\lambda}_u - \boldsymbol{\lambda}_l \\ -(\mathbf{u}_k^\top \otimes \mathbf{I}_p)\boldsymbol{\lambda}_u - (\mathbf{u}_k^\top \otimes \mathbf{I}_p)\boldsymbol{\lambda}_l + \tau\mathbf{u}_p \\ \mathbf{u}_k \otimes \mathbf{z} - \mathbf{b} - \mathbf{s}_u \\ \mathbf{u}_k \otimes \mathbf{z} + \mathbf{b} - \mathbf{s}_l \\ t - \mathbf{u}_p^\top \mathbf{z} - \zeta \\ \boldsymbol{\Lambda}_u \mathbf{S}_u \mathbf{u}_{pk} - \sigma\mu\mathbf{u}_{pk} + \hat{\mathbf{r}}_{hi} \\ \boldsymbol{\Lambda}_l \mathbf{S}_l \mathbf{u}_{pk} - \sigma\mu\mathbf{u}_{pk} + \hat{\mathbf{r}}_{lo} \\ \tau\zeta - \sigma\mu + \hat{r}_{\tau\zeta} \end{pmatrix}. \quad (19)
$$

The modified Newton system is then

$$
\begin{pmatrix} \mathbf{Q} & & \mathbf{I}_{pk} & -\mathbf{I}_{pk} & & & & \\ & & -\mathbf{u}_k^\top \otimes \mathbf{I}_p & -\mathbf{u}_k^\top \otimes \mathbf{I}_p & \mathbf{u}_p & & & \\ -\mathbf{I}_{pk} & \mathbf{u}_k \otimes \mathbf{I}_p & & & & -\mathbf{I}_{pk} & & \\ \mathbf{I}_{pk} & \mathbf{u}_k \otimes \mathbf{I}_p & & & & & -\mathbf{I}_{pk} & \\ & -\mathbf{u}_p^\top & & & & & & -1 \\ & & \mathbf{S}_u & & & \boldsymbol{\Lambda}_u & & \\ & & & \mathbf{S}_l & & & \boldsymbol{\Lambda}_l & \\ & & & & \zeta & & & \tau \end{pmatrix}
$$

$$
\times \begin{pmatrix} \Delta\mathbf{b} \\ \Delta\mathbf{z} \\ \Delta\boldsymbol{\lambda}_u \\ \Delta\boldsymbol{\lambda}_l \\ \Delta\tau \\ \Delta\mathbf{s}_u \\ \Delta\mathbf{s}_l \\ \Delta\zeta \end{pmatrix} = - \begin{pmatrix} \mathbf{r}_b \\ \mathbf{r}_z \\ \mathbf{r}_u \\ \mathbf{r}_l \\ r_t \\ \mathbf{r}_{hi} \\ \mathbf{r}_{lo} \\ r_{\tau\zeta} \end{pmatrix}. \quad (20)
$$

Thus, at each iteration of our interior-point algorithm, we have to solve systems of equations of the form (20). Although this system of equations is very large, it is also highly structured, and some block elimination greatly reduces its dimension. Details of these algebraic manipulations are given in the technical report (available from the authors on request) on which this article is based.

We continue with some details of our implementation of the Mehrotra algorithm, which actually solves two systems of the form (20) at each iteration, with the same coefficient matrices but different right sides (19). In the first of these systems, we

set $\sigma = 0$ and $\hat{\mathbf{r}}_{hi} = \mathbf{0}$, $\hat{\mathbf{r}}_{lo} = \mathbf{0}$, and $\hat{r}_{\tau\zeta} = 0$, to obtain the *affine-scaling direction*. This direction, which we denote by

$$
(\Delta\mathbf{b}^{\text{aff}}, \Delta\mathbf{z}^{\text{aff}}, \Delta\boldsymbol{\lambda}_u^{\text{aff}}, \Delta\boldsymbol{\lambda}_l^{\text{aff}}, \Delta\tau^{\text{aff}}, \Delta\mathbf{s}_u^{\text{aff}}, \Delta\mathbf{s}_l^{\text{aff}}, \Delta\zeta^{\text{aff}}), \quad (21)
$$

is simply the pure Newton direction for the system of equations $\mathbf{F}$ given in (15). We then find the largest step length $\alpha_{\text{aff}} \in (0, 1]$ such that a step of length $\alpha_{\text{aff}}$ along this direction from the current iterate satisfies the conditions (14g). We then calculate the value $\mu_{\text{aff}}$ from (18) that would occur if we actually took this step, and set

$$
\sigma = \left(\frac{\mu_{\text{aff}}}{\mu}\right)^3, \quad (22)
$$

where $\mu$ is calculated using the current iterate. We use this value of $\sigma$ to form the right side for the second system, and also use the components of (21) to define the residual modifications

$$
\begin{aligned} \hat{\mathbf{r}}_{hi} &= \Delta\boldsymbol{\Lambda}_u^{\text{aff}} \Delta\mathbf{S}_u^{\text{aff}} \mathbf{u}_{pk}, \\ \hat{\mathbf{r}}_{lo} &= \Delta\boldsymbol{\Lambda}_l^{\text{aff}} \Delta\mathbf{S}_l^{\text{aff}} \mathbf{u}_{pk}, \quad (23) \\ \hat{r}_{\tau\zeta} &= \Delta\tau^{\text{aff}} \Delta\zeta^{\text{aff}}. \end{aligned}
$$

We then solve (20) with the new right side to obtain the actual search direction. The heuristic for $\sigma$ in (22) yields a value close to 0 when the pure Newton direction appears to be a profitable search direction. Thus the calculated step will not be much different from the pure Newton direction, and will move quite aggressively to reduce the value of $\mu$ on this iteration. When the affine-scaling direction does not make much progress in reducing $\mu$, the heuristic yields a conservative choice of $\sigma$, closer to 1.

