

Orthogonalizing Penalized Regression

Shifeng XIONG¹, Bin DAI², and Peter Z. G. QIAN²

¹ Academy of Mathematics and Systems Science

Chinese Academy of Sciences, Beijing 100190

² Department of Statistics

University of Wisconsin-Madison, Madison, WI 53706

Abstract

Since the penalized likelihood function of the smoothly clipped absolute deviation (SCAD) penalty is highly non-linear and has many local optima, finding a local solution to achieve the so-called oracle property is an open problem. We propose an iterative algorithm, called the OEM algorithm, to fill this gap. The development of the algorithm draws direct impetus from a missing-data problem arising in design of experiments with an orthogonal complete matrix. In each iteration, the algorithm imputes the missing data based on the current estimates of the parameters and updates a closed-form solution associated with the complete data. By introducing a procedure called active orthogonalization, we make the algorithm broadly applicable to problems with arbitrary regression matrices. In addition to the SCAD penalty, the proposed algorithm works for other penalties like the MCP, lasso and nonnegative garrote. Convergence and convergence rate of the algorithm are examined. The algorithm has several unique theoretical properties. For the SCAD and MCP penalties, an OEM sequence can achieve the oracle property after sufficient iterations. For various penalties, an OEM sequence converges to a point having grouping coherence for fully aliased regression matrices. For computing the ordinary least squares estimator with a singular regression matrix, an OEM sequence converges to the Moore-Penrose generalized inverse-based least squares estimator.

KEY WORDS: Design of experiments; MCP; Missing data; Optimization; Oracle property; Orthogonal design; SCAD; The EM algorithm; The Lasso.

²Corresponding author: Peter Z. G. Qian. Email: peterq@stat.wisc.edu

1 INTRODUCTION

Fan and Li (2001) proposed the smoothly clipped absolute deviation (SCAD) penalty to achieve simultaneous estimation and variable selection. Consider a linear model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (1)$$

where $\mathbf{X} = (x_{ij})$ is the $n \times p$ regression matrix, $\mathbf{Y} \in \mathbb{R}^n$ is the response vector, $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ is the vector of regression coefficients and the distribution of the vector of random error $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$ is $N(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)$ with $\mathbf{0}_n$ being the n th zero vector and \mathbf{I}_n being the $n \times n$ identity matrix. Throughout, let $\|\cdot\|$ denote the Euclidean norm. A *regularized least squares estimator* of $\boldsymbol{\beta}$ with this penalty is given by solving

$$\min_{\boldsymbol{\beta}} \left[\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + 2 \sum_{j=1}^p P_{\lambda}(|\beta_j|) \right], \quad (2)$$

where for $\theta > 0$,

$$P'_{\lambda}(\theta) = \lambda I(\theta \leq \lambda) + (a\lambda - \theta)_+ I(\theta > \lambda) / (a - 1), \quad (3)$$

$a > 2$, $\lambda > 0$ is the tuning parameter and I is the indicator function. In order to apply the penalty P_{λ} equally on all the variables, \mathbf{X} can be standardized so that

$$\sum_{i=1}^n x_{ij}^2 = 1, \text{ for } j = 1, \dots, p. \quad (4)$$

Both theory and computation of the estimator in (2) have been actively studied. On the theoretical side, Fan and Li (2001) introduced an important concept, called the *oracle* property. An estimator of $\boldsymbol{\beta}$ having this property can not only select the correct submodel asymptotically, but also estimate the nonzero coefficients as efficiently as if the correct submodel were known in advance. On the computational side, existing algorithms for solving this optimization problem include local quadratic approximation (Fan and Li 2001; Hunter and Li 2005), local linear approximation (Zou and Li 2008), the coordinate descent algorithm (Tseng 2001;

Tseng and Yun 2009; Breheny and Huang 2010; Mazumder, Friedman, and Hastie 2010) and the minimization by iterative soft thresholding (MIST) algorithm (Schifano, Strawderman, and Wells 2010), among others.

Departing from the existing work, we study the SCAD penalty from a new perspective, targeting on the *interface* between theory and computing. Fan and Li (2001) proved that there exists a local solution to (2) with the oracle property. From the optimization viewpoint, (2) can have many local minima (Huo and Chen 2010) and it is very challenging to find one of them to achieve the oracle property. To the best of our knowledge, no theoretical results are available to show that any existing algorithm can provide such a local minimum. We propose an iterative algorithm, called orthogonalizing EM (OEM), to fill this gap. We will show in Section 4 that the OEM solution to (2) can indeed achieve the oracle property under regularity conditions. OEM draws its direct impetus from a missing data problem with a complete orthogonal design arising in design of experiments. Throughout, a matrix is orthogonal if its columns are orthogonal. In each iteration, the algorithm imputes the missing data based on the current estimate of $\boldsymbol{\beta}$ and updates a closed-form solution to (2) associated with the complete data. Much beyond this orthogonal design formulation, the OEM algorithm applies to general data structures by *actively orthogonalizing* arbitrary regression matrices.

Though the inspiration of the OEM algorithm stems from the SCAD penalty, it, not surprisingly, works for the general penalized regression problem:

$$\min_{\boldsymbol{\beta} \in \Theta} [\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + P(\boldsymbol{\beta}; \lambda)], \quad (5)$$

where $\boldsymbol{\beta} \in \Theta$, Θ is a subset of \mathbb{R}^p and λ is the vector of tuning parameters. Besides the SCAD penalty, choices for $P(\boldsymbol{\beta}; \lambda)$ include the ridge regression (Hoerl and Kennard 1970), the nonnegative garrote (Breiman 1995), the lasso (Tibshirani 1996) and the MCP (Zhang 2010). Algorithms for solving the problem in (5) include those developed in Fu (1998), Grandvalet (1998), Osborne, Presnell, and Turlach (2000), the LARS algorithm introduced in Efron, Hastie, Johnstone, and Tibshirani (2004) and the coordinate descent algorithm

(Tseng 2001; Friedman, Hastie, Hofling and Tibshirani 2007; Wu and Lange 2008; Tseng and Yun 2009), and are available in R packages like `lars` (Hastie and Efron 2011), `glmnet` (Friedman, Hastie, and Tibshirani 2011) and `scout` (Witten and Tibshirani 2011).

In addition to achieving the oracle property for the SCAD and MCP penalties, the OEM algorithm has several other unique theoretical features. 1. *Having grouping coherence*: An estimator $\hat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$ in (1) is said to have grouping coherence if it has the same coefficient for full aliased columns in \mathbf{X} (Zou and Hastie 2005). For the lasso, SCAD and MCP, an OEM sequence converges to a point having grouping coherence. 2. *Convergence in singular case*: When \mathbf{X} in (1) is singular, the ordinary least squares estimator given by (5) without any penalty is not unique. For this singular case, an OEM solution, or essentially the Healy-Westmacott estimator (Healy and Westmacott 1956), converges to the Moore-Penrose generalized inverse-based least squares estimator.

The remainder of the article will unfold as follows. Section 2 derives the OEM algorithm for a missing data problem with a complete orthogonal design. Section 3 significantly broadens the applicability of the algorithm by introducing an idea for actively expanding any regression matrix to an orthogonal matrix. Section 4 establishes the oracle property of the OEM solution for the SCAD and MCP. Section 5 provides convergence properties of the OEM algorithm. Section 6 shows that for a regression matrix with full aliased columns, an OEM sequence for the lasso, SCAD or MCP converges to a solution with grouping coherence and illustrates how to use the OEM algorithm to compute the ordinary least squares estimator for a singular regression matrix. Section 7 provides some discussion.

2 THE OEM ALGORITHM AS A PENALIZED HEALY-WESTMACOTT PROCEDURE

Orthogonal designs are widely used in science and engineering. Such designs have been intensively studied in different branches of statistics including design of experiments, information theory (MacWilliams and Sloane 1977), linear models, sampling survey and computer

experiments. Popular classes of orthogonal designs include orthogonal arrays (Hedayat, Sloane, and Stufken 1999), orthogonal main-effect plans (Addelman 1962; Wu and Hamada 2009) and orthogonal Latin hypercube designs (Ye 1998; Steinberg and Lin 2006; Bingham, Sitter, and Tang 2009; Lin, Mukerjee, and Tang 2009; Pang, Liu, and Lin, 2009; Sun, Liu, and Lin 2009; Lin, Bingham, Sitter, and Tang 2010) from the computer experiments literature (Santner, Williams, and Notz 2003; Fang, Li, and Sudjianto 2005). In this section, we motivate the OEM algorithm by using a missing data problem with an orthogonal complete design. Suppose that the matrix \mathbf{X} in (1) for this problem is a submatrix of an $m \times p$ *complete* orthogonal matrix

$$\mathbf{X}_c = (\mathbf{X}' \ \mathbf{\Delta}')$$
 (6)

where $\mathbf{\Delta}$ is the $(m - n) \times p$ *missing* matrix. Let

$$\mathbf{Y}_c = (\mathbf{Y}', \mathbf{Y}'_{\text{miss}})'$$
 (7)

define the vector of complete observations with \mathbf{Y}_{miss} corresponding to $\mathbf{\Delta}$. If \mathbf{Y}_{miss} were observable, then the ordinary least square estimator of $\boldsymbol{\beta}$ based on the complete data $(\mathbf{X}_c, \mathbf{Y}_c)$ has a closed form as \mathbf{X}_c is orthogonal. In light of this fact, Healy and Westmacott (1956) proposed an iterative procedure to compute the ordinary least squares estimator $\boldsymbol{\beta}_{OLS}$ of $\boldsymbol{\beta}$. In each iteration, their procedure imputes the values of \mathbf{Y}_{miss} and updates the closed-form ordinary least squares estimator associated with the complete data. The OEM algorithm follows the same idea but solves (2) with the SCAD penalty. If \mathbf{Y}_{miss} were observable, then \mathbf{X} in (2) and \mathbf{Y} can be replaced by \mathbf{X}_c and \mathbf{Y}_c , yielding a closed-form solution to (2). Much beyond this orthogonal design formulation, we will significantly broaden the applicability of the algorithm in Section 3 by introducing an idea, called active orthogonalization, to actively expand any regression matrix into an orthogonal matrix.

Define

$$\mathbf{A} = \mathbf{\Delta}'\mathbf{\Delta}.$$
 (8)

Let (d_1, \dots, d_p) denote the diagonal elements of $\mathbf{X}'_c\mathbf{X}_c$. The OEM algorithm for solving the

optimization problem in (2) proceeds as follows. Let $\boldsymbol{\beta}^{(0)}$ be an initial estimate of $\boldsymbol{\beta}$. For $k = 0, 1, \dots$, impute \mathbf{Y}_{miss} as $\mathbf{Y}_I = \boldsymbol{\Delta}\boldsymbol{\beta}^{(k)}$, let $\mathbf{Y}_c = (\mathbf{Y}', \mathbf{Y}'_I)'$, and solve

$$\boldsymbol{\beta}^{(k+1)} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left[\|\mathbf{Y}_c - \mathbf{X}_c\boldsymbol{\beta}\|^2 + 2 \sum_{j=1}^p P_\lambda(|\beta_j|) \right] \quad (9)$$

until $\{\boldsymbol{\beta}^{(k)}\}$ converges. Letting

$$\mathbf{u} = (u_1, \dots, u_p)' = \mathbf{X}'\mathbf{Y} + \mathbf{A}\boldsymbol{\beta}^{(k)}, \quad (10)$$

(9) becomes

$$\boldsymbol{\beta}^{(k+1)} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left[\sum_{j=1}^p (d_j\beta_j^2 - 2u_j\beta_j) + 2 \sum_{j=1}^p P_\lambda(|\beta_j|) \right], \quad (11)$$

which is *separable* in the dimensions of $\boldsymbol{\beta}$. If \mathbf{X} in (1) is standardized as in (4) with $d_j \geq 1$ for all j , (11) has a closed-form

$$\beta_j^{(k+1)} = \begin{cases} \operatorname{sign}(u_j)(|u_j| - \lambda)_+/d_j, & \text{when } |u_j| \leq (d_j + 1)\lambda, \\ \operatorname{sign}(u_j)[(a - 1)|u_j| - a\lambda]/[(a - 1)d_j - 1], & \text{when } (d_j + 1)\lambda < |u_j| \leq a\lambda d_j, \\ u_j/d_j, & \text{when } |u_j| > a\lambda d_j. \end{cases} \quad (12)$$

As pointed out in Dempster, Laird, and Rubin (1977) that the Healy-Westmacott procedure is essentially an EM algorithm, OEM is an EM algorithm as well. The complete data $\mathbf{Y}_c = (\mathbf{Y}', \mathbf{Y}'_{\text{miss}})'$ in (7) follow a regression model $\mathbf{Y}_c = \mathbf{X}_c\boldsymbol{\beta} + \boldsymbol{\varepsilon}_c$, where $\boldsymbol{\varepsilon}_c$ is from $N(\mathbf{0}_m, \mathbf{I}_m)$. Let $\boldsymbol{\beta}_{\text{SCAD}}$ be a solution to (2), where $\boldsymbol{\beta}_{\text{SCAD}} = \operatorname{argmax}_{\boldsymbol{\beta}} L(\boldsymbol{\beta} | \mathbf{Y})$ and the penalized likelihood function $L(\boldsymbol{\beta} | \mathbf{Y})$ is

$$(2\pi)^{-n/2} \exp\left(-\frac{1}{2}\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2\right) \exp\left[-\sum_{j=1}^p P_\lambda(|\beta_j|)\right].$$

Given $\boldsymbol{\beta}^{(k)}$, the E-step of the OEM algorithm for the SCAD is

$$\begin{aligned}
& E[\log\{L(\boldsymbol{\beta}|\mathbf{Y}_c)\} \mid \mathbf{Y}, \boldsymbol{\beta}^{(k)}] \\
&= -C\{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + E(\|\mathbf{Y}_{\text{miss}} - \mathbf{X}\boldsymbol{\beta}\|^2 \mid \boldsymbol{\beta}^{(k)}) + 2\sum_{j=1}^p P_\lambda(|\beta_j|)\} \\
&= -C\{n + \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \|\Delta\boldsymbol{\beta}^{(k)} - \Delta\boldsymbol{\beta}\|^2 + 2\sum_{j=1}^p P_\lambda(|\beta_j|)\}
\end{aligned}$$

for some constant $C > 0$. Define

$$Q_{\text{SCAD}}(\boldsymbol{\beta} \mid \boldsymbol{\beta}^{(k)}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \|\Delta\boldsymbol{\beta}^{(k)} - \Delta\boldsymbol{\beta}\|^2 + 2\sum_{j=1}^p P_\lambda(|\beta_j|). \quad (13)$$

The M-step of the OEM algorithm for the SCAD is

$$\boldsymbol{\beta}^{(k+1)} = \arg \min_{\boldsymbol{\beta}} Q_{\text{SCAD}}(\boldsymbol{\beta} \mid \boldsymbol{\beta}^{(k)}),$$

which is equivalent to (11).

