



Penalized composite quasi-likelihood for ultrahigh dimensional variable selection

Jelena Bradic and Jianqing Fan

Princeton University, USA

and Weiwei Wang

University of Texas Health Science Center, Houston, USA

[Received December 2009. Revised October 2010]

Summary. In high dimensional model selection problems, penalized least square approaches have been extensively used. The paper addresses the question of both robustness and efficiency of penalized model selection methods and proposes a data-driven weighted linear combination of convex loss functions, together with weighted L_1 -penalty. It is completely data adaptive and does not require prior knowledge of the error distribution. The weighted L_1 -penalty is used both to ensure the convexity of the penalty term and to ameliorate the bias that is caused by the L_1 -penalty. In the setting with dimensionality much larger than the sample size, we establish a strong oracle property of the method proposed that has both the model selection consistency and estimation efficiency for the true non-zero coefficients. As specific examples, we introduce a robust method of composite L_1 – L_2 , and an optimal composite quantile method and evaluate their performance in both simulated and real data examples.

Keywords: Composite quasi-maximum likelihood estimation; Lasso; Model selection; Non-polynomial dimensionality; Oracle property; Robust statistics; Smoothly clipped absolute deviation

1. Introduction

Feature extraction and model selection are important for sparse high dimensional data analysis in many research areas such as genomics, genetics and machine learning. Motivated by the need for a robust and efficient high dimensional model selection method, we introduce a new penalized quasi-likelihood estimator for linear models with high dimensional parameter space.

Consider the estimation of the unknown parameter β in the linear regression model

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

where $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ is an n -vector of response, $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^T$ is an $n \times p$ matrix of independent variables with \mathbf{X}_i^T being its i th row, $\beta = (\beta_1, \dots, \beta_p)^T$ is a p -vector of unknown parameters and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$ is an n -vector of independent identically distributed random errors with mean 0, independent of \mathbf{X} . When the dimension p is high it is commonly assumed that only a small number of predictors actually contribute to the response vector \mathbf{Y} , which leads to the sparsity pattern in the unknown parameters and thus makes variable selection crucial. In many applications such as genetic association studies and disease classifications using high throughput data such as microarrays with gene–gene interactions, the number of

Address for correspondence: Jelena Bradic, Department of Operations Research and Financial Engineering, Sherrerd Hall, Princeton University, Princeton, NJ 08544, USA.
E-mail: jbradic@princeton.edu

variables p can be much larger than the sample size n . We shall refer to such problems as ultra-high dimensional problems and model them by assuming that $\log(p) = O(n^\delta)$ for some $\delta \in (0, 1)$. Following Fan and Lv (2010), we shall refer to p as a non-polynomial order or non-polynomial dimensionality.

Popular approaches such as the lasso (Tibshirani, 1996), smoothly clipped absolute deviation (SCAD) (Fan and Li, 2001), the adaptive lasso (Zou, 2006) and the elastic net (Zou and Zhang, 2009) use penalized least square regression:

$$\hat{\beta} = \arg \min_{\beta} \left\{ \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \beta)^2 + n \sum_{j=1}^p p_{\lambda}(|\beta_j|) \right\}. \quad (2)$$

where $p_{\lambda}(\cdot)$ is a specific penalty function. Quadratic loss is popular for its mathematical beauty but is not robust to non-normal errors and outliers. Robust regressions such as least absolute deviation and quantile regressions have recently been used in variable selection techniques when p is finite (Wu and Liu, 2009; Zou and Yuan, 2008; Li and Zhu, 2008). Other possible choices of robust loss functions include Huber's loss (Huber, 1964), Tukey's bisquare and Hampel's ψ , among others. Each of these loss functions performs well under a certain class of error distributions: quadratic loss is suitable for normal distributions, least absolute deviation is suitable for heavy-tail distributions and is the most efficient for double-exponential (DE) distributions and Huber's loss performs well for contaminated normal distributions. However, none of them is universally better than all the others. How to construct an adaptive loss function that is applicable to a large collection of error distributions is the question.

