



## Sequential Lasso Cum EBIC for Feature Selection With Ultra-High Dimensional Feature Space

Shan Luo & Zehua Chen

**To cite this article:** Shan Luo & Zehua Chen (2014) Sequential Lasso Cum EBIC for Feature Selection With Ultra-High Dimensional Feature Space, Journal of the American Statistical Association, 109:507, 1229-1240, DOI: [10.1080/01621459.2013.877275](https://doi.org/10.1080/01621459.2013.877275)

**To link to this article:** <https://doi.org/10.1080/01621459.2013.877275>



View supplementary material [↗](#)



Published online: 02 Oct 2014.



Submit your article to this journal [↗](#)



Article views: 1368



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 14 View citing articles [↗](#)

# Sequential Lasso Cum EBIC for Feature Selection With Ultra-High Dimensional Feature Space

Shan LUO and Zehua CHEN

In this article, we propose a method called sequential Lasso (SLasso) for feature selection in sparse high-dimensional linear models. The SLasso selects features by sequentially solving partially penalized least squares problems where the features selected in earlier steps are not penalized. The SLasso uses extended BIC (EBIC) as the stopping rule. The procedure stops when EBIC reaches a minimum. The asymptotic properties of SLasso are considered when the dimension of the feature space is ultra high and the number of relevant feature diverges. We show that, with probability converging to 1, the SLasso first selects all the relevant features before any irrelevant features can be selected, and that the EBIC decreases until it attains the minimum at the model consisting of exactly all the relevant features and then begins to increase. These results establish the selection consistency of SLasso. The SLasso estimators of the final model are ordinary least squares estimators. The selection consistency implies the oracle property of SLasso. The asymptotic distribution of the SLasso estimators with diverging number of relevant features is provided. The SLasso is compared with other methods by simulation studies, which demonstrates that SLasso is a desirable approach having an edge over the other methods. The SLasso together with the other methods are applied to a microarray data for mapping disease genes. Supplementary materials for this article are available online.

**KEY WORDS:** Extended BIC; Oracle property; Selection consistency; Sparse high-dimensional linear models.

## 1. INTRODUCTION

Sparse high-dimensional regression (SHR) models arise in many important contemporary scientific fields. A SHR model is

$$y_i = \beta_0 + \sum_{j=1}^p \beta_j x_{ij} + \epsilon_i, \quad i = 1, \dots, n, \quad (1.1)$$

where the number of features  $p$  is much larger than the sample size  $n$ , and only a relatively small number of the  $\beta_j$ 's are nonzero. Feature selection is crucial in the analysis of SHR models. There are usually two goals of feature selection: (i) to build a model with desirable prediction properties and (ii) to identify the features with nonzero coefficients (for convenience, such features are referred to as relevant features in this article). These two goals are intertwined but are not the same.

Regularized regression approaches to the analysis of SHR models have attracted considerable attention of the researchers. A regularized regression approach selects the features and estimates the coefficients simultaneously by minimizing a penalized sum of squares of the form

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \sum_{j=1}^p p_\lambda(|\beta_j|), \quad (1.2)$$

where  $\lambda$  is a regulating parameter and  $p_\lambda$  is a penalty function such that the number of fitted nonzero coefficients can be regulated by  $\lambda$ ; that is, only a certain number of  $\beta_j$ 's are estimated nonzero when  $\lambda$  is set at a certain value. Various penalty functions have been proposed and studied, including

Lasso (Tibshirani 1996):  $p_\lambda(|\beta_j|) = \lambda|\beta_j|$ , SCAD (Fan and Li 2001), which smoothly clips a  $L_1$  penalty (for small  $|\beta_j|$ ) and a constant penalty (for large  $|\beta_j|$ ), adaptive Lasso (Zou 2006):  $p_\lambda(|\beta_j|) = \lambda w_j |\beta_j|$ , where  $w_j$  are given weights, and MCP (Zhang 2010), which smoothly approaches the  $L_1$  penalty from a constant penalty (for large  $|\beta_j|$ 's) by an asymptote.

A so-called oracle property is of major concern for any feature selection method. The oracle property refers to two asymptotic natures: (i) selection consistency, that is, the sparse relevant features can be exactly selected with probability converging to 1, and (ii) the effects of relevant features can be consistently estimated the same as they would be, were they obtained by knowing the relevant features in advance. For fixed  $p$ , it was shown in Knight and Fu (2000), and Leng, Lin, and Wahba (2004) that Lasso is consistent in estimating the regression coefficients but, in general, it does not have the oracle property. A condition on the feature matrix was provided in Zou (2006) for Lasso to possess the oracle property. The condition was also discovered in Meinshausen and Bühlmann (2006), and Zhao and Yu (2006) and was dubbed as irreducibility condition in Zhao and Yu (2006). When  $p$  is allowed to diverge to infinity faster than  $n$  (but not too fast), the selection consistency of Lasso under the irreducibility condition was established in Meinshausen and Bühlmann (2006), and Zhao and Yu (2006). To relax the irreducibility condition, adaptive Lasso was considered in Zou (2006) for fixed  $p$  using the ordinary least-square estimates as the weights in the penalty, and adaptive Lasso was shown to have the oracle property. For diverging  $p$ , Huang, Ma, and Zhang (2008) showed that Adaptive Lasso with marginal least-square estimates as the weights has the oracle property if a partial orthogonality condition holds. The properties of SCAD were studied in Fan and Li (2001, 2004), Fan and Peng (2004), and Xie and Huang (2009). In these papers, the oracle property of SCAD was established for various models when  $p$  is fixed

Shan Luo, Department of Mathematics, Shanghai Jiao Tong University (E-mail: [sluo2012@gmail.com](mailto:sluo2012@gmail.com)). Zehua Chen, Department of Statistics and Applied Probability, National University of Singapore (E-mail: [stachen@nus.edu.sg](mailto:stachen@nus.edu.sg)). The authors express their gratitude to the associate editor and the two anonymous referees for their valuable comments and suggestions which led to a great improvement of the original manuscript of this article. They also thank the co-editor, Prof. Jun Liu, for his patience and encouragement during the long review process. The research of the authors, which led to this article, was supported by Singapore Ministry of Education grants: R-155-000-091-112 and R-155-000-125-112.

or diverging to infinity not too fast. The MCP penalty is similar to the SCAD penalty. The asymptotic properties of the MCP penalty were studied in Zhang (2010). To realize the oracle property of the various regularized regression methods in finite samples, a proper choice of the regulating parameter has to be made. A multifold cross-validation (CV) is commonly used in these methods for the choice of the regulating parameter.

Sequential methods have also received attention in recent years for feature selection in SHR models. The traditional sequential procedures such as forward stepwise regression (FSR) were criticized for their greedy nature. However, it was discovered recently that the greedy nature is indeed a good one if the goal is to identify relevant features, see Tropp (2004), and Tropp and Gilbert (2007), especially, in the presence of high spurious correlations due to extremely high dimensionality of the feature space. In many practical problems, the identification of the relevant features is of primary interest. For example, in genetic quantitative trait loci (QTL) mapping and disease gene mapping, of interest are the markers which are either QTL or disease gene themselves or are in linkage disequilibrium with QTL or disease genes. This revived interests in sequential approaches.

The properties of FSR for feature selection in SHR models were re-examined in Wang (2009). It was shown that FSR has a so-called sure screening property when the procedure is carried out until a certain step before the number of steps reaches the sample size. The sure screening property means that the selected set contains the set of relevant features with probability converging to 1; see Fan and Lv (2008).

A different version of forward regression referred to as forward selection in Weisberg (1980) was reconsidered and dubbed as orthogonal matching pursuit (OMP) in Pati, Rezaiifar, and Krishnaprasad (1993). At each step of OMP, the response vector is projected onto the space spanned by the currently selected features, and the next feature is selected to maximize the correlation with the current residual. The procedure stops when the residual is reduced below a certain level specified by a stopping rule. The properties of OMP have been studied under quite strict conditions in, for example, Cai and Wang (2011), Tropp (2004), and Tropp and Gilbert (2007). A thorough investigation of the properties of OMP is still lacking.

An adaptive forward-backward greedy algorithm (FoBa) was considered in Zhang (2011). The FoBa is a variant of OMP. At the forward step of FoBa, the same mechanism of OMP is used to select new features. A new feature is selected if the amount of decrease in the residual exceeds a specified threshold. Following each forward step, a backward step is carried out if the amount of increase in the residual caused by deleting one of the selected features is less than half of the amount of increase at the forward step. The procedure stops when the amount of decrease in the residual at the forward step falls below the specified threshold. Some oracle inequalities have been derived under a so-called restricted isometry condition in Zhang (2011), however, whether or not FoBa is selection consistent is left untouched.

A nonsequential but closely related procedure called compressive sampling matching pursuit (CoSaMP) was proposed in Needell and Tropp (2009). The procedure of CoSaMP assumes the knowledge of sparsity level  $k$ , that is, the number of relevant features. For a given  $k$ , the CoSaMP starts with a sparse set consisting of the  $k$  features having highest correlations with

the response vector, then updates the sparse set by iteration. At each iteration, the response vector is projected onto the space spanned by the sparse set,  $2k$  additional features having highest correlations with the residual of the projection together with the sparse set are fitted to a least square regression model, the updated sparse set consists of the  $k$  features having the largest absolute fitted coefficients in the regression model. The number of iterations is either fixed or determined by a certain rule. Eventually, the true sparsity level is chosen by a certain method. Though the procedure of CoSaMP is appealing, its properties are not fully investigated.

A sequential procedure of a different nature called least angle regression (LAR) was proposed in Efron et al. (2004). The LAR continuously updates the estimate of the expected responses along a direction having equal angle with the features already selected and selects new features having the largest absolute correlation with the updated current residuals. The fitted regression coefficients at each step are shrunk. The LAR algorithm has been modified to compute the solution path of Lasso. There are also variants of LAR, for example, the forward Lasso adaptive shrinkage (FLASH) considered in Radchenko and James (2011). The fitted regression coefficients at each step of FLASH are not fully shrunk as in LAR or Lasso.

Sequential approach has also been considered for models other than SHR models. For example, a sequential procedure called correlation pursuit has been proposed in Zhong et al. (2012) for selecting predictors and estimating index coefficients simultaneously for index models.

Besides their desirable theoretical properties, sequential approaches are computationally more appealing. They are more stable and less affected by the dimensionality of the feature space.