For the general penalized regression problem in (5), the M-step of the OEM algorithm becomes

$$\boldsymbol{\beta}^{(k+1)} = \arg \min_{\boldsymbol{\beta} \in \Theta} Q(\boldsymbol{\beta} \mid \boldsymbol{\beta}^{(k)}), \quad (14)$$

where Q replaces Q_{SCAD} in (13) for the corresponding penalty function. If Θ and P in (5) are *decomposable* as $\Theta = \prod_{j=1}^p \Theta_j$ and $P(\boldsymbol{\beta}; \lambda) = \sum_{j=1}^p P_j(\beta_j; \lambda)$, similarly to (11), (14) reduces to p one-dimensional problems

$$\beta_j^{(k+1)} = \arg \min_{\beta_j \in \Theta_j} [d_j \beta_j^2 - 2u_j \beta_j + P_j(\beta_j; \lambda)], \text{ for } j = 1, \dots, p, \quad (15)$$

with $\mathbf{u} = (u_1, \dots, u_p)'$ defined in the same way as in (10). This shortcut applies to the following penalties:

1. The lasso (Tibshirani 1996), where $\Theta_j = \mathbb{R}$,

$$P_j(\beta_j; \lambda) = 2\lambda|\beta_j|, \quad (16)$$

and (15) becomes

$$\beta_j^{(k+1)} = \text{sign}(u_j) \left(\frac{|u_j| - \lambda}{d_j} \right)_+. \quad (17)$$

Here, for $a \in \mathbb{R}$, $(a)_+$ denotes $\max\{a, 0\}$.

2. The nonnegative garrote (Breiman 1995), where $\Theta_j = \{x : x\hat{\beta}_j \geq 0\}$, $P_j(\beta_j; \lambda) = 2\lambda\beta_j/\hat{\beta}_j$, $\hat{\beta}_j$ is the ordinary least squares estimator of β_j , and (15) becomes

$$\beta_j^{(k+1)} = \left(\frac{u_j\hat{\beta}_j - \lambda}{d_j\hat{\beta}_j^2} \right)_+ \hat{\beta}_j.$$

3. The elastic-net (Zou and Hastie 2005), where $\Theta_j = \mathbb{R}$,

$$P_j(\beta_j; \lambda) = 2\lambda_1|\beta_j| + \lambda_2\beta_j^2. \quad (18)$$

and (15) becomes

$$\beta_j^{(k+1)} = \text{sign}(u_j) \left(\frac{|u_j| - \lambda_1}{d_j + \lambda_2} \right)_+. \quad (19)$$

5. The MCP (Zhang 2010), where $\Theta_j = \mathbb{R}$, $P_j(\beta_j; \lambda) = 2P_\lambda(|\beta_j|)$, and

$$P'_\lambda(\theta) = (\lambda - \theta/a)I(\theta \leq a\lambda) \quad (20)$$

with $a > 1$ and $\theta > 0$. If \mathbf{X} in (1) is standardized as in (4) with $d_j \geq 1$ for all j , (15) becomes

$$\beta_j^{(k+1)} = \begin{cases} \text{sign}(u_j)a(|u_j| - \lambda)_+/(ad_j - 1), & \text{when } |u_j| \leq a\lambda d_j, \\ u_j/d_j, & \text{when } |u_j| > a\lambda d_j \end{cases} \quad (21)$$

6. The “Berhu” penalty (Owen 2006), where $\Theta_j = \mathbb{R}$, $P_j(\beta_j; \lambda) = 2\lambda\{|\beta_j|I(|\beta_j| < \delta) + (\beta_j^2 + \delta^2)I(|\beta_j| \geq \delta)/(2\delta)\}$ for some $\delta > 0$, and (15) becomes

$$\beta_j^{(k+1)} = \begin{cases} \text{sign}(u_j)(|u_j| - \lambda)_+/d_j, & \text{when } |u_j| < \lambda + d_j\delta, \\ u_j\delta/(\lambda + d_j\delta), & \text{when } |u_j| \geq \lambda + d_j\delta. \end{cases}$$

Obviously, if the penalty on $\boldsymbol{\beta}$ disappears, the OEM algorithm reduces to the Healy-Westmacott procedure. Quite interestingly, Theorem 6 in Section 5 shows that, for the same \mathbf{X} and \mathbf{Y} in (1), the OEM algorithm for the elastic-net and lasso numerically converges faster than the Healy-Westmacott procedure.

Example 1. For the model in (1), let the complete matrix \mathbf{X}_c be a fractional factorial design from Xu (2009) with 4096 runs in 30 factors. Clearly, \mathbf{X}_c is an orthogonal design. Let \mathbf{X} in (1) be the submatrix of \mathbf{X}_c consisting of the first 3000 rows and let \mathbf{Y} be generated with $\sigma = 1$ and

$$\beta_j = (-1)^j \exp[-2(j-1)/20] \text{ for } j = 1, \dots, p. \quad (22)$$

Here, $p = 30$ and $n = 3000$. Assume the response values corresponding to the last 1096 rows of \mathbf{X}_c are missing. We used the OEM algorithm to solve the optimization problem in (2) with an initial value $\boldsymbol{\beta}^{(0)} = \mathbf{0}$ and a criterion to stop when relative changes in all coefficients are less than 10^{-6} . For $\lambda = 1$ and $a = 3.7$ in (3), Figure 1 plots values of the objective function in (2) of the OEM sequence against iteration numbers, where the algorithm converges at iteration 13.

3 THE GENERAL FORMULATION WITH ACTIVE ORTHOGONALIZATION

The OEM algorithm in Section 2 was derived for a missing-data problem where \mathbf{X} in (1) is imbedded in a *pre-specified* orthogonal matrix. We drop this assumption in this section

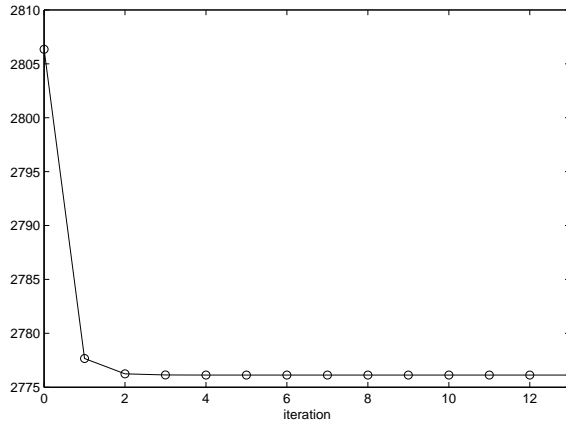


Figure 1. Values of the objective function of an OEM sequence for the SCAD against iterations for Example 1.

and further develop the algorithm for general data structures by introducing a procedure to *actively* expand an arbitrary matrix to an orthogonal matrix. The general idea of augmenting extra data has been used for EM problems before. For example, for a covariance estimation problem in Rubin and Szatrowski (1982), extra data are added elaborately to make the maximum likelihood estimator of the expanded patterned covariance matrices have an explicit form. To facilitate the use of the OEM algorithm in Section 2, the contribution here is to develop a scheme to *orthogonalize* the matrix \mathbf{X} with an arbitrary structure.

Take \mathbf{S} to be a $p \times p$ diagonal matrix with non-zero diagonal elements s_1, \dots, s_p . Define

$$\mathbf{Z} = \mathbf{X}\mathbf{S}^{-1}. \quad (23)$$

The eigenvalue decomposition of $\mathbf{Z}'\mathbf{Z}$ (Wilkinson 1965) is

$$\mathbf{V}'\mathbf{\Gamma}\mathbf{V},$$

where \mathbf{V} is an orthogonal matrix and $\mathbf{\Gamma}$ is a diagonal matrix whose diagonal elements, $\gamma_1 \geq \dots \geq \gamma_p$, are the nonnegative eigenvalues of $\mathbf{Z}'\mathbf{Z}$. Let

$$t = \#\{j : \gamma_j = \gamma_1, j = 1, \dots, p\} \quad (24)$$

denote the number of the γ_j equal to γ_1 . For example, if $\gamma_1 = \gamma_2$ and $\gamma_1 > \gamma_j$ for $j = 3, \dots, p$, then $t = 2$. Define

$$\mathbf{B} = \text{diag}(\gamma_1 - \gamma_{t+1}, \dots, \gamma_1 - \gamma_p) \quad (25)$$

and

$$\mathbf{\Delta} = \mathbf{B}^{1/2} \mathbf{V}_1 \mathbf{S}, \quad (26)$$

where \mathbf{V}_1 is the submatrix of \mathbf{V} consisting of the last $p - t$ rows. Let \mathbf{X}_c be the augmented matrix of $\mathbf{\Delta}$ and \mathbf{X} .

Lemma 1. The matrix \mathbf{X}_c constructed above is orthogonal.

Proof. Note that

$$\mathbf{X}'_c \mathbf{X}_c = \mathbf{X}' \mathbf{X} + \mathbf{\Delta}' \mathbf{\Delta},$$

which, by plugging (25) and (26), becomes

$$\mathbf{S}[\mathbf{Z}' \mathbf{Z} + \mathbf{V}'(\gamma_1 \mathbf{I}_p - \mathbf{\Gamma}) \mathbf{V}] \mathbf{S} = \gamma_1 \mathbf{S}^2. \quad (27)$$

Now, because

$$\gamma_1 \mathbf{I}_p - \mathbf{\Gamma} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{pmatrix},$$

(27) is orthogonal, which completes the proof. \square

The underlying geometry of the active orthogonalization in Lemma 1 can be described as follows. For a vector $\mathbf{x} \in \mathbb{R}^m$, let $P_\omega \mathbf{x}$ denote its projection onto a subspace ω of \mathbb{R}^m . This lemma implies that for the column vectors of \mathbf{X} in (1), $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$, there exists a set of mutually orthogonal vectors $\mathbf{x}_{c1}, \dots, \mathbf{x}_{cp} \in \mathbb{R}^{n+p-t}$, essentially the column vectors of \mathbf{X}_c in (6), satisfy the condition that $P_{\mathbb{R}^n} \mathbf{x}_{ci} = \mathbf{x}_i$, for $j = 1, \dots, p$. Proposition 1 makes this precise.

Proposition 1. Let ω be an n -dimensional subspace of \mathbb{R}^m with $n \leq m$. If $p \leq m - n + 1$, then for any p vectors $\mathbf{x}_1, \dots, \mathbf{x}_p \in \omega$, there exist p vectors $\mathbf{x}_{c1}, \dots, \mathbf{x}_{cp} \in \mathbb{R}^m$ such that $P_\omega \mathbf{x}_{ci} = \mathbf{x}_i$ for $j = 1, \dots, p$ and $\mathbf{x}'_{ci} \mathbf{x}_{cj} = 0$ for $i \neq j$.

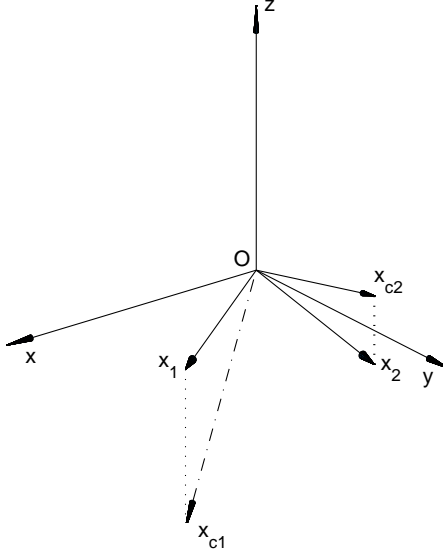


Figure 2. Expand two two-dimensional vectors \mathbf{x}_1 and \mathbf{x}_2 to two three-dimensional vectors \mathbf{x}_{c1} and \mathbf{x}_{c2} with $\mathbf{x}'_{c1}\mathbf{x}_{c2} = 0$.

Figure 2 expands two vectors \mathbf{x}_1 and \mathbf{x}_2 in \mathbb{R}^2 to two orthogonal vectors \mathbf{x}_{c1} and \mathbf{x}_{c2} in \mathbb{R}^3 .

Now, if \mathbf{X}_c from Lemma 1 is treated as the complete matrix defined in (6), the OEM algorithm in Section 2 follows through immediately.

When using the OEM algorithm to solve (5), in (10) instead of computing $\mathbf{\Delta}$ in (26), one may compute $\mathbf{A} = \mathbf{\Delta}'\mathbf{\Delta}$ and the diagonal entries d_1, \dots, d_p of $\mathbf{X}'_c\mathbf{X}_c$. Note that

$$\mathbf{A} = \gamma_1\mathbf{S}^2 - \mathbf{X}'\mathbf{X} \quad (28)$$

and

$$d_j = \gamma_1 s_j^2 \text{ for } j = 1, \dots, p, \quad (29)$$

where \mathbf{S} and \mathbf{Z} are defined in (23) and γ_1 is the largest eigenvalue of $\mathbf{Z}'\mathbf{Z} = \mathbf{S}^{-1}\mathbf{X}'\mathbf{X}\mathbf{S}^{-1}$. One way to compute γ_1 is to use the power method (Wilkinson 1965) described below. Given a nonzero initial vector $\mathbf{a}^{(0)} \in \mathbb{R}^p$, let $\gamma_1^{(0)} = \|\mathbf{a}^{(0)}\|$. For $k = 0, 1, \dots$, compute $\mathbf{a}^{(k+1)} = \mathbf{X}'\mathbf{X}\mathbf{a}^{(k)}/\gamma_1^{(k)}$ and $\gamma_1^{(k+1)} = \|\mathbf{a}^{(k+1)}\|$ until convergence. If $\mathbf{a}^{(0)}$ is not an eigenvector

of any γ_j that does not equal γ_1 , then $\gamma_1^{(k)}$ converges to γ_1 . For t defined in (24), the convergence rate is linear (Watkins 2002) specified by

$$\lim_{k \rightarrow \infty} \frac{\|\gamma_1^{(k+1)} - \gamma_1\|}{\|\gamma_1^{(k)} - \gamma_1\|} = \frac{\gamma_{t+1}}{\gamma_1}.$$

An easy way to make $\mathbf{A} = \mathbf{\Delta}'\mathbf{\Delta}$ in (28) positive definite is to replace \mathbf{B} in (25) by

$$\mathbf{B} = \text{diag}(d - \gamma_{t+1}, \dots, d - \gamma_p)$$

with $d \geq \gamma_1$, which changes (28) and (29) to

$$\mathbf{A} = d\mathbf{S}^2 - \mathbf{X}'\mathbf{X} \quad (30)$$

and

$$d_j = ds_j^2, \text{ for } j = 1, \dots, p, \quad (31)$$

respectively. If $d > \gamma_1$, then $\mathbf{A} = \mathbf{\Delta}'\mathbf{\Delta}$ is positive definite.