We propose a simple and yet effective quasi-likelihood function, which replaces quadratic loss by a weighted linear combination of convex loss functions:

$$\rho_{\mathbf{w}} = \sum_{k=1}^K w_k \rho_k, \quad (3)$$

where ρ_1, \dots, ρ_K are convex loss functions and w_1, \dots, w_K are positive constants chosen to minimize the asymptotic variance of the resulting estimator. From the point of view of non-parametric statistics, the functions $\{\rho_1, \dots, \rho_K\}$ can be viewed as a set of basis functions, which are not necessarily orthogonal, used to approximate the unknown log-likelihood function of the error distribution. When the set of loss functions is large, the quasi-likelihood function can well approximate the log-likelihood function and therefore yield a nearly efficient method. This kind of idea has appeared already in traditional statistical inference with finite dimensionality (Koenker, 1984; Bai *et al.*, 1992). We shall extend it to sparse statistical inference with non-polynomial dimensionality.

The quasi-likelihood function (3) can be directly used together with any penalty function such as the L_p -penalty with $0 < p < 1$ (Frank and Friedman, 1993), the lasso, i.e. the L_1 -penalty (Tibshirani, 1996), SCAD (Fan and Li, 2001) or the hierarchical penalty (Bickel *et al.*, 2008), resulting in the penalized composite quasi-likelihood problem:

$$\min_{\beta} \left\{ \sum_{i=1}^n \rho_{\mathbf{w}}(Y_i - \mathbf{X}_i^T \beta) + n \sum_{j=1}^p p_{\lambda}(|\beta_j|) \right\}. \quad (4)$$

Instead of using folded concave penalty functions, we use the weighted L_1 -penalty of the form

$$n \sum_{j=1}^p \gamma_{\lambda}(|\beta_j^{(0)}|) |\beta_j|$$

for some function γ_{λ} and initial estimator $\beta^{(0)}$, to ameliorate the bias in L_1 -penalization (Fan and Li, 2001; Zou, 2006; Fan and Lv, 2010) and to maintain the convexity of the problem. This leads to the following convex optimization problem:

$$\hat{\beta}_{\mathbf{w}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \rho_{\mathbf{w}}(Y_i - \mathbf{X}_i^T \beta) + n \sum_{j=1}^p \gamma_{\lambda}(|\beta_j^{(0)}|) |\beta_j| \right\}. \quad (5)$$

When $\gamma_{\lambda}(\cdot) = p'_{\lambda}(\cdot)$, the derivative of the penalty function, problem (5) can be regarded as the local linear approximation to problem (4) (Zou and Li, 2008). In particular, the lasso (Tibshirani, 1996) corresponds to $\gamma_{\lambda}(x) = \lambda$, SCAD reduces to (Fan and Li, 2001)

$$\gamma_{\lambda}(x) = \lambda \left\{ I(x \leq \lambda) + \frac{(a\lambda - x)_+}{(a-1)\lambda} I(x > \lambda) \right\} \quad (6)$$

and the adaptive lasso (Zou, 2006) takes the form $\gamma_{\lambda}(x) = \lambda |x|^{-a}$ where $a > 0$.

There is a rich literature in establishing the oracle property for penalized regression methods, mostly for large but fixed p (Fan and Li, 2001; Zou, 2006; Yuan and Lin, 2007; Zou and Yuan, 2008). One of the early references on diverging p is Fan and Peng (2004) under conditions of $p = O(n^{1/5})$. More recent works of a similar kind include Huang *et al.* (2008), Zou and Zhang (2009) and Xie and Huang (2009), which assume that the number of non-sparse elements s is finite. When the dimensionality p is of polynomial order, Kim *et al.* (2008) recently gave the conditions under which the SCAD estimator is an oracle estimator. We would like to address this problem further when $\log(p) = O(n^{\delta})$ with $\delta \in (0, 1)$ and $s = O(n^{\alpha_0})$ for $\alpha_0 \in (0, 1)$, i.e. when the dimensionality is of exponential order.