In this article, we propose a sequential procedure called sequential Lasso (SLasso) with the emphasis on the goal of identifying relevant features. We give a conceptual description of SLasso in the following. Its computation algorithm is given in Section 2. In summary, SLasso solves a sequence of partially penalized least squares problems. The features selected in an earlier step are not penalized in the subsequent steps. Let the vectors  $\mathbf{y} = (y_1, \dots, y_n)^T$ ,  $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})^T$ , be standardized such that they have length  $\sqrt{n}$  and are orthogonal to the vector with all elements 1. Thus in model (1.1), the intercept  $\beta_0$  can be omitted. At the initial step, SLasso minimizes the following penalized sum of squares:

$$l_1 = \left\| \mathbf{y} - \sum_{j=1}^p \beta_j \mathbf{x}_j \right\|^2 + \lambda_1 \sum_{j=1}^p |\beta_j|,$$

where  $\|\cdot\|$  is the  $L_2$ -norm, and  $\lambda_1$  is the largest value of the penalty parameter such that at least one of the  $\beta_j$ 's will be estimated nonzero. The features with nonzero estimated coefficients are selected and the set of their indices is denoted by  $s_{*1}$ . For  $k \geq 1$ , let  $s_{*k}$  be the index set of the features selected until step  $k$ . At step  $k+1$ , SLasso minimizes the following partially penalized sum of squares:

$$l_{k+1} = \left\| \mathbf{y} - \sum_{j=1}^p \beta_j \mathbf{x}_j \right\|^2 + \lambda_{k+1} \sum_{j \notin s_{*k}} |\beta_j|,$$

where no penalty is imposed on the  $\beta_j$ 's for  $j \in s_{*k}$  and  $\lambda_{k+1}$  is the largest value of the penalty parameter such that at least one of the  $\beta_j$ 's,  $j \notin s_{*k}$ , will be estimated nonzero. The selected set is then updated to  $s_{*k+1}$ . The EBIC proposed in Chen and Chen (2008) is used as the stopping rule. For each  $s_{*k}$ , the EBIC of the model with features in  $s_{*k}$  is computed. The procedure continues, if the EBIC keeps decreasing. If the EBIC attains a minimum at step  $k^*$ , the procedure stops and the set  $s_{*k^*}$  is taken as the final selected set.

The minimization of  $l_{k+1}$  is equivalent to the minimization of

$$\|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\tilde{\boldsymbol{\beta}}\|^2 + \lambda_{k+1} \sum_{j \notin s_{*k}} |\beta_j|, \quad (1.3)$$

where  $\tilde{\mathbf{y}}$  is the residual of  $\mathbf{y}$  projected on the space spanned by the  $\mathbf{x}_j$ 's with  $j \in s_{*k}$  and  $\tilde{\mathbf{X}}$  is the residual matrix of the  $\mathbf{x}_j$ 's,  $j \notin s_{*k}$ , projected on the same space; see Proposition 2. The active features  $\mathbf{x}_j$  in the minimization of (1.3) must attain  $\max_{j' \notin s_{*k}} |\tilde{\mathbf{y}}^\tau \mathbf{x}_{j'}|$ . Thus, the minimization of (1.3) further reduces to the minimization of

$$\|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}_{\text{TEMP}}\tilde{\boldsymbol{\beta}}_{\text{TEMP}}\|^2 + \lambda_{k+1} \sum_{j \in s_{\text{TEMP}}} |\beta_j|, \quad (1.4)$$

where  $s_{\text{TEMP}} = \{j : |\tilde{\mathbf{y}}^\tau \mathbf{x}_j| = \max_{j' \notin s_{*k}} |\tilde{\mathbf{y}}^\tau \mathbf{x}_{j'}|\}$ ,  $\tilde{\mathbf{X}}_{\text{TEMP}}$  and  $\tilde{\boldsymbol{\beta}}_{\text{TEMP}}$  are, respectively, the corresponding projected residual matrix and the coefficient vector. If a partial positive cone condition (condition A2 in Section 3) is satisfied then  $s_{\text{TEMP}}$  is exactly the index set of the active  $\mathbf{x}_j$ 's. When  $s_{\text{TEMP}}$  is a singleton, the partial positive cone condition is automatically satisfied. For these results, see the proof of Theorem 1. The nonsingleton case rarely occurs. Therefore, the minimization of (1.4) is rarely called. If the need for the minimization of (1.4) does arise, the active  $\mathbf{x}_j$ 's can be easily obtained by applying the R function `glmPath` developed by Park and Hastie (2007) to  $\tilde{\mathbf{y}}$  and  $\tilde{\mathbf{X}}_{\text{TEMP}}$  and extracting the first feature (or features) with nonzero coefficient in the solution path. The results discussed above give rise to an efficient computation algorithm which is provided in Section 2.

We consider the properties of SLasso cum EBIC in the scenario that  $p = \exp(cn^\kappa)$ ,  $0 < \kappa < 1$ , and the number of relevant features  $p_0$  is also diverging to infinity at a proper rate. We establish the following properties. Let  $s_{*1}, s_{*2}, \dots, s_{*k}, \dots$  be the sequence generated by SLasso. Under reasonable conditions, there is a  $k = k^*$  such that  $s_{*k^*} = s_0$  with probability converging to 1 as  $n$  goes to infinity, where  $s_0$  is the exact index set of the relevant features (Theorems 1 and 2). Further, with probability converging to 1 uniformly for all  $k < k^*$ ,  $\text{EBIC}(s_{*k}) > \text{EBIC}(s_{*k+1})$  and  $\text{EBIC}(s) > \text{EBIC}(s_0)$  for all  $s$  such that  $p_0 < |s| \leq k_0 p_0$  with any fixed  $k_0 > 1$ ,  $|s|$  denoting the number of features in  $s$ , (Theorem 3). These results imply the selection consistency of the SLasso cum EBIC procedure. The asymptotic distribution of the SLasso estimators with diverging  $p_0$  is given in Theorem 4, which justifies the second part of the oracle property.

The remainder of the article is arranged as follows. The basic properties of SLasso cum EBIC and its computation algorithm are given in Section 2. The theoretical properties of SLasso cum EBIC are studied in Section 3. Simulation studies comparing SLasso cum EBIC with various other methods are reported in Section 4. A real data analysis is provided in Section 5. The article is concluded by a discussion of the similarities and dif-

ferences between SLasso cum EBIC and other related methods in Section 6. Some technical details are provided in a supplementary document.

## 2. BASIC PROPERTIES AND COMPUTATION ALGORITHM

We consider the scenario that both the total number of features and the number of relevant features diverge. We also allow the set of relevant features and their effects vary as  $n$  varies. For the sake of clarity, we do not index these quantities explicitly by  $n$ , but their dependence on  $n$  should be kept in mind. Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)$  be the design matrix. Let  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^\tau$ ,  $\mathbf{y} = (y_1, \dots, y_n)^\tau$ , and  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^\tau$ . In matrix notation, model (1.1) is expressed as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

Let  $S$  denote the set of indices  $\{1, 2, \dots, p\}$ . Let  $s$  be any subset of  $S$ . Denote by  $\mathbf{X}(s)$  the matrix consisting of the columns of  $\mathbf{X}$  with indices in  $s$ . Similarly, let  $\boldsymbol{\beta}(s)$  denote the vector consisting of the corresponding components of  $\boldsymbol{\beta}$ . Let  $\mathcal{R}(s)$  be the linear space spanned by the columns of  $\mathbf{X}(s)$  and  $\mathbf{H}(s)$  its corresponding projection matrix, i.e.,  $\mathbf{H}(s) = \mathbf{X}(s)[\mathbf{X}^\tau(s)\mathbf{X}(s)]^{-1}\mathbf{X}^\tau(s)$ .

*Proposition 1.* Let  $s_{*k}$  denote the index set of the features selected at the  $k$ th step of SLasso. For  $k \geq 1$  and any  $\tilde{j} \in s_{*k}^c$ , if  $\mathbf{X}(\{\tilde{j}\}) \in \mathcal{R}(s_{*k})$  then  $\tilde{j} \notin s_{*k+1}$ .

*Proof.* If  $\mathbf{X}(\{\tilde{j}\}) \in \mathcal{R}(s_{*k})$  then there exists an  $\mathbf{a}_k$  such that  $\mathbf{X}(\{\tilde{j}\}) = \mathbf{X}(s_{*k})\mathbf{a}_k$  and hence

$$\begin{aligned} l_{k+1} &= \|\mathbf{y} - \mathbf{X}(s_{*k})(\boldsymbol{\beta}(s_{*k}) + \beta_{\tilde{j}}\mathbf{a}_k) - \mathbf{X}(s_{*k}^c/\{\tilde{j}\})\boldsymbol{\beta}(s_{*k}^c/\{\tilde{j}\})\|_2^2 \\ &\quad + \lambda \left( |\beta_{\tilde{j}}| + \sum_{j \in s_{*k}^c/\{\tilde{j}\}} |\beta_j| \right) \\ &= \|\mathbf{y} - \mathbf{X}(s_{*k})\tilde{\boldsymbol{\beta}}(s_{*k}) - \mathbf{X}(s_{*k}^c/\{\tilde{j}\})\boldsymbol{\beta}(s_{*k}^c/\{\tilde{j}\})\|_2^2 \\ &\quad + \lambda \left( |\beta_{\tilde{j}}| + \sum_{j \in s_{*k}^c/\{\tilde{j}\}} |\beta_j| \right) \\ &\geq \|\mathbf{y} - \mathbf{X}(s_{*k})\tilde{\boldsymbol{\beta}}(s_{*k}) - \mathbf{X}(s_{*k}^c/\{\tilde{j}\})\boldsymbol{\beta}(s_{*k}^c/\{\tilde{j}\})\|_2^2 \\ &\quad + \lambda \sum_{j \in s_{*k}^c/\{\tilde{j}\}} |\beta_j|. \end{aligned}$$

Thus when  $l_{k+1}$  is minimized  $\beta_{\tilde{j}}$  must be 0, that is,  $\tilde{j} \notin s_{*k+1}$ .  $\square$

Proposition 1 implies that, for any  $k$ , the matrix  $\mathbf{X}(s_{*k})$  is of full column rank. It also suggests that, in the SLasso procedure, any feature that is highly correlated with the features selected already will have little chance to be selected subsequently. This nature of SLasso is favorable when it is used for feature selection in ultra-high dimensional feature space where high spurious correlations present, see Fan and Lv (2008).

*Proposition 2.* For  $k \geq 1$ , the minimization of  $l_{k+1}$  is equivalent to the minimization of

$$\tilde{l}_{k+1} = \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\tilde{\boldsymbol{\beta}}\|^2 + \lambda_{k+1} \sum_{j \in s_{*k}^c} |\beta_j|, \quad (2.5)$$

where  $\tilde{\mathbf{y}} = [\mathbf{I} - \mathbf{H}(s_{*k})]\mathbf{y}$ ,  $\tilde{\mathbf{X}} = [\mathbf{I} - \mathbf{H}(s_{*k})]\mathbf{X}(s_{*k}^c)$ ,  $\tilde{\boldsymbol{\beta}} = \boldsymbol{\beta}(s_{*k}^c)$ .