The step length along the search direction is chosen by means of a heuristic due to Mehrotra (1992, sec. 6) and described by Wright (1997, p. 205). The heuristic is modified in an obvious way to account for the fact that our objective function is quadratic rather than linear. This choice of step ensures that the *strict* inequalities (17) are satisfied by the new iterate.

We terminate the algorithm when $\mu$ and the residuals in (19) become sufficiently small. (In our code, we apply the simple test $\mu < 10^{-8}$.) We obtain a starting point by simply setting $\mathbf{z} = \mathbf{0}$ and $\mathbf{b} = \mathbf{0}$, whereas the components of $\boldsymbol{\lambda}_u$, $\boldsymbol{\lambda}_l$, $\mathbf{s}_u$, $\mathbf{s}_l$, $\tau$, and $\zeta$ are all set to some value (in our code, $10^5$). An outline of the overall algorithm is as follows:

1. Choose a starting point.
2. Calculate $\mu$ from (18). If $\mu$ is small enough, then stop.
3. Calculate the affine-scaling direction (21) by solving the system (19)–(20) with $\sigma = 0$ and zero residual modifications.
4. Use the affine-scaling direction to calculate $\sigma$ according to (22) and the residual modifications according to (23).
5. Solve (19)–(20) with the new right side to obtain the actual step.
6. Calculate the step length and take the step.
7. Return to step 2 and iterate.

When the algorithm terminates, the final iterate is usually close to the solution of (2), but has all its components nonzero. We use a heuristic to determine which of these components

represent indices of the variables that should be in the model. Specifically, we set

$$\mathcal{I} = \left\{ l : \left\| \mathbf{b}_{(l)} \right\|_\infty > t10^{-4}, \; l = 1, \ldots, p \right\},$$

and set all $b_{lj}$'s with $l \notin \mathcal{I}$ to 0. For the case where $k = 1$, there is some evidence that this heuristic is too liberal, in the sense that coefficients that are 0 at the exact solution are some distance from 0 at the final interior-point iterate. However, this occurs only for some values of $t$—namely, some of those at which variables enter or drop out of the model. This behavior is thus not of great concern, because as argued earlier, we regard the methodology as exploratory only.

## 4. THE INFRARED SPECTROMETRY DATA REVISITED

We implemented the algorithm described in Section 3.3 in C and applied it to the infrared spectrometry data discussed in Section 1. Our hope is that, by using all response variables simultaneously to select a single set of regressor variables, we can avoid problems of overfitting and high variability.

As remarked earlier, Figure 2 indicates that one observation in NLeco—namely, the 10th—is suspicious and possibly an outlier. Hence we ran our analysis twice, once using all observations and once with the 10th observation removed from *all* of the variables. In this way we also hoped to get some insight into the robustness of the proposed methodology with respect to outliers.

We used several values for the tuning parameter $t$: $t = .1$, .2, .25, .3, .4, .5, .6, .7, .75, .8, .9, and 1.0. The results are summarized in Table 1 for the complete dataset and in Table 2 for the dataset with the 10th observation missing. Each table lists all chosen regressor variables for various values of $t$. The horizontal lines link the values of $t$ for which each coefficient remains nonzero in the solution.

Note that the tables are quite similar, showing essentially the same group of variables selected over the range of $t$ values. Only X691 from Table 1 is replaced by X690 in Table 2 and X87 is added. Thus, at least for this extreme example, the method appears to be fairly robust with respect to the outlier in the 10th observation of NLeco.

*Table 1. Selected Variables Using the Complete Dataset*

| | .1 | .2 | .25 | .3 | .4 | .5 | .6 | .7 | .75 | .8 | .9 | 1.0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| X81 | | | | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X82 | | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X84 | ——— | ——— | ——— | ——— | ——— | ——— | ——— | | | | | |
| X85 | | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X86 | | | | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X112 | | | | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X113 | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X114 | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X118 | | | | | | | | | ——— | ——— | ——— | |
| X220 | | | | | | | | | | | ——— | |
| X622 | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X623 | | | | | | | ——— | ——— | ——— | ——— | ——— | ——— |
| X662 | | | | ——— | ——— | ——— | | | | | | |
| X667 | | | | | | | ——— | ——— | ——— | ——— | ——— | ——— |
| X672 | | | | | ——— | ——— | ——— | ——— | | | | |
| X673 | | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X691 | | | | | | | | | | | | ——— |

*Table 2. Selected Variables Without the 10th Observation*

| | .1 | .2 | .25 | .3 | .4 | .5 | .6 | .7 | .75 | .8 | .9 | 1.0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| X81 | | | | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X82 | | | | | | ——— | ——— | ——— | | | | |
| X84 | ——— | ——— | ——— | ——— | ——— | ——— | | | | | | |
| X85 | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X86 | | | | | | | | | ——— | ——— | ——— | ——— |
| X87 | | | | | | | | | | | ——— | ——— |
| X112 | | | | | | | ——— | ——— | ——— | ——— | ——— | ——— |
| X113 | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X114 | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X118 | | | | | | | | | | ——— | ——— | |
| X220 | | | | | | | | | | | ——— | |
| X622 | | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X623 | | | | | | | | | | | ——— | |
| X662 | | | | ——— | ——— | ——— | ——— | | | | | |
| X667 | | | | | | | ——— | ——— | ——— | ——— | ——— | ——— |
| X672 | | | | | ——— | ——— | ——— | ——— | ——— | | | |
| X673 | | | | | | ——— | ——— | ——— | ——— | ——— | ——— | ——— |
| X690 | | | | | | | | | | | | ——— |