Remark 1. The matrix \mathbf{S} in (23) can be chosen flexibly. One possibility is to use $\mathbf{S} = \mathbf{I}_p$ so that

$$\mathbf{X}'\mathbf{X} + \mathbf{\Delta}'\mathbf{\Delta} = d\mathbf{I}_p \quad (32)$$

with $d \geq \gamma_1$, and \mathbf{X}_c/\sqrt{d} is standardized as in (4).

Example 2. Suppose that \mathbf{X} in (1) is orthogonal. Take

$$\mathbf{S} = \text{diag} \left[\left(\sum_{i=1}^n x_{i1}^2 \right)^{1/2}, \dots, \left(\sum_{i=1}^n x_{ip}^2 \right)^{1/2} \right]. \quad (33)$$

Since $t = p$, $\mathbf{\Delta}$ in (26) is empty. This result indicates the active orthogonalization procedure will not overshoot: if \mathbf{X} is orthogonal already, it adds no row.

Example 3. Let

$$\mathbf{X} = \begin{pmatrix} 0 & 0 & 3/2 \\ -4/3 & -2/3 & 1/6 \\ 2/3 & 4/3 & 1/6 \\ -2/3 & 2/3 & -7/6 \end{pmatrix}.$$

If $\mathbf{S} = \mathbf{I}_3$, (26) gives $\mathbf{\Delta} = (-2/\sqrt{3}, 2/\sqrt{3}, 1/\sqrt{3})$.

Example 4. Consider a two-level design in three factors, A , B and C :

$$\begin{pmatrix} -1 & -1 & -1 \\ -1 & 1 & 1 \\ 1 & -1 & 1 \\ 1 & 1 & -1 \end{pmatrix}.$$

The regression matrix including all main effects and two-way interactions is

$$\mathbf{X} = \begin{pmatrix} -1 & -1 & -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & 1 & -1 & -1 \end{pmatrix},$$

where the last three columns for the interactions are fully aliased with the first three columns for the main effects. For $\mathbf{S} = \mathbf{I}_3$, (26) gives

$$\mathbf{\Delta} = \begin{pmatrix} 0 & -2 & 0 & 0 & -2 & 0 \\ 0 & 0 & -2 & -2 & 0 & 0 \\ -2 & 0 & 0 & 0 & 0 & -2 \end{pmatrix}.$$

The elements in $\mathbf{\Delta}$ are chosen flexibly, not restricted to ± 1 .

Example 5. Consider a 1000×10 random matrix $\mathbf{X} = (x_{ij})$ with entries independently

drawn from the uniform distribution on $[0, 1)$. Using \mathbf{S} in (33), (26) gives

$$\Delta = \begin{pmatrix} -7.99 & 16.06 & -6.39 & -18.26 & 12.91 & -8.67 & 7.56 & 34.08 & -17.04 & -11.81 \\ 26.83 & -12.09 & 7.91 & 1.02 & -22.75 & -6.90 & -19.98 & 26.10 & -0.86 & 0.88 \\ -4.01 & 1.48 & 9.51 & -21.99 & 19.46 & -10.27 & -25.12 & -3.39 & 7.29 & 27.90 \\ 21.77 & 10.72 & -0.61 & -6.46 & 28.00 & 1.28 & -6.86 & -7.04 & 11.13 & -30.64 \\ -15.78 & 5.60 & -15.26 & -7.67 & -9.76 & 23.93 & -14.71 & 12.25 & 29.45 & -7.89 \\ 16.34 & 10.61 & -41.82 & 11.82 & 6.49 & -7.38 & -6.14 & -1.82 & -1.86 & 13.09 \\ -8.15 & 24.97 & 12.11 & 24.35 & 3.66 & -2.59 & -27.84 & -3.45 & -9.40 & -13.72 \\ -5.35 & -21.70 & -4.16 & 7.42 & 13.98 & 29.84 & -10.26 & 7.60 & -25.13 & 7.78 \\ -19.62 & -22.43 & -2.61 & 22.58 & 11.80 & -22.08 & 1.25 & 15.87 & 14.94 & 0.31 \end{pmatrix}.$$

Only nine rows need to be added to make this large \mathbf{X} matrix orthogonal.

4 ACHIEVING THE ORACLE PROPERTY WITH NONCONVEX PENALTIES

Fan and Li (2001) introduced an important concept called the oracle property and showed that there exists one local minimum of (2) with this property. However, because the optimization problem in (2) has an exponential number of local optima (Huo and Ni 2007; Huo and Chen 2010), no theoretical results in the literature claim that an existing algorithm can provide such a local minimum. In this section, we prove that an OEM sequence for the SCAD and MCP can indeed achieve this property. The theoretical results in this and the following sections work for the OEM algorithm in both Sections 2 and 3.

First, we describe the oracle concept. A penalized least squares estimator of $\boldsymbol{\beta}$ in (1) has this property if it can not only select the correct submodel asymptotically, but also estimate the nonzero coefficients $\boldsymbol{\beta}_1$ in (34) as efficiently as if the correct submodel were known in advance. Suppose that the number of nonzero coefficients of $\boldsymbol{\beta}$ in (1) is p_1 (with $p_1 \leq p$) and

partition $\boldsymbol{\beta}$ as

$$\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)', \quad (34)$$

where $\boldsymbol{\beta}_2 = \mathbf{0}$ and no component of $\boldsymbol{\beta}_1$ is zero. Divide the regression matrix \mathbf{X} in (1) to $(\mathbf{X}_1 \ \mathbf{X}_2)$ with \mathbf{X}_1 corresponding to $\boldsymbol{\beta}_1$. If all the variables that influence the response in (1) are known, an *oracle* estimator can be given as $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}'_1, \hat{\boldsymbol{\beta}}'_2)'$ with $\hat{\boldsymbol{\beta}}_2 = \mathbf{0}$, where

$$\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_1 \sim N(\mathbf{0}, \sigma^2(\mathbf{X}'_1 \mathbf{X}_1)^{-1}).$$

We need several assumptions.

Assumption 1. As $n \rightarrow \infty$,

$$\frac{\mathbf{X}'\mathbf{X}}{n} \rightarrow \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_1 & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_2 \end{pmatrix},$$

where $\boldsymbol{\Sigma}$ is positive definite and $\boldsymbol{\Sigma}_1$ is $p_1 \times p_1$. Furthermore, \mathbf{X}/\sqrt{n} is standardized such that each entry on the diagonal of $\mathbf{X}'\mathbf{X}/n$ is 1, and $\mathbf{X}'\mathbf{X}/n + \boldsymbol{\Delta}'\boldsymbol{\Delta} = d\mathbf{I}_p$ with $d \geq \gamma_1$, where $d = O(1)$ and γ_1 is the largest eigenvalue of $\mathbf{X}'\mathbf{X}/n$.

Assumption 2. The tuning parameter $\lambda = \lambda_n$ in (3) satisfies the condition that, as $n \rightarrow \infty$, $\lambda_n/n \rightarrow 0$ and $\lambda_n/\sqrt{n} \rightarrow \infty$.

Let $\{\boldsymbol{\beta}^{(k)}, k = 0, 1, \dots, \}$ be the OEM sequence from (12) for the SCAD with a fixed $a > 2$ in (3). We need an assumption on $k = k_n$. Let η be the largest eigenvalue of $\mathbf{I}_{p_1} - \mathbf{X}'_1 \mathbf{X}_1/(nd)$. Under Assumption 1, η tends to a limit lying between 0 and 1 as $n \rightarrow \infty$.

Assumption 3. As $n \rightarrow \infty$, $\eta^{k_n} \lambda_n/\sqrt{n} \rightarrow 0$ and $k_n^2 \exp(-c(\lambda_n/\sqrt{n})^2) \rightarrow 0$ for any $c > 0$.

One choice for k_n to satisfy Assumption 3 is

$$k_n = \left(\frac{\lambda_n}{\sqrt{n}} \right)^\nu \text{ for some } \nu > 0.$$

Under Assumption 2, $k_n \rightarrow \infty$ as $n \rightarrow \infty$.

Under the above assumptions, Theorem 1 shows that $\boldsymbol{\beta}^{(k)} = (\boldsymbol{\beta}_1^{(k)'}, \boldsymbol{\beta}_2^{(k)'})'$ can achieve the oracle property.

Theorem 1. Suppose that Assumption 1-3 hold. If $\boldsymbol{\beta}^{(0)} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$, then as $n \rightarrow \infty$,

- (i) $P(\boldsymbol{\beta}_2^{(k)} = \mathbf{0}) \rightarrow 1$;
- (ii) $\sqrt{n}(\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1) \rightarrow N(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma}_1^{-1})$ in distribution.

The proof of Theorem 1 is deferred to the Appendix.

Remark 2. From (60) in the proof of Theorem 1, for $k = 1, 2, \dots$, $\boldsymbol{\beta}^{(k)}$ is *consistent* in variable selection. That is, $P(\beta_j^{(k)} \neq 0 \text{ for } j = 1, \dots, p_1) \rightarrow 1$ and $P(\boldsymbol{\beta}_2^{(k)} = \mathbf{0}) \rightarrow 1$ as $n \rightarrow \infty$.

Remark 3. The proof of Theorem 1 uses the convergence rates of $P(A_k)$, $P(B_k)$ and $P(C_k^h)$. If an OEM sequence satisfies the condition that $\beta_j^{(k+1)} = 0$ when $|u_j| < \lambda$ and $\beta_j^{(k+1)} = u_j/d$ when $|u_j| > c\lambda$ for some $c = O(1)$, then $P(A_{k+1}) = P(|u_j| < \lambda)$ and $P(B_{k+1}) = P(|u_j| > c\lambda)$. Since an OEM sequence for the MCP satisfies the above condition, an argument very similar to the proof in the Appendix shows that the convergence rates of $P(A_k)$, $P(B_k)$ and $P(C_k^h)$ for the MCP are the same as those with the SCAD. Thus, under Assumption 1-3, Theorem 1 holds for the MCP with a fixed $a > 1$ in (20).

Huo and Chen (2010) showed that, for the SCAD penalty, solving the global minimum of (5) leads to an NP-hard problem. Theorem 1 indicates that as far as the oracle property is concerned, the local solution given by OEM will suffice.

5 CONVERGENCE OF THE OEM ALGORITHM

In this section, we derive theoretical results on convergence properties of the OEM algorithm and compare the convergence rates of OEM for the ordinary least squares estimator and the elastic-net and lasso. The general penalty in (5) is considered here. Our derivations employ the main tool in Wu (1983) in conjunction with special properties of the penalties mentioned in Section 2.

We make several assumptions for Θ and $P(\boldsymbol{\beta}; \lambda)$ in (5).

Assumption 4. The parameter space Θ is a closed convex subset of \mathbb{R}^p .

Assumption 5. For a fixed λ , the penalty $P(\boldsymbol{\beta}; \lambda) \rightarrow +\infty$ as $\|\boldsymbol{\beta}\| \rightarrow +\infty$.

Assumption 6. For a fixed λ , the penalty $P(\boldsymbol{\beta}; \lambda)$ is continuous with respect to $\boldsymbol{\beta} \in \Theta$.

All penalties discussed in Section 2 satisfy these assumptions. The assumptions cover the case in which the iterative sequence $\{\boldsymbol{\beta}^{(k)}\}$ defined in (14) may fall on the boundary of Θ (Nettleton 1999), like the nonnegative garrote (Breiman 1995) and the nonnegative lasso (Efron et al. 2004). The bridge penalty (Frank and Friedman 1993) in (37) also satisfies the above assumptions.

For the model in (1), denote the objective function in (5) by

$$l(\boldsymbol{\beta}) = \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + P(\boldsymbol{\beta}; \lambda). \quad (35)$$

For some penalizes like the bridge, it may be numerically infeasible to perform the M-step in (14). For this situation, following the generalized EM algorithm in Dempster, Laird, and Rubin (1977), we define a *generalized OEM* algorithm to be an iterative scheme

$$\boldsymbol{\beta}^{(k)} \rightarrow \boldsymbol{\beta}^{(k+1)} \in \mathcal{M}(\boldsymbol{\beta}^{(k)}), \quad (36)$$

where $\boldsymbol{\beta} \rightarrow \mathcal{M}(\boldsymbol{\beta}) \subset \Theta$ is a point-to-set map such that

$$Q(\boldsymbol{\phi} \mid \boldsymbol{\beta}) \leq Q(\boldsymbol{\beta} \mid \boldsymbol{\beta}), \quad \text{for all } \boldsymbol{\phi} \in \mathcal{M}(\boldsymbol{\beta}).$$

Here, Q is given by replacing the SCAD with $P(\boldsymbol{\beta}; \lambda)$ in (13). The OEM sequence defined by (14) is a special case of (36). For example, the generalized OEM algorithm can be used for the bridge penalty, where $\Theta_j = \mathbb{R}$ and

$$P_j(\beta_j; \lambda) = \lambda|\beta_j|^a \quad (37)$$

for some $a \in (0, 1)$ in (5). Since the solution to (15) with the bridge penalty has no closed form, one may use one-dimensional search to compute $\beta_j^{(k+1)}$ that satisfies (36). By Assumption 1, $\{\boldsymbol{\beta} \in \Theta : l(\boldsymbol{\beta}) \leq l(\boldsymbol{\beta}^{(0)})\}$ is compact for any $l(\boldsymbol{\beta}^{(0)}) > -\infty$. By Assumption 6, \mathcal{M} is a closed point-to-set map (Zangwill 1969; Wu 1983).