*Proof.* Differentiating  $l_{k+1}$  with respect to  $\beta(s_{*k})$ , we have

$$\frac{\partial l_{k+1}}{\partial \beta(s_{*k})} = -2X^\tau(s_{*k})\mathbf{y} + 2X^\tau(s_{*k})X(s_{*k})\beta(s_{*k}) + 2X^\tau(s_{*k})X(s_{*k}^c)\beta(s_{*k}^c).$$

Setting the above derivative to zero, we obtain

$$\hat{\beta}(s_{*k}) = [X^\tau(s_{*k})X(s_{*k})]^{-1}X^\tau(s_{*k})[\mathbf{y} - X(s_{*k}^c)\beta(s_{*k}^c)]. \quad (2.6)$$

Substituting (2.6) into  $\|\mathbf{y} - X\beta\|^2$  we have

$$\begin{aligned} l_{k+1} &= \|\mathbf{y} - X(s_{*k})\beta(s_{*k}) - X(s_{*k}^c)\beta(s_{*k}^c)\|^2 + \lambda_{k+1} \sum_{j \in s_{*k}^c} |\beta_j| \\ &= \|\mathbf{y} - X(s_{*k}^c)\beta(s_{*k}^c) - X(s_{*k})[X^\tau(s_{*k})X(s_{*k})]^{-1}X^\tau(s_{*k})[\mathbf{y} \\ &\quad - X(s_{*k}^c)\beta(s_{*k}^c)]\|^2 + \lambda_{k+1} \sum_{j \in s_{*k}^c} |\beta_j| \\ &= \|[I - H(s_{*k})][\mathbf{y} - X(s_{*k}^c)\beta(s_{*k}^c)]\|^2 + \lambda_{k+1} \sum_{j \in s_{*k}^c} |\beta_j| \\ &= \|\tilde{\mathbf{y}} - \tilde{X}\tilde{\beta}\|^2 + \lambda_{k+1} \sum_{j \in s_{*k}^c} |\beta_j|. \end{aligned}$$

□

As a byproduct of the above proof, the components of  $\hat{\beta}(s_{*k})$  are almost surely nonzero since  $\mathbf{y}$  is a vector of continuous random variables. This implies that  $s_{*1} \subset s_{*2} \subset \dots \subset s_{*k} \subset \dots$ ; that is, the feature sets selected in the sequential steps are nested.

**Proposition 3.** Let  $s_{\text{TEMP}} = \{j : j \in s_{*k}^c, |\tilde{\mathbf{y}}^\tau \mathbf{x}_j| = \max_{l \in s_{*k}^c} |\tilde{\mathbf{y}}^\tau \mathbf{x}_l|\}$ . If  $s_{\text{TEMP}}$  is a singleton, then the  $\mathbf{x}_j$  with  $j \in s_{\text{TEMP}}$  is the only feature with nonzero estimated coefficient in the minimization of (2.5); otherwise, the minimization of (2.5) is equivalent to the minimization of

$$\|\tilde{\mathbf{y}} - \tilde{X}_{\text{TEMP}}\tilde{\beta}_{\text{TEMP}}\|^2 + \lambda_{k+1} \sum_{j \in s_{\text{TEMP}}} |\beta_j|,$$

where  $\tilde{X}_{\text{TEMP}}$  consists of the  $\tilde{\mathbf{x}}_j$  with  $j \in s_{\text{TEMP}}$ ,  $\tilde{\beta}_{\text{TEMP}}$  is the corresponding coefficient vector.

This proposition follows from the proof of Theorem 1. Proposition 3 gives rise to the following computation algorithm.

*SLasso cum EBIC algorithm:*

- Initial Step: Standardize  $\mathbf{y}$ ,  $\mathbf{x}_j$ ,  $j = 1, \dots, p$ , such that  $\mathbf{y}^\tau \mathbf{1} = 0$ ,  $\mathbf{x}_j^\tau \mathbf{1} = 0$  and  $\mathbf{y}^\tau \mathbf{y} = n$ ,  $\mathbf{x}_j^\tau \mathbf{x}_j = n$ . Compute  $\mathbf{x}_j^\tau \mathbf{y}$  for  $j \in S$ . Let

$$s_{\text{TEMP}} = \{j : |\mathbf{x}_j^\tau \mathbf{y}| = \max_{j' \in S} |\mathbf{x}_{j'}^\tau \mathbf{y}|\}.$$

If  $s_{\text{TEMP}}$  is a singleton, let  $s_{*1} = s_{\text{TEMP}}$ , otherwise, apply `glmPath` to  $\mathbf{y}$  and  $X(s_{\text{TEMP}})$  and extract the first feature with nonzero coefficient in the solution path, and let  $s_{*1}$  be the active set. Compute  $I - H(s_{*1})$  and  $\text{EBIC}(s_{*1})$ .

- General Step: For  $k \geq 1$ , compute  $\tilde{\mathbf{x}}_j^\tau \tilde{\mathbf{y}}$  for  $j \in s_{*k}^c$ , where  $\tilde{\mathbf{y}} = [I - H(s_{*k})]\mathbf{y}$ ,  $\tilde{\mathbf{x}}_j = [I - H(s_{*k})]\mathbf{x}_j$ . Let

$$s_{\text{TEMP}} = \{j : |\tilde{\mathbf{x}}_j^\tau \tilde{\mathbf{y}}| = \max_{j' \in s_{*k}^c} |\tilde{\mathbf{x}}_{j'}^\tau \tilde{\mathbf{y}}|\}.$$

If  $s_{\text{TEMP}}$  is a singleton, let  $s_{*k+1} = s_{*k} \cup s_{\text{TEMP}}$ , otherwise, apply `glmPath` to  $\tilde{\mathbf{y}}$  and  $\tilde{X}(s_{\text{TEMP}})$  and extract the first feature with nonzero coefficient in the solution path, and let  $s_{*k+1}$  be  $s_{*k}$  union the active set. Compute  $\text{EBIC}(s_{*k+1})$ . If

$\text{EBIC}(s_{*k+1}) > \text{EBIC}(s_{*k})$ , stop; otherwise, compute  $I - H(s_{*k+1})$  and continue.

- When the process stops, the parameters in the selected model are estimated by their least-square estimates.

The EBIC for  $s_{*k}$ ,  $k = 1, 2, \dots$ , in the above algorithm is given by

$$\begin{aligned} \text{EBIC}(s_{*k}) &= n \ln \left( \frac{\|[I - H(s_{*k})]\mathbf{y}\|_2^2}{n} \right) \\ &\quad + |s_{*k}| \ln n + 2 \left( 1 - \frac{\ln n}{r \ln p} \right) \ln \left( \frac{p}{|s_{*k}|} \right), \end{aligned}$$

where  $r$  is a positive number slightly bigger than 2, say  $r = 2.1$ . For more details on EBIC, see Section 3.3. The matrix  $I - H(s_{*k+1})$  can be updated from  $I - H(s_{*k})$  recursively. Suppose there are  $K$  active features with indices  $\{j_m : m = 1, \dots, K\}$  at step  $k + 1$ . Denote by  $J_m = \{j_1, \dots, j_m\}$ . Let  $J_0 = \emptyset$ . The recursive formula is given by

$$\begin{aligned} I - H(s_{*k} \cup J_m) &= [I - H(s_{*k} \cup J_{m-1})] \\ &\quad \times \left\{ I - \frac{\mathbf{x}_{j_m} \mathbf{x}_{j_m}^\tau [I - H(s_{*k} \cup J_{m-1})]}{\mathbf{x}_{j_m}^\tau [I - H(s_{*k} \cup J_{m-1})] \mathbf{x}_{j_m}} \right\}, \end{aligned} \quad (2.7)$$

The amount of computation in the above algorithm is minimal. The computation of the projection matrices does not involve any matrix inversion. The call for `glmPath` is in fact seldom invoked. As mentioned earlier that, at steps where the partial positive cone condition holds, the set  $s_{\text{TEMP}}$  is indeed the index set of the active features at those steps. Thus, the SLasso selects features at those steps by maximizing  $|\tilde{\mathbf{x}}_j^\tau \tilde{\mathbf{y}}|$ , which is the same as OMP.

### 3. THE ORACLE PROPERTY OF SLASSO CUM EBIC

We assume in model (1.1) that the  $\epsilon_i$ 's are iid normal random variables with mean zero and variance  $\sigma^2$ . We consider the design matrix  $X$  either as a deterministic or a random matrix. Let  $s_0 = \{j : \beta_j \neq 0, j = 1, \dots, p\}$ . Assume  $\ln p = O(n^\kappa)$  for some  $\kappa > 0$  and  $p_0 = |s_0| = O(n^c)$  for some  $0 < c < 1$ . We establish the oracle property of the SLasso cum EBIC procedure in this section. We first present the result that, with probability converging to 1, the SLasso selects all the relevant features before any irrelevant feature can be selected. Then we present the result that, with probability converging to 1, the SLasso procedure using the EBIC as the stopping rule stops exactly at the step when all the relevant features are selected. The asymptotic distribution of the SLasso estimator is also provided. The proofs of these results are given in the supplementary document. Some special cases are discussed at the end of this section.

#### 3.1 The Case of Deterministic Design Matrix

In the case of deterministic design matrix, suppose the columns of  $X$  are standardized. We now introduce some notations. For  $s \subset S$ , let  $s^- = s^c \cap s_0$ . If  $s \subset s_0$  then  $s^-$  is the complement of  $s$  in  $s_0$ . For  $s \subset s_0$ , define

$$\gamma_n(j, s, \beta) = \frac{1}{n} \mathbf{x}_j^\tau [I - H(s)] X \beta.$$

In fact,  $\gamma_n(j, s, \beta)$  only depends on  $\beta(s^c)$ . But for the ease of notation,  $\beta$  and  $\beta(s^c)$  will be used interchangeably. Unless otherwise stated,  $\beta$  also denotes the unknown true value of the parameter vector. We make the following assumptions.

- A1.  $\max_{j \in s_0^c} |\gamma_n(j, s, \beta)| < q \max_{j \in s^-} |\gamma_n(j, s, \beta)|$ ,  $0 < q < 1$ .  
 A2. (Partial positive cone condition). If  $s^- \neq \emptyset$ , let

$$\mathcal{A}_s = \{\tilde{j} : \tilde{j} \in s^-, |\gamma_n(\tilde{j}, s, \beta)| = \max_{j \in s^c} |\gamma_n(j, s, \beta)|\},$$

and  $\tilde{X}(\mathcal{A}_s) = [I - H(s)]X(\mathcal{A}_s)$ . Then  $[\tilde{X}^T(\mathcal{A}_s)\tilde{X}(\mathcal{A}_s)]^{-1} \mathbf{1} > 0$ , where  $\mathbf{1}$  is the vector with all components 1.

- A3.  $\frac{\sqrt{n}}{\ln p} \lambda_{\min}[\frac{1}{n} X^T(s_0)X(s_0)] \min_{j \in s_0} |\beta_j| \rightarrow +\infty$ , as  $n \rightarrow \infty$ , where  $\lambda_{\min}$  denotes the smallest eigenvalue.