The paper is organized as follows. Section 2 introduces an easy-to-implement two-step computation procedure. Section 3 proves the strong oracle property of the weighted L_1 -penalized quasi-likelihood approach with discussion on the choice of weights and corrections for convexity. Section 4 defines two specific instances of the approach proposed and compares their asymptotic efficiencies. Section 5 provides a comprehensive simulation study as well as a real data example of single-nucleotide polymorphism (SNP) selection for Down syndrome. Section 6 is devoted to discussion. To facilitate the readability, all the proofs are relegated to Appendices A–C.

2. Penalized adaptive composite quasi-likelihood

We would like to describe the two-step adaptive computation procedure proposed and to defer the justification of the appropriate choice of the weight vector \mathbf{w} to Section 3.

In the first step, we obtain the initial estimate $\hat{\beta}^{(0)}$ by using the lasso procedure, i.e.

$$\hat{\beta}^{(0)} = \arg \min_{\beta} \left\{ \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \beta)^2 + n \lambda \sum_{j=1}^p |\beta_j| \right\},$$

and estimate the residual vector $\varepsilon^0 = \mathbf{Y} - \mathbf{X} \hat{\beta}^{(0)}$ (for justification see the discussion following condition 2 in Appendix A). The matrix \mathbf{M} and vector \mathbf{a} are calculated as follows:

$$\begin{aligned} \mathbf{M}_{kl} &= \frac{1}{n} \sum_{i=1}^n \psi_k(\varepsilon_i^0) \psi_l(\varepsilon_i^0), \\ a_k &= \frac{1}{n} \sum_{i=1}^n \partial \psi_k(\varepsilon_i^0), \quad k, l = 1, \dots, K, \end{aligned}$$

where $\psi_k(t)$ is a choice of the subgradient of $\rho_k(t)$, ε_i^0 is the i th component of ε^0 and a_k should be considered as a consistent estimator of $E\{\partial \psi_k(\varepsilon)\}$, which is the derivative of $E\{\psi_k(\varepsilon + c)\}$ at $c = 0$. For example, when $\psi_k(x) = \text{sgn}(x)$, then $E\{\psi_k(\varepsilon + c)\} = 1 - 2F_{\varepsilon}(-c)$ and $E\{\partial \psi_k(\varepsilon)\} = 2f_{\varepsilon}(0)$. The optimal weight is then determined as

$$\mathbf{w}_{\text{opt}} = \arg \min_{\mathbf{w} \geq 0, \mathbf{a}^T \mathbf{w} = 1} (\mathbf{w}^T \mathbf{M} \mathbf{w}). \quad (7)$$

In the second step, we calculate the quasi-maximum-likelihood estimator by using weights \mathbf{w}_{opt} as

$$\hat{\beta}^{\mathbf{a}} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \rho_{\mathbf{w}_{\text{opt}}}(Y_i - \mathbf{X}_i^T \beta) + n \sum_{j=1}^p \gamma_{\lambda}(|\hat{\beta}_j^{(0)}|) |\beta_j| \right\}. \quad (8)$$

Remark 1. Zero is not an absorbing state in the minimization problem (8). Those elements that are estimated as 0 in the initial estimate $\beta^{(0)}$ have a chance to escape from zero, whereas those non-vanishing elements can be estimated as 0 in problem (8).

Remark 2. The number of loss functions K is typically small or moderate in practice. Problem (7) can be easily solved by using a quadratic programming algorithm. The resulting vector \mathbf{w}_{opt} can have vanishing components, automatically eliminating inefficient loss functions in the second step (8) and hence learning the best approximation of the unknown log-likelihood function. This can lead to considerable computational gains. See Section 4 for additional details.

Remark 3. Problem (8) is a convex optimization problem when the ρ_k s are all convex and $\gamma_{\lambda}(|\hat{\beta}_j^{(0)}|)$ are all non-negative. This class of problems can be solved with fast and efficient computational algorithms such as pathwise co-ordinate optimization (Friedman *et al.*, 2008) and least angle regression (Efron *et al.*, 2004).