Although there still appears to be some variation, and several variables are identified as being nonzero for only a few values of $t$, the method consistently picks regressor variables from only three separate ranges of the spectrum. Roughly speaking, these ranges are the 81st–87th, the 112th–118th, and the 622nd–691st frequencies. Within each of these ranges, the regressor variables are highly correlated. The minimum correlation in the group X81, X82, X84, X85, X86, and X87 is larger than .9989, whereas the minimum correlation in the group X112, X113, X114, and X118 is larger than .9993. For the last group X622, X623, X662, X667, X672, X673, X690, and X691, the minimum correlation is larger than .9431. Because this group spans a wider range of frequencies, it is not surprising that its correlation is slightly smaller than the others.

Given the high correlations, there are essentially two ways in which one could proceed: either select one regressor variable from each group or average over the variables in each group. Here, for illustrative purposes we use the first method and choose those variables that are selected for most values of $t$, that is, X85, X114, and X622. Both Table 1 and Table 2 suggest using this set of explanatory variables.

To investigate how well this selection of variables performs, we fit linear regression models to each of the response variables using these three regressor variables. The resulting fits are shown in Figure 3, for the case in which all observations are used and in Figure 6 for the case in which the 10th observation is removed. Figures 4 and 7 show the corresponding plots of the jackknifed residuals, whereas Figures 5 and 8 show the normal quantile plots based on these jackknifed residuals.

From these figures, it seems that a linear regression model using these three predictors is satisfactory, at least in most cases. Of course, in the figures that are produced from the complete dataset, the outlier in NLeco is clearly visible. The residual plot for Na shows a lot of structure, but this is due to the granularity of this response variable. Modeling of Na clearly is a difficult task, because this response variable takes only 5 distinct values, with 18 replications of the smallest value and 3 replications of the median.