Assumption A1 is implied by the following condition

$$\|\tilde{x}_j^T \tilde{X}(s^-) [\tilde{X}^T(s^-) \tilde{X}(s^-)]^{-1} \mathbf{1}\| < 1 - \eta, \forall j \in s_0^c, \quad (3.8)$$

where  $\tilde{x}_j = [I - H(s)]x_j$  and  $0 < \eta < 1$ . The claim above follows because

$$\begin{aligned} |\gamma_n(j, s, \beta)| &= \frac{1}{n} |x_j^T [I - H(s)]\mu| \\ &= |\tilde{x}_j^T \tilde{X}(s^-) [\tilde{X}^T(s^-) \tilde{X}(s^-)]^{-1} \frac{1}{n} \tilde{X}^T(s^-) [I - H(s)]\mu| \\ &\leq \|\tilde{x}_j^T \tilde{X}(s^-) [\tilde{X}^T(s^-) \tilde{X}(s^-)]^{-1} \mathbf{1}\| \frac{1}{n} \|\tilde{X}^T(s^-) [I - H(s)]\mu\| \\ &< (1 - \eta) \frac{1}{n} \|\tilde{X}^T(s^-) [I - H(s)]\mu\| \\ &= (1 - \eta) \frac{1}{n} \max_{j \in s^-} |x_j^T [I - H(s)]\mu| \\ &= (1 - \eta) \max_{j \in s^-} |\gamma_n(j, s, \beta)|, \end{aligned}$$

where the strict inequality holds by (3.8).

Under Assumption A1, the  $\mathcal{A}_s$  in A2 is a subset of  $s_0$ . Assumption A2 holds if and only if

$$\tilde{x}_j^T \tilde{X}(\mathcal{A}_s \setminus \{j\}) [\tilde{X}^T(\mathcal{A}_s \setminus \{j\}) \tilde{X}(\mathcal{A}_s \setminus \{j\})]^{-1} \mathbf{1} < 1, \forall j \in \mathcal{A}_s. \quad (3.9)$$

We establish the equivalence of A2 and (3.9) below. Let  $A = \tilde{X}(\mathcal{A}_s \setminus \{j\})$  and  $b = \tilde{x}_j$ . Since a permutation of the rows and columns does not change the sum of the rows, it suffices to verify that the sum of the last row of  $\begin{pmatrix} A^T A & A^T b \\ b^T A & b^T b \end{pmatrix}^{-1}$  is positive if and only if  $b^T A(A^T A)^{-1} \mathbf{1} < 1$ . Let  $E = I - A(A^T A)^{-1} A^T$  and  $F = I - b(b^T b)^{-1} b^T$ . By the formula for the inverse of blocked matrices, we have

$$\begin{pmatrix} A^T A & A^T b \\ b^T A & b^T b \end{pmatrix}^{-1} = \begin{pmatrix} (A^T F A)^{-1} & -(A^T A)^{-1} A^T b (b^T E b)^{-1} \\ -(b^T b)^{-1} b^T A (A^T F A)^{-1} & (b^T E b)^{-1} \end{pmatrix}.$$

and

$$\begin{aligned} (A^T F A)^{-1} &= [A^T A - A^T b (b^T b)^{-1} b^T A]^{-1} \\ &= (A^T A)^{-1} + (A^T A)^{-1} A^T (b^T E b)^{-1} b^T A (A^T A)^{-1}. \end{aligned}$$

Substituting the expression of  $(A^T F A)^{-1}$  into the first block of the last row of the above matrix, we obtain

$$-(b^T b)^{-1} b^T A (A^T F A)^{-1} = -(b^T E b)^{-1} b^T A (A^T A)^{-1}.$$

Thus the sum of the last row becomes

$$\begin{aligned} (b^T E b)^{-1} - (b^T E b)^{-1} b^T A (A^T A)^{-1} \mathbf{1} \\ = (b^T E b)^{-1} [1 - b^T A (A^T A)^{-1} \mathbf{1}] \end{aligned}$$

which is greater than 0 if and only if  $b^T A (A^T A)^{-1} \mathbf{1} < 1$ .

Condition (3.8) is a conditional version of the exact recovery condition (ERC) assumed in Tropp (2004) while conditioning on the subset  $s$  of the relevant features. Condition (3.9) is similar to but much weaker than the *irrepresentability condition*. The above arguments suggest that Conditions A1 and A2 might be weaker than the ERC and the *irrepresentability condition*. This is indeed the case. We will demonstrate this by special cases where the conditions for the selection consistency of the SLasso hold but the ERC and the *irrepresentability condition* are not satisfied. If  $\lambda_{\min}(\frac{1}{n} X^T(s_0)X(s_0))$  is bounded away from zero, which is a common assumption in the case of ultra-high dimensional feature space, then Condition A3 is equivalent to  $\frac{\sqrt{n}}{\ln p} \min_{j \in s_0} |\beta_j| \rightarrow \infty$ . If  $\ln p = O(n^\kappa)$  with  $\kappa < 1/2$  and  $\min_{j \in s_0} |\beta_j| \geq C n^{-\delta}$  for some constant  $C$  and  $\delta < 1/2 - \kappa$ , A3 is then satisfied.

**Theorem 1.** Let  $s_{*1}, s_{*2}, \dots, s_{*k}, \dots$  be the sequence generated by the SLasso procedure. Suppose that assumptions A1–A3 hold. Let  $\ln p = O(n^\kappa)$ , where  $\kappa < 1/2$ . Then, there is a  $k^*$  such that

$$\Pr(s_{*k^*} = s_0) \rightarrow 1, \text{ as } n \rightarrow \infty,$$

where  $s_0$  is the exact index set of the relevant features.

### 3.2 The Case of Random Design Matrix

Assume  $x_i = (x_{i1}, \dots, x_{ip})^T$ ,  $i = 1, \dots, n$ , are iid copies of a random vector  $z = (z_1, \dots, z_p)^T$ . Without loss of generality, assume that  $Ez = 0$  and  $\text{var}(z) = \Sigma$  with diagonal elements 1 and off-diagonal elements independent of  $n$ . Assume that

- a1. The off-diagonal elements of  $\Sigma$  are bounded by a constant less than 1; that is, the correlation between any two features are bounded below from  $-1$  and above from  $1$ .
- a2.  $\sigma_{\max} \equiv \max_{1 \leq j, k \leq p} \sigma(z_j z_k) < \infty$  where  $\sigma(z_j z_k)$  denotes the standard deviation of  $z_j z_k$ .
- a3.  $\max_{1 \leq j, k \leq p} E \exp(t z_j z_k)$  and  $\max_{1 \leq j \leq p} E \exp(t z_j)$  are finite for  $t$  in a neighborhood of zero.

For any  $s, \tilde{s} \subset S$ , denote by  $\Sigma_{s\tilde{s}}$  the sub matrix of  $\Sigma$  with row indices in  $s$  and column indices in  $\tilde{s}$ . Define

$$\Gamma(j, s, \beta) = (\Sigma_{jS} - \Sigma_{j\tilde{s}} \Sigma_{\tilde{s}\tilde{s}}^{-1} \Sigma_{\tilde{s}S})\beta.$$

The following assumptions are imposed:

- A1'. For any  $s \subset s_0$ ,  $s \neq s_0$ ,  $\max_{j \in s_0^c} |\Gamma(j, s, \beta)| < \max_{j \in s^-} |\Gamma(j, s, \beta)|$ .  
 A2'. Let  $\mathcal{A}_s = \{j : j \in s^c, |\Gamma(j, s, \beta)| = \max_{l \in s^c} |\Gamma(l, s, \beta)|\}$ . Then

$$(\Sigma_{\mathcal{A}_s \mathcal{A}_s} - \Sigma_{\mathcal{A}_s s} \Sigma_{ss}^{-1} \Sigma_{s \mathcal{A}_s})^{-1} \mathbf{1} > 0.$$

- A3'.  $\frac{n^{1/2}}{\ln p} \lambda_{\min}(\Sigma_{s_0 s_0}) (\min_{j \in s_0} |\beta_j|) \rightarrow +\infty$  as  $n \rightarrow +\infty$ .

The Assumptions A1'–A3' are in fact the Assumptions A1–A3 with the empirical variances and covariances of the features replaced by their theoretical counterparts. To establish the selection consistency of SLasso in the case of random feature matrix, we need to pass from assumptions A1'–A3' to assumptions A1–A3. The following lemma ensures that if A1'–A3' hold then A1–A3 hold with probability converging to 1 as  $n$  goes to infinity.

**Lemma 1.** Under assumptions a1–a3,

- (i)  $P(\max_{1 \leq j, k \leq p} \left| \frac{1}{n} \sum_{i=1}^n x_{ij}x_{ik} - \Sigma_{jk} \right| > n^{-\frac{1}{3}} \sigma_{\max}) \rightarrow 0$ .
- (ii)  $P(\max_{1 \leq j \leq p} \left| \frac{1}{n} \sum_{i=1}^n x_{ij}\epsilon_i \right| > n^{-\frac{1}{3}} \sigma) \rightarrow 0$ .
- (iii) Let  $\Sigma_{jl|s} = \Sigma_{jl} - \Sigma_{js}\Sigma_{ss}^{-1}\Sigma_{sl}$  and  $\hat{\Sigma}_{jl|s} = \mathbf{x}_j^\tau [\mathbf{I} - \mathbf{H}(s)]\mathbf{x}_l/n$ . Then

$$\max_{1 \leq j, l \leq p} \max_{s: |s| \leq p_0} |\hat{\Sigma}_{jl|s} - \Sigma_{jl|s}| = o_p(1).$$

The proof of the lemma is given in the supplementary document.

**Theorem 2.** Let  $\ln p = O(n^\kappa)$ ,  $\kappa < 1/3$ , and  $p_0 = O(n^c)$ ,  $\kappa/2 < c < 1/6$ . Assume that conditions a1–a3 and A1'–A3' are satisfied. Then, there is a  $k^*$  such that

$$Pr(s_{*k^*} = s_0) \rightarrow 1, \text{ as } n \rightarrow \infty.$$

The theorem is in fact a corollary of Lemma 1. It follows from the lemma immediately that if a1–a3 and A1'–A3' are satisfied then A1–A3 hold with probability converging to 1. Thus the selection consistency of SLasso with random feature matrix is established.

### 3.3 Property of the Stopping Rule and the Asymptotic Distribution of the SLasso Estimators

For a linear model with features in  $s$ , the EBIC proposed in Chen and Chen (2008) is defined as

$$\text{EBIC}_\gamma(s) = n \ln \left( \frac{\|[\mathbf{I} - \mathbf{H}(s)]\mathbf{y}\|_2^2}{n} \right) + |s| \ln n + 2\gamma \ln \left( \frac{p}{|s|} \right), \gamma \geq 0.$$

The properties of EBIC for sparse high-dimensional linear models are investigated in Chen and Chen (2008) and Luo and Chen (2011). It is shown that, if  $\gamma > 1 - \ln n/(2 \ln p)$ , EBIC is selection consistent in the sense that

$$P(\min_{|s| \leq k_0 p_0} \text{EBIC}_\gamma(s) > \text{EBIC}_\gamma(s_0)) \rightarrow 1, \text{ as } n \rightarrow \infty,$$

where  $k_0 > 1$  is any fixed number.