The objective functions of the EM algorithms in the literature like those discussed in Wu (1983), Green (1990) and McLachlan and Krishnan (2008) are typically continuously differentiable. This condition does not hold for the objective function in (5) with the lasso and other penalties. A more general definition of stationary points is needed here. We call $\boldsymbol{\beta} \in \Theta$ a stationary point of l if

$$\liminf_{t \rightarrow 0_+} \frac{l((1-t)\boldsymbol{\beta} + t\boldsymbol{\phi}) - l(\boldsymbol{\beta})}{t} \geq 0 \quad \text{for all } \boldsymbol{\phi} \in \Theta.$$

Let S denote the set of stationary points of l . Analogous to Theorem 1 in Wu (1983) on the global convergence of the EM algorithm, we have the following result.

Theorem 2. Let $\{\boldsymbol{\beta}^{(k)}\}$ be a generalized OEM sequence generated by (36). Suppose that

$$l(\boldsymbol{\beta}^{(k+1)}) < l(\boldsymbol{\beta}^{(k)}) \quad \text{for all } \boldsymbol{\beta}^{(k)} \in \Theta \setminus S. \quad (38)$$

Then all limit points of $\{\boldsymbol{\beta}^{(k)}\}$ are elements of S and $l(\boldsymbol{\beta}^{(k)})$ converges monotonically to $l^* = l(\boldsymbol{\beta}^*)$ for some $\boldsymbol{\beta}^* \in S$.

Theorem 3. If $\boldsymbol{\beta}^*$ is a local minimum of $Q(\boldsymbol{\beta} \mid \boldsymbol{\beta}^*)$, then $\boldsymbol{\beta}^* \in S$.

This theorem follows from the fact that $l(\boldsymbol{\beta}) - Q(\boldsymbol{\beta} \mid \boldsymbol{\beta}^*)$ is differentiable and

$$\left. \frac{\partial [l(\boldsymbol{\beta}) - Q(\boldsymbol{\beta} \mid \boldsymbol{\beta}^*)]}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\boldsymbol{\beta}^*} = 0.$$

Remark 4. By Theorem 3, if $\boldsymbol{\beta}^{(k)} \notin S$, then $\boldsymbol{\beta}^{(k)}$ cannot be a local minimum of $Q(\boldsymbol{\beta} \mid \boldsymbol{\beta}^{(k)})$. Thus, there exists at least one point $\boldsymbol{\beta}^{(k+1)} \in \mathcal{M}(\boldsymbol{\beta}^{(k)})$ such that $Q(\boldsymbol{\beta}^{(k+1)} \mid \boldsymbol{\beta}^{(k)}) < Q(\boldsymbol{\beta}^{(k)} \mid \boldsymbol{\beta}^{(k)})$ and therefore satisfies the condition in (38). As a special case, an OEM sequence generated by (14) satisfies (38) in Theorem 2.

Next, we consider the convergence of a generalized OEM sequence $\{\boldsymbol{\beta}^{(k)}\}$ in (36). By Theorem 3, such results will automatically hold for an OEM sequence as well. If the penalty function $P(\boldsymbol{\beta}; \lambda)$ is convex and $l(\boldsymbol{\beta})$ has a unique minimum, Theorem 4 shows that $\{\boldsymbol{\beta}^{(k)}\}$ converges to the global minimum.

Theorem 4. Let $\{\boldsymbol{\beta}^{(k)}\}$ be defined in Theorem 2. Suppose that $l(\boldsymbol{\beta})$ in (35) is a convex function on Θ with a unique minimum $\boldsymbol{\beta}^*$ and that (38) holds for $\{\boldsymbol{\beta}^{(k)}\}$. Then $\boldsymbol{\beta}^{(k)} \rightarrow \boldsymbol{\beta}^*$ as $k \rightarrow \infty$.

Proof. We only need to show that $S = \{\boldsymbol{\beta}^*\}$. For $\boldsymbol{\phi} \in \Theta$ with $\boldsymbol{\phi} \neq \boldsymbol{\beta}^*$ and $t > 0$, we have

$$\frac{l((1-t)\boldsymbol{\phi} + t\boldsymbol{\beta}^*) - l(\boldsymbol{\beta}^*)}{t} \leq \frac{tl(\boldsymbol{\beta}^*) + (1-t)l(\boldsymbol{\phi}) - l(\boldsymbol{\phi})}{t} = l(\boldsymbol{\beta}^*) - l(\boldsymbol{\phi}) < 0.$$

This implies $\boldsymbol{\phi} \notin S$. □

Theorem 5 discusses the convergence of an OEM sequence $\{\boldsymbol{\beta}^{(k)}\}$ for more general penalties. For $a \in \mathbb{R}$, define $S(a) = \{\boldsymbol{\phi} \in S : l(\boldsymbol{\phi}) = a\}$. From Theorem 2, all limit points of an OEM sequence are in $S(l^*)$, where l^* is the limit of $l(\boldsymbol{\beta}^{(k)})$ in Theorem 2. Theorem 5 states that the limit point is unique under certain conditions.

Theorem 5. Let $\{\boldsymbol{\beta}^{(k)}\}$ be a generalized OEM sequence generated by (36) with $\boldsymbol{\Delta}'\boldsymbol{\Delta} > 0$. If (38) holds, then all limit points of $\{\boldsymbol{\beta}^{(k)}\}$ are in a connected and compact subset of $S(l^*)$. In particular, if the set $S(l^*)$ is discrete in that its only connected components are singletons, then $\boldsymbol{\beta}^{(k)}$ converges to some $\boldsymbol{\beta}^*$ in $S(l^*)$ as $k \rightarrow \infty$.

Proof. Note that $Q(\boldsymbol{\beta}^{(k+1)} \mid \boldsymbol{\beta}^{(k)}) = l(\boldsymbol{\beta}^{(k+1)}) + \|\boldsymbol{\Delta}\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\Delta}\boldsymbol{\beta}^{(k)}\|^2 \leq Q(\boldsymbol{\beta}^{(k)} \mid \boldsymbol{\beta}^{(k)}) = l(\boldsymbol{\beta}^{(k)})$. By Theorem 2, $\|\boldsymbol{\Delta}\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\Delta}\boldsymbol{\beta}^{(k)}\|^2 \leq l(\boldsymbol{\beta}^{(k)}) - l(\boldsymbol{\beta}^{(k+1)}) \rightarrow 0$ as $k \rightarrow \infty$. Thus, $\|\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\beta}^{(k)}\| \rightarrow 0$. This theorem now follows immediately from Theorem 5 of Wu (1983). □

Since the bridge, SCAD and MCP penalties all satisfy the condition that $S(l^*)$ is discrete, an OEM sequence for any of them converges to the stationary points of l . Theorem 5 is obtained under the condition $\boldsymbol{\Delta}'\boldsymbol{\Delta} > 0$. Since the error $\boldsymbol{\varepsilon}$ in (1) has a continuous distribution,

it is easy to show that Theorem 5 holds with probability one if $\Delta' \Delta$ is singular when d defined in (30) and (31) equals γ_1 .

We now derive the convergence rate of the OEM sequence in (14). Following Dempster, Laird, and Rubin (1977), write

$$\boldsymbol{\beta}^{(k+1)} = \mathbf{M}(\boldsymbol{\beta}^{(k)}),$$

where the map $\mathbf{M}(\boldsymbol{\beta}) = (M_1(\boldsymbol{\beta}), \dots, M_p(\boldsymbol{\beta}))'$ is defined by (14). We capture the convergence rate of the OEM sequence $\{\boldsymbol{\beta}^{(k)}\}$ through \mathbf{M} . Assume that (32) holds for $d \geq \gamma_1$, where γ_1 is the largest eigenvalue of $\mathbf{X}'\mathbf{X}$. For the active orthogolization in (30) and (31), taking $\mathbf{S} = \mathbf{I}_p$ satisfies this assumption; see Remark 1.

Let $\boldsymbol{\beta}^*$ be the limit of the OEM sequence $\{\boldsymbol{\beta}^{(k)}\}$. As in Meng (1994), we call

$$R = \limsup_{k \rightarrow \infty} \frac{\|\boldsymbol{\beta}^{(k+1)} - \boldsymbol{\beta}^*\|}{\|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*\|} = \limsup_{k \rightarrow \infty} \frac{\|\mathbf{M}(\boldsymbol{\beta}_k) - \mathbf{M}(\boldsymbol{\beta}^*)\|}{\|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*\|} \quad (39)$$

the global rate of convergence for the OEM sequence. If there is no penalty in (5), i.e., computing the ordinary least squares estimator, the global rate of convergence R in (39) becomes the largest eigenvalue of $\mathbf{J}(\boldsymbol{\beta}^*)$, denoted by R_0 , where $\mathbf{J}(\boldsymbol{\phi})$ is the $p \times p$ Jacobian matrix for $\mathbf{M}(\boldsymbol{\phi})$ having (i, j) th entry $\partial M_i(\boldsymbol{\phi}) / \partial \phi_j$. If (32) holds, then $\mathbf{J}(\boldsymbol{\beta}^*) = \mathbf{A}/d$ with $\mathbf{A} = \Delta' \Delta$. Thus,

$$R_0 = \frac{d - \gamma_p}{d}. \quad (40)$$

For (5), the penalty function $P(\boldsymbol{\beta}; \lambda)$ typically is not sufficiently smooth and R in (39) does not have an analytic form. Theorem 6 gives an upper bound of R_{net} , the value of R for the elastic-net penalty in (18) with $\lambda_1, \lambda_2 \geq 0$.

Theorem 6. For Δ from (6), if (32) holds, then $R_{\text{NET}} \leq R_0$.

Proof. Let \mathbf{x}_j denote the j th column of \mathbf{X} and \mathbf{a}_j denote the j th column of $\mathbf{A} = \Delta' \Delta$, respectively. For an OEM sequence for the elastic-net, by (19),

$$M_j(\boldsymbol{\beta}) = f(\mathbf{x}'_j \mathbf{Y} + \mathbf{a}'_j \boldsymbol{\beta}), \text{ for } j = 1, \dots, p,$$

where

$$f(u) = \text{sign}(u) \left(\frac{|u| - \lambda_1}{d + \lambda_2} \right)_+.$$

For $j = 1, \dots, p$, observe that

$$\begin{aligned} \frac{|M_j(\boldsymbol{\beta}^{(k)}) - M_j(\boldsymbol{\beta}^*)|}{\|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*\|} &= \frac{|f(\mathbf{x}'_j \mathbf{Y} + \mathbf{a}'_j \boldsymbol{\beta}^{(k)}) - f(\mathbf{x}'_j \mathbf{Y} + \mathbf{a}'_j \boldsymbol{\beta}^*)|}{|(\mathbf{x}'_j \mathbf{Y} + \mathbf{a}'_j \boldsymbol{\beta}^{(k)}) - (\mathbf{x}'_j \mathbf{Y} + \mathbf{a}'_j \boldsymbol{\beta}^*)|} \\ &\quad \cdot \frac{|(\mathbf{x}'_j \mathbf{Y} + \mathbf{a}'_j \boldsymbol{\beta}^{(k)}) - (\mathbf{x}'_j \mathbf{Y} + \mathbf{a}'_j \boldsymbol{\beta}^*)|}{\|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*\|} \\ &\leq \frac{1}{d} \cdot \frac{|\mathbf{a}'_j(\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*)|}{\|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*\|}. \end{aligned}$$

Thus,

$$\frac{\|M(\boldsymbol{\beta}^{(k)}) - M(\boldsymbol{\beta}^*)\|}{\|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*\|} \leq \frac{1}{d} \cdot \frac{\|A(\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*)\|}{\|\boldsymbol{\beta}^{(k)} - \boldsymbol{\beta}^*\|} \leq \frac{d - \gamma_p}{d}.$$

This completes the proof. \square

Remark 5. Theorem 6 indicates that, for the same \mathbf{X} and \mathbf{Y} in (1), the OEM solution for the elastic-net converges faster than its counterpart for the ordinary least squares. Since the lasso is a special case of the elastic-net with $\lambda_2 = 0$ in (18), this theorem holds for the lasso as well.

Remark 6. From (40) and Theorem 6, the convergence rate of the OEM algorithm depends on the ratio of the smallest eigenvalue, γ_p , and the largest eigenvalue, γ_1 , of $\mathbf{X}'\mathbf{X}$. This rate is the fastest when $\gamma_1 = \gamma_p$, i.e., if \mathbf{X} is orthogonal and standardized. This result suggests that OEM converges faster if \mathbf{X} is orthogonal or nearly orthogonal like from a supersaturated design or a nearly orthogonal Latin hypercube design (Owen 1994; Tang 1998). This result is in agreement with the recent finding in the design of experiments community that the use of orthogonal or nearly orthogonal designs can significantly improve the accuracy of penalized variable selection methods (Phoa, Pan, and Xu 2009; Deng, Lin, and Qian 2010; Zhu 2011).