The results of significance tests for each parameter in each of the linear models are summarized in Table 3. An entry of " *** " in this table means that in the linear model for the response
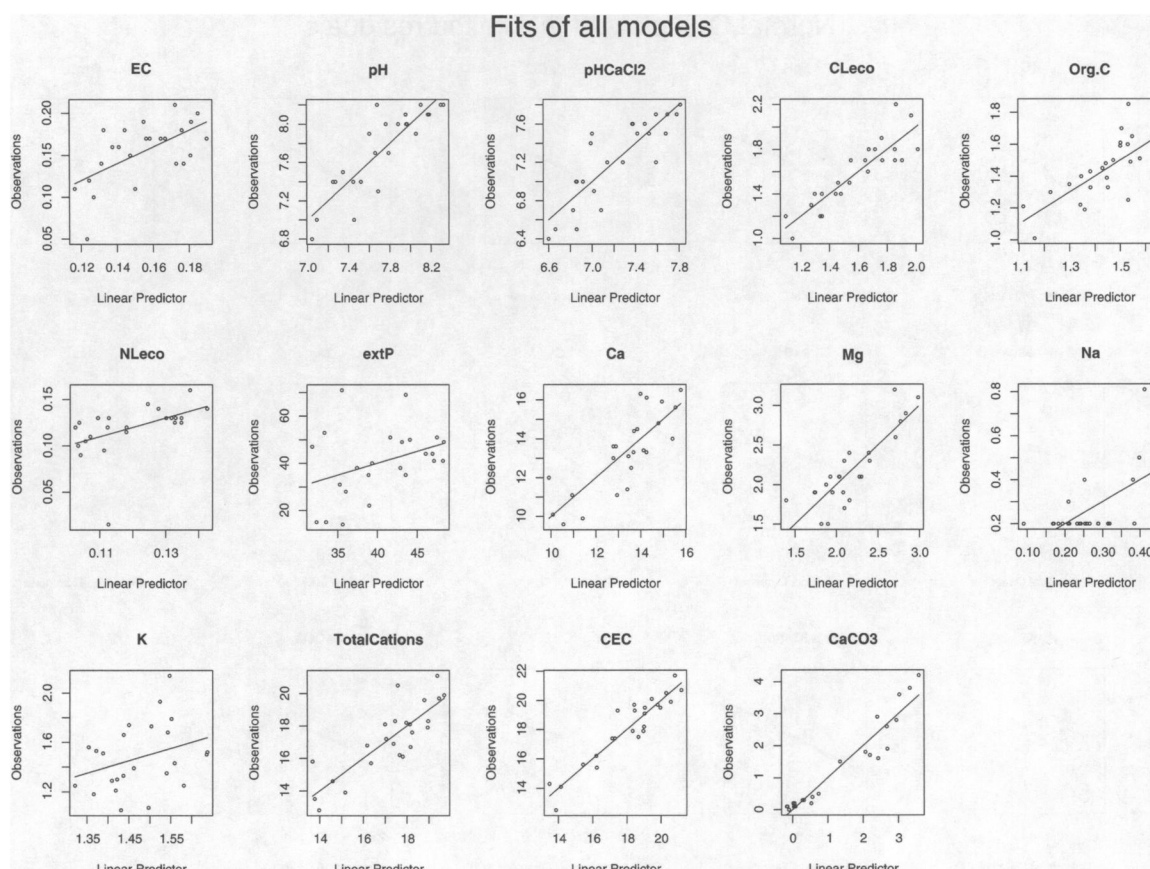
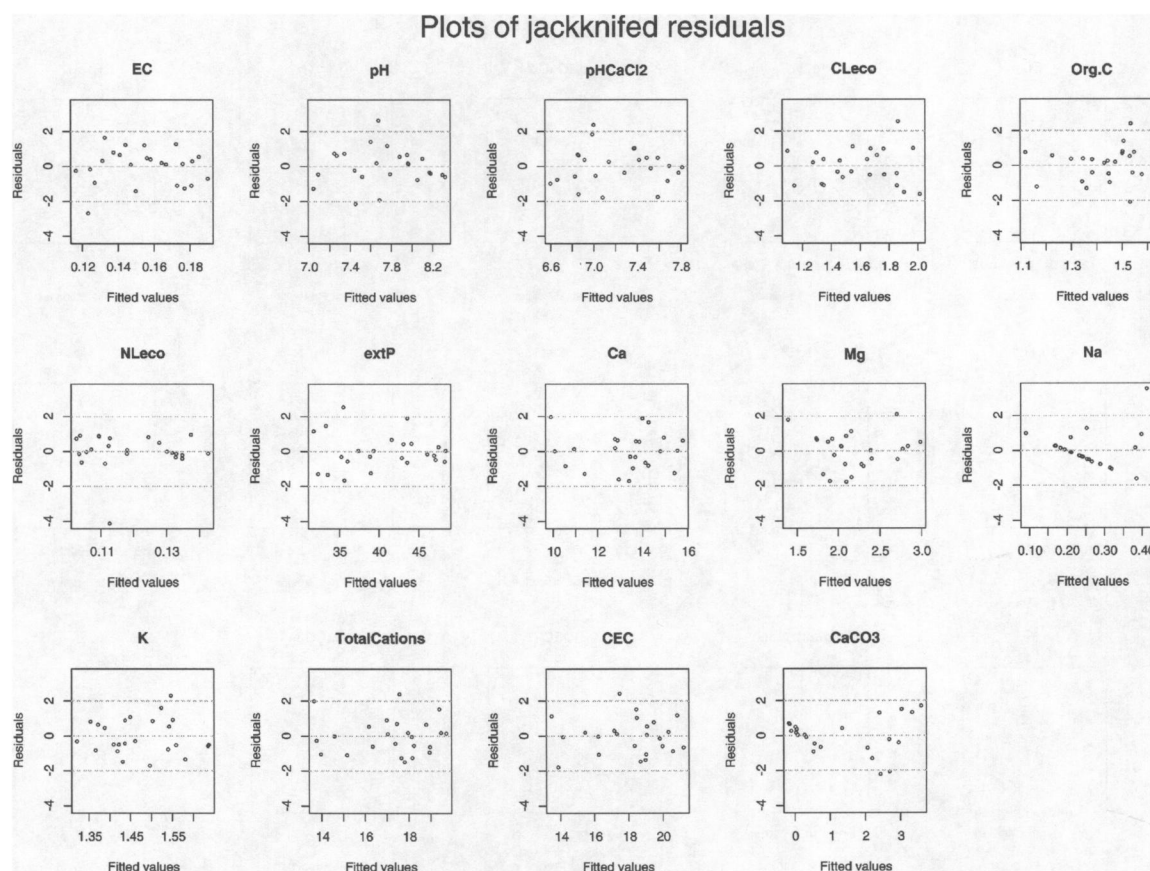Figure 3. Linear Fits Using Three Explanatory Variables (all observations).



Figure 4. Residual Plots (all observations).
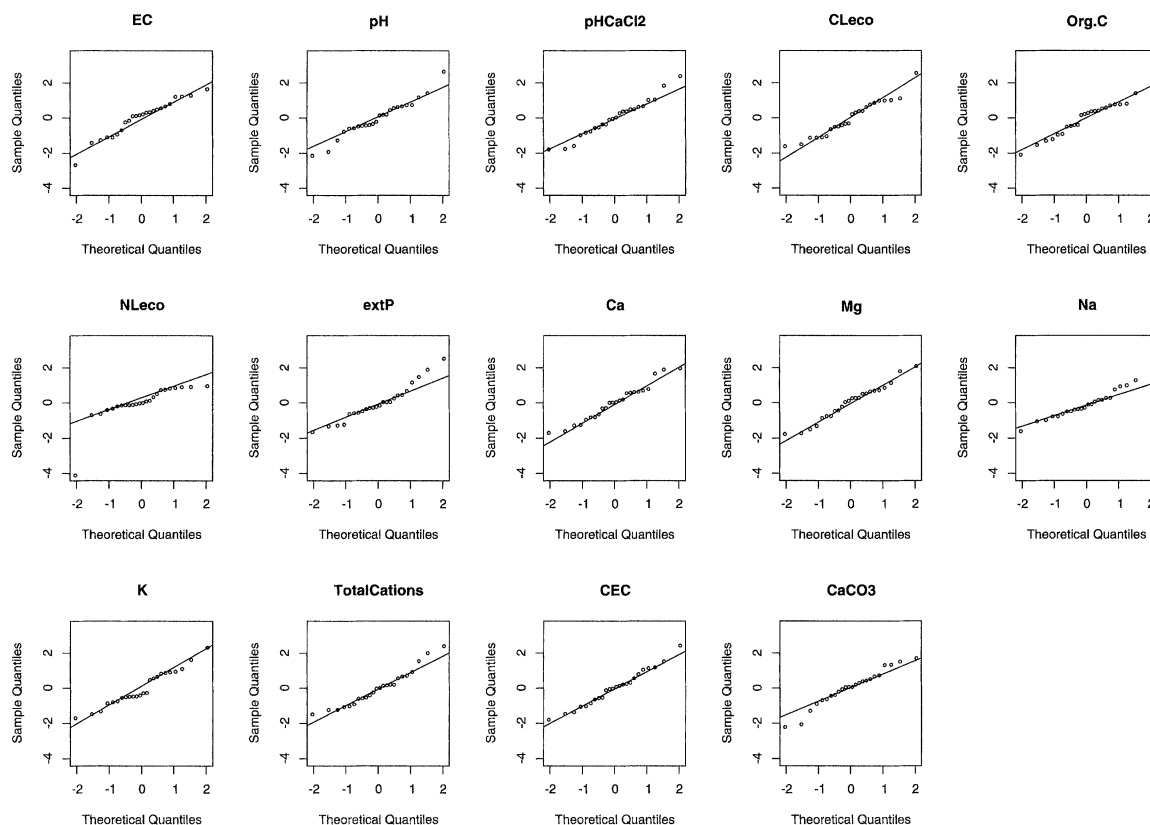
## Normal Q-Q plots of jackknifed residuals



