# Regression Coefficient and Autoregressive Order Shrinkage and Selection via Lasso

Hansheng Wang

*Guanghua School of Management, Peking University*

hansheng@gsm.pku.edu.cn

Guodong Li

*Department of Statistics and Actuarial Science, University of Hong Kong*

ligd@hkusua.hku.hk

and Chih-Ling Tsai

*Graduate School of Management, University of California - Davis*

cltsai@ucdavis.edu

### Abstract

The *least absolute shrinkage and selection operator* (lasso) has been widely used in regression shrinkage and selection. In this article, we extend its application to the REGression model with AutoRegressive errors (REGAR). Two types of lasso estimators are carefully studied. The first is similar to the traditional lasso estimator with only two tuning parameters (one for regression coefficients and the other for autoregression coefficients). These tuning parameters can be easily calculated via a data driven method, but the resulting lasso estimator may not be fully efficient (Fan and Li, 2001). In order to overcome this limitation, we propose a second lasso estimator which uses different tuning parameters for each coefficient. We show that this modified lasso is able to produce the estimator as efficiently as the *oracle*. Moreover, we propose an algorithm for tuning parameter estimates to obtain the modified lasso estimator. Simulation studies demonstrate that the modified estimator is superior to the traditional one. One empirical example is also presented to illustrate the usefulness of lasso estimators. The extension of lasso to the autoregression with exogenous variables (ARX) model is briefly discussed.

*Key words:* ARX; Lasso; Oracle Estimator; REGAR

## 1    Introduction

The linear regression model is a commonly used statistical tool for analysis of the relationships between response and explanatory variables. One of its standard assumptions is that different observations are independent. However, significant serial correlation might occur when the data are collected sequentially in time. In this case, the linear REGression with AutoRegressive errors (REGAR) model is often considered, as it takes into account the autocorrelated structure in regression analysis (Shumway and Stoffer, 2000; Tsay, 1984; Harvey, 1981).

In model building, it is known that making the model unnecessarily complex can degrade the efficiency of the resulting parameter estimator and yield less accurate predictions. Hence, two heuristic selection criteria, AIC (Akaike, 1973) and BIC (Schwarz, 1978), are often applied to select regression variables. In the context of time series, both criteria are also employed to choose the order

of the autoregressive (AR) process (Brockwell and Davis, 1991; Choi, 1992; McQuarrie and Tsai, 1998; Shumway and Stoffer, 2000). Moreover, Ramanathan (1989) extends the application of AIC and BIC to the linear REGAR model. However, as noted by researchers, the statistical performance of AIC and BIC can be unstable (Breiman, 1996), and selection bias may cause inference problems (Hurvich and Tsai, 1990).

To amend the deficiencies of classical linear model selections, Tibshirani (1996) developed the *least absolute shrinkage and selection operator* (lasso), which selects variables and estimates parameters simultaneously. This motivated us to obtain the shrinkage estimator in the AR process. To this end, we employ the lasso-type penalty not only on the regression coefficients but also on the autoregression coefficients. Consequently, a direct extension of lasso to the REGAR model involves two regularization parameters (i.e., one for regression coefficients and the other for autoregression coefficients), which can be easily tuned via a data driven method (e.g., cross-validation). We show that the resulting lasso estimator satisfies a Knight & Fu - type asymptotic property (Knight and Fu, 2000). However, it suffers an appreciable bias (Fan and Li, 2001). Hence, the traditional lasso estimator cannot achieve the same efficiency as the *oracle*, i.e., the estimator obtained based on the true model (Fan and Li, 2001).

To improve the utility of the traditional lasso approach to the REGAR model, we modify the penalty function so that different tuning parameters can be used for each coefficient. As a result, large amounts of shrinkage can be used for the insignificant variables, while small amounts of shrinkage can be used for the significant variables. We show that the resulting modified lasso estimator shares the same asymptotic distribution as the *oracle*. In practice, however, simultaneously tuning many regularization parameters is not realistic. Therefore, we further propose the tuning parameter algorithm via the unpenalized REGAR estimator. Simulation studies indicate that the resulting lasso estimator outperforms the traditional lasso estimator.

The rest of the paper is organized as follows. Section 2 introduces the REGAR model and the two lasso estimators. The asymptotic theories of the two lasso estimators are established in Section 3. The practical implementations of these two estimators are developed in Section 4, and numerical studies are presented in Section 5. Section 6 concludes the article with a brief discussion.

## 2 Least Absolute Shrinkage and Selection Operators

Consider the linear regression with autoregressive errors (REGAR) model

$$y_t = x_t'\beta + e_t, \qquad (t = 1, \cdots, n_0), \tag{1}$$

where $x_t = (x_{t1}, \cdots, x_{tp})'$ is the $p$-dimensional regression covariate and $\beta = (\beta_1, \cdots, \beta_p)'$ is the associated regression coefficient. In addition, the variable $e_t$ follows the autoregressive process

$$e_t = \phi_1 e_{t-1} + \phi_2 e_{t-2} + \cdots + \phi_q e_{t-q} + \epsilon_t, \tag{2}$$

where $\phi = (\phi_1, \cdots, \phi_q)'$ is the autoregression coefficient and $\epsilon_t$ are independent and identically distributed random variables with mean 0 and variance $\sigma^2$. Moreover, we define the regression and autoregressive parameters as $\theta = (\beta', \phi')'$. For practical implementation, it is a common practice to standardize the predictor $x_{tj}$ so that it has zero mean and unit variance (Tibshirani, 1996).

Analogously, the response $y_t$ is scaled by dividing it with the estimate of $[\mathrm{var}(e_t)]^{1/2}$.

Suppose that $\epsilon_t$ in (2) follows a normal distribution and the first $q$ observations are fixed. Then the conditional likelihood function of the remaining $n_0 - q$ observations, $(y_{q+1}, \cdots, y_{n_0})'$, is

$$
\left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left\{-\frac{1}{2\sigma^2} \sum_{t=q+1}^{n_0} \left[(y_t - x_t'\beta) - \sum_{j=1}^{q} \phi_j(y_{t-j} - x_{t-j}'\beta)\right]^2\right\},
$$

where $n = n_0 - q$ is the effective sample size. Maximizing the above likelihood function yields a conditional maximum likelihood estimator (MLE) of $\theta$. This estimator can also be obtained by minimizing the following least squares type objective function,

$$
L_n(\theta) = \sum_{t=q+1}^{n_0} \left[(y_t - x_t'\beta) - \sum_{j=1}^{q} \phi_j(y_{t-j} - x_{t-j}'\beta)\right]^2, \tag{3}
$$

where $L_n(\theta)$ is an extension of the method proposed by Cochrane and Orcutt (1949) for $q = 1$; see Harvey (1981) and Hamilton (1994).

In order to shrink unnecessary coefficients to zero, we next adapt Tibshirani's (1996) approach for obtaining the estimator by minimizing the following lasso criterion

$$
Q_n(\theta) = \sum_{t=q+1}^{n_0} \left[(y_t - x_t'\beta) - \sum_{j=1}^{q} \phi_j(y_{t-j} - x_{t-j}'\beta)\right]^2 + n\sum_{j=1}^{p} \lambda|\beta_j| + n\sum_{j=1}^{q} \gamma|\phi_j|. \tag{4}
$$

Because lasso uses the same tuning parameters $\lambda$ and $\gamma$ for the regression and autoregressive coefficients, respectively, the resulting estimator, $\hat{\theta} = (\hat{\beta}', \hat{\phi}')'$, may suffer an appreciable bias. This is mainly due to the fact that all the regression (or autoregression) coefficients share the same amount of shrinkage (see Fan and Li, 2001). To overcome this limitation, we propose the following modified lasso criterion, lasso*,

$$
Q_n^*(\theta) = \sum_{t=q+1}^{n_0} \left[(y_t - x_t'\beta) - \sum_{j=1}^{q} \phi_j(y_{t-j} - x_{t-j}'\beta)\right]^2 + n\sum_{j=1}^{p} \lambda_j^*|\beta_j| + n\sum_{j=1}^{q} \gamma_j^*|\phi_j|, \tag{5}
$$

which allows for different tuning parameters $\lambda_j^*$ and $\gamma_j^*$ for different coefficients. As a result, a larger amount of shrinkage can be applied for the insignificant coefficients, while a smaller amount of shrinkage can be employed to the significant coefficients. Hence, the resulting estimator, $\hat{\theta}^* = (\hat{\beta}^{*'}, \hat{\phi}^{*'})'$, is expected to have a smaller bias than $\hat{\theta}$. The detailed investigations of these two estimators are given in the next section.