For the sequence  $s_{*1}, s_{*2}, \dots, s_{*k}, \dots$  selected by the procedure of SLasso, we have shown that  $s_{*1} \subset s_{*2} \subset \dots \subset s_{*k} \subset \dots$  and that, with probability converging to 1, there is a  $k^*$  such that  $s_{*k^*} = s_0$ . In this section, we provide the result that, with probability converging to 1,  $\text{EBIC}(s_{*k})$  decreases when  $k < k^*$  and reaches its minimum at step  $k^*$ , and that  $\text{EBIC}(s_{*k}) > \text{EBIC}(s_{*k^*})$  for any  $k > k^*$ . The result is given in Theorem 3. This result implies that, with probability converging

to 1, the procedure of SLasso cum EBIC stops at step  $k^*$ . The parameters in the selected model are estimated by their least-square estimates. Theorems 1–3 imply that the estimators are obtained as if  $s_0$  were known in advance. This implies that the SLasso cum EBIC procedure possesses the oracle property. Since  $p_0$  diverges, the asymptotic theory on ordinary least squares estimators with fixed  $p_0$  does not apply. We derive the asymptotic distribution of the SLasso estimator of  $\beta(s_0)$  in Theorem 4.

**Theorem 3.** Assume conditions A1 and A2. Suppose that  $\ln p_n = O(n^\kappa)$ ,  $\kappa < 1/3$ ,  $p_0 = O(n^c)$ ,  $c < 1/6$ , and there is a constant  $C$  such that  $\lambda_{\min}(\frac{1}{n} \mathbf{X}(s_0)^\tau \mathbf{X}(s_0)) \min_{j \in s_0} |\beta_j| \geq Cn^{-1/6+\delta}$ , where  $\delta$  is an arbitrarily small positive number. Let  $s_{*1} \subset s_{*2} \subset \dots \subset s_{*k} \subset \dots$  be the sets generated by the procedure of SLasso. Let  $k^*$  be as given in Theorems 1 and 2. Then

- (i) Uniformly, for  $k < k^*$ ,

$$P(\text{EBIC}_\gamma(s_{*k+1}) < \text{EBIC}_\gamma(s_{*k})) \rightarrow 1, \text{ when } \gamma > 0.$$

- (ii)  $P(\min_{p_0 < |s_{*k}| \leq k_0 p_0} \text{EBIC}_\gamma(s_{*k}) > \text{EBIC}_\gamma(s_0)) \rightarrow 1$ , when  $\gamma > 1 - \frac{\ln n}{2 \ln p}$ , where  $k_0 > 1$  is an arbitrarily fixed constant.

The proof of the theorem is given in the supplementary document.

In the stopping rule,  $\gamma$  is taken as  $1 - \frac{\ln n}{r \ln p}$  where  $r$  is slightly bigger than 2. This choice of  $\gamma$  is to keep the EBIC selection consistent at one hand and to achieve the largest power for the identification of relevant feature at another hand. A brief justification is given as follows. For a sample of size  $n$ , define the positive discovery rate ( $\text{PDR}_n$ ) and the false discovery rate ( $\text{FDR}_n$ ) as follows:

$$\text{PDR}_n = \frac{|s_{*k^*} \cap s_0|}{|s_0|}, \quad \text{FDR}_n = \frac{|s_{*k^*} \cap s_0^c|}{|s_{*k^*}|}. \quad (3.10)$$

The asymptotic property  $P(s_{*k^*} = s_0) \rightarrow 1$  is equivalent to that  $\text{FDR}_n \rightarrow 0$  and  $\text{PDR}_n \rightarrow 1$  simultaneously. For any  $\gamma > 1 - \frac{\ln n}{2 \ln p}$ , the above convergences are guaranteed, but the convergence rates are different for different  $\gamma$  values. For a bigger  $\gamma$ , both  $\text{FDR}_n$  and  $\text{PDR}_n$  are smaller. By choosing the  $\gamma$  as small as possible in its consistent range, the  $\text{PDR}_n$  is maximized while the  $\text{FDR}_n$  still converges to zero.

Let  $s^*$  be the set selected by SLasso cum EBIC and  $\hat{\beta}(s^*)$  the SLasso estimator of  $\beta(s^*)$  (which is indeed the least squares estimator). Let  $\mathbf{a} = (a_1, a_2, \dots)$  be an infinite sequence of constants. For any index set  $s$ , let  $\mathbf{a}(s)$  denote the vector with components  $a_j$ ,  $j \in s$ . We have the following theorem.

**Theorem 4.** Let  $\mathbf{z}_i^\tau$  be the  $i$ th row vector of  $\mathbf{X}(s_0)$ ,  $i = 1, \dots, n$ . Assume that

$$\lim_{n \rightarrow \infty} \max_{1 \leq i \leq n} \mathbf{z}_i^\tau [\mathbf{X}(s_0)^\tau \mathbf{X}(s_0)]^{-1} \mathbf{z}_i \rightarrow 0. \quad (3.11)$$

Then, for any fixed sequence  $\mathbf{a}$ ,

$$\frac{\mathbf{a}(s^*)^\tau [\hat{\beta}(s^*) - \beta(s^*)]}{\sqrt{\mathbf{a}(s^*)^\tau [\mathbf{X}(s^*)^\tau \mathbf{X}(s^*)]^{-1} \mathbf{a}(s^*)}} \rightarrow_d N(0, \sigma^2),$$

where  $\sigma^2 = \text{var}(Y_i)$ .

Let  $\hat{\beta}(s_0) = [X(s_0)^\tau X(s_0)]^{-1} X(s_0)^\tau y$ . Then under (3.11),

$$V_n = \frac{a(s_0)^\tau [\hat{\beta}(s_0) - \beta(s_0)]}{\sqrt{a(s_0)^\tau [X(s_0)^\tau X(s_0)]^{-1} a(s_0)}} \rightarrow_d N(0, \sigma^2),$$

which follows from the Linderberg's central limit theorem; see Corollary 1.3 in Shao (2003). Its proof is by checking the validity of the conditions for this corollary, which is straightforward and is omitted here. Let  $U_n = \frac{a(s^*)^\tau [\hat{\beta}(s^*) - \beta(s^*)]}{\sqrt{a(s^*)^\tau [X(s^*)^\tau X(s^*)]^{-1} a(s^*)}}$ . Since  $P(U_n \neq V_n) = P(s^* \neq s_0) \rightarrow 0$ , we have that  $U_n - V_n \rightarrow 0$  in probability. Thus, by Slutsky's theorem,  $U_n = V_n + (U_n - V_n) \rightarrow_d N(0, \sigma^2)$ . Theorem 4 implies that any fixed dimensional subvector of  $\hat{\beta}(s^*)$  has an asymptotic multivariate normal distribution.

### 3.4 Special Cases

We give two special cases which demonstrate that the conditions required for the oracle property of SLasso are weaker than the well-known *irrepresentability condition*.

*Special case I:* Let the correlation matrix of  $z$  be given by

$$\Sigma = (1 - \rho)I + \rho \mathbf{1}\mathbf{1}^\tau,$$

where  $I$  is the identity matrix of dimension  $p$ ,  $\mathbf{1}$  is a  $p$ -vector of all elements 1, and  $0 < \rho \leq \rho_0 < 1$ . In this case, for the *irrepresentability condition* to be satisfied, some restriction must be imposed. But such restriction is not needed for sequential Lasso.

*Special case II.* Without loss of generality, let  $s_0 = \{1, \dots, p_0\}$ . Assume that

- (i)  $|\beta_1| > |\beta_2| > \dots > |\beta_{p_0}| = Cn^{-1/2+\delta}$  for some constant  $C$  and an arbitrarily small positive  $\delta$ ;
- (ii) The correlation matrix  $\Sigma$  has the following structure:

$$\Sigma_{s_0 s_0} = I, \quad \Sigma_{j s_0} = \frac{1}{p_0} \text{sign} \beta(s_0)^\tau, \text{ for } j \in s_0^c.$$

In this case, the *irrepresentability condition* is violated but the conditions for the oracle property of SLasso are satisfied.

The verification for the claims on the two special cases are provided in the supplementary document.

## 4. SIMULATION STUDY

In our simulation study, we compare SLasso with adaptive Lasso (ALasso; Huang, Ma, and Zhang 2008), SCAD (Kim, Choi, and Oh 2008; Xie and Huang 2009), SIS+SCAD (Fan and Lv 2008), forward stepwise regression (FSR; Wang 2009), and a forward-backward greedy algorithm (FoBa; Zhang 2011). The first three competing methods have also been shown to have the oracle property (or selection consistency) under certain conditions. The FSR and FoBa are included in the comparison because of their close relationship with SLasso. We do not include CoSaMP in the simulation study since no concrete approach for the selection of models was given in Needell and Tropp (2009). For ALasso, SCAD, SIS+SCAD and FoBa, cross-validation is used to determine the final selected model as in the cited references. For SLasso and FSR, EBIC is used as the stopping rule. The R packages *parcor* (Kraemer and Schaefer 2011), *ncvreg* (Breheny 2011), *SIS* (Fan et al. 2010), and *foba* are used for

the computation of ALasso, SCAD, SIS+SCAD, and FoBa, respectively.

We take two groups of settings: group A and group B. In group A, we consider the diverging pattern  $(n, p, p_0) = (n, [4n^{0.16}], [5e^{n^{0.3}}])$  for  $n = 100, 200$ , which is in consistence with the assumption in the theory of SLasso. The coefficients are generated as independent random variables distributed as  $(-1)^u(4n^{-0.15} + |z|)$ , where  $u \sim P$  Bernoulli (0.4) and  $z$  is a normal random variable with mean 0 and satisfies  $P(|z| \geq 0.1) = 0.25$ . The absolute values of the coefficients are roughly of order  $O(n^{-0.15})$ . The variance of the error term in the linear model is determined by

$$h = \frac{\beta^\tau \Sigma \beta}{\beta^\tau \Sigma \beta + \sigma^2} = 0.8,$$

where  $\Sigma$  is the variance-covariance matrix of relevant features. Five settings of the covariance structure for the design matrix  $X$  are considered. They are named GA1, GA2, ..., GA5. In group B, three settings named GB1, GB2, and GB3, which are adapted from the cited references, are considered. In these settings, the triplet  $(n, p, p_0)$  does not follow the diverging pattern. The details of the above simulation settings are provided in the supplementary document.

The methods are compared in terms of PDR, FDR, model size (Msize), and prediction mean square error (PMSE). The definition of PDR and FDR are given in (3.10). The MSize is the number of features selected. The PMSE is the average squared differences between the observations in a sample and their predicted values obtained using the model built from another independent sample. Thus, for the computation of PMSE, we generate two independent samples with the same sample size  $n$  in each replicate of the settings. One sample is used for the selection of features and the estimation of the coefficients, and the other one is used to compute the PMSE. For the settings in group A, these quantities are averaged over 200 replicates, for those in group B, they are averaged over 500 replicates. The results for group A are reported in Tables 1 and 2. The results for group B are reported in Table 3.