Figure 5. Normal Quantiles Plots for the Residuals (all observations).
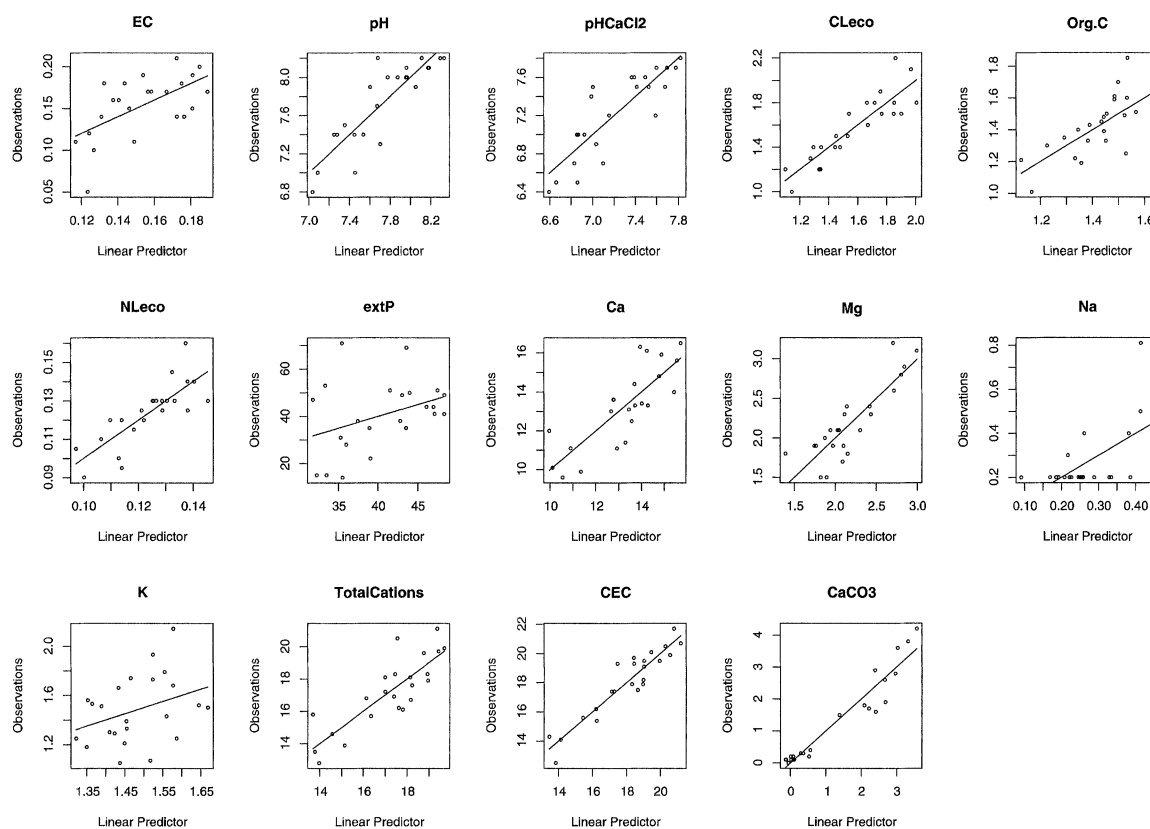
## Fits of all models