Example 6. We generate \mathbf{X} from p dimensional Gaussian distribution $N(\mathbf{0}, \mathbf{V})$ with n independent observations, where the (i, j) th entry of \mathbf{V} is 1 for $i = j$ and ρ for $i \neq j$. Values

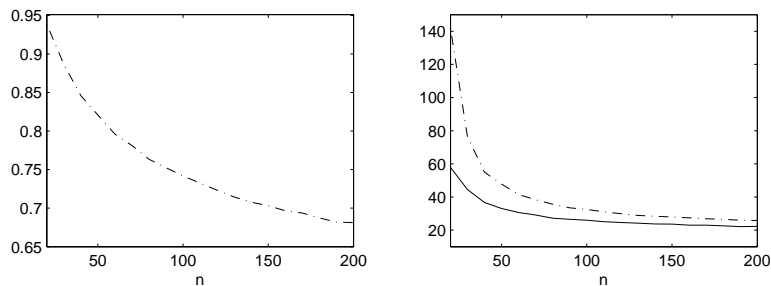


Figure 3. (Left) the average values of R_0 in (40) against increasing n for Example 6; (right) the average iteration numbers against increasing n for Example 6, where the dashed and solid lines denote the ordinary least squares estimator and the lasso, respectively.

of \mathbf{Y} and $\boldsymbol{\beta}$ are generated by (1) and (22). The same setup was used in Friedman, Hastie, and Tibshirani (2009). For $p = 10$, $\rho = 0.1$, $\lambda = 0.5$ and increasing n , the left panel of Figure 3 depicts the average values of R_0 in (40) against increasing n and the right panel of the figure depicts the average iteration numbers against increasing n , with the dashed and solid lines corresponding to the ordinary least squares estimator and the lasso, respectively. Quite strikingly, this figure indicates that OEM requires *fewer* iterations as n becomes larger, which makes OEM particularly attractive for situations with massive data (SAMSI 2012). It is important to note that here the OEM sequence for the lasso requires fewer iterations than its counterpart for the ordinary least squares, empirically validating Theorem 6.

6 POSSESSING GROUPING COHERENCE

In this section, we consider the convergence of the OEM algorithm when the regression matrix \mathbf{X} in (1) is singular due to fully aliased columns or other conditions. Let \mathbf{X} be standardized as in (4) with columns $\mathbf{x}_1, \dots, \mathbf{x}_p$. If \mathbf{x}_i and \mathbf{x}_j are fully aliased, i.e., $|\mathbf{x}_i| = |\mathbf{x}_j|$, then the objective function in (5) for the lasso is not strictly convex and has many minima (Zou and Hastie 2005). Data with fully aliasing structures commonly appear in observational studies and various classes of experimental designs like supersaturated designs (Wu 1993; Lin 1993; Tang and Wu 1997; Li and Lin 2002; Xu and Wu 2005) and factorial designs (Dey

and Mukerjee 1999; Mukerjee and Wu 2006).

Zou and Hastie (2005) states that if some columns of \mathbf{X} are identical, it is desirable to have grouping coherence by assigning the same coefficient to them. Definition 1 makes this precise.

Definition 1. An estimator $\hat{\boldsymbol{\beta}} = (\hat{\beta}_1, \dots, \hat{\beta}_p)'$ of $\boldsymbol{\beta}$ in (1) has *grouping coherence* if $\mathbf{x}_i = \mathbf{x}_j$ implies $\hat{\beta}_i = \hat{\beta}_j$ and $\mathbf{x}_i = -\mathbf{x}_j$ implies $\hat{\beta}_i = -\hat{\beta}_j$.

Let $\mathbf{0}_p$ denote the zero vector in \mathbb{R}^p . Let \mathbf{e}_{ij}^+ be the vector obtained by replacing the i th and j th entries of $\mathbf{0}_p$ with 1. Let \mathbf{e}_{ij}^- be the vector obtained by replacing the i th and j th entries of $\mathbf{0}_p$ with 1 and -1 , respectively. Let \mathcal{E} denote the set of all \mathbf{e}_{ij}^+ and \mathbf{e}_{ij}^- . By Definition 1, an estimator $\hat{\boldsymbol{\beta}}$ has grouping coherence if and only if for any $\boldsymbol{\alpha} \in \mathcal{E}$ with $\mathbf{X}\boldsymbol{\alpha} = \mathbf{0}$, $\boldsymbol{\alpha}'\hat{\boldsymbol{\beta}} = 0$.

Lemma 2. Suppose that (32) holds. For the OEM sequence $\{\boldsymbol{\beta}^{(k)}\}$ of the lasso, SCAD or MCP, if $\mathbf{X}\boldsymbol{\alpha} = \mathbf{0}$ and $\boldsymbol{\alpha}'\boldsymbol{\beta}^{(k)} = 0$ for $\boldsymbol{\alpha} \in \mathcal{E}$, then $\boldsymbol{\alpha}'\boldsymbol{\beta}^{(k+1)} = 0$.

Proof. For \mathbf{u} defined in (10), we have that $\boldsymbol{\alpha}'\mathbf{u} = \boldsymbol{\alpha}'\mathbf{X}'\mathbf{Y} + \boldsymbol{\alpha}'(d\mathbf{I}_p - \mathbf{X}'\mathbf{X})\boldsymbol{\beta}^{(k)} = 0$ for any $\boldsymbol{\alpha} \in \mathcal{E}$ with $\mathbf{X}\boldsymbol{\alpha} = \mathbf{0}$ and $\boldsymbol{\alpha}'\boldsymbol{\beta}^{(k)} = 0$. Then by (17), (12) and (21), an OEM sequence of the lasso, SCAD or MCP satisfies the condition that if $\boldsymbol{\alpha}'\mathbf{u} = 0$, then $\boldsymbol{\alpha}'\boldsymbol{\beta}^{(k+1)} = 0$ for $\boldsymbol{\alpha} \in \mathcal{E}$. This completes the proof. \square

Remark 7. Lemma 2 implies that, for $k = 1, 2, \dots$, $\boldsymbol{\beta}^{(k)}$ has grouping coherence if $\boldsymbol{\beta}^{(0)}$ has grouping coherence. Thus, if $\{\boldsymbol{\beta}^{(k)}\}$ converges, then its limit has grouping coherence. By Theorem 5, if $d > \lambda_1$ in (32), then an OEM sequence for the SCAD or MCP converges to a point with grouping coherence.

When \mathbf{X} in (1) has fully aliased columns, the objective function in (5) for the lasso has many minima and hence the condition in Theorem 4 does not hold. Theorem 7 shows that, even with full aliasing, an OEM sequence (17) for the lasso converges to a point with grouping coherence.

Theorem 7. Suppose that (32) holds. If $\boldsymbol{\beta}^{(0)}$ has grouping coherence, then as $k \rightarrow \infty$, the OEM sequence $\{\boldsymbol{\beta}^{(k)}\}$ of the lasso converges to a limit that has grouping coherence.

Proof. Partition the matrix \mathbf{X} in (1) as $(\mathbf{X}_1 \ \mathbf{X}_2)$, where no two columns of \mathbf{X}_2 are fully aliased and any column of \mathbf{X}_1 is fully aliased with at least one column of \mathbf{X}_2 . Let J denote the number of columns in \mathbf{X}_1 . Partition $\boldsymbol{\beta}$ as $(\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)'$ and $\boldsymbol{\beta}^{(k)}$ as $(\boldsymbol{\beta}^{(k)'}_1, \boldsymbol{\beta}^{(k)'}_2)'$, corresponding to \mathbf{X}_1 and \mathbf{X}_2 , respectively. For $j = 1, \dots, p$, let

$$\omega(j) = \#\{i = 1, \dots, p : |\mathbf{x}_i| = |\mathbf{x}_j|\}.$$

By Lemma 2, for $j = 1, \dots, J$, $\beta_j^{(k)} = \beta_{j'}^{(k)}$ if $\mathbf{x}_j = \mathbf{x}_{j'}$ and $\beta_j^{(k)} = -\beta_{j'}^{(k)}$ otherwise, where $j' \in \{J+1, \dots, p\}$. It follows that $\{\boldsymbol{\beta}_2^{(k)}\}$ can be viewed as an OEM sequence for solving

$$\min_{\boldsymbol{\theta}} \|\mathbf{Y} - \tilde{\mathbf{X}}\boldsymbol{\theta}\|^2 + 2 \sum_{j=1}^{p-J} |\theta_j|, \quad (41)$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_{p-J})'$, and the columns of $\tilde{\mathbf{X}}$ are $\omega(J+1)\mathbf{x}_{J+1}, \dots, \omega(p)\mathbf{x}_p$. Because the objective function in (41) is strictly convex, by Theorem 4, $\{\boldsymbol{\beta}_2^{(k)}\}$ converges to a limit with grouping coherence. This completes the proof. \square

Example 7. Consider \mathbf{X} in Example 4. Let $\mathbf{Y} = (2, 1, -4, 1.5)'$. Using an initial point $\boldsymbol{\beta}^{(0)} = 0$, the OEM sequence of the lasso with $\lambda = 1$ converges to

$$\hat{\boldsymbol{\beta}} = (-0.5625, 0.4375, -0.6875, 0.6875, -0.4375, 0.5625)',$$

which has grouping coherence. For the same data and the same initial point, the coordinate descent sequence converges to $2(-0.5625, 0.4375, -0.6875, 0, 0, 0)'$, which does not have grouping coherence.

Theorem 7 shows that, if the initial point has grouping coherence, then every limit point of the OEM sequence for the lasso inherits this property. It is now tempting to ask whether such a result holds for the OEM sequence with $\lambda = 0$ in (16), i.e., the Healy and Westmacott procedure. Since full aliasing is just one possible culprit for making the matrix \mathbf{X} in (1) lack full column rank and hence $\mathbf{X}'\mathbf{X}$ become singular, Theorem 8 provides an answer to this question for the general singular situation.

Let r denote the rank of \mathbf{X} . When $r < p$, the singular value decomposition (Wilkinson 1965) of \mathbf{X} is

$$\mathbf{X} = \mathbf{U}' \begin{pmatrix} \mathbf{\Gamma}_0^{1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{V},$$

where \mathbf{U} is an $n \times n$ matrix, \mathbf{V} is a $p \times p$ orthogonal matrix and $\mathbf{\Gamma}_0$ is a diagonal matrix whose diagonal elements are the positive eigenvalues, $\gamma_1 \geq \dots \geq \gamma_r$, of $\mathbf{X}'\mathbf{X}$. Define

$$\hat{\boldsymbol{\beta}}^* = (\mathbf{X}'\mathbf{X})^+ \mathbf{X}'\mathbf{Y}, \quad (42)$$

where $^+$ denotes the Moore-Penrose generalized inverse (Ben-Israel and Greville 2003).

Theorem 8. Suppose that $\mathbf{X}'\mathbf{X} + \mathbf{\Delta}'\mathbf{\Delta} = \gamma_1 \mathbf{I}_p$. If $\boldsymbol{\beta}^{(0)}$ lies in the linear space spanned by the first r columns of \mathbf{V}' , then as $k \rightarrow \infty$, for the OEM sequence $\{\boldsymbol{\beta}^{(k)}\}$ of the ordinary least squares, $\boldsymbol{\beta}^{(k)} \rightarrow \hat{\boldsymbol{\beta}}^*$.

Proof. Denote $\mathbf{D} = \mathbf{I}_p - \gamma_1^{-1} \mathbf{X}'\mathbf{X}$. Note that $\boldsymbol{\beta}^{(k+1)} = \gamma_1^{-1} \mathbf{X}'\mathbf{Y} + \mathbf{D}\boldsymbol{\beta}^{(k)}$. By induction,

$$\begin{aligned} \boldsymbol{\beta}^{(k)} &= \gamma_1^{-1} (\mathbf{I}_p + \mathbf{D} + \dots + \mathbf{D}^{k-1}) \mathbf{X}'\mathbf{Y} + \mathbf{D}^k \boldsymbol{\beta}^{(0)} \\ &= \gamma_1^{-1} \mathbf{V}' \left\{ \mathbf{I}_p + \begin{pmatrix} \mathbf{I}_r - \gamma_1^{-1} \mathbf{\Gamma}_0 & \mathbf{0} \\ \mathbf{0} & -\mathbf{I}_{p-r} \end{pmatrix} + \dots + \begin{pmatrix} (\mathbf{I}_r - \gamma_1^{-1} \mathbf{\Gamma}_0)^{k-1} & \mathbf{0} \\ \mathbf{0} & (-1)^{k-1} \mathbf{I}_{p-r} \end{pmatrix} \right\} \mathbf{V} \\ &\quad \cdot \mathbf{V}' \begin{pmatrix} \mathbf{\Gamma}_0^{1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{U}\mathbf{Y} + \mathbf{D}^k \boldsymbol{\beta}^{(0)} \\ &= \gamma_1^{-1} \mathbf{V}' \begin{pmatrix} \{\mathbf{I}_r + (\mathbf{I}_r - \gamma_1^{-1} \mathbf{\Gamma}_0) + \dots + (\mathbf{I}_r - \gamma_1^{-1} \mathbf{\Gamma}_0)^{k-1}\} \mathbf{\Gamma}_0^{1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{U}\mathbf{Y} + \mathbf{D}^k \boldsymbol{\beta}^{(0)}. \end{aligned}$$

As $k \rightarrow \infty$, we have that

$$\mathbf{D}^k \rightarrow \mathbf{V}' \begin{pmatrix} \mathbf{0} \\ \mathbf{I}_{p-r} \end{pmatrix} \mathbf{V}$$

and $\mathbf{D}^k \boldsymbol{\beta}^{(0)} \rightarrow 0$, which implies that

$$\boldsymbol{\beta}^{(k)} \rightarrow \mathbf{V}' \begin{pmatrix} \boldsymbol{\Gamma}_0^{-1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{U}\mathbf{Y} = \hat{\boldsymbol{\beta}}^*.$$

This completes the proof. □

The condition $\mathbf{X}'\mathbf{X} + \boldsymbol{\Delta}'\boldsymbol{\Delta} = \gamma_1 \mathbf{I}_p$ holds if $\mathbf{S} = \mathbf{I}_p$ in (26).

Remark 8. Computing the Moore-Penrose generalized inverse $\hat{\boldsymbol{\beta}}^*$ in (42) is a difficult problem. Theorem 8 says that the OEM algorithm provides an efficient solution to this problem. When $r < p$, the limiting matrix $\hat{\boldsymbol{\beta}}^*$ given by an OEM sequence has the following properties. First, it has the minimal Euclidean norm among the least squares estimators $(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}$ (Ben-Israel and Greville 2003). Second, its model error has a simple form, $E[(\hat{\boldsymbol{\beta}}^* - \boldsymbol{\beta})'(\mathbf{X}'\mathbf{X})(\hat{\boldsymbol{\beta}}^* - \boldsymbol{\beta})] = r\sigma^2$. Third, $\mathbf{X}\boldsymbol{\alpha} = \mathbf{0}$ implies $\boldsymbol{\alpha}'\hat{\boldsymbol{\beta}}^* = 0$ for any vector $\boldsymbol{\alpha}$, which immediately implies that $\hat{\boldsymbol{\beta}}^*$ has grouping coherence.