## 3   Theoretical Properties

To study the theoretical properties of the two lasso estimators, we assume that there is a correct model with the regression and autoregression coefficients $\theta^0 = (\beta^{0'}, \phi^{0'})' = (\beta_1^0, \cdots, \beta_p^0, \phi_1^0, \cdots, \phi_q^0)'$. Furthermore, we assume that there are a total of $p_0 \leq p$ non-zero regression coefficients and $q_0 \leq q$

non-zero autoregression coefficients. For the sake of convenience, we define $\mathcal{S}_1 = \{1 \leq j \leq p : \beta_j^0 \neq 0\}$ and $\mathcal{S}_2 = \{1 \leq j \leq q : \phi_j^0 \neq 0\}$. Then, the set $\mathcal{S}_1$ ($\mathcal{S}_2$) contains the indices of the significant regression (autoregression) coefficients, while its complement $\mathcal{S}_1^c$ ($\mathcal{S}_2^c$) contains the indices of the insignificant regression (autoregression) coefficients. Next, let $\beta_{\mathcal{S}_1}$ denote the $p_0 \times 1$ significant regression coefficient vector with $\hat{\beta}_{\mathcal{S}_1}$ as its associated lasso estimator. Moreover, other related parameters and their corresponding estimators are analogously defined (e.g., $\beta_{\mathcal{S}_1^c}$, $\hat{\beta}_{\mathcal{S}_1^c}$, $\hat{\beta}_{\mathcal{S}_1}^*$, $\phi_{\mathcal{S}_2}$, $\hat{\phi}_{\mathcal{S}_2}$, etc). Finally, let $\theta_1^0 = (\beta_{\mathcal{S}_1}^{0\prime}, \phi_{\mathcal{S}_2}^{0\prime})'$ and $\theta_2^0 = (\beta_{\mathcal{S}_1^c}^{0\prime}, \phi_{\mathcal{S}_2^c}^{0\prime})'$. Then, $\hat{\theta}_k$ and $\hat{\theta}_k^*$ ($k = 1, 2$) are the associated lasso and lasso* estimators, respectively. To investigate the theoretical properties of $\hat{\theta}$ and $\hat{\theta}^*$, we introduce the following conditions.

(C.1) The sequence $\{x_t\}$ is independent of $\{\epsilon_t\}$.

(C.2) All roots of the polynomial $1 - \sum_{j=1}^q \phi_j^0 z^j$ are outside the unit circle.

(C.3) The error $\epsilon_t$ has the finite fourth order moment, i.e., $\mathrm{E}(\epsilon_t^4) < \infty$.

(C.4) The covariate $x_t$ is strictly stationary and ergodic with finite second order moment (i.e., $\mathrm{E}\|x_t\|^2 < \infty$). Furthermore, the following matrix is positive definite,

$$B = \mathrm{E}\left\{ \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right) \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right)' \right\}.$$

The technical conditions above are typically used to assure the $\sqrt{n}$-consistency and asymptotic normality of the unpenalized least square estimator.

## 3.1 The Lasso Estimator

In this subsection, we study the property of the traditional lasso estimator given below.

**Theorem 1.** *Assume that $\sqrt{n}\lambda_n \to \lambda_0$ and $\sqrt{n}\gamma_n \to \gamma_0$ for some $\lambda_0 \geq 0$ and $\gamma_0 \geq 0$. Then, under the conditions (C.1) - (C.4), we have $\sqrt{n}(\hat{\theta} - \theta^0) \to_d argmin\{\kappa(\delta)\}$, where*

$$\kappa(\delta) = -2\delta'w + \delta'\Sigma\delta + \lambda_0 \sum_{j=1}^p \left\{ u_j sgn\{\beta_j^0\}I(\beta_j^0 \neq 0) + |u_j|I(\beta_j^0 = 0) \right\}$$

$$+ \gamma_0 \sum_{j=1}^q \left\{ v_j sgn\{\phi_j^0\}I(\phi_j^0 \neq 0) + |v_j|I(\phi_j^0 = 0) \right\},$$

$\delta = (u', v')'$, $w \sim N(0, \sigma^2\Sigma)$, $\Sigma = diag\{B, C\}$, $C = (\xi_{|i-j|})$, *and* $\xi_k = E(e_t e_{t+k})$.

The proof is given in Appendix A. Theorem 1 shows that the lasso estimator possesses a Knight & Fu - type asymptotic property (Knight and Fu, 2000). This implies that the tuning parameters used in the traditional lasso estimator cannot shrink to 0 at a speed faster than $n^{-1/2}$. Otherwise, both $\lambda_0$ and $\gamma_0$ degenerate to zero and $\kappa(\delta)$ becomes a standard quadratic function,

$$\kappa(\delta) = \kappa(u, v) = -2(u', v')w + (u', v')\Sigma(u', v')',$$

which is unable to produce sparse solutions. Therefore, Theorem 1 suggests that $\lambda_0 > 0$ and $\gamma_0 > 0$ are needed for obtaining the traditional lasso estimator.

*Remark 1.* In a standard regression model with independent observations, Fan and Li (2001) noticed that the traditional lasso estimator may suffer an appreciable bias. Therefore, it is of interest to investigate whether the traditional lasso estimator for the REGAR model encounters the same problem. For the sake of illustration, we consider a special case with $\beta_j^0 > 0$ for $1 \le j \le p$ and $\phi_j^0 = 0$ for $1 \le j \le q$. If the minimizer of $\kappa(\delta)$ can correctly identify the true model, then $u \ne 0$ but $v = 0$. In addition, $\kappa(\delta)$ satisfies the following equation

$$\frac{\partial \kappa(u,0)}{\partial u} = -2w_1 + 2u'B + \lambda_0 \mathbf{1} = 0,$$

where $w_1$ consists of the first $p$ components of $w$ and $\mathbf{1}$ is a $p \times 1$ vector with the elements of ones. As a result, $\sqrt{n}(\hat{\beta} - \beta^0) \to_d u = B^{-1}(w_1 - 0.5\lambda_0 \mathbf{1})$, which is distributed as $N(-0.5\lambda_0 B^{-1}\mathbf{1}, B^{-1})$. Because $\lambda_0 > 0$, Theorem 1 indicates that the traditional lasso estimator is asymptotically biased. Thus, it is not as efficient as the *oracle* estimator, whose asymptotic distribution is $N(0, B^{-1})$.

## 3.2   The Lasso* Estimator

In this subsection, we focus on the modified lasso (lasso*) estimator. To facilitate studying the properties of this estimator, we introduce the following notation:

$$a_n = \max\{\lambda_{j_1}^*, \gamma_{j_2}^*, j_1 \in \mathcal{S}_1, j_2 \in \mathcal{S}_2\} \quad \text{and} \quad b_n = \min\{\lambda_{j_1}^*, \gamma_{j_2}^*, j_1 \in \mathcal{S}_1^c, j_2 \in \mathcal{S}_2^c\},$$

where $\lambda_{j_1}^*$ and $\gamma_{j_2}^*$ are functions of $n$. We first investigate the consistency of the lasso* estimator.