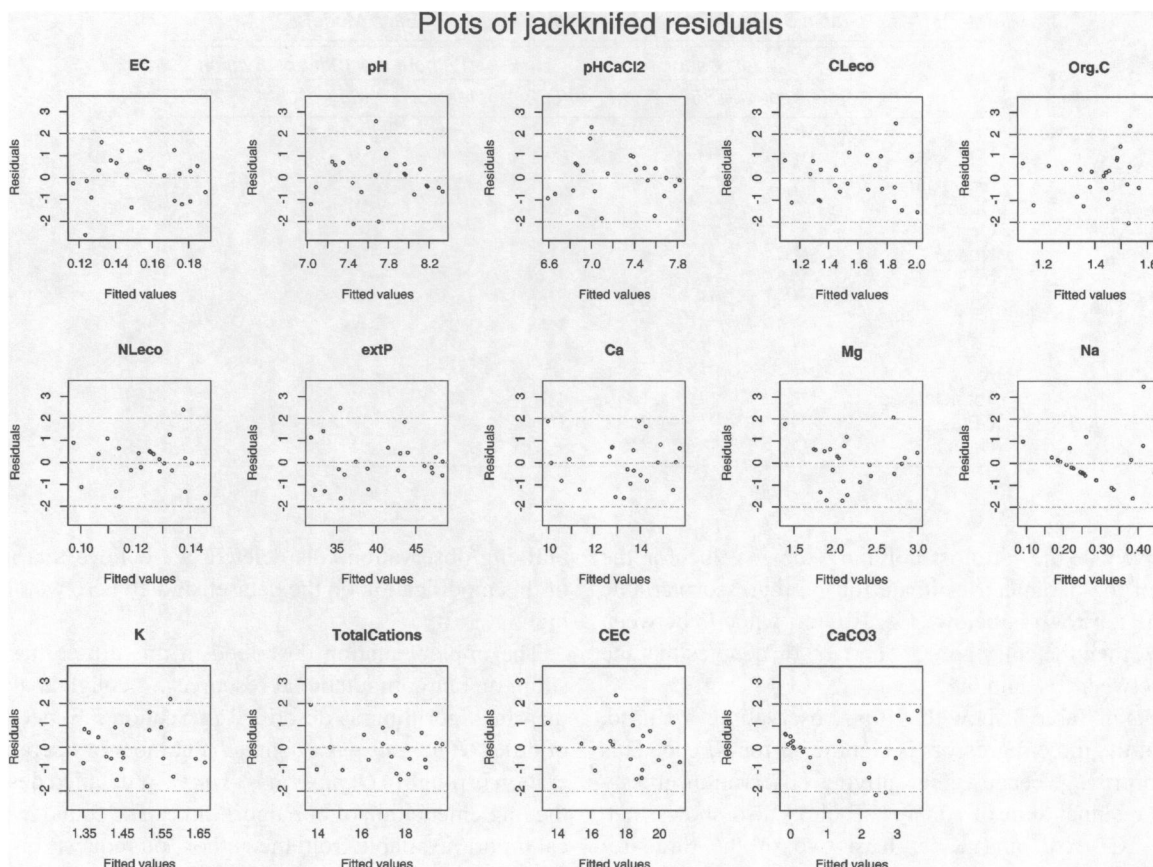Figure 6. Linear Fits Using Three Explanatory Variables (without the 10th observation).
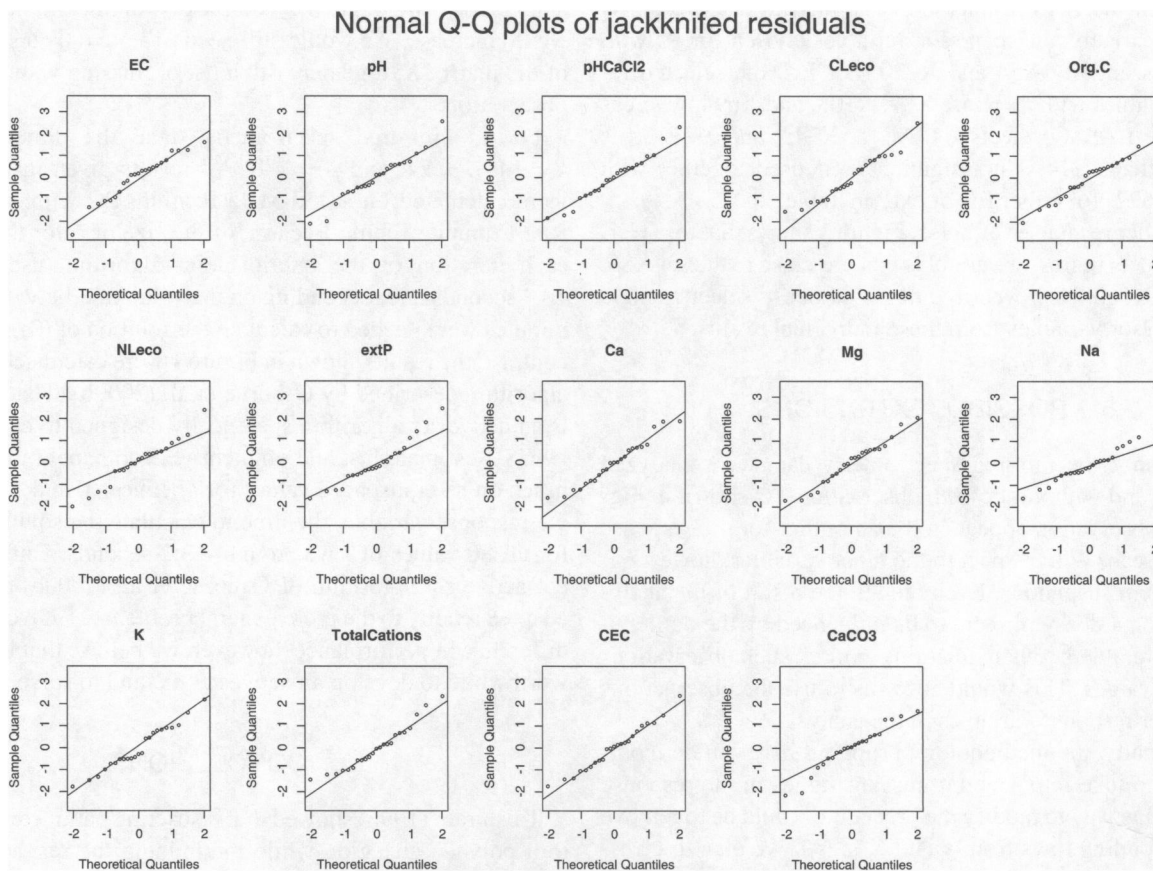
Figure 7. Residual Plots (without the 10th observation).