Example 8. Use the same data and the same initial point in Example 7. The OEM sequence of the ordinary least squares converges to

$$\hat{\boldsymbol{\beta}}^* = (-0.6875, 0.5625, -0.8125, 0.8125, -0.5625, 0.6875)',$$

which has grouping coherence.

7 DISCUSSION

For the regularized least squares method with the SCAD penalty, finding a local solution to achieve the oracle property is a well-known open problem. We have proposed an algorithm, called the OEM algorithm, to fill this gap. For the SCAD and MCP penalties, this algorithm can provide a local solution with the oracle property. The discovery of the algorithm is quite accidental, drawing direct impetus from a missing-data problem arising in design of experiments. The introduction of the active orthogonization procedure in Section 3

makes the algorithm applicable to very general data structures from observational studies and experiments. Recent years have witnessed an explosive interest in both the theoretical and computational aspects of penalized methods. Our introduction of an algorithm that not only has desirable numerical convergence properties but also possesses an important theoretical property suggests a new interface between these two aspects. Subsequent work in this direction is expected. The active orthogonization idea is general and may have potential applications beyond the scope of the OEM algorithm, such as other EM algorithms (Meng and Rubin 1991; Meng and Rubin 1993; Meng 2007), data augmentation (Tanner and Wong 1987), Markov chain Monte Carlo (Liu 2001), smoothing splines (Wahba and Luo 1997; Luo 1998), mesh-free methods in numerical analysis (Fasshauer 2007) and parallel computing (Kumar, Grama, Gupta, and Karypis 2003). This procedure may also has intriguing mathematical connection with complementary theory in design of experiments (Tang and Wu 1996; Chen and Hedayat 1996; Xu and Wu 2005).

The result on the oracle property in Section 4 uses the assumption that the sample size n goes to infinity. This result is appealing for practical situations with massive data (SAMSI 2012), such as the data deluge in astronomy, the Internet and marketing (the Economist 2010), large-scale industrial experiments (Xu 2009) and modern simulations in engineering (NAE 2008), to just name a few. For applications like micro-array and image analysis, one might be interested in extending the result to the small n and large p case, like in Fan and Peng (2004). Such an extension, however, poses significant challenges. Even for a fixed p , the penalized likelihood function for the SCAD can have a large number of local minima (Huo and Chen 2010). When p goes to infinity, that number can be prohibitively large, which makes it very difficult to sort out a local minima with the oracle property.

In addition to achieving the oracle property for nonconvex penalties, an OEM sequence has other unique theoretical properties, including convergence to a point having grouping coherence for the lasso, SCAD or MCP and convergence to the Moore-Penrose generalized inverse-based least squares estimator for singular regression matrices. These theoretical results together with the active orthogonization scheme form the main contribution of the article.

A computer package for distributing the OEM algorithm to the general audience is under development and will be released. We now remark on the acceleration issue and directions for future work. The algorithm can be speeded up by using various methods from the EM literature (McLachlan and Krishnan 2008). For example, following the idea in Varadhan and Roland (2008), one can replace the OEM iteration in (14) by

$$\boldsymbol{\beta}^{(k+1)} = \boldsymbol{\beta}^{(k)} - 2\gamma\mathbf{r} + \gamma^2\mathbf{v},$$

where $\mathbf{r} = \mathbf{M}(\boldsymbol{\beta}^{(k)}) - \boldsymbol{\beta}^{(k)}$, $\mathbf{v} = \mathbf{M}(\mathbf{M}(\boldsymbol{\beta}^{(k)})) - \mathbf{M}(\boldsymbol{\beta}^{(k)}) - \mathbf{r}$ and $\gamma = -\|\mathbf{r}\|/\|\mathbf{v}\|$. This scheme is found to lead to significant reduction of the running time in several examples. For problems with very large p , one may consider a hybrid algorithm to combine the OEM and coordinate descent ideas. It partitions $\boldsymbol{\beta}$ in (1) into G groups and in each iteration, it minimizes the objective function l in (35) by using the OEM algorithm with respect to one group while holding the other groups fixed. Here are some details. Group $\boldsymbol{\beta}$ as $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \dots, \boldsymbol{\beta}'_G)'$. For $k = 0, 1, \dots$, solve

$$\boldsymbol{\beta}_g^{(k+1)} = \arg \min_{\boldsymbol{\beta}_g} l(\boldsymbol{\beta}_1^{(k+1)}, \dots, \boldsymbol{\beta}_{g-1}^{(k+1)}, \boldsymbol{\beta}_g, \boldsymbol{\beta}_{g+1}^{(k)}, \dots, \boldsymbol{\beta}_G^{(k)}) \text{ for } g = 1, \dots, G \quad (43)$$

by OEM until convergence. Note that (43) has a much lower dimension than the iteration in (14). For $G = 1$, the hybrid algorithm reduces to the OEM algorithm and for $G = p$, it becomes the coordinate descent algorithm. Theoretical properties of this hybrid algorithm will be studied and reported elsewhere.

Extension of the OEM algorithm can be made by imposing special structures on regression matrices, such as grouped variables (Yuan and Lin 2006; Zhou and Zhu 2008; Huang, Ma, Xie, and Zhang 2009; Zhao, Rocha, and Yu 2009; Wang, Chen, and Li 2009; Xiong 2010), mixtures (Khalili and Chen 2007) and heredity constraints (Yuan, Joseph, and Lin 2007; Yuan, Joseph, and Zou 2009; Choi, Li, and Zhu 2010), among many other possibilities.

APPENDIX: PROOF OF THEOREM 1

Proof. We first give some definitions and notation. Let Φ be the cumulative distribution function of the standard normal random variable. For $a > 2$ and λ in (3) and $d \geq \gamma_1$ in Assumption 1, define

$$s(u; \lambda) = \begin{cases} \text{sign}(u)(|u| - \lambda)_+/d, & \text{when } |u| \leq (d+1)\lambda, \\ \text{sign}(u)\{(a-1)|u| - a\lambda\}/\{(a-1)d - 1\}, & \text{when } (d+1)\lambda < |u| \leq ad\lambda, \\ u/d, & \text{when } |u| > ad\lambda, \end{cases}$$

Let $\mathbf{s}(\mathbf{u}; \lambda) = [s(u_1; \lambda), \dots, s(u_p; \lambda)]'$. The OEM sequence from (12) satisfies the condition that $\boldsymbol{\beta}^{(k+1)} = \mathbf{s}(\mathbf{u}^{(k)}; \lambda_n/n)$, where

$$\mathbf{u}^{(k)} = (\mathbf{u}_1^{(k)'}, \mathbf{u}_2^{(k)'})' = \frac{\mathbf{X}'\mathbf{Y}}{n} + \left(d\mathbf{I}_d - \frac{\mathbf{X}'\mathbf{Y}}{n} \right) \boldsymbol{\beta}^{(k)}. \quad (44)$$

For $k = 1, 2, \dots$, define two sequences of events $A_k = \{\boldsymbol{\beta}_2^{(k)} = \mathbf{0}\}$ and $B_k = \{\boldsymbol{\beta}_1^{(k)} = \mathbf{u}_1^{(k-1)}/d\}$. For $h > 0$ and $k = 0, 1, \dots$, let $C_k^h = \{\|\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1\| \leq h\lambda_n/n\}$. The flow of the proof is to first show that $P(A_k)$, $P(B_k)$ and $P(C_k^h)$ all tend to one at exponential rates as n goes to infinity, thus establishing Theorem 1 (i), then show that $P(\cap_{i=1}^k A_i)$ and $P(\cap_{i=1}^{k+1} B_i)$ tend to one and finally establish Theorem 1(ii) by noting that the asymptotic normality of $\boldsymbol{\beta}_1^{(k)}$ follows when $\cap_{i=1}^k A_i$ and $\cap_{i=1}^{k+1} B_i$ both occur.

Step 1. Let $\mathbf{G} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$ with columns $\mathbf{g}_1, \dots, \mathbf{g}_p$. Let \mathbf{G}_1 denote its submatrix with the first p_1 columns. Let $\tau \geq 0$ denote the largest eigenvalue of $\mathbf{X}'_1\mathbf{X}_2\mathbf{X}'_2\mathbf{X}_1/n^2$. Define

$$h_A = \begin{cases} 1/(2\tau), & \text{when } \tau > 0, \\ +\infty, & \text{when } \tau = 0. \end{cases} \quad (45)$$

Let $(\mathbf{v}_1, \dots, \mathbf{v}_{p_1}) = d\mathbf{I}_{p_1} - n^{-1}\mathbf{X}'_1\mathbf{X}_1$ and $b = \max\{\|\mathbf{v}_1\|, \dots, \|\mathbf{v}_{p_1}\|\}$. Define

$$h_B = \frac{ad}{b}, \quad (46)$$

with a and d given in (3) and Assumption 1, respectively. For $h > 0$, define

$$h_C = \frac{h}{2\eta}, \quad (47)$$

where η , used in Assumption 3, is the largest eigenvalue of $\mathbf{I}_{p_1} - \mathbf{X}'_1 \mathbf{X}_1 / (nd)$.

For C_0^h , since $\boldsymbol{\beta}^{(0)} = \boldsymbol{\beta} + \mathbf{G}' \boldsymbol{\varepsilon}$, we have that

$$\begin{aligned} \mathrm{P}(C_0^h) &= \mathrm{P}(\|\mathbf{G}'_1 \boldsymbol{\varepsilon}\| \leq h\lambda_n/n) \\ &\geq \mathrm{P}(|\mathbf{g}'_j \boldsymbol{\varepsilon}| \leq h\lambda_n/(n\sqrt{p_1}) \text{ for } j = 1, \dots, p_1) \\ &\geq 1 - \sum_{j=1}^{p_1} \left[1 - \mathrm{P}(|\mathbf{g}'_j \boldsymbol{\varepsilon}| \leq h\lambda_n/(n\sqrt{p_1})) \right] \\ &= 1 - 2 \sum_{j=1}^{p_1} \left[1 - \Phi \left(\frac{h\lambda_n}{n\sqrt{p_1}\sigma\|\mathbf{g}_j\|} \right) \right]. \end{aligned} \quad (48)$$

For A_1 , note that

$$\begin{aligned} \mathrm{P}(A_1) &= \mathrm{P}(|\beta_j^{(0)}| \leq \lambda_n/(nd) \text{ for } j = p_1 + 1, \dots, p) \\ &\geq 1 - \sum_{j=p_1+1}^p \left[1 - \mathrm{P}(|\mathbf{g}'_j \boldsymbol{\varepsilon}| \leq \lambda_n/(nd)) \right] \\ &= 1 - 2 \sum_{j=p_1+1}^p \left[1 - \Phi \left(\frac{\lambda_n}{nd\sigma\|\mathbf{g}_j\|} \right) \right]. \end{aligned} \quad (49)$$

For B_1 , note that

$$\begin{aligned} \mathrm{P}(B_1) &= \mathrm{P}(|\beta_j^{(0)}| \geq a\lambda_n/n \text{ for } j = 1, \dots, p_1) \\ &\geq 1 - \sum_{j=1}^{p_1} \left[1 - \mathrm{P}(|\beta_j + \mathbf{g}'_j \boldsymbol{\varepsilon}| \geq a\lambda_n/n) \right] \\ &= 1 - \sum_{j=1}^{p_1} \left[\Phi \left(\frac{n\beta_j + a\lambda_n}{n\sigma\|\mathbf{g}_j\|} \right) - \Phi \left(\frac{n\beta_j - a\lambda_n}{n\sigma\|\mathbf{g}_j\|} \right) \right]. \end{aligned} \quad (50)$$

For any $h > 0$, by (48),

$$\begin{aligned}
P(C_1^h) &\geq P(C_1^h \cap B_1) = P(C_0^h \cap B_1) \\
&\geq 1 - 2 \sum_{j=1}^{p_1} \left[1 - \Phi \left(\frac{h\lambda_n}{n\sqrt{p_1}\sigma\|\mathbf{g}_j\|} \right) \right] \\
&\quad - \sum_{j=1}^{p_1} \left[\Phi \left(\frac{n\beta_j + a\lambda_n}{n\sigma\|\mathbf{g}_j\|} \right) - \Phi \left(\frac{n\beta_j - a\lambda_n}{n\sigma\|\mathbf{g}_j\|} \right) \right].
\end{aligned} \tag{51}$$

Next, consider A_k , B_k and C_k^h , for $k = 2, 3, \dots$. If A_{k-1} occurs, then by (44),

$$\mathbf{u}^{(k-1)} = \frac{\mathbf{X}'\mathbf{X}\boldsymbol{\beta}}{n} + \frac{\mathbf{X}'\boldsymbol{\varepsilon}}{n} + \begin{pmatrix} d\mathbf{I}_{p_1} - \mathbf{X}'_1\mathbf{X}_1/n & -\mathbf{X}'_1\mathbf{X}_2/n \\ -\mathbf{X}'_2\mathbf{X}_1/n & d\mathbf{I}_{p-p_1} - \mathbf{X}'_2\mathbf{X}_2/n \end{pmatrix} \begin{pmatrix} \boldsymbol{\beta}_1^{(k-1)} \\ \mathbf{0} \end{pmatrix}.$$

Thus,

$$\mathbf{u}_1^{(k-1)} = d\boldsymbol{\beta}_1^{(k-1)} + \frac{\mathbf{X}'_1\mathbf{X}_1}{n}[\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^{(k-1)}] + \frac{\mathbf{X}'_1\boldsymbol{\varepsilon}}{n}, \tag{52}$$

$$\mathbf{u}_2^{(k-1)} = \frac{\mathbf{X}'_2\mathbf{X}_1}{n}[\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^{(k-1)}] + \frac{\mathbf{X}'_2\boldsymbol{\varepsilon}}{n}. \tag{53}$$

By (53),

$$\begin{aligned}
P(A_k) &\geq P(A_k \cap A_{k-1} \cap C_{k-1}^{h_A}) \\
&= P(\{|u_j^{(k-1)}| \leq \lambda_n/n \text{ for } j = p_1 + 1, \dots, p\} \cap A_{k-1} \cap C_{k-1}^{h_A}) \\
&\geq P(\{\|\mathbf{u}_2^{(k-1)}\| \leq \lambda_n/n\} \cap A_{k-1} \cap C_{k-1}^{h_A}) \\
&= P(\{\|n^{-1}\mathbf{X}'_2\mathbf{X}_1(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^{(k-1)}) + n^{-1}\mathbf{X}'_2\boldsymbol{\varepsilon}\| \leq \lambda_n/n\} \cap A_{k-1} \cap C_{k-1}^{h_A}) \\
&\geq P(\{\|n^{-1}\mathbf{X}'_2\boldsymbol{\varepsilon}\| \leq \lambda_n/n - \tau\|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^{(k-1)}\|\} \cap A_{k-1} \cap C_{k-1}^{h_A}) \\
&= P(\{\|n^{-1}\mathbf{X}'_2\boldsymbol{\varepsilon}\| \leq \lambda_n/(2n)\} \cap A_{k-1} \cap C_{k-1}^{h_A}) \\
&\geq P(\{|x'_j\boldsymbol{\varepsilon}| \leq \lambda_n/(2\sqrt{p-p_1}) \text{ for } j = p_1 + 1, \dots, p\} \cap A_{k-1} \cap C_{k-1}^{h_A}) \\
&\geq 1 - 2 \sum_{j=p_1+1}^p \left[1 - \Phi \left(\frac{\lambda_n}{2\sqrt{p-p_1}\sigma\|\mathbf{x}_j\|} \right) \right] \\
&\quad - [1 - P(A_{k-1})] - [1 - P(C_{k-1}^{h_A})],
\end{aligned} \tag{54}$$

where h_A is defined in (45).