The findings of the simulation study are summarized as follows. First, consider the performance in prediction. For settings of Group A, ALasso, SCAD, FSR, and SLasso have comparable PMSE, however, FSR and SLasso have smaller (in GA1 and GA2, much smaller) MSize. The other two, SIS+SCAD and FoBa, have much larger PMSE. SIS+SCAD has smaller MSize than all the others and FoBa has much larger MSize than all the others. FSR and SLasso always have about the same MSize, but SLasso has smaller PMSE than FSR except in setting GA2. For settings of Group B, in GB1 and GB2, the three methods, SCAD, FSR, and SLasso, have comparable PMSE which are much smaller than the other three methods. SCAD has larger MSize than FSR and SLasso whose MSize are about the same. In GB3, ALasso, FoBa, and FSR have much smaller PMSE but also have much larger MSize than SCAD, SIS+SCAD, and SLasso.

Now, consider the performance in the identification of relevant features. First, let our attention be drawn to the performance of FoBa. In all the settings, FoBa has extremely high FDR which are much higher than all the other methods. In settings of Group A with sample size 100 and 200, its minimum FDR is,



Table 1. Comparison of SLasso with ALasso, SCAD, SIS+SCAD, and FSR in terms of PDR, FDR, PMSE, and model size (MSize) averaged over 200 simulation replicates in Group A settings (sample size  $n = 100$ ; numbers in parentheses are standard deviations)

Setting	Methods	MSize	PDR	FDR	PMSE
GA1	ALasso	37.8 (12.5)	0.999 (0.009)	0.767 (0.072)	13.714 (2.837)
	SCAD	14.9 (3.8)	0.998 (0.022)	0.430 (0.142)	10.523 (2.943)
	SIS+SCAD	5.0 (0.1)	0.488 (0.098)	0.219 (0.157)	25.135 (5.091)
	FoBa	87.3 (1.3)	0.998 (0.020)	0.909 (0.002)	32.583 (6.759)
	FSR	8.3 (1.4)	0.963 (0.149)	0.063 (0.098)	11.250 (5.782)
GA2	SLasso	8.4 (1.3)	0.968 (0.134)	0.071 (0.098)	11.129 (5.233)
	ALasso	26.8 (15.8)	0.861 (0.112)	0.674 (0.142)	17.07 (3.334)
	SCAD	13.7 (6.1)	0.724 (0.189)	0.513 (0.179)	18.178 (4.162)
	SIS+SCAD	4.9 (0.3)	0.484 (0.066)	0.210 (0.086)	16.68 (3.142)
	FoBa	87.2 (1.6)	0.802 (0.169)	0.926 (0.015)	48.407 (13.017)
GA3	FSR	5.0 (1.7)	0.579 (0.201)	0.057 (0.104)	17.554 (3.995)
	SLasso	4.6 (1.6)	0.511 (0.176)	0.091 (0.133)	19.236 (4.323)
	ALasso	15.1 (8.1)	0.999 (0.009)	0.356 (0.242)	4.676 (0.82)
	SCAD	8.1 (0.7)	0.979 (0.099)	0.029 (0.099)	4.45 (0.933)
	SIS+SCAD	5.0 (0.3)	0.543 (0.068)	0.126 (0.100)	17.033 (3.178)
GA4	FoBa	87.2 (1.4)	0.862 (0.042)	0.921 (0.004)	15.121 (3.794)
	FSR	8.2 (0.8)	0.745 (0.190)	0.270 (0.169)	5.658 (1.290)
	SLasso	8.2 (1.0)	0.932 (0.153)	0.092 (0.143)	4.892 (1.266)
	ALasso	11.2 (5.2)	1.00 (0.000)	0.188 (0.225)	2.312 (0.419)
	SCAD	8.1 (0.1)	0.999 (0.009)	0.000 (0.000)	2.582 (1.049)
GA5	SIS+SCAD	4.9 (0.4)	0.527 (0.076)	0.132 (0.115)	10.02 (2.652)
	FoBa	86.9 (1.4)	0.873 (0.032)	0.920 (0.003)	7.679 (2.419)
	FSR	6.7 (3.4)	0.779 (0.398)	0.049 (0.086)	5.766 (6.555)
	SLasso	6.5 (3.5)	0.784 (0.409)	0.028 (0.059)	5.745 (6.525)
	ALasso	11.2 (5.2)	1.00 (0.000)	0.190 (0.226)	5.196 (0.842)
	SCAD	8.1 (0.1)	0.999 (0.012)	0.000 (0.000)	5.147 (1.031)
	SIS+SCAD	5.1 (0.1)	0.504 (0.034)	0.190 (0.046)	10.324 (1.601)
	FoBa	87.1 (1.5)	0.871 (0.021)	0.920 (0.002)	17.665 (4.927)
	FSR	7.3 (2.0)	0.782 (0.167)	0.124 (0.104)	7.500 (2.734)
	SLasso	7.5 (1.7)	0.911 (0.188)	0.027 (0.057)	6.440 (2.768)

respectively, 0.909 and 0.949. In the settings of Group B, its minimum FDR is 0.852. Whatever high PDR it might achieve cannot be justified with such high FDRs. Let alone the fact that its PDR is even lower than at least one of the other methods in all the settings except in GA2 with sample size 200 and GB3. FoBa fails the goal of identifying relevant features. The poor performance of FoBa is not surprising. Though it has essentially the same mechanism for the selection of features as SLasso, it has an improper stopping rule. FoBa stops at a forward step when the decrease in the residual,  $\frac{1}{n} \|\tilde{y}\|_2^2$ , falls below a threshold of order  $O(\sigma^2 \ln p/n)$  without penalizing the increase in the number of features. The backward step is hardly activated since it is in force only when the increase in the residual by deleting one of the selected features is less than half of the amount of decrease at the forward step. See Sections 1 and 6. In a SHR model,  $\frac{1}{n} \|\tilde{y}\|_2^2$  can be reduced to zero, it is easy for a feature, relevant or not, to reduce  $\frac{1}{n} \|\tilde{y}\|_2^2$  by an amount larger than the threshold. Thus, features with high spurious correlation with the response can be easily selected. Note that the order of threshold,  $O(\sigma^2 \ln p/n)$ , decreases as  $n$  increases, which suggests that, the larger the sample size, the more features can be selected and hence the higher the FDR. This is in fact demonstrated in the simulation study, when the sample size goes from 100 to 200 in the settings of Group A, the minimum FDR of FoBa goes from 0.909 to 0.949. It is interestingly contrasted with

SLasso cum EBIC whose FDR decreases as  $n$  increases, which is also demonstrated in the simulation study with the settings of Group A.

In what follows, we compare the performance of other methods. Under settings GA1 and GA2, which are common settings in many simulation studies, SLasso and FSR are better than the other methods. They have high PDR and very low FDR. SLasso and FSR are comparable while FSR is slightly better. ALasso and SCAD have higher PDR than SLasso and FSR, but their FDR are too much higher. Averaged over the four simulations, FSR, SLasso, ALasso, and SCAD have averaged PDR 0.847, 0.823, 0.950, and 0.913, respectively, and averaged FDR 0.042, 0.067, 0.741, and 0.459, respectively. The difference in PDR between FSR, SLasso and ALasso, SCAD, Foba is not too much, but the difference in FDR is strikingly large. FSR and SLasso are absolutely much better than SIS+SCAD in terms of both PDR and FDR.

Under settings GA3–GA5, SCAD is absolutely better than all the other methods. The performance of SLasso is close to SCAD. SLasso is absolutely better than SIS+SCAD and FSR. Though ALasso has a slightly higher PDR in a few cases, its FDR is too high to be acceptable in terms of the identification of relevant features. It is not surprising that SLasso is absolutely better than FSR. In settings GA3–GA5, all the irrelevant features are equally and highly correlated with the relevant features.

Table 2. Comparison of SLasso with ALasso, SCAD, SIS+SCAD, and FSR in terms of PDR, FDR, PMSE, and model size (MSize) averaged over 200 simulation replicates in Group A settings (sample size  $n = 200$ , numbers in parentheses are standard deviations)

Setting	Methods	MSize	PDR	FDR	PMSE
GA1	ALasso	49.0 (18.0)	1.00 (0.000)	0.791 (0.077)	10.937 (1.463)
	SCAD	13.6 (4.8)	1.00 (0.000)	0.283 (0.183)	8.638 (0.897)
	SIS+SCAD	8.7 (0.5)	0.793 (0.077)	0.181 (0.086)	12.355 (3.309)
	FoBa	176.2 (1.7)	1.00 (0.000)	0.949 (0.000)	28.439 (4.261)
	FSR	9.4 (0.7)	1.00 (0.000)	0.035 (0.060)	8.688 (1.024)
	SLasso	9.4 (0.7)	1.00 (0.000)	0.035 (0.061)	8.683 (1.025)
GA2	ALasso	40.6 (19.8)	0.941 (0.072)	0.735 (0.140)	15.297 (2.03)
	SCAD	23.7 (7.3)	0.931 (0.110)	0.612 (0.127)	14.159 (2.928)
	SIS+SCAD	8.1 (0.8)	0.661 (0.028)	0.255 (0.076)	14.715 (1.715)
	FoBa	176.2 (1.6)	0.943 (0.102)	0.952 (0.005)	41.228 (8.005)
	FSR	7.9 (1.7)	0.846 (0.179)	0.035 (0.070)	13.541 (2.915)
	SLasso	7.8 (2.1)	0.796 (0.190)	0.073 (0.100)	14.462 (3.207)
GA3	ALasso	25.5 (15.9)	0.956 (0.071)	0.507 (0.283)	4.205 (0.539)
	SCAD	9.1 (1.1)	0.972 (0.121)	0.031 (0.124)	3.963 (0.62)
	SIS+SCAD	8.9 (0.4)	0.864 (0.064)	0.128 (0.046)	4.498 (1.987)
	FoBa	176.3 (1.5)	0.882 (0.034)	0.955 (0.002)	12.544 (2.252)
	FSR	9.2 (0.9)	0.708 (0.206)	0.311 (0.183)	4.688 (0.672)
	SLasso	9.2 (1.0)	0.873 (0.209)	0.148 (0.190)	4.272 (0.714)
GA4	ALasso	13.3 (6.4)	1.00 (0.000)	0.215 (0.242)	2.186 (0.267)
	SCAD	9.0 (0.0)	1.00 (0.000)	0.000 (0.000)	2.320 (0.753)
	SIS+SCAD	8.7 (0.7)	0.449 (0.064)	0.535 (0.061)	3.327 (1.679)
	FoBa	174.6 (1.5)	0.888 (0.011)	0.954 (0.001)	6.875 (1.229)
	FSR	9.3 (0.6)	0.993 (0.043)	0.037 (0.074)	2.207 (0.297)
	SLasso	9.2 (0.6)	1.00 (0.000)	0.023 (0.052)	2.183 (0.268)
GA5	ALasso	15.7 (9.5)	0.986 (0.044)	0.276 (0.284)	5.303 (0.622)
	SCAD	9.0 (0.1)	0.999 (0.011)	0.000 (0.000)	5.199 (0.71)
	SIS+SCAD	7.8 (0.8)	0.681 (0.066)	0.206 (0.070)	7.975 (1.258)
	FoBa	175.1 (1.6)	0.886 (0.017)	0.954 (0.001)	16.644 (2.932)
	FSR	9.4 (0.6)	0.943 (0.086)	0.091 (0.100)	5.545 (0.799)
	SLasso	9.3 (0.6)	1.00 (0.000)	0.024 (0.054)	5.241 (0.608)

In these situations, FSR is more prone to error compared with SLasso, since FSR tends to select features which are highly correlated with the features already selected though they might have less correlation with the current residuals; see the discussion in Section 6.