**Lemma 1.** *Assume that $a_n = o(1)$ as $n \to \infty$ . Then under the conditions (C.1) - (C.4), there exists a local minimizer $\hat{\theta}^*$ of $Q_n^*(\theta)$ such that $\hat{\theta}^* - \theta^0 = O_p(n^{-1/2} + a_n)$.*

The proof is given in Appendix B. Lemma 1 implies that if the tuning parameters associated with the significant regression variables and autoregressive orders converge to 0 at a speed faster than $n^{-1/2}$, then there exists a local minimizer of $Q_n^*(\theta)$, which is $\sqrt{n}$-consistent.

We next show that if the tuning parameters associated with the non-significant regression and autoregressive variables shrink to 0 slower than $n^{-1/2}$, then their regression and autoregression coefficients can be estimated exactly as 0 with probability tending to one.

**Theorem 2.** *Assume that $\sqrt{n}b_n \to \infty$ and $\|\hat{\theta}^* - \theta^0\| = O_p(n^{-1/2})$. Then*

$$P(\hat{\beta}_{\mathcal{S}_1^c}^* = 0) \to 1 \quad and \quad P(\hat{\phi}_{\mathcal{S}_2^c}^* = 0) \to 1.$$

The proof is in Appendix C. Theorem 2 shows that lasso* has the ability to consistently produce a sparse solution for insignificant regression and autoregression coefficients. Furthermore, this theorem, together with Lemma 1, indicates that the $\sqrt{n}$-consistent estimator $\hat{\theta}^*$ must satisfy $P(\hat{\theta}_2^* = 0) \to 1$ when the tuning parameters fulfill the appropriate conditions (e.g., $\lambda_j$ and $\gamma_j$ are defined as in equations (7) of the next section). Finally, we obtain the asymptotic distribution of the lasso* estimator.

**Theorem 3.** *Assume that $\sqrt{n}a_n \to 0$ and $\sqrt{n}b_n \to \infty$. Then, under the conditions (C.1) - (C.4), the component $\hat{\theta}_1^*$ of the local minimizer $\hat{\theta}^*$ given in Lemma 1 satisfies*

$$\sqrt{n}(\hat{\theta}_1^* - \theta_1^0) \to_d N(0, \sigma^2 \Sigma_0^{-1}),$$

*where $\Sigma_0$ is the submatrix of $\Sigma$ corresponding to $\theta_1^0$.*

The proof is given in Appendix D. Theorem 3 implies that if the tuning parameters satisfy the conditions $\sqrt{n}a_n \to 0$ and $\sqrt{n}b_n \to \infty$, then, asymptotically, the resulting lasso* estimator can be as efficient as the *oracle* estimator.

# 4    Algorithm

After arriving at an understanding of the properties of the two lasso estimators, it is natural to implement them for real applications. To this end, we propose the following algorithm to obtain the local minimizers for lasso estimators $\hat{\theta}$ and $\hat{\theta}^*$. In addition, we provide an approach to simultaneously estimate a total of $(p + q)$ tuning parameters for the lasso* estimator.

## 4.1    The Iterative Process

The objective function $Q_n^*(\theta)$ contains $Q_n(\theta)$ as a special case (i.e., $\lambda_j = \lambda$ and $\gamma_j = \gamma$). Therefore, we focus mainly on the optimization problem of $Q_n^*(\theta)$ in the rest of this section. Because Equation (5) contains both regression and autoregression parameters, it is sensible to optimize the objective function $Q_n^*(\theta)$ iteratively by minimizing the following two lasso-type objective functions:

$$\sum_{t=q+1}^{n_0} \left[ (y_t - x_t'\beta) - \sum_{j=1}^q \phi_j (y_{t-j} - x_{t-j}'\beta) \right]^2 + n \sum_{j=1}^p \lambda_j |\beta_j| \text{ with a fixed } \phi,$$

and

$$\sum_{t=q+1}^{n_0} \left[ (y_t - x_t'\beta) - \sum_{j=1}^q \phi_j (y_{t-j} - x_{t-j}'\beta) \right]^2 + n \sum_{j=1}^q \gamma_j |\phi_j| \text{ with a fixed } \beta.$$

As a result, many well developed procedures can be used to find the solution for the above non-concave penalized functions. For example, quadratic programming (Tibshirani, 1996), the shooting algorithm (Fu, 1998), local quadratic approximation (Fan and Li, 2001), and, most recently, the least angle regression method (Efron *et al.*, 2004). For the sake of simplicity, we adapt the local quadratic approximation procedure, which was first developed by Fan and Li (2001) and has been used extensively in the literature (e.g., see Fan and Li (2002), Fan and Peng (2004), and Cai *et al.* (2005)). Our simulation studies indicate that this procedure converges with a reasonable degree of speed and accuracy.

*Remark 2.* The solution of the local quadratic approximation does not yield a sparse solution. However, the small parameter estimate produced by the local quadratic approximation can be arbitrarily close to 0, as long as a sufficiently small threshold for its tolerance of accuracy is set

up. For the sake of illustration, the ordinary linear regression is considered. In this case, the local quadratic approximation produces the one step ahead estimate, $\beta^{(m+1)}$, by minimizing

$$||Y - X\beta^{(m+1)}||^2 + n\sum_{j=1}^{p}\lambda_j\frac{(\beta_j^{(m+1)})^2}{|\beta_j^{(m)}|},$$

where $Y = (y_1, \cdots, y_{n_0})'$ and $X = (x_1, \cdots, x_{n_0})'$. If one of the coefficients (e.g., $\beta_1^{(m)}$) is very small (but not sparse), then the ridge effect induced by $\beta_1^{(m)}$, $\lambda_1/|\beta_1^{(m)}|$, can be very large. As a result, the value of $|\beta_1^{(m+1)}|$ is forced to be even smaller. Because this is an iterative process, the value of $|\beta_1^{(m)}|$ can be arbitrarily close to 0 as long as we can have a sufficiently small threshold for accuracy. Therefore, it is possible to set up an arbitrarily small thresholding value to shrink small estimates to be exactly 0. By doing so, the sparse solution is obtained. In simulation studies, we use the thresholding value of $10^{-9}$ so that any coefficient whose the absolute value is smaller than $10^{-9}$, is shrunk to be exactly 0.

## 4.2   Local Convexity

Although the proposed iterative process is easy to implement, one cannot be assured that the resulting estimator converges to the global minimizer. This is because the least squares term, $L_n(\theta)$, in the objective function $Q_n^*(\theta)$ is not a convex function. This motivates us to develop the following theorem, which shows that there is a sufficiently small but fixed local region containing the true parameter in which $L_n(\theta)$ is almost surely guaranteed to be convex.

**Theorem 4.** *There is a probability null set $\mathcal{N}_0$ and a sufficiently small but fixed $\delta > 0$, such that for any $\omega \notin \mathcal{N}_0$, there is an integer $n_\omega$, such that for any $n > n_\omega$, $L_n(\theta)$ is convex in $\theta \in \mathcal{B}_\delta$, where $\mathcal{B}_\delta = \{\theta : \|\theta - \theta^0\| < \delta\}$ is a ball containing the true value $\theta^0$.*

The proof of the above theorem can be obtained upon request from the authors. Theorem 4 indicates that, with probability tending to 1, there will be at most one local minimizer located in $\mathcal{B}_\delta$. According to Lemma 1, $\hat{\theta}^*$ exists and is consistent in probability. Hence, Theorem 4 together with Lemma 1 implies that, with probability tending to one, $\hat{\theta}^*$ is the unique local minimizer in $\mathcal{B}_\delta$. As a result, the desired local minimizer $\hat{\theta}^*$ can be obtained by finding the unique local minimizer in $\mathcal{B}_\delta$.