By (52),

$$\begin{aligned}
\mathbb{P}(B_k) &\geq \mathbb{P}(B_k \cap A_{k-1} \cap C_{k-1}^{h_B}) \\
&= \mathbb{P}(\{|d\beta_j + \mathbf{u}'_j(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^{(k-1)}) + n^{-1}\mathbf{x}'_j\boldsymbol{\varepsilon}| \geq ad\lambda_n/n \text{ for } j = 1, \dots, p_1\} \cap A_{k-1} \cap C_{k-1}^{h_B}) \\
&\geq \mathbb{P}(\{|d\beta_j + n^{-1}\mathbf{x}'_j\boldsymbol{\varepsilon}| \geq |\mathbf{u}'_j(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^{(k-1)})| + ad\lambda_n/n \text{ for } j = 1, \dots, p_1\} \cap A_{k-1} \cap C_{k-1}^{h_B}) \\
&\geq \mathbb{P}(\{|d\beta_j + n^{-1}\mathbf{x}'_j\boldsymbol{\varepsilon}| \geq b\|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_1^{(k-1)}\| + ad\lambda_n/n \text{ for } j = 1, \dots, p_1\} \cap A_{k-1} \cap C_{k-1}^{h_B}) \\
&\geq \mathbb{P}(\{|d\beta_j + n^{-1}\mathbf{x}'_j\boldsymbol{\varepsilon}| \geq 2ad\lambda_n/n \text{ for } j = 1, \dots, p_1\} \cap A_{k-1} \cap C_{k-1}^{h_B}) \\
&\geq 1 - \sum_{j=1}^{p_1} \left[\Phi\left(\frac{nd\beta_j + 2ad\lambda_n}{\sigma\|\mathbf{x}_j\|}\right) - \Phi\left(\frac{nd\beta_j - 2ad\lambda_n}{\sigma\|\mathbf{x}_j\|}\right) \right] - [1 - \mathbb{P}(A_{k-1})] \\
&\quad - [1 - \mathbb{P}(C_{k-1}^{h_B})], \tag{55}
\end{aligned}$$

where h_B is defined in (46).

For any $h > 0$, we have that

$$\begin{aligned}
\mathbb{P}(C_k^h) &\geq \mathbb{P}(C_k^h \cap B_k \cap A_{k-1} \cap C_{k-1}^{h_C}) \\
&= \mathbb{P}(\{\|(\mathbf{I}_{p_1} - \mathbf{X}'_1\mathbf{X}_1/(nd))(\boldsymbol{\beta}_1^{(k-1)} - \boldsymbol{\beta}_1) + \mathbf{X}'_1\boldsymbol{\varepsilon}/(nd)\| \leq h\lambda_n/n\} \cap B_k \cap A_{k-1} \cap C_{k-1}^{h_C}) \\
&\geq \mathbb{P}(\{\eta\|\boldsymbol{\beta}_1^{(k-1)} - \boldsymbol{\beta}_1\| + \|\mathbf{X}'_1\boldsymbol{\varepsilon}/(nd)\| \leq h\lambda_n/n\} \cap B_k \cap A_{k-1} \cap C_{k-1}^{h_C}) \\
&\geq \mathbb{P}(\{\|\mathbf{X}'_1\boldsymbol{\varepsilon}/(nd)\| \leq h\lambda_n/(2n)\} \cap B_k \cap A_{k-1} \cap C_{k-1}^{h_C}) \\
&\geq \mathbb{P}(\{|\mathbf{x}'_j\boldsymbol{\varepsilon}/(nd)| \leq h\lambda_n/(2n\sqrt{p_1}) \text{ for } j = 1, \dots, p_1\} \cap B_k \cap A_{k-1} \cap C_{k-1}^{h_C}) \\
&\geq 1 - 2 \sum_{j=1}^{p_1} \left[1 - \Phi\left(\frac{h\lambda_n}{2n\sqrt{p_1}\sigma\|\mathbf{x}_j\|}\right) \right] - [1 - \mathbb{P}(A_{k-1})] - [1 - \mathbb{P}(B_k)] \\
&\quad - [1 - \mathbb{P}(C_{k-1}^{h_C})], \tag{56}
\end{aligned}$$

where h_C is defined in (47).

Since $1 - \Phi(x) = o(\exp(-x^2/2))$ as $x \rightarrow +\infty$, (49), (50) and (51) imply that

$$\begin{aligned} 1 - P(A_1) &= o(\exp[-c_1(\lambda_n/\sqrt{n})^2]), \\ 1 - P(B_1) &= o(\exp(-c_2n)) = o(\exp[-c_2(\lambda_n/\sqrt{n})^2]), \\ 1 - P(C_1^h) &= o(\exp[-c_3(\lambda_n/\sqrt{n})^2]), \end{aligned}$$

where c_1 , c_2 and c_3 are positive constants. By induction, it now follows from (54), (55) and (56) that

$$\begin{aligned} 1 - P(A_k) &= [1 - P(A_{k-1})] + [1 - P(C_{k-1}^h)] + o(\exp[-c_4(\lambda_n/\sqrt{n})^2]) \\ &= o(k \exp[-c_4(\lambda_n/\sqrt{n})^2]), \end{aligned} \tag{57}$$

where c_4 and c_5 are positive constants. Similarly,

$$1 - P(B_k) = o(k \exp[-c_6(\lambda_n/\sqrt{n})^2]), \tag{58}$$

$$1 - P(C_k^h) = o(k \exp[-c_7(\lambda_n/\sqrt{n})^2]), \tag{59}$$

where c_6 and c_7 are positive constants. By (57), a sufficient condition for $P(A_k) \rightarrow 1$ is

$$k \exp[-c(\lambda_n/\sqrt{n})^2] \rightarrow 0 \tag{60}$$

for any $c > 0$, which is covered by Assumption 3.

Step 2. Consider the asymptotic normality of $\boldsymbol{\beta}_1^{(k)}$. When the events A_{k-1} , A_k , B_k and B_{k+1} all occur, by (52),

$$\frac{\mathbf{X}'_1 \mathbf{X}_1}{n} (\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1) = \frac{\mathbf{X}'_1 \boldsymbol{\varepsilon}}{n} + d(\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1^{(k+1)}), \tag{61}$$

$$\|\boldsymbol{\beta}_1^{(k+1)} - \boldsymbol{\beta}_1^{(k)}\| = \left\| \left(\mathbf{I}_{p_1} - \frac{\mathbf{X}'_1 \mathbf{X}_1}{nd} \right) (\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1^{(k-1)}) \right\| \leq \eta \|\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1^{(k-1)}\|. \tag{62}$$

When the events $C_1^{1/2}$, $C_2^{1/2}$, $A^{(k)} = \bigcap_{i=1}^k A_i$ and $B^{(k+1)} = \bigcap_{i=1}^{k+1} B_i$ all occur, by (62),

$$\|\boldsymbol{\beta}_1^{(k+1)} - \boldsymbol{\beta}_1^{(k)}\| \leq \eta^{k-1} \|\boldsymbol{\beta}_1^{(2)} - \boldsymbol{\beta}_1^{(1)}\| \leq \eta^{k-1} \lambda_n / n. \quad (63)$$

For any $\boldsymbol{\alpha} \in \mathbb{R}^p$ with $\boldsymbol{\alpha} \neq \mathbf{0}$, by (61) and (63),

$$\begin{aligned} & |\sqrt{n} \boldsymbol{\alpha}' (\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1) - \sqrt{n} \boldsymbol{\alpha}' (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \boldsymbol{\varepsilon}| \\ &= dn^{3/2} |\boldsymbol{\alpha}' (\mathbf{X}'_1 \mathbf{X}_1)^{-1} (\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1^{(k+1)})| \\ &\leq n^{-1/2} \lambda_n \eta^{k-1} d \|\boldsymbol{\alpha}\| \|n (\mathbf{X}'_1 \mathbf{X}_1)^{-1}\|. \end{aligned}$$

For any $x \in \mathbb{R}$, note that

$$\begin{aligned} & \left| \mathbb{P}(\sqrt{n} \boldsymbol{\alpha}' (\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1) \leq x) - \Phi \left(\frac{x}{\sigma(\boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\alpha})^{1/2}} \right) \right| \\ &\leq \left| \mathbb{P}(\{\sqrt{n} \boldsymbol{\alpha}' (\boldsymbol{\beta}_1^{(k)} - \boldsymbol{\beta}_1) \leq x\} \cap C_1^{1/2} \cap C_2^{1/2} \cap A^{(k)} \cap B^{(k+1)}) - \Phi \left(\frac{x}{\sigma(\boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\alpha})^{1/2}} \right) \right| \\ &\quad + [1 - \mathbb{P}(C_1^{1/2})] + [1 - \mathbb{P}(C_2^{1/2})] + [1 - \mathbb{P}(A^{(k)})] + [1 - \mathbb{P}(B^{(k+1)})] \\ &\leq \left| \mathbb{P}(\sqrt{n} \boldsymbol{\alpha}' (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \boldsymbol{\varepsilon} \leq x - n^{-1/2} \lambda_n \eta^{k-1} d \|\boldsymbol{\alpha}\| \|n (\mathbf{X}'_1 \mathbf{X}_1)^{-1}\|) - \Phi \left(\frac{x}{\sigma(\boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\alpha})^{1/2}} \right) \right| \\ &\quad + \left| \mathbb{P}(\sqrt{n} \boldsymbol{\alpha}' (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \mathbf{X}'_1 \boldsymbol{\varepsilon} \leq x + n^{-1/2} \lambda_n \eta^{k-1} d \|\boldsymbol{\alpha}\| \|n (\mathbf{X}'_1 \mathbf{X}_1)^{-1}\|) - \Phi \left(\frac{x}{\sigma(\boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\alpha})^{1/2}} \right) \right| \\ &\quad + 3[1 - \mathbb{P}(C_1^{1/2})] + 3[1 - \mathbb{P}(C_2^{1/2})] + 3[1 - \mathbb{P}(A^{(k)})] + 3[1 - \mathbb{P}(B^{(k+1)})] \\ &= \left| \Phi \left(\frac{x - n^{-1/2} \lambda_n \eta^{k-1} d \|\boldsymbol{\alpha}\| \|n (\mathbf{X}'_1 \mathbf{X}_1)^{-1}\|}{\sigma[\boldsymbol{\alpha}' n (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \boldsymbol{\alpha}]^{1/2}} \right) - \Phi \left(\frac{x}{\sigma(\boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\alpha})^{1/2}} \right) \right| \\ &\quad + \left| \Phi \left(\frac{x + n^{-1/2} \lambda_n \eta^{k-1} d \|\boldsymbol{\alpha}\| \|n (\mathbf{X}'_1 \mathbf{X}_1)^{-1}\|}{\sigma[\boldsymbol{\alpha}' n (\mathbf{X}'_1 \mathbf{X}_1)^{-1} \boldsymbol{\alpha}]^{1/2}} \right) - \Phi \left(\frac{x}{\sigma(\boldsymbol{\alpha}' \boldsymbol{\Sigma}_1^{-1} \boldsymbol{\alpha})^{1/2}} \right) \right| \\ &\quad + 6[1 - \mathbb{P}(C_1^{1/2})] + 3[1 - \mathbb{P}(A^{(k)})] + 3[1 - \mathbb{P}(B^{(k+1)})]. \end{aligned} \quad (64)$$

Now, under Assumption 3, by (57), (58) and (59), (64) converges to zero as $n \rightarrow \infty$. This completes the proof. \square

ACKNOWLEDGEMENTS

Xiong is supported by grant 10801130 of the National Natural Science Foundation of China. Dai is partially support by Grace Wahba through NIH grant R01 EY009946, ONR grant N00014-09-1-0655 and NSF grant DMS-0906818. Qian is partially supported by NSF grant CMMI 0969616, an IBM Faculty Award and an NSF Faculty Early Career Development Award DMS 1055214. The authors thank Xiao-Li Meng and Grace Wahba for comments and suggestions.