One particular example is the combination of L_1 - and L_2 -regressions, in which $K = 2$, $\rho_1(t) = |t - b_0|$ and $\rho_2(t) = t^2$. Here b_0 denotes the median of error distribution ε . If the error distribution is symmetric, then $b_0 = 0$. If the error distribution is completely unknown, b_0 is unknown and can be estimated from the residual vector $\{\varepsilon_i^0\}$ or be regarded as an additional parameter and optimized together with β in problem (8). Another example is the combination of multiple quantile check functions, i.e.

$$\rho_k(t) = \tau_k(t - b_k)_+ + (1 - \tau_k)(t - b_k)_-,$$

where $\tau_k \in (0, 1)$ is a preselected quantile and b_k is the τ_k -quantile of the error distribution. Again, when b_k s are unknown, they can be estimated by using the sample quantiles τ_k of the estimated residuals ε^0 or along with β in problem (8). See Section 4 for additional discussion.

3. Sampling properties and their applications

In this section, we plan to establish the sampling properties of estimator (5) under the assumption that the number of parameters (true dimensionality) p and the number of non-vanishing components (effective dimensionality) $s = \|\beta^*\|_0$ satisfy $\log(p) = O(n^\delta)$ and $s = O(n^{\alpha_0})$ for some $\delta \in (0, 1)$ and $\alpha_0 \in (0, 1)$. Particular focus will be given to the oracle property of Fan and Li (2001), but we shall strengthen it and prove that estimator (5) is an oracle estimator with overwhelming probability. Fan and Lv (2010) were among the first to discuss the oracle properties with non-polynomial dimensionality by using the full likelihood function in generalized linear models with a class of folded concave penalties. We work on a quasi-likelihood function and a class of weighted convex penalties.

3.1. Asymptotic properties

To facilitate presentation, we relegate technical conditions and the details of proofs to the appendices. We consider more generally the weighted L_1 -penalized estimator with non-negative weights d_1, \dots, d_p . Let

$$L_n(\beta) = \sum_{i=1}^n \rho_{\mathbf{w}}(Y_i - \mathbf{X}_i^T \beta) + n \lambda_n \sum_{j=1}^p d_j |\beta_j| \quad (9)$$

denote the penalized quasi-likelihood function. Estimator (5) is a particular case of equation (9) and corresponds to the case with $d_j = \gamma_\lambda(|\beta_j^{(0)}|)/\lambda_n$.

Without loss of generality, assume that parameter β^* can be arranged in the form of $\beta^* = (\beta_1^{*T}, \mathbf{0}^T)^T$, with $\beta_1^* \in R^s$ a vector of non-vanishing elements of β^* . Let us call $\hat{\beta}^o = (\hat{\beta}_1^{oT}, \mathbf{0}^T)^T \in R^p$ the biased oracle estimator, where $\hat{\beta}_1^o$ is the minimizer of $L_n(\beta_1, \mathbf{0})$ in R^s and $\mathbf{0}$ is the vector of all 0s in R^{p-s} . Here, we suppress the dependence of $\hat{\beta}^o$ on \mathbf{w} and $\mathbf{d} = (d_1, \dots, d_p)^T$. The estimator $\hat{\beta}^o$ is called the biased oracle estimator, since the oracle knows the true submodel $\mathcal{M}_* = \{j: \beta_j^* \neq 0\}$, but nevertheless applies a penalized method to estimate the non-vanishing regression coefficients. The bias becomes negligible when the weights in the first part are 0 or uniformly small (see theorem 2). When the design matrix \mathbf{S} is non-degenerate, the function $L_n(\beta_1, \mathbf{0})$ is strictly convex and the biased oracle estimator is unique, where \mathbf{S} is a submatrix of \mathbf{X} such that $\mathbf{X} = (\mathbf{S}, \mathbf{Q})$ with \mathbf{S} and \mathbf{Q} being $n \times s$ and $n \times (p-s)$ submatrices of \mathbf{X} respectively.

The following theorem shows that $\hat{\beta}^o$ is the unique minimizer of $L_n(\beta)$ on the whole space R^p with an overwhelming probability. As a consequence, $\hat{\beta}_{\mathbf{w}}$ becomes the biased oracle. We establish the following theorem under conditions on the non-stochastic vector \mathbf{d} (see condition 2). It is also applicable to stochastic penalty weights as in problem (8); see the remark following condition 2 in Appendix A.