*Remark 3.* Theorem 4 is applicable not only for the modified lasso estimator $\hat{\theta}^*$, but also for the traditional lasso estimator $\hat{\theta}$. Specifically, Theorem 4 together with Theorem 1 implies that $\hat{\theta}$ can be obtained by finding the unique local minimizer in $\mathcal{B}_\delta$. In practice, however, it is not necessary to know $\mathcal{B}_\delta$ exactly. This is because if the initial estimator is consistent, then it must be located within $\mathcal{B}_\delta$ with a probability tending to 1. As a result, the proposed iterative process (with aforementioned initial estimator) leads to the local minimizer (i.e., $\hat{\theta}^*$ or $\hat{\theta}$) in $\mathcal{B}_\delta$ with probability tending to 1.

## 4.3 Initial Estimator

To obtain the consistent estimator for the suggested iterative process, we consider the following ordinary least squares estimator as an initial estimator for the regression coefficient $\beta^0$:

$$\hat{\beta}^{(0)} = (X'X)^{-1}(X'Y).$$

Using the fact that $\epsilon_t$ is independent of $x_t$ (see Condition (C.1)), it can be shown that $\hat{\beta}^{(0)}$ is a consistent estimator of $\beta^0$ under classical regularity conditions. Then, computing the ordinary residual $\hat{e}_t = y_t - x_t'\hat{\beta}^{(0)}$ and employing the least squares approach by fitting $\hat{e}_t$ versus $(\hat{e}_{t-1}, \cdots, \hat{e}_{t-q})$, we obtain the initial estimator for the autoregressive coefficient $\phi^0$ given below.

$$\hat{\phi}^{(0)} = (W'W)^{-1}(W'V),$$

where $V = (\hat{e}_{q+1}, \cdots, \hat{e}_{n_0})'$ and $W$ is a $n \times q$ matrix with the $t$'th row given by $(\hat{e}_{t+q-1}, \cdots, \hat{e}_t)$. It can also be shown that $\hat{\phi}^{(0)}$ is a consistent estimator of $\phi^0$ under classical regularity conditions.

## 4.4 Tuning Parameters

After obtaining the initial estimator, we need to select the tuning parameters in the iterative process to complete the whole algorithm. The traditional lasso estimator contains only two tuning parameters (i.e., $\lambda$ and $\gamma$). Hence, one can directly apply the commonly used cross-validation (CV) method to select the optimal tuning parameters. Due to the time series structure, we use the first half of the data for model training and the rest for model testing. In the classical linear regression setting, however, Shao (1997) indicated that BIC would perform better than CV if the true model has a finite dimension and is among the candidate models. This motivates us to adapt Zou *et al*'s (2004) BIC-type tuning parameter selector

$$\text{BIC} = \log(\hat{\sigma}^2) + \hat{df} \times \log n / n, \tag{6}$$

where $\hat{\sigma}^2 = n^{-1} \sum_{t=q+1}^{n_0} \left[(y_t - x_t'\hat{\beta}) - \sum_{j=1}^q \hat{\phi}_j(y_{t-j} - x_{t-j}'\hat{\beta})\right]^2$ and $\hat{df}$ is the number of nonzero coefficients of $\hat{\theta}$.

As for the modified lasso estimator, it becomes a challenging task since there are $(p + q)$ regularization parameters that need to be tuned. Following an anonymous referee's suggestion, we propose the adaptive estimators:

$$\lambda_j^* = \lambda^* \frac{\log(n)}{n|\tilde{\beta}_j|} \quad \text{and} \quad \gamma_j^* = \gamma^* \frac{\log(n)}{n|\tilde{\phi}_j|}, \tag{7}$$

where $\tilde{\theta} = (\tilde{\beta}', \tilde{\phi}')'$ is the unpenalized least square estimator by assuming that $\lambda = \gamma = 0$ in Equation (4). In addition, both $\lambda^*$ and $\gamma^*$ are positive constants and estimated from the data. The advantage of (7) is that it converts the original $(p+q)$-dimensional tuning problem for finding $\lambda_j$ and $\gamma_j$ into a two dimensional task for searching $\lambda^*$ and $\gamma^*$, which can be easily determined by using either CV or BIC.

According to Theorem 1, $\tilde{\theta}$ is an $\sqrt{n}$-consistent estimator of $\theta^0$. Hence, for any $\beta_j^0 \neq 0$ and

$\phi_j^0 \neq 0$, we have $\lambda_j^* = O_p(\log n / n) = o_p(n^{-1/2})$ and $\gamma_j^* = O_p(\log n / n) = o_p(n^{-1/2})$. Consequently, both $\lambda_j^*$ and $\gamma_j^*$ satisfy the condition $\sqrt{n} a_n \to 0$, where $a_n$ is defined in Section 3.2. In contrast, for any $\beta_j^0 = 0$ and $\phi_j^0 = 0$, Theorem 1 implies that $\tilde{\beta}_j = O_p(n^{-1/2})$ and $\tilde{\phi}_j = O_p(n^{-1/2})$. Therefore,

$$\sqrt{n}\lambda_j^* = \lambda^* \frac{\log(n)}{\sqrt{n}\tilde{\beta}_j} \quad \text{and} \quad \sqrt{n}\gamma_j^* = \gamma^* \frac{\log(n)}{\sqrt{n}\tilde{\phi}_j},$$

where the denominators of the above equations are $O_p(1)$ and the numerators go to infinity as $n \to \infty$. As a result, $\sqrt{n}\lambda_j^* \to_p \infty$ and $\sqrt{n}\gamma_j^* \to_p \infty$, which imply that both of them satisfy the condition $\sqrt{n}b_n \to \infty$, where $b_n$ is defined in Section 3.2. In sum, the proposed tuning parameters $\lambda_j^*$ and $\gamma_j^*$ are able to produce the modified lasso estimator $\hat{\theta}^*$, which is as efficient as the *oracle* estimator asymptotically.

# 5 Simulation and Example

## 5.1 Simulation Results

We present Monte Carlo simulations to evaluate the finite sample performance of the lasso estimators. They consist of the traditional and modified lasso estimators with the tuning parameters selected by CV and BIC, respectively. For the traditional lasso estimator, we adapt Zou and Hastie's (2005) approach to select the optimal tuning parameters, $\hat{\lambda}$ and $\hat{\gamma}$, from the grid points $\{0, 0.01, 0.1, 1.0, 10, 100\}$. For the lasso* estimator, the optimal tuning parameter $\hat{\tau}$ is selected from one of 6 equally spaced grid points from 0 to 0.5 (i.e., 0, 0.1, 0.2, $\cdots$, 0.5). Our simulation experience seems to suggest that such a search region and spacing work out satisfactorily. In addition, the estimation algorithm stops if $\sum_j |\theta_j^{(m)} - \theta_j^{(m+1)}| < 10^{-12}$, where $\theta^{(m)} = (\theta_1^{(m)}, \cdots, \theta_{p+q}^{(m)})'$ is the estimator of $\theta$ at the $m$'th iteration, $\theta_j^{(m)} = \beta_j^{(m)}$ for $j = 1, \cdots, p$, and $\theta_j^{(m)} = \phi_{j-p}^{(m)}$ for $j = p+1, \cdots, p+q$. When the convergence is obtained, any parameter estimator whose absolute value is less than $10^{-9}$ is shrunk to 0. Based on our extensive simulation studies, the above proposed stopping and shrinking rules lead to a reasonable convergence speed.