REFERENCES

- Addelman, S. (1962), “Orthogonal Main-Effect Plans for Asymmetrical Factorial Experiments,” *Technometrics*, 4, 21-46.
- Ben-Israel, A. and Greville, T. N. E. (2003), *Generalized Inverses, Theory and Applications*, 2nd ed., New York: Springer.
- Bertsekas, D. P. (1999), *Nonlinear Programming*, 2nd ed., Belmont, MA: Athena Scientific.
- Bingham, D., Sitter, R. R. and Tang, B. (2009). “Orthogonal and Nearly Orthogonal Designs for Computer Experiments,” *Biometrika*, 96, 51–65.
- Breiman, L. (1995), “Better Subset Regression Using the Nonnegative Garrote,” *Technometrics*, 37, 373–384.
- Breheny, P. and Huang, J. (2011), “Coordinate Descent Algorithms for Nonconvex Penalized Regression, With Applications to Biological Feature Selection,” *The Annals of Applied Statistics*, in press.
- Chen, H. and Hedayat, A. S. (1996), “ 2^{n-l} designs with Weak Minimum Aberration,” *The Annals of Statistics*, 24, 2536–2548.
- Choi, N. H., Li W. and Zhu J. (2010), “Variable Selection With the Strong Heredity Constraint and Its Oracle Property,” *Journal of the American Statistical Association*, 105, 354–364.

- Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977), “Maximum Likelihood from Incomplete Data via the EM Algorithm,” *Journal of the Royal Statistical Society, Ser. B*, 39, 1–38.
- Deng, X., Lin, C. D. and Qian, P. Z. G. (2010), “Designs for the Lasso,” *Technical Report*.
- Dey, A. and Mukerjee, R. (1999), *Fractional Factorial Plans*, New York: Wiley.
- Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004), “Least Angle Regression,” *The Annals of Statistics*, 32, 407–451.
- Fan, J. and Li, R. (2001), “Variable Selection via Nonconcave Penalized Likelihood and Its Oracle Properties,” *Journal of the American Statistical Association*, 96, 1348–1360.
- Fan, J. and Peng, H. (2004), “On Non-concave Penalized Likelihood With Diverging Number of Parameters,” *The Annals of Statistics*, 32, 928–961.
- Fang, K. T., Li, R. and Sudjianto, A. (2005), *Design and Modeling for Computer Experiments*, New York: Chapman and Hall/CRC Press.
- Fasshauer, G. F. (2007), *Meshfree Approximation Methods With MATLAB*, Singapore: World Scientific Publishing Company.
- Frank, L. E. and Friedman, J. (1993), “A Statistical View of Some Chemometrics Regression Tools,” *Technometrics*, 35, 109–135.
- Friedman, J., Hastie, T., Hofling, H. and Tibshirani, R. (2007), “Pathwise Coordinate Optimization,” *The Annals of Applied Statistics*, 1, 302–332.
- Friedman, J., Hastie, T. and Tibshirani, R. (2009), “Regularization Paths for Generalized Linear Models via Coordinate Descent,” *Journal of Statistical Software*, 33, 1–22.
- Friedman, J., Hastie, T. and Tibshirani, R. (2011), “Glmnet,” R package.
- Fu, W. J. (1998), “Penalized Regressions: The Bridge versus the LASSO,” *Journal of Computational and Graphical Statistics*, 7, 397–416.

- Grandvalet, Y. (1998), “Least Absolute Shrinkage Is Equivalent to Quadratic Penalization,”
*In: Niklasson, L., Bodén, M., Ziemcke, T. (eds.), ICANN’98. Vol. 1 of Perspectives
in Neural Computing, Springer, 201–206.*
- Green, P. J. (1990), “On Use of the EM Algorithm for Penalized Likelihood Estimation,”
Journal of the Royal Statistical Society, Ser. B, 52, 443–452.
- Hastie, T. and Efron, B. (2011), “Lars,” R package.
- Healy, M. J. R. and Westmacott, M. H. (1956), “Missing Values in Experiments Analysed
on Automatic Computers,” *Journal of the Royal Statistical Society, Ser. C, 5, 203–206.*
- Hedayat, A. S., Sloane, N. J. A. and Stufken, J. (1999), *Orthogonal Arrays: Theory and
Applications*, New York: Springer.
- Hoerl, A. E. and Kennard, R. W. (1970), “Ridge Regression: Biased Estimation for
Nonorthogonal Problems,” *Technometrics, 12, 55–67.*
- Huang, J., Ma, S., Xie, H. and Zhang, C-H. (2009), “A Group Bridge Approach for Variable
Selection,” *Biometrika, 96, 339–355.*
- Hunter, D. R. and Li, R. (2005), “Variable Selection Using MM Algorithms,” *The Annals
of Statistics, 33, 1617–1642.*
- Huo, X. and Chen, J. (2010), “Complexity of Penalized Likelihood Estimation,” *Journal of
Statistical Computation and Simulation, 80, 747–759 .*
- Huo, X. and Ni, X. L. (2007), “When Do Stepwise Algorithms Meet Subset Selection
Criteria?” *The Annals of Statistics, 35, 870–887.*
- Khalili, A. and Chen, J. (2007), “Variable Selection in Finite Mixture of Regression Mod-
els,” *Journal of the American Statistical Association, 102, 1025–1038.*
- Kumar, V., Grama, A., Gupta, A. and Karypis, G. (2003), *Introduction to Parallel Com-
puting*, 2nd ed., Boston, MA: Addison Wesley.

- Lin, C. D., Bingham, D., Sitter, R. R. and Tang, B. (2010), “A New and Flexible Method for Constructing Designs for Computer Experiments,” *The Annals of Statistics*, 38, 1460–1477.
- Lin, C. D., Mukerjee, R. and Tang, B. (2009), “Construction of Orthogonal and Nearly Orthogonal Latin Hypercubes,” *Biometrika*, 96, 243–247.
- Li, R. and Lin, D. K. J. (2002), “Data Analysis in Supersaturated Designs,” *Statistics and Probability Letters*, 59, 135–144.
- Lin, D. K. J. (1993), “A New Class of Supersaturated Designs,” *Technometrics*, 35, 28–31.
- Liu, J. S. (2001), *Monte Carlo Strategies in Scientific Computing*, New York: Springer.
- Luo, Z. (1998), “Backfitting in Smoothing Spline ANOVA,” *The Annals of Statistics*, 26, 1733–1759.
- MacWilliams, F. J. and Sloane, N. J. A. (1977), *The Theory of Error Correcting Codes*, New York: North-Holland.
- Mazumder, R., Friedman, J. and Hastie, T. (2010), “SparseNet: Coordinate Descent with Non-Convex Penalties,” *Technical Report*.
- McLachlan, G. and Krishnan, T. (2008), *The EM Algorithm and Extensions*, 2nd ed., New York: Wiley.
- Meng, X. L. (1994), “On the Rate of Convergence of the ECM Algorithm,” *The Annals of Statistics*, 22, 326–339.
- Meng, X. L. (2007), “Thirty Years of EM and Much More,” *Statistica Sinica*, 17, 839–840.
- Meng, X. L. and Rubin, D. (1991), “Using EM to Obtain Asymptotic Variance-Covariance Matrices: The SEM Algorithm,” *Journal of the American Statistical Association*, 86, 899–909.

- Meng, X. L. and Rubin, D. (1993), “Maximum Likelihood Estimation via the ECM Algorithm: A General Framework,” *Biometrika*, 80, 267–278.
- Mukerjee, R. and Wu, C. F. J. (2006), *A Modern Theory of Factorial Design*, New York: Springer.
- NAE (2008), “Grand Challenges for Engineering,” *Technical report*, The National Academy of Engineering.
- Nettleton, D. (1999), “Convergence Properties of the EM Algorithm in Constrained Parameter Spaces,” *Canadian Journal of Statistics*, 27, 639–648.
- Osborne, M. R., Presnell, B. and Turlach, B. (2000), “On the LASSO and Its Dual,” *Journal of Computational and Graphical Statistics*, 9, 319–337.
- Owen, A. B. (1994), “Controlling Correlations in Latin Hypercube Samples,” *Journal of the American Statistical Association*, 89, 1517–1522.
- Owen, A. B. (2006), “A Robust Hybrid of Lasso and Ridge Regression,” *Technical Report*.
- Pang, F., Liu, M. Q. and Lin, D. K. J. (2009), “A Construction Method for Orthogonal Latin Hypercube Designs with Prime Power Levels,” *Statistica Sinica*, 19, 1721–1728.
- Phoa, F. K. H., Pan, Y-H. and Xu, H. (2009), “Analysis of Supersaturated Designs via The Dantzig Selector,” *Journal of Statistical Planning and Inference*, 139, 2362–2372.
- Rubin, D. B. and Szatrowski, T. H. (1982), “Finding Maximum Likelihood Estimates of Patterned Covariance Matrices by the EM Algorithm,” *Biometrika*, 69, 657–660.
- SAMSI (2012), Program on Statistical and Computational Methodology for Massive Datasets, 2012-2013.
- Santner, T. J., Williams, B. J. and Notz, W. I. (2003), *The Design and Analysis of Computer Experiments*, New York: Springer.

- Schifano, E. D., Strawderman, R. and Wells, M. T. (2010), “Majorization-Minimization Algorithms for Nonsmoothly Penalized Objective Functions,” *Electronic Journal of Statistics*, 23, 1258–1299.
- Steinberg, D. M. and Lin, D. K. J. (2006), “A Construction Method for Orthogonal Latin Hypercube Designs,” *Biometrika*, 93, 279–288.
- Sun, F., Liu, M. and Lin, D. K. J. (2009), “Construction of Orthogonal Latin Hypercube Designs,” *Biometrika*, 96, 971–974.
- Tang, B. (1998), “Selecting Latin Hypercubes Using Correlation Criteria,” *Statistica Sinica*, 8, 965–977.
- Tang, B. and Wu, C. F. J. (1996), “Characterization of Minimum Aberration 2^{n-k} Designs in terms of Their Complementary Designs,” *The Annals of Statistics*, 24, 2549–2559.
- Tang, B. and Wu, C. F. J. (1997), “A Method for Constructing Supersaturated Designs and Its $E(s^2)$ Optimality,” *Canadian Journal of Statistics*, 25, 191–201.
- Tanner, M. A. and Wong, W. H. (1987). “The Calculation of Posterior Distributions by Data Augmentation,” *Journal of the American Statistical Association*, 82, 528–540.
- The Economist (2010), “Special Report on the Data Deluge, (February 27),” *The Economist*, 394, 3–18.
- Tibshirani, R. (1996), “Regression Shrinkage and Selection via the Lasso,” *Journal of the Royal Statistical Society, Ser. B*, 58, 267–288.
- Tseng, P. (2001), “Convergence of a Block Coordinate Descent Method for Nondifferentiable Minimization,” *Journal of Optimization Theory and Applications*, 109, 475–494.
- Tseng, P. and Yun, S. (2009), “A Coordinate Gradient Descent Method for Nonsmooth Separable Minimization,” *Mathematical Programming B*, 117, 387–423.

- Varadhan, R. and Roland, C. (2008), “Simple and Globally Convergent Methods for Accelerating the Convergence of Any EM Algorithm,” *Scandinavian Journal of Statistics*, 35, 335–353.
- Wahba, G. and Luo, Z. (1997), “Smoothing Spline ANOVA Fits for Very Large, Nearly Regular Data Sets, with Application to Historical Global Climate Data,” *Annals of Numerical Mathematics*, 4579-4598. (Festschrift in Honor of Ted Rivlin, C. Micchelli, Ed.)
- Wang, L., Chen, G. and Li, H. (2007), “Group SCAD Regression Analysis for Microarray Time Course Gene Expression Data,” *Bioinformatics*, 23, 1486–1494.
- Watkins, D. S. (2002), *Fundamentals of Matrix Computations*, 2nd ed., New York: Wiley.
- Wilkinson, J. H. (1965), *The Algebraic Eigenvalue Problem*, New York: Oxford University Press.
- Witten, D. M. and Tibshirani, R. (2011), “Scout,” R package.
- Wu, C. F. J. (1983), “On the Convergence Properties of the EM Algorithm,” *The Annals of Statistics*, 11, 95–103.
- Wu, C. F. J. (1993), “Construction of Supersaturated Designs through Partially Aliased Interactions,” *Biometrika*, 80, 661–669.
- Wu, C. F. J. and Hamada M. (2009), *Experiments: Planning, Analysis, and Optimization*, 2nd ed., New York: Wiley.
- Wu, T. and Lange, K. (2008), “Coordinate Descent Algorithm for Lasso Penalized Regression,” *The Annals of Applied Statistics*, 2, 224–244.
- Xiong, S. (2010), “Some Notes on the Nonnegative Garrote,” *Technometrics*, 52, 349–361.
- Xu, H. (2009), “Algorithmic Construction of Efficient Fractional Factorial Designs With Large Run Sizes,” *Technometrics*, 51, 262–277.

- Xu, H. and Cheng, C-S. (2008), “A Complementary Design Theory for Doubling,” *The Annals of Statistics*, 36, 445–457.
- Xu, H. and Wu, C. F. J. (2005), “Construction of Optimal Multi-Level Supersaturated Designs,” *The Annals of Statistics*, 33, 2811–2836.
- Ye, K. Q. (1998). “Orthogonal Column Latin Hypercubes and Their Applications in Computer Experiments,” *Journal of the American Statistical Association*, 93, 1430–1439.
- Yuan, M. and Lin, Y. (2006), “Model Selection and Estimation in Regression with Grouped Variables,” *Journal of the Royal Statistical Society, Ser. B*, 68, 49–68.
- Yuan, M., Joseph, V. R. and Lin, Y. (2007), “An Efficient Variable Selection Approach for Analyzing Designed Experiments,” *Technometrics*, 49, 430–439.
- Yuan, M., Joseph, V. R. and Zou, H. (2009), “Structured Variable Selection and Estimation,” *Annals of Applied Statistics*, 3, 1738–1757.
- Zangwill, W. I. (1969), *Nonlinear Programming: A Unified Approach*, Englewood Cliffs, New Jersey: Prentice Hall.
- Zhang, C-H. (2010), “Nearly Unbiased Variable Selection under Minimax Concave Penalty,” *The Annals of Statistics*, 38, 894–942.
- Zhao, P., Rocha, G. and Yu, B. (2009), “The Composite Absolute Penalties Family for Grouped and Hierarchical Variable Selection,” *The Annals of Statistics*, 37, 3468–3497.
- Zhou, N. and Zhu, J. (2008), “Group Variable Selection via a Hierarchical Lasso and Its Oracle Property,” *Technical Report*.
- Zhu, Y. (2011), *Personal communication*.
- 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.

Zou, H. and Li, R. (2008), “One-step Sparse Estimates in Nonconcave Penalized Likelihood Models,” *The Annals of Statistics*, 36, 1509–1533.