Under settings GB1 and GB2, the pattern is similar to that under settings GA1 and GA2. Under setting GB3, though the condition (which is sufficient but not necessarily necessary) for the selection consistency of SLasso is not satisfied, SLasso performs better than all the other methods; it has comparable or higher PDR than other methods and has the lowest (much lower than the others) FDR.

The simulation study demonstrates that (i) the SLasso cum EBIC method is one of the best feature selection methods for the purpose of prediction, and (ii) in terms of the identification of relevant features the performance of SLasso is satisfactory and robust; it always has a very low FDR and it is always close to the best, though it is not the best over all the simulation settings. On the contrast, the performance of SCAD and FSR are erratic over the settings. They are the best in certain settings but perform much worse in other settings.

## 5. REAL DATA ANALYSIS

The data, which were reported in Scheetz et al. (2006), consist of the expression levels of over 31,042 different probes from 120

$F_2$  male rats generated from an intercross experiment. A cross of SR/JrHsd male rats and SHRSP female rats was performed to generate  $F_1$  and the  $F_1$  rats were intercrossed to generate the  $F_2$  rats. The probes that were not expressed in the eye or that lacked sufficient variation were excluded. A probe was considered expressed if its maximum expression value observed among the 120  $F_2$  rats was greater than the 25th percentile of the entire set of RMA (robust multichip averaging) expression values. A probe was considered “sufficiently variable” if it exhibited at least two-fold variation in expression level among the 120  $F_2$  rats. A total of 18,976 probes that met these criteria were retained. Among the 18,976 probes, there is one, 1389163\_at, from gene TRIM32. This gene was found to cause Bardet-Biedl syndrom Chiang et al. (2006). Of interest is to find the probes among the remaining 18,975 probes that are most related to TRIM32. This has been studied by using different methods in the literature; see Huang, Ma, and Zhang (2008), Kim, Choi, and Oh (2008), Hwang, Zhang, and Ghosal (2009), Fan, Feng, and Song (2011), and Sun and Zhang (2011). In this section, we apply the five methods considered in our simulation study, that is, ALasso, SCAD, SIS+SCAD, FSR and SLasso, to the above problem. The response variable is the expression level of probe 1389163\_at. The features are the expression levels of the remaining 18,975 probes. The expression levels are standardized to have mean 0 and standard deviation 1 in the analysis.

Table 3. Comparison of SLasso with ALasso, SCAD, SIS+SCAD, and FSR in terms of PDR, FDR, PMSE, and model size (MSize) averaged over 500 simulation replicates in Group B settings ( numbers in parentheses are standard deviations)

Setting	Methods	MSize	PDR	FDR	PMSE
GB1	ALasso	45.3 (13.4)	0.8583 (0.072)	0.6953 (0.079)	5.5143 (0.853)
	SCAD	13.03 (2.8)	0.7313 (0.078)	0.1353 (0.115)	4.2373 (0.893)
	SIS+SCAD	3.03 (0.7)	0.1983 (0.042)	0.0163 (0.062)	21.5813 (1.881)
	FoBa	87.03 (1.6)	0.8583 (0.076)	0.8523 (0.013)	9.7503 (2.282)
	FSR	11.23 (1.3)	0.7053 (0.066)	0.0503 (0.067)	4.0093 (0.973)
GB2	SLasso	11.23 (2.1)	0.6773 (0.109)	0.0833 (0.082)	4.5033 (2.654)
	ALasso	56.53 (14.7)	0.8813 (0.077)	0.7523 (0.059)	6.8023 (1.371)
	SCAD	17.23 (4.0)	0.7923 (0.065)	0.2773 (0.136)	3.8583 (0.997)
	SIS+SCAD	4.73 (0.6)	0.2593 (0.041)	0.1753 (0.127)	15.3273 (2.527)
	FoBa	86.83 (1.6)	0.8553 (0.056)	0.8523 (0.010)	12.4693 (3.468)
GB3	FSR	11.03 (2.2)	0.6963 (0.128)	0.0403 (0.062)	4.3493 (1.625)
	SLasso	10.53 (2.6)	0.6603 (0.155)	0.0533 (0.073)	4.9663 (2.630)
	ALasso	69.63 (5.9)	0.8523 (0.051)	0.8773 (0.012)	7.8933 (4.293)
	SCAD	8.8 (2.6)	0.583 (0.104)	0.308 (0.125)	28.737 (10.868)
	SIS+SCAD	4.3 (0.7)	0.000 (0.000)	1.00 (0.000)	58.334 (10.717)
	FoBa	82.5 (1.7)	0.997 (0.018)	0.879 (0.004)	7.701 (3.389)
	FSR	18.2 (3.0)	0.785 (0.122)	0.561 (0.075)	8.638 (7.522)
	SLasso	9.8 (3.5)	0.754 (0.31)	0.262 (0.146)	19.470 (20.808)

Following the same strategy of Huang, Ma, and Zhang (2008), the probes are first screened according to their variances and the top 3000 probes with the largest variances are retained for further selection. But, unlike in Huang, Ma, and Zhang (2008) where these 3000 probes were further reduced to 200 probes that are marginally most correlated with TRIM32, the concerned five methods are directly applied to the 3000 probes. The numbers of probes selected from these 3000 probes by ALasso, SCAD, SIS+SCAD, FSR, and SLasso are 21, 28, 5, 3, and 2, respectively. The ID of the selected probes are reported in Table 4. The two probes selected by SLasso, that is, 1383110\_at and 1392692\_at, are also selected by FSR and ALasso. But they are not selected by SCAD and SIS+SCAD. The additional probe selected by FSR, 1389584\_at, is also selected by SCAD and ALasso. There is an intersection of seven probes selected by ALasso and SCAD. There is no intersection of the probes selected by SIS+SCAD with any other methods.

It is interesting to note that one of the probes selected by SLasso, that is, 1383110\_at, is also detected by other methods (Lasso, scaled Lasso, scaled MC) and the other one, 1392692\_at, is also detected by Lasso, as reported in Sun and Zhang (2011). Combining all these findings together and taking into account the low FDR of SLasso evidenced in the simulation studies, we have a strong belief that the two probes selected by SLasso are associated with TRIM32.

## 6. DISCUSSION

The properties of SLasso show that SLasso and the orthogonal matching pursuit (OMP) differ only at steps where the partial positive cone condition is violated. When the partial positive cone condition is satisfied at each step, SLasso is equivalent to OMP. Since the set  $S_{TEMP}$  at the steps of SLasso is rarely nonsingleton, the procedure of SLasso is essentially the same

as OMP, and hence, as a byproduct of the article, we reveal new

Table 4. The ID of the probes selected by ALasso, SCAD, SIS+SCAD, FSR, and SLasso in the analysis of the rat data

Method	Probes ID				
ALasso	1387060_at	1388538_at	1380070_at	1370052_at	1382452_at
	1379079_at	1397489_at	1374131_at	1383110_at	1389584_at
	1392692_at	1379971_at	1385687_at	1369353_at	1374106_at
	1383673_at	1379495_at	1383749_at	1382835_at	1395415_at
	1383996_at				
SCAD	1394689_at	1370434_a_at	1375724_at	1378765_at	1375139_at
	1388538_at	1370052_at	1382452_at	1377781_at	1383841_at
	1380311_at	1379460_at	1385921_at	1384886_at	1384136_at
	1387111_at	1390789_at	1376693_at	1389584_at	1389231_at
	1390788_a_at	1367741_at	1374106_at	1387455_a_at	1383749_at
SIS+SCAD	1379803_at	1383996_at	1382633_at		
	1377546_at	1396809_at	1381430_at	1393543_at	1372481_at
FSR	1383110_at	1392692_at	1389584_at		
SLasso	1383110_at	1392692_at			

properties of OMP other than those discovered in Cai and Wang (2011), Tropp (2004), and Tropp and Gilbert (2007) under much weaker conditions.

Since the mechanism for the selection of features in FoBa and CoSaMP are essentially the same as OMP, our results also reveal that the backward steps in FoBa are indeed not needed and that the iterative procedure of CoSaMP to get the best sparse set at a given sparsity level is not really necessary. The crucial issue in all these procedures is actually the stopping rule. Some ad hoc stopping rules are adopted for OMP in Cai and Wang (2011). These rules compare the norm ( $L_2$  or  $L_\infty$ ) of the residual  $\tilde{y}$  with an upper bound of the norm of error  $\epsilon$ . The procedure continues until the norm of the residual falls below the upper bound. For example, when  $\epsilon \sim N(0, \sigma^2 I)$ ,  $\|\tilde{y}\|_2$  is compared with  $\sigma\sqrt{n+2\sqrt{n}\ln n}$ . With this rule, the OMP can correctly select  $s_0$  with probability  $1 - 1/n$ , if  $\min_{j \in s_0} |\beta_j| \geq \frac{2\sigma\sqrt{n+2\sqrt{n}\ln n}}{1-(2p_0-1)MI}$  and  $MI < \frac{1}{2p_0-1}$ , where  $MI$  is the mutual incoherence defined by  $MI = \max_{i \neq j} |\text{CORR}(\mathbf{x}_i, \mathbf{x}_j)|$ . The limitations of this rule are obvious. First, its effectiveness relies on the strict condition  $MI < \frac{1}{2p_0-1}$  which essentially imposes mutual independence among the features when  $p_0 \rightarrow \infty$ , as assumed in our setting. Second, the lower bound required for  $\min_{j \in s_0} |\beta_j|$  is too large compared with the requirement in SLasso that  $\min_{j \in s_0} |\beta_j| > Cn^{-\delta}$  where  $\delta < 1/2 - \kappa$  for some  $\kappa < 1/2$ ; see condition A3 and the remark that follows in Section 3.1. Third, it requires an accurate estimate of  $\sigma$  which cannot be easily obtained without knowing the true model. The other ad hoc stopping rules considered in Cai and Wang (2011) have the same limitations. A similar stopping rule is considered in Zhang (2011) for FoBa. The rule compares the decrease in  $\frac{1}{n}\|\tilde{y}\|_2^2$  with a threshold of order  $O(\sigma^2 \ln p/n)$  at the forward steps. Similar comments can be made on this stopping rule. A common nature of the above stopping rules is that only the contribution of the features to the decrease of the residual is taken into account, the contribution is not penalized by the increase in the number of features. This common nature is a crucial drawback. It has the potential to select more irrelevant features, which is demonstrated for FoBa in the simulation studies. As a contrast, by using EBIC as the stopping rule, a feature can be selected only when its contribution to the decrease of the residual is large enough to compensate the increase in the number of features. This is perhaps the main reason why SLasso cum EBIC is selection consistent under much weaker conditions.