We generated the data from the following REGAR model

$$y_t = 3.0 x_{t1} + 1.5 x_{t2} + 2.0 x_{t5} + e_t, \tag{8}$$

where

$$e_t = 0.5 e_{t-1} - 0.70 e_{t-3} + \sigma \epsilon_t, \tag{9}$$

and $\epsilon_t$ were independent and identically standard normal random variables for $t = 1, \cdots, n_0$. The regression and autocorrelated coefficients are $\beta^0 = (3, 1.5, 0, 0, 2, 0, 0, 0)'$ and $\phi^0 = (0.50, 0, -0.70, 0, 0)'$, respectively. In addition, the covariate $x_t = (x_{t1}, \cdots, x_{t8})'$ were independently generated from the multivariate normal distribution with mean $\mathbf{0}_{8 \times 1}$, and the pairwise correlation between $x_{tj_1}$ and $x_{tj_2}$ is $\rho^{|j_1 - j_2|}$. Note that the regression model (8) is adapted from Tibshirani (1996) and has been used in other simulation studies (e.g., see Fan and Li, 2001; Zou *et al.*, 2005; Leng *et al.*, 2005), while the autoregression model (9) is modified from Shi and Tsai (2004).

In this study, we consider three sample sizes ($n_0 = 50$, 100, and 300) and two standard deviations

($\sigma = 3.0$ and $\sigma = 0.5$). In addition, the correlation coefficients ($\rho$ values) are 0.75, 0.50, and 0.25, which represent the high, moderate, and low linear correlations between the covariates. For each setting, a total of 1,000 realizations were carried out, and the percentage of correctly (under, over) estimated numbers of regression variables, the percentage of correctly (under, over) estimated numbers of autoregressive orders, and the percentage of the correct model identified by two lasso estimators were computed.

When $\rho = 0.5$, Table 1 shows that lasso performs poorly across various sample sizes and noises. This is because lasso's tuning parameter is fixed, and therefore is not able to effectively shrink non-significant coefficients to zero. As a result, it tends to overfit in both regression and autoregression variable selection. In contrast, the lasso* with CV selector (lasso*-CV) demonstrates a considerable improved finite sample performance. Furthermore, the lasso* with BIC selector (lasso*-BIC) performs the best in correct model identifications across various sample sizes and noise levels. Moreover, as the sample size increases, the correct model percentage approaches 100% rapidly. In sum, we recommend employing lasso*-BIC to jointly choose variables and estimate coefficients.

In addition to the correct model identification, an anonymous referee suggested comparing the prediction accuracies of four lasso estimates in terms of their mean squared prediction error (MSPE). To this end, we generated an additional 10,000 independent testing samples within each realization, which are used to evaluate the prediction accuracy. Analogous to the correct model selection results, Table 1 shows that lasso-CV performs the worst, while lasso*-BIC outperforms the rest of lasso estimates. Similar patterns (not presented here) are also found when $\rho = 0.25$ and $\rho = 0.75$.

## 5.2 Electricity Demand Study

We consider a dataset taken from Ramanathan (1989), which studies the electricity consumption of residential customers served by San Diego Gas and Electric Company. The data contains a total of 53 quarterly observations, running from the first quarter of 1970 to the first quarter of 1983. The response variable is the electricity consumption, which is measured by the log-transformed electricity consumption per residential customer in Millions of Kilowatt-Hours (LKWH). The five explanatory variables are the logarithm of per capita real income (LY), the logarithm of real average price of residential electricity in dollars per Kilowatt-Hour (LELP), the logarithm of real ex-post average price of residential gas in Dollars per Therm (LGSP), the cooling degree days per quarter (CDD), and the heating degree days per quarter (HDD).

We first fit the data with the classical multiple regression model, and the resulting estimated equation is $LK\hat{W}H = -8.988 + 0.819LY + 0.154LELP - 0.159LGSP + 0.00012CDD + 0.00042HDD$. The signs of the parameter estimates of variables LY, CDD, and HDD meet our expectations. In other words, an increase in real income (LY), the cooling degree days (CDD), or the heating degree days (HDD) yields more demand for heating. However, the variables LELP and LGSP have unexpected signs since the higher electricity price (LELP) and the higher gas price (LGSP) result in more and less electricity consumption, respectively. Because this is a time series data, the unexpected signs may occur as a result of ignoring the autocorrelation structure. Hence, Ramanathan (1989) naturally recommended the regression model with autoregressive errors.

10

Following Ramanathan's suggestion, we employ lasso and lasso* with CV and BIC to jointly shrink both the regression and autoregression coefficients. The $p$ and $q$ of the candidate models are 5 and 4, respectively, and the maximum autoregressive order 4 is naturally chosen for the quarterly data. Table 2 indicates that lasso with CV and BIC yields the most complicated model, which is consistent with the simulation findings. This overfitted model also leads to an unexpected sign on the variable LGSP. In contrast, both lasso* with CV and lasso* with BIC select the same yet simpler model with variable LELP, CDD, HDD and four lags. It is noteworthy to point out that the two important temperature variables (CDD and HDD) are successfully identified by lasso*. In addition to that the sign of LELP is corrected as compared with the full model regression estimate. To check the adequacy of the model fitting, the $\chi^2$ test statistics for assessing the autocorrelation of residuals (see Box, Jenkins and Reinsel, 1994, p. 314) are computed (see the last line of Table 2). No statistically significant serial correlation is detected in the residuals. In sum, the lasso* estimator with either CV or BIC produces the same simple, interpretable, yet adequate model fitting to the electricity demand data.

# 6    Discussion

In regression with autoregressive errors (REGAR) models, we propose the lasso approach to jointly shrink regression and autoregression coefficients. In contrast to the REGAR model, the autoregression with exogenous variables (ARX) model (Harvey, 1981, and Shumway and Stoffer, 2000) provides an alternative approach to explicitly take into account serial dependency via the lagged variables. Specifically, the ARX model is

$$y_t = x_t'\beta + \sum_{j=1}^{q} \phi_{t-j} y_{t-j} + \epsilon_t.$$

To simultaneously shrink the regression and lagged coefficients, we consider the following lasso criterion:

$$\sum_{t=q+1}^{n_0} \left( y_t - x_t'\beta - \sum_{j=1}^{q} \phi_j y_{t-j} \right)^2 + n \sum_{j=1}^{p} \lambda_j^* |\beta_j| + n \sum_{j=1}^{q} \gamma_j^* |\phi_j|.$$

Analogous to the REGAR model, it can be shown that the lasso approach produces a sparse solution not only for exogenous variables but also for lagged dependent variables. Moreover, the resulting lasso estimator enjoys the *oracle* property when the tuning parameters satisfy the proper conditions. Extensive simulation studies (not presented here) also indicate their satisfactorily finite sample performance.

Finally, we identify three research areas for further study. The first is extending the application of lasso to both the dynamic regression model (Greene, 2003) and the regression model with seasonal autoregressive errors. The second is to obtain the lasso estimator for the regression model with autoregressive conditional heteroscedastic (ARCH) errors (Gouriéroux, 1997) and the autoregressive and moving average with exogenous variables (ARMAX) model (Shumway and Stoffer, 2000). The third is to investigate autoregressive shrinkage and selection by compressing the partial autocorrelations sequentially. We believe that these efforts would further enhance the usefulness of the lasso estimators in real data analysis.