Figure 8. Normal Quantile Plots for the Residuals (without the 10th observation).

Table 3. Significance of Selected Variables in Each Model

| | All observations | | | | Without the 10th observation | | | |
|---|---|---|---|---|---|---|---|---|
| | Intercept | X85 | X114 | X622 | Intercept | X85 | X114 | X622 |
| EC | *** | * | | | *** | * | | |
| pH | *** | *** | * | ** | *** | *** | * | * |
| pHCaCl2 | *** | *** | ** | | *** | *** | ** | |
| CLeco | *** | *** | ** | * | *** | *** | ** | * |
| Org.C | *** | *** | | * | *** | ** | | * |
| NLeco | *** | | | | *** | *** | | * |
| extP | *** | | | | *** | | | |
| Ca | *** | *** | | *** | *** | ** | | *** |
| Mg | *** | | *** | *** | *** | | *** | *** |
| Na | *** | | * | * | *** | | * | * |
| K | *** | | | | *** | | | |
| TotalCations | *** | *** | | *** | *** | ** | | *** |
| CEC | *** | *** | *** | *** | *** | *** | *** | *** |
| CaCO3 | *** | ** | *** | | *** | ** | *** | |

variable (given in the leftmost column), the $p$ value for the $t$-statistics of the parameter estimate for the regressor variable (given in the top row) is below .1%. If the $p$ value is between .1% and 1%, then the entry is " ** " and a " * " denotes that the $p$ value is between 1% and 5%.

The results in Table 3 show that if all observations are used, then none of the three regressors is significant for NLeco. This is hardly surprising, because the outlying observation grossly inflates the residual sum of squares. Table 3 also shows that, except for EC, extP, and K, at least two of the three regressor variables that we have chosen are significant for each response variable. In the case of K, this observation is not surprising, given the results of the LASSO summarized in Figure 1. When the LASSO method is applied to each response variable separately, the regressor variables chosen for K, with $t = 1$ are essentially X54 and X220. For EC, the selected regressor variables with $t = 1$ are X87, X108, and X694, whereas for extP the LASSO selects X1, X118, X345, and X732. It is surprising that X114 is not significant (if used together with X85 and X622) for either EC or extP in Table 3.

For all other response variables, Figure 1 shows that for $t = 1$, LASSO selects regressor variables that are close to the set X85, X114, and X662. However, it is much harder to select a single set of regressor variables from these individual results.

## 5. POSSIBLE EXTENSIONS

In Section 4 the infrared spectrometry data were analyzed twice, with and without the 10th observation. As shown in Figure 2, this observation appears to be an outlier for the response variable NLeco. With respect to the other variables, however, it does not seem suspicious. Even though the result of the analysis in Section 4 does not seem to be influenced by the outlier, it may be preferable in other situations to mark such observations as missing values. This would allow us to use the observations for the other response variables in the analysis.

Theoretically, the methodology proposed here can be modified easily to take into account missing values in the response variables. One way to modify the procedure would be to remove the corresponding rows from $\mathbf{y}$ and $\tilde{\mathbf{X}}$ in (4). We plan to investigate this modification in a future project. However, given the similarity of the results in Section 4, regardless of whether the

outlying observation was deleted, we believe that application of this modification on the dataset studied here would not be of high interest.

The implementation described in this article has put some strain on our computational resources. A rough analysis shows that the algorithm, as described in Section 3.3, needs memory of order $O(p^2 + np + kn^2)$, and the number of operations per iteration is roughly $O(p^3 + kn^3 + kpn^2)$. A detailed description of the implementation of our algorithm can be found in the technical report (available from the authors on request) on which this article is based.

Implementation of the modification with missing rows could be performed without increasing the computational demands appreciably, although the complexity of the implementation would increase. We would still need to store only a single copy of the matrix $\mathbf{X}$, together with a list of missing values for each observation.

For the infrared spectrometry data, the dimensions are $k = 14$, $n = 24$, and $p = 770$. All results given in this article were calculated on a 450-MHz Pentium PC with 128 MB of RAM running Linux. Because of the size of $p$ for this dataset, each iteration of the interior-point algorithm used roughly 10.4 seconds and, depending on the value of $t$, between 2 and 5 minutes were needed to calculate the solution of (2). By way of contrast, the results shown in Figure 1 were calculated using the algorithm described by Osborne et al. (2000b). That algorithm is an active set algorithm specifically designed to calculate the LASSO estimate fast and efficiently. Each panel in Figure 1 is based on 80 equispaced values for $t$ between 0 and 1. For a single response variable, the time to calculate the solutions of (1) for all 80 values of $t$ was roughly 3.2 seconds. Unfortunately, the active set algorithm of Osborne et al. (2000b) can not be adapted readily to the more general problem (2). Given the large difference in performance, however, we believe that it would be worthwhile to develop an active set method to solve (2).