Let  $g_1(j) = |\mathbf{x}_j^T[\mathbf{I} - \mathbf{H}(s_{*k})]\mathbf{y}|$ , where  $\mathbf{H}(s_{*k})$  is the projection matrix of the space spanned by the features in  $s_{*k}$ . SLasso selects the next features among the features that maximize  $g_1(j)$  after the submodel  $s_{*k}$  is selected. This is to be compared with FSR that selects the next feature by minimizing  $\text{RSS}(j) = \mathbf{y}^T[\mathbf{I} - \mathbf{H}(s_{*k} \cup \{j\})]\mathbf{y}$ , which is equivalent to maximizing  $g_2(j) = \frac{|\mathbf{x}_j^T[\mathbf{I} - \mathbf{H}(s_{*k})]\mathbf{y}|}{\sqrt{\mathbf{x}_j^T[\mathbf{I} - \mathbf{H}(s_{*k})]\mathbf{x}_j}}$ . The equivalence follows from (2.7). SLasso selects the next feature that has the highest correlation with the current residual but the FSR selects the next feature that has the highest inflated correlation with the current residual by an inflating factor  $[\mathbf{x}_j^T[\mathbf{I} - \mathbf{H}(s_{*k})]\mathbf{x}_j]^{-1/2}$ . The more correlated the  $\mathbf{x}_j$  is with the features in  $s_{*k}$ , the larger the inflating factor. If two features have the same absolute correlation with the current residual, the FSR will select the one

that is more correlated with the features in  $s_{*k}$ . If one feature has a lower correlation with the current residual but is more correlated with the features in  $s_{*k}$  than another feature, it might turn out that this feature has a higher inflated correlation and is selected by FSR. Obviously, this is a disadvantage of FSR in terms of the identification of relevant features, especially when high spurious correlations present.

Like SLasso, solution path of Lasso, LAR, and variants of LAR also select the next feature that has the highest correlation with the current residual. But, in these methods, the current residual is obtained from a shrunk estimate of  $E\mathbf{y}$ , that is, they select  $\mathbf{x}_j$  that maximizes  $g_3(j) = |\mathbf{x}_j^T[\mathbf{y} - \mathbf{X}(s_{*k})\tilde{\boldsymbol{\beta}}(s_{*k})]|$  where  $\tilde{\boldsymbol{\beta}}(s_{*k})$  is a shrunk estimate. In the shrunk estimate, the effects on  $\mathbf{y}$  of the features in  $s_{*k}$  are not fully counted. This leaves more chance for those features that have high spurious correlations with the features in  $s_{*k}$  to be selected in subsequent steps than in the case of SLasso. This is a potential disadvantage for the identification of relevant features.

## 7. SUPPLEMENTARY MATERIALS

The supplementary materials contained in the online document include

1. Technical proofs.
  - 1.1 Proof of Theorem 1.
  - 1.2 Proof of Lemma 1.
  - 1.3 Proof of Theorem 3.
2. The details of the two special cases given in Section 3.4.
3. Covariance structure of the design matrix in the simulation studies.

[Received July 2011. Revised October 2013.]

## REFERENCES

- Breheny, P. (2011), "ncvreg: Regularization Paths for SCAD- and MCP-Penalized Regression Models," R package version 2.3-2. Available at <http://cran.r-project.org/web/packages/ncvreg/index.html> [1235]
- Cai, T. T., and Wang, L. (2011), "Orthogonal Matching Pursuit for Sparse Signal Recovery With Noise," *IEEE Transactions on Information Theory*, 57, 4680–4688. [1230,1239]
- Chen, J. H., and Chen, Z. H. (2008), "Extended Bayesian Information Criteria for Model Selection With Large Model Spaces," *Biometrika*, 95, 759–771. [1231,1234]
- Chiang, A. P., Beck, J. S., Yen, H.-J., Tayeh, M. K., Scheetz, T. E., Swiderski, R. E., Nishimura, D. Y., Braun, T. A., Kim, K.-Y., Huang, J., Elbedour, K., Carmi, R., Slusarski, D. C., Casavant, T. L., Stone, E. M., and Sheffield, V. C. (2006), "Homozygosity Mapping With SNP Arrays Identifies TRIM32, an E3 Ubiquitin Ligase, as a Bardet-Biedl Syndrome Gene (BBS11)," *Proceedings of the National Academy of Sciences*, 103, 6287–6292. [1237]
- Huang, J., Ma, S., and Zhang, C.-H. (2008), "Adaptive Lasso for Sparse High-Dimensional Regression Models," *Statistica Sinica*, 18, 1603–1618. [1229,1235,1237]
- Hwang, W. Y., Zhang, H. H., and Ghosal, S. (2009), "FIRST: Combining Forward Iterative Selection and Shrinkage in High Dimensional Sparse Linear Regression," *Statistics and Its Interface*, 2, 341–348. [1237]
- Efron, B., Hastie, T., Johnstone, I., and Tibshirani, R. (2004), "Least Angle Regression" (with discussion), *The Annals of Statistics*, 32, 407–499. [1230]
- Fan, J., Feng, Y., Samworth, R., and Wu, Y. (2010), "SIS: Sure Independence Screening," R package version 0.6. Available at <http://cran.r-project.org/web/packages/SIS/index.html> [1235]



- Fan, J. Q., Feng, Y., and Song, R. (2011), "Nonparametric Independence Screening in Sparse Ultra-High Dimensional Additive Models," *Journal of the American Statistical Association*, 116, 544–557. [1237]
- Fan, J., and Li, R. (2001), "Variable Selection Via Non-Concave Penalized Likelihood and Its Oracle Properties," *Journal of the American Statistical Association*, 96, 1348–1360. [1229]
- (2004), "New Estimation and Model Selection Procedures for Semiparametric Modeling in Longitudinal Data Analysis," *Journal of the American Statistical Association*, 99, 710–723. [1229]
- Fan, J., and Lv, J. (2008), "Sure Independence Screening for Ultra-High Dimensional Feature Space," *Journal of the Royal Statistical Society, Series B*, 70, 849–911. [1230,1231,1235]
- Fan, J., and Peng, H. (2004), "Nonconcave Penalized Likelihood With a Diverging Number of Parameters," *The Annals of Statistics*, 32, 928–961. [1229]
- Kim, Y. D., Choi, H., and Oh, H.-S. (2008), "Smoothly Clipped Absolute Deviation on High Dimensions," *Journal of the American Statistical Association*, 103, 1665–1673. [1235,1237]
- Knight, K., and Fu, W. (2000), "Asymptotics for Lasso-Type Estimators," *The Annals of Statistics*, 28, 1356–1378. [1229]
- Kraemer, N., and Schaefer, J. (2010), "parcor: Regularized Estimation of Partial Correlation Matrices," R package version 0.2-2. Available at <http://cran.r-project.org/web/packages/parcor/index.html> [1235]
- Leng, C., Lin, Y., and Wahba, G. (2004), "A Note on the Lasso and Related Procedures in Model Selection," *Statistica Sinica*, 16, 1273–1284. [1229]
- Luo, S., and Chen, Z. (2011), "Extended BIC for Linear Regression Models With Diverging Number of Relevant Features and High or Ultra-High Feature Spaces," *Journal of Statistical Planning and Inference*, 143, 494–504. [1234]
- Meinshausen, N., and Bühlmann, P. (2006), "High Dimensional Graphs and Variable Selection With the Lasso," *The Annals of Statistics*, 34, 1436–1462. [1229]
- Needell, D., and Tropp, J. A. (2009), "CoSaMP: Iterative Signal Recovery From Incomplete and Inaccurate Samples," *Applied and Computational Harmonic Analysis*, 26, 301–321. [1230,1235]
- Park, M. Y., and Hastie, T. (2007), " $L_1$ -Regularization Path Algorithm for Generalized Linear Models," *Journal of the Royal Statistical Society, Series B*, 69, 659–677. [1231]
- Radchenko, P., and James, G. (2011), "Improved Variable Selection With Forward-LASSO Adaptive Shrinkage," *The Annals of Applied Statistics*, 5, 427–448. [1230]
- Scheetz, T. E., Kim, K.-Y. A., Swiderski, R. E., Philip, A. R., Braun, T. A., Knudtson, K. L., Dorrance, A. M., DiBona, G. F., Huang, J., Casavant, T. L., Sheffield, V. C., and Stone, E. M. (2006), "Regularization of Gene Expression in the Mammalian Eye and Its Relevance to Eye Disease," *Proceedings of the National Academy of Sciences*, 103, 14429–14434. [1237]
- Shao, J. (2003), *Mathematical Statistics* (2nd ed.), New York: Springer. [1235]
- Sun, T. N., and Zhang, C. H. (2011), "Scaled Sparse Linear Regression," arXiv:1104.4595v1. [1237,1238]
- Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society, Series B*, 58, 267–288. [1229]
- Tropp, J. A. (2004), "Greed is Good: Algorithmic Results for Sparse Approximation," *IEEE Transactions on Information Theory*, 50, 1–21. [1230,1233,1239]
- Tropp, J. A., and Gilbert, A. C. (2007), "Signal Recovery From Random Measurements via Orthogonal Matching Pursuit," *IEEE Transactions on Information Theory*, 53, 4655–4666. [1230,1239]
- Wang, H. (2009), "Forward Regression for Ultra-High Dimensional Variable Screening," *Journal of the American Statistical Association*, 104, 1512–1524. [1230,1235]
- Weisberg, S. (1980), *Applied Linear Regression*, New York: Wiley. [1230]
- Xie, H., and Huang, J. (2009), "Scad-Penalized Regression in High-Dimensional Partially Linear Models," *The Annals of Statistics*, 37, 673–696. [1229,1235]
- Pati, Y. C., Rezaiifar, R., and Krishnaprasad, P. S. (1993), "Orthogonal Matching Pursuit: Recursive Function Approximation With Applications to Wavelet Decomposition," in *1993 Conference Record of The Twenty-Seventh Asilomar Conference on Signals, Systems and Computers*, pp. 40–44. [1230]
- Zhang, C.-H. (2010), "Nearly Unbiased Variable Selection Under Minimax Concave Penalty," *The Annals of Statistics*, 38, 894–942. [1229]
- Zhang, T. (2011), "Adaptive Forward-Backward Greedy Algorithm for Learning Sparse Representations," *IEEE Transactions on Information Theory*, 57, 4689–4708. [1230,1235,1239]
- Zhao, P., and Yu, B. (2006), "On Model Selection Consistency of Lasso," *Journal of Machine Learning Research*, 7, 2541–2567. [1229]
- Zhong, W., Zhang, T., Zhu, Y., and Liu, J. (2012), "Correlation Pursuit: Forward Stepwise Variable Selection for Index Models," *Journal of the Royal Statistical Society, Series B*, 74, 849–870. [1230]
- Zou, H. (2006), "The Adaptive Lasso and Its Oracle Properties," *Journal of the American Statistical Association*, 101, 1418–1429. [1229]