# Appendix

## Appendix A: Proof of Theorem 1

Let $\delta = (u', v')'$, $u = (u_1, \cdots, u_p)'$, and $v = (v_1, \cdots, v_q)'$, and then define

$$
\begin{aligned}
\kappa_n(\delta) &= Q_n(\theta^0 + n^{-1/2}\delta) - Q_n(\theta^0) \\
&= \left\{ L_n(\theta^0 + n^{-1/2}\delta) - L_n(\theta^0) \right\} + n\lambda_n \sum_{j=1}^{p} \left\{ |\beta_j^0 + u_j n^{-1/2}| - |\beta_j^0| \right\} \\
&\quad + n\gamma_n \sum_{j=1}^{q} \left\{ |\phi_j^0 + v_j n^{-1/2}| - |\phi_j^0| \right\}.
\end{aligned}
$$

Adopting Knight and Fu's (2000) approach, we have

$$
n\lambda_n \sum_{j=1}^{p} \left\{ |\beta_j^0 + u_j n^{-1/2}| - |\beta_j^0| \right\} \to \lambda_0 \sum_{j=1}^{p} \left\{ u_j \mathrm{sgn}\{\beta_j^0\} I(\beta_j^0 \neq 0) + |u_j| I(\beta_j^0 = 0) \right\}
$$

$$
n\gamma_n \sum_{j=1}^{q} \left\{ |\phi_j^0 + v_j n^{-1/2}| - |\phi_j^0| \right\} \to \gamma_0 \sum_{j=1}^{q} \left\{ v_j \mathrm{sgn}\{\phi_j^0\} I(\phi_j^0 \neq 0) + |v_j| I(\phi_j^0 = 0) \right\}.
$$

Furthermore,

$$
\begin{aligned}
& L_n(\theta^0 + n^{-1/2}\delta) - L_n(\theta^0) \\
&= \sum_t \left\{ [y_t - x_t'(\beta_0 + n^{-1/2}u)] - \sum_{j=1}^{p}(\phi_j^0 + n^{-1/2}v_j)[y_{t-j} - x_{t-j}'(\beta_0 + n^{-1/2}u)] \right\}^2 - \sum_t \epsilon_t^2 \\
&= \sum_t \left\{ e_t - \sum_{j=1}^{q}(\phi_j^0 + n^{-1/2}v_j)e_{t-j} - n^{-1/2}u'\left[ x_t - \sum_{j=1}^{q}(\phi_j^0 + n^{-1/2}v_j)x_{t-j} \right] \right\}^2 - \sum_t \epsilon_t^2 \\
&= \sum_t \left\{ \epsilon_t - n^{-1/2}\sum_{j=1}^{q}v_j e_{t-j} - n^{-1/2}u'\left( x_t - \sum_{j=1}^{q}\phi_j^0 x_{t-j} \right) + n^{-1}u'\sum_{j=1}^{q}v_j x_{t-j} \right\}^2 - \sum_t \epsilon_t^2 \\
&= R_1 + R_2 + R_3 + R_4 + R_5,
\end{aligned}
$$

where

$$
\begin{aligned}
R_1 &= -2n^{-1/2}\sum_t \left(\epsilon_t \sum_{j=1}^{q} v_j e_{t-j}\right) - 2n^{-1/2}u' \sum_t \left[\epsilon_t \left(x_t - \sum_{j=1}^{q} \phi_j^0 x_{t-j}\right)\right] \\
R_2 &= 2n^{-1}u' \sum_t \left\{\left(\sum_{j=1}^{q} v_j e_{t-j}\right)\left(x_t - \sum_{j=1}^{q} \phi_j^0 x_{t-j}\right)\right\} \\
R_3 &= n^{-1}\sum_t \left(\sum_{j=1}^{q} v_j e_{t-j}\right)^2 + n^{-1}u' \sum_t \left(x_t - \sum_{j=1}^{q} \phi_j^0 x_{t-j}\right)\left(x_t - \sum_{j=1}^{q} \phi_j^0 x_{t-j}\right)' u \\
R_4 &= 2n^{-1}\sum_t \left(u' \sum_{j=1}^{q} v_j x_{t-j}\right)\left\{\epsilon_t - n^{-1/2}\sum_{j=1}^{q} v_j e_{t-j} - n^{-1/2}u' \left(x_t - \sum_{j=1}^{q} \phi_j^0 x_{t-j}\right)\right\} \\
R_5 &= n^{-2}u' \sum_t \left(\sum_{j=1}^{q} v_j x_{t-j}\right)\left(\sum_{j=1}^{q} v_j x_{t-j}\right)' u.
\end{aligned}
$$

Employing the martingale central limit theorem and the ergodic theorem, we are able to show that $R_1 \to_d -2\delta' w$, $R_2 = o_p(1)$, $R_3 \to_p \delta'\Sigma\delta$, $R_4 = o_p(1)$, and $R_5 = o_p(1)$. Consequently,

$$
L_n(\theta^0 + n^{-1/2}\delta) - L_n(\theta^0) \to_d -2\delta' w + \delta'\Sigma\delta.
$$

In order to show that $\operatorname{argmin}\{\kappa_n(\delta)\} \to_d \operatorname{argmin}\{\kappa(\delta)\}$, we have to prove that $\operatorname{argmin}\{\kappa_n(\delta)\} = O_p(1)$. Note that

$$
\begin{aligned}
\kappa_n(\delta) &\geq \sum_t \left\{\left[\epsilon_t - n^{-1/2}\sum_{j=1}^{q} v_j e_{t-j} - n^{-1/2}u' \left(x_t - \sum_{j=1}^{q} \phi_j^0 x_{t-j}\right) + n^{-1}u' \sum_{j=1}^{q} v_j x_{t-j}\right]^2 - \epsilon_t^2\right\} \\
&\quad - n\lambda_n \sum_{j=1}^{p} |u_j n^{-1/2}| - n\gamma_n \sum_{j=1}^{q} |v_j n^{-1/2}| \\
&\geq \sum_t \left\{\left[\epsilon_t - n^{-1/2}\sum_{j=1}^{q} v_j e_{t-j} - n^{-1/2}u' \left(x_t - \sum_{j=1}^{q} \phi_j^0 x_{t-j}\right)\right]^2 - \epsilon_t^2\right\} \\
&\quad - (\lambda_0 + \epsilon_0)\sum_{j=1}^{p} |u_j| - (\gamma_0 + \epsilon_0)\sum_{j=1}^{q} |v_j| + \xi_n(\delta) \doteq \tilde{\kappa}_n(\delta),
\end{aligned}
$$

where $\epsilon_0 > 0$ is some positive constant. In addition, $\kappa_n(0) = \tilde{\kappa}_n(0)$ and $\xi_n(\delta) = o_p(1)$. Moreover, for all $\delta$ and sufficiently large $n$, the quadratic terms in $\tilde{\kappa}_n(\delta)$ grow faster than the $|u_j|$ and $|v_j|$. As a result, $\operatorname{argmin}\{\tilde{\kappa}_n(\delta)\} = O_p(1)$ and $\operatorname{argmin}\{\kappa_n(\delta)\} = O_p(1)$. Because $\operatorname{argmin}\{\kappa(\delta)\}$ is unique with probability 1, the proof is completed.

## Appendix B: Proof of Lemma 1

Let $\alpha_n = n^{-1/2} + a_n$ and $\{\theta^0 + \alpha_n \delta : \|\delta\| \leq d\}$ be the ball around $\theta^0$. Then, for $\|\delta\| = d$, we have

$$
\begin{aligned}
D_n(\delta) &\doteq Q_n^*(\theta^0 + \alpha_n \delta) - Q_n^*(\theta^0) \\
&\geq L_n(\theta^0 + \alpha_n \delta) - L_n(\theta^0) + n \sum_{j \in \mathcal{S}_1} \lambda_j \left( |\beta_j^0 + \alpha_n u_j| - |\beta_j^0| \right) + n \sum_{j \in \mathcal{S}_2} \gamma_j \left( |\phi_j^0 + \alpha_n v_j| - |\phi_j^0| \right) \\
&\geq L_n(\theta^0 + \alpha_n \delta) - L_n(\theta^0) - n\alpha_n \sum_{j \in \mathcal{S}_1} \lambda_j |u_j| - n\alpha_n \sum_{j \in \mathcal{S}_2} \gamma_j |v_j| \\
&\geq L_n(\theta^0 + \alpha_n \delta) - L_n(\theta^0) - n\alpha_n^2 p_0 d - n\alpha_n^2 q_0 d \\
&= L_n(\theta^0 + \alpha_n \delta) - L_n(\theta^0) - n\alpha_n^2 (p_0 + q_0) d. 
\end{aligned} \tag{A.1}
$$