## 6. CONCLUSIONS

Tibshirani (1996) showed that restricting parameter estimates to a polyhedral region while minimizing the residual sum of squares yields a method that successfully combines elements of ridge regression and subset selection. In this article we have

extended Tibshirani's idea to the situation in which we seek a subset of regressor variables that is useful for several response variables simultaneously. This extension leads again to a convex programming problem; we describe an interior-point algorithm for solving this problem.

Application of this method to the infrared spectrometry data shows that this method can be quite useful in identifying a single subset for simultaneous modeling purposes. The example chosen to motivate and illustrate much of what we have done here was extreme in the sense that the number of regressors is huge, the number of responses is moderate, and the number of observations is almost unrealistically small. We contend that in less extreme cases the computational load will be more manageable, and the methodology will be just as useful. Testing the method on more datasets would be enlightening.

## ACKNOWLEDGMENTS

## REFERENCES

Bakin, S. (1999), "Adaptive Regression and Model Selection in Data Mining Problems," unpublished doctoral thesis, Australian National University.

Breiman, L. (1995), "Better Subset Regression Using the Nonnegative Garrote," *Technometrics*, 37, 373–384.

Breiman, L., and Friedman, J. H. (1997), "Predicting Multivariate Responses in Multiple Linear Regression" (with discussion), *Journal of the Royal Statistical Society*, Ser. B, 59, 3–54.

Brown, P. J. (1993), *Measurement, Regression, and Calibration*, Oxford, U.K.: Clarendon Press.

Brown, P. J., Fearn, T., and Vannucci, M. (1999), "The Choice of Variables in Multivariate Regression: A Non-Conjugate Bayesian Decision Approach," *Biometrika*, 86, 635–648.

Brown, P. J., Vannucci, M., and Fearn, T. (1998), "Multivariate Bayesian Variable Selection and Prediction," *Journal of the Royal Statistical Society*, Ser. B, 60, 627–641.

—— (2002), "Bayes Model Averaging With Selection of Regressors," *Journal of the Royal Statistical Society*, Ser. B, 64, 519–536.

Burnham, K. P., and Anderson, D. A. (1998), *Model Selection and Inference: A Practical Information Theoretic Approach*, New York: Springer-Verlag.

Clark, D. I., and Osborne, M. R. (1988), "On Linear Restricted and Interval Least-Squares Problems," *IMA Journal of Numerical Analysis*, 8, 23–36.

Clarke, F. H. (1990), *Optimization and Nonsmooth Analysis*, Philadelphia: SIAM.

Draper, N. R., and Smith, H. (1998), *Applied Regression Analysis* (3rd. ed.), New York: Wiley.

Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004), "Least Angle Regression" (with discussion), *The Annals of Statistics*, 32, 407–499.

Frank, I. E., and Friedman, J. H. (1993), "A Statistical View of Some Chemometrics Regression Tools" (with discussion), *Technometrics*, 35, 109–148.

Hastie, T. J., and Tibshirani, R. J. (1990), *Generalized Additive Models*, London: Chapman & Hall.

Hocking, R. R. (1996), *Methods and Applications of Linear Models: Regression and the Analysis of Variance*, New York: Wiley.

Leamer, E. E. (1978), "Regression Selection Strategies and Revealed Priors," *Journal of the American Statistical Association*, 73, 580–587.

Lobo, M. S., Vandenberghe, L., Boyd, S., and Lebret, H. (1998), "Applications of Second-Order Cone Programming," *Linear Algebra and Its Applications*, 284, 193–228.

Martens, H., and Naes, T. (1989), *Multivariate Calibration*, Chichester, U.K.: Wiley.

McCullagh, P., and Nelder, J. A. (1989), *Generalized Linear Models* (2nd ed.), London: Chapman & Hall.

Mehrotra, S. (1992), "On the Implementation of a Primal–Dual Interior Point Method," *SIAM Journal on Optimization*, 2, 575–601.

Miller, A. J. (1990), *Subset Selection in Regression*, London: Chapman & Hall.

Osborne, M. R. (1985), *Finite Algorithms in Optimization and Data Analysis*, Chichester, U.K.: Wiley.

—— (1992), "An Effective Method for Computing Regression Quantiles," *IMA Journal of Numerical Analysis*, 12, 151–166.

Osborne, M. R., Presnell, B., and Turlach, B. A. (2000a), "A New Approach to Variable Selection in Least Squares Problems," *IMA Journal of Numerical Analysis*, 20, 389–403.

—— (2000b), "On the LASSO and Its Dual," *Journal of Computational and Graphical Statistics*, 9, 319–337.

Rockafellar, R. T. (1970), *Convex Analysis*, Princeton, NJ: Princeton University Press.

Roos, C., Terlaky, T., and Vial, J. P. (1997), *Theory and Algorithms for Linear Optimization: An Interior Point Approach*, New York: Wiley.

Sturm, J. F. (1999), "Using SEDUMI 1.02: A Matlab Toolbox for Optimization Over Symmetric Cones," *Optimization Methods and Software*, 11/12, 625–653.

Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society*, Ser. B, 58, 267–288.

Wold, H. (1984), "PLS Regression," in *Encyclopedia of Statistical Sciences*, Vol. 6, eds. N. L. Johnson and S. Kotz, New York: Wiley, pp. 581–591.

Wright, S. J. (1997), *Primal–Dual Interior-Point Methods*, Philadelphia: SIAM.

Ye, Y. (1997), *Interior Point Algorithms: Theory and Analysis*, New York: Wiley.

Zou, H., and Hastie, T. (2005), "Regularization and Variable Selection via the Elastic Net," *Journal of the Royal Statistical Society*, Ser. B, 67, 301–320.