Furthermore,

$$
\begin{aligned}
& L_n(\theta^0 + \alpha_n \delta) - L_n(\theta^0) \\
&= \sum_t \left\{ \epsilon_t - \alpha_n \sum_{j=1}^q v_j e_{t-j} - \alpha_n u' \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right) + \alpha_n^2 u' \sum_{j=1}^q v_j x_{t-j} \right\}^2 - \sum_t \epsilon_t^2 \\
&= A_1 + A_2 + A_3 + A_4 + A_5, 
\end{aligned} \tag{A.2}
$$

where

$$
\begin{aligned}
A_1 &= \alpha_n^2 \sum_t \left\{ \left( \sum_{j=1}^q v_j e_{t-j} \right)^2 + u' \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right) \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right)' u \right\} \\
A_2 &= -2\alpha_n \sum_t \epsilon_t \left[ \sum_{j=1}^q v_j e_{t-j} + u' \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right) \right] \\
A_3 &= 2\alpha_n^2 \sum_t \left( \sum_{j=1}^q v_j e_{t-j} \right) u' \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right) \\
A_4 &= \alpha_n^3 \sum_t \left( u' \sum_{j=1}^q v_j x_{t-j} \right) \left[ \alpha_n u' \sum_{j=1}^q v_j x_{t-j} - 2 u' \left( x_t - \sum_{j=1}^q \phi_j^0 x_{t-j} \right) - 2 \sum_{j=1}^q v_j e_{t-j} \right] \\
A_5 &= 2\alpha_n^2 \sum_t \epsilon_t \left( u' \sum_{j=1}^q v_j x_{t-j} \right).
\end{aligned}
$$

Moreover, we have

$$
\begin{aligned}
A_1 &= n\alpha_n^2 \times \left\{ \delta' \Sigma \delta + o_p(1) \right\} \\
A_2 &= \delta' O_p(n\alpha_n^2) \\
A_3 &= n\alpha_n^2 \times o_p(1) = o_p(n\alpha_n^2)
\end{aligned}
$$

14

$$
\begin{aligned}
A_4 &= n\alpha_n^3 \times O_p(1) = n\alpha_n^2 o_p(1) = o_p(n\alpha_n^2) \\
A_5 &= n\alpha_n^2 \times o_p(1) = o_p(n\alpha_n^2).
\end{aligned}
$$

Because $A_1$ dominates the rest of four terms in (A.2) and also $n\alpha_n^2(p_0 + q_0)d$ in (A.1). Hence, for any given $\epsilon > 0$, there exists a large constant $d$ such that

$$
P\left\{\inf_{\|\delta\|=d} Q_n^*(\theta^0 + \alpha_n\delta) > Q_n^*(\theta^0)\right\} \geq 1 - \epsilon.
$$

This implies that with probability at least $1 - \epsilon$, there exists a local minimizer in the ball $\{\theta^0 + \alpha_n\delta : \|\delta\| \leq d\}$ (Fan and Li, 2001). Consequently, there exists a local minimizer of $Q_n^*(\theta)$ such that $\|\hat{\theta}^* - \theta^0\| = O_p(\alpha_n)$. This completes the proof.

## Appendix C: Proof of Theorem 2

It follows from the fact that the local minimizer $\hat{\theta}^*$ must satisfy the following equation,

$$
\frac{\partial Q_n^*(\hat{\theta}^*)}{\partial\beta_j} = \frac{\partial L_n(\hat{\theta}^*)}{\partial\beta_j} - n\lambda_j \mathrm{sgn}(\hat{\beta}_j^*)
$$

$$
= \frac{\partial L_n(\theta^0)}{\partial\beta_j} + n\Sigma_j(\hat{\theta}^* - \theta^0)\{1 + o_p(1)\} - n\lambda_j \mathrm{sgn}(\hat{\beta}_j^*), \tag{A.3}
$$

where $\Sigma_j$ denotes the $j$'th row of $\Sigma$ and $j \in \mathcal{S}_1^c$. Employing the central limit theorem, the first term in (A.3) is of the order $O_p(n^{1/2})$. Furthermore, the condition in Theorem 2 implies that its second term is also of the order $O_p(n^{1/2})$. Both of them are dominated by $n\lambda_j$ since $\sqrt{n}b_n \to \infty$. Therefore, the sign of (A.3) is dominated by the sign of $\hat{\beta}_j^*$. Consequently, we must have $\hat{\beta}_j^* = 0$ in probability. Analogously, we can show that $P(\hat{\phi}_{\mathcal{S}_2^c}^* = 0) \to 1$. This completes the proof.

## Appendix D: Proof of Theorem 3

Applying Lemma 1 and Theorem 2, we have $P(\hat{\theta}_2^* = 0) \to 1$. Hence, the minimizer of $Q_n^*(\theta)$ is the same as that of $Q_n^*(\theta_1)$ with probability tending to one. This implies that the lasso estimator, $\hat{\theta}_1^*$, satisfies the following equation

$$
\left.\frac{\partial Q_n^*(\theta_1)}{\partial\theta_1}\right|_{\theta_1 = \hat{\theta}_1^*} = 0. \tag{A.4}
$$

According to Lemma 1, $\hat{\theta}_1^*$ is a $\sqrt{n}$-consistent estimator. Thus, the Taylor's expansion of (A.4) yields

$$
\begin{aligned}
0 &= \frac{1}{\sqrt{n}} \times \frac{\partial L_n(\hat{\theta}_1^*)}{\partial\theta_1} + \sqrt{n}P(\hat{\theta}_1^*) = \frac{1}{\sqrt{n}} \times \frac{\partial L_n(\theta_1^0)}{\partial\theta_1} + \sqrt{n}P(\theta_1^0) \\
&\quad + \Sigma_0\sqrt{n}(\hat{\theta}_1^* - \theta_1^0) + o_p(1),
\end{aligned}
$$

where $P$ is the first order derivative of the penalty function

$$\sum_{j \in \mathcal{S}_1} \lambda_j |\beta_j| + \sum_{j \in \mathcal{S}_2} \gamma_j |\phi_j|,$$

and $P(\hat{\theta}_1^*) = P(\theta_1^0)$ as $n$ large enough. Furthermore, it can be easily shown that $\sqrt{n}P(\theta_1^0) = o_p(1)$, which implies that

$$\sqrt{n}(\hat{\theta}_1^* - \theta_1^0) = \frac{\Sigma_0^{-1}}{\sqrt{n}} \times \frac{\partial L_n(\theta_1^0)}{\partial \theta_1} + o_p(1) \xrightarrow{d} N(0, \sigma^2 \Sigma_0^{-1}).$$

This completes the proof.

## Acknowledgments

## References

Akaike (1973) Information theory and an extension of the maximum likelihood principle. *In 2nd International Symposium on Information Theory, Ed. B. N. Petrov & F. Csaki*, 267–281. Budapest: Akademia Kiado.

Box, G. E. P., Jenkins, G. M. and Reinsel, G. C. (1994) *Time Series Analysis - Forecasting and Control*. New York: Prentice Hall, 3 edn.

Breiman, L. (1996) Heuristics of instability and stabilization in model selection. *Ann. Statist.*, **24**, 2350–2383.

Brockwell, P. J. and Davis, R. A. (1991) *Time Series: Theory and Methods*. New York: Springer, 2 edn.

Cai, J., Fan, J., Li, R. and Zhou, H. (2005) Model selection for multivariate failure time data. *Biometrika*, **92**, 303–316.

Choi, B. S. (1992) *ARMA Model Identification*. New York: Springer.

Cochrane, D. and Orcutt, G. H. (1949) Application of least squares regression to relationships containing autocorrelated error terms. *J. Amer. Statist. Assoc.*, **44**, 32–61.

Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004) Least angle regression. *Ann. Statist.*, **32**, 407–489.

Fan, J. and Li, R. (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.*, **96**, 1348–1360.

— (2002) Variable selection for cox's proportional hazards model and frailty model. *Ann. Statist.*, **30**, 74–99.

Fan, J. and Peng, H. (2004) On non-concave penalized likelihood with diverging number of parameters. *Ann. Statist.*, **32**, 928–961.

Fu, W. J. (1998) Penalized regression: the bridge versus the lasso. *J. Comput. and Graph. Statist.*, **7**, 397–416.

Gouriéroux, C. (1997) *ARCH Models and Financial Applications.* New York: Springer.

Greene, W. H. (2003) *Econometric Analysis (Fifth Edition).* Prentice Hall.

Hamilton, J. D. (1994) *Time Series Analysis.* Princeton University Press, Princeton, NJ.

Harvey, A. C. (1981) *The Econometric Analysis of Time Series.* Wiley, New York, NY.

Hurvich, C. M. and Tsai, C. L. (1990) The impact of model selection on inference in linear regression. *Amer. Statist.*, **44**, 214–217.

Knight, K. and Fu, W. (2000) Asymptotics for lasso-type estimators. *Ann. Statist.*, **28**, 1356–1378.

Leng, C., Lin, Y. and Wahba, G. (2005) A note on lasso and related procedures in model selection. *Statistica Sinica, To appear.*

McQuarrie, D. R. and Tsai, C. L. (1998) *Regression and Time Series Model Selection.* World Scientific, Singapore.

Ramanathan, R. (1989) *Introductory Econometrics with Applications.* Harcourt Brace Jovanovich, Orlando.

Schwarz, G. (1978) Estimating the dimension of a model. *Ann. Statist.*, **6**, 461–464.

Shao, J. (1997) An asymptotic theory for linear model selection (with discussion). *Statistica Sinica*, **7**, 221–264.

Shi, P. and Tsai, C. L. (2004) A joint regression variable and autoregressive order selection criterion. *J. Time Series*, **25**, 923–941.

Shumway, R. H. and Stoffer, D. S. (2000) *Time Series Analysis and Its Application.* New York: Springer.

Tibshirani, R. J. (1996) Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B*, **58**, 267–288.

Tsay, R. S. (1984) Regression models with time series errors. *J. Amer. Statist. Assoc.*, **79**, 118–124.

Zou, H. and Hastie, T. (2005) Regression shrinkage and selection via the elastic net with application to microarrays. *J. Roy. Statist. Soc. Ser. B*, **67**, 301–320.

Zou, H., Hastie, T. and Tibshirani, R. (2004) On the 'degrees of freedom' of lasso. *Technical Report, Statistics Department, Standford University.*

Table 1: Simulation Results with $\rho = 0.50$

| Estimator | Tuning Method | Regression Variable | | | Autoregressive Order | | | Correctly fitted model | Median of MSPE |
|---|---|---|---|---|---|---|---|---|---|
| | | under fitted | correctly fitted | over fitted | under fitted | correctly fitted | over fitted | | |
| $\sigma = 3.0, n = 50$ | | | | | | | | | |
| LASSO | CV | 0.019 | 0.174 | 0.807 | 0.043 | 0.142 | 0.815 | 0.026 | 17.430 |
| | BIC | 0.017 | 0.223 | 0.760 | 0.018 | 0.206 | 0.776 | 0.045 | 16.715 |
| LASSO* | CV | 0.078 | 0.412 | 0.510 | 0.063 | 0.585 | 0.352 | 0.245 | 16.228 |
| | BIC | 0.101 | 0.578 | 0.321 | 0.074 | 0.752 | 0.174 | 0.455 | 15.382 |
| $\sigma = 3.0, n = 100$ | | | | | | | | | |
| LASSO | CV | 0.001 | 0.235 | 0.764 | 0.001 | 0.126 | 0.873 | 0.020 | 14.713 |
| | BIC | 0.001 | 0.367 | 0.632 | 0.000 | 0.176 | 0.824 | 0.054 | 14.392 |
| LASSO* | CV | 0.003 | 0.572 | 0.425 | 0.002 | 0.654 | 0.344 | 0.376 | 13.826 |
| | BIC | 0.003 | 0.852 | 0.145 | 0.003 | 0.932 | 0.065 | 0.796 | 13.504 |
| $\sigma = 3.0, n = 300$ | | | | | | | | | |
| LASSO | CV | 0.000 | 0.144 | 0.856 | 0.000 | 0.133 | 0.867 | 0.011 | 13.194 |
| | BIC | 0.000 | 0.167 | 0.833 | 0.000 | 0.233 | 0.767 | 0.035 | 13.111 |
| LASSO* | CV | 0.000 | 0.683 | 0.317 | 0.000 | 0.678 | 0.322 | 0.449 | 12.900 |
| | BIC | 0.000 | 0.946 | 0.054 | 0.000 | 0.971 | 0.029 | 0.919 | 12.862 |
| $\sigma = 0.5, n = 50$ | | | | | | | | | |
| LASSO | CV | 0.000 | 0.174 | 0.826 | 0.047 | 0.138 | 0.815 | 0.026 | 1.530 |
| | BIC | 0.000 | 0.228 | 0.772 | 0.017 | 0.207 | 0.776 | 0.045 | 1.461 |
| LASSO* | CV | 0.000 | 0.566 | 0.434 | 0.056 | 0.579 | 0.365 | 0.340 | 1.320 |
| | BIC | 0.000 | 0.802 | 0.198 | 0.071 | 0.758 | 0.171 | 0.636 | 1.275 |
| $\sigma = 0.5, n = 100$ | | | | | | | | | |
| LASSO | CV | 0.000 | 0.234 | 0.766 | 0.001 | 0.126 | 0.873 | 0.020 | 1.289 |
| | BIC | 0.000 | 0.370 | 0.630 | 0.000 | 0.176 | 0.824 | 0.056 | 1.260 |
| LASSO* | CV | 0.000 | 0.623 | 0.377 | 0.002 | 0.650 | 0.348 | 0.416 | 1.189 |
| | BIC | 0.000 | 0.941 | 0.059 | 0.003 | 0.930 | 0.067 | 0.877 | 1.165 |
| $\sigma = 0.5, n = 300$ | | | | | | | | | |
| LASSO | CV | 0.000 | 0.144 | 0.856 | 0.000 | 0.133 | 0.867 | 0.011 | 1.156 |
| | BIC | 0.000 | 0.168 | 0.832 | 0.000 | 0.233 | 0.767 | 0.035 | 1.148 |
| LASSO* | CV | 0.000 | 0.685 | 0.315 | 0.000 | 0.675 | 0.325 | 0.452 | 1.132 |
| | BIC | 0.000 | 0.969 | 0.031 | 0.000 | 0.972 | 0.028 | 0.943 | 1.124 |

Table 2: Three Models Selected by lasso and lasso* for the Electricity Demand Study

| Variable | LASSO | | LASSO* | |
|---|---|---|---|---|
| | CV | BIC | CV | BIC |
| LY | 0.117796 | 0.196611 | - | - |
| LELP | -0.150010 | -0.154907 | -0.168853 | -0.168853 |
| LGSP | -0.035808 | -0.057948 | - | - |
| CDD | 0.000237 | 0.000246 | 0.000226 | 0.000226 |
| HDD | 0.000216 | 0.000231 | 0.000222 | 0.000222 |
| LAG1 | 0.608868 | 0.598635 | 0.627451 | 0.627451 |
| LAG2 | -0.705199 | -0.689985 | -0.713343 | -0.713343 |
| LAG3 | 0.590666 | 0.581069 | 0.604899 | 0.604899 |
| LAG4 | 0.253515 | 0.271175 | 0.225005 | 0.225005 |
| $\chi^2$-Test | 0.039428 | 0.005632 | 0.305896 | 0.305896 |