Theorem 1. Under conditions 1–4 in Appendix A, the estimators $\hat{\beta}^o$ and $\hat{\beta}_{\mathbf{w}}$ exist and are unique on a set with probability tending to 1. Furthermore,

$$P(\hat{\beta}_{\mathbf{w}} = \hat{\beta}^o) \geq 1 - (p-s) \exp(-cn^{(\alpha_0 - 2\alpha_1)_+ + 2\alpha_2})$$

for a positive constant c .

For theorem 1 to be non-trivial, we need to impose the dimensionality restriction $\delta < (\alpha_0 - 2\alpha_1)_+ + 2\alpha_2$, where α_1 controls the rate of growth of the correlation coefficients between the matrices \mathbf{S} and \mathbf{Q} , the important predictors and unimportant predictors (see condition 5) and $\alpha_2 \in [0, \frac{1}{2})$ is a non-negative constant, related to the maximum absolute value of the design matrix (see condition 4). It can be taken as 0 and is introduced to deal with the situation where $(\alpha_0 - 2\alpha_1)_+$ is small or 0 so that the result is trivial. The larger α_2 is, the more stringent restriction is imposed on the choice of λ_n . When the above conditions hold, the penalized composite quasi-likelihood estimator $\hat{\beta}_{\mathbf{w}}$ is equal to the biased oracle estimator $\hat{\beta}^o$, with probability tending to 1 exponentially fast.

Remark 4. The result of theorem 1 is stronger than the oracle property that was defined in Fan and Li (2001) once the properties of $\hat{\beta}^o$ have been established (see theorem 2). It was formulated by Kim *et al.* (2008) for the SCAD estimator with polynomial dimensionality p . It implies not only the model selection consistency but also sign consistency (Zhao and Yu, 2006; Bickel *et al.*, 2008, 2009):

$$P\{\text{sgn}(\hat{\beta}_{\mathbf{w}}) = \text{sgn}(\beta^*)\} = P\{\text{sgn}(\hat{\beta}^o) = \text{sgn}(\beta^*)\} \rightarrow 1.$$

In this way, the result of theorem 1 nicely unifies the two approaches in discussing the oracle property in high dimensional spaces.

Let $\hat{\beta}_{\mathbf{w}1}$ and $\hat{\beta}_{\mathbf{w}2}$ be the first s components and the remaining $p-s$ components of $\hat{\beta}_{\mathbf{w}}$ respectively. According to theorem 1, we have $\hat{\beta}_{\mathbf{w}2} = \mathbf{0}$ with probability tending to 1. Hence, we only need to establish the properties of $\hat{\beta}_{\mathbf{w}1}$.

Theorem 2. Under conditions 1–5 in Appendix A, the asymptotic bias of non-vanishing component $\hat{\beta}_{\mathbf{w}1}$ is controlled by $D_n = \max\{d_j : j \in \mathcal{M}_*\}$ with

$$\|\hat{\beta}_{\mathbf{w}1} - \beta_1^*\|_2 = O_P\{(\lambda_n D_n + n^{-1/2})\sqrt{s}\}.$$

Furthermore, when $0 \leq \alpha_0 < \frac{2}{3}$, $\hat{\beta}_{\mathbf{w}1}$ has asymptotic normality:

$$\mathbf{b}^T (\mathbf{S}^T \mathbf{S})^{1/2} (\hat{\beta}_{\mathbf{w}1} - \beta_1^*) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma_{\mathbf{w}}^2) \quad (10)$$

where \mathbf{b} is a unit vector in \mathbb{R}^s and

$$\sigma_{\mathbf{w}}^2 = \sum_{k,l=1}^K w_k w_l E\{\psi_k(\varepsilon) \psi_l(\varepsilon)\} / \left[\sum_{k=1}^K w_k E\{\partial \psi_k(\varepsilon)\} \right]^2. \quad (11)$$

Since the dimensionality s depends on n , the asymptotic normality of $\hat{\beta}_{\mathbf{w}1}$ is not well defined in the conventional probability sense. The arbitrary linear combination $\mathbf{b}^T \hat{\beta}_{\mathbf{w}1}$ is used to overcome the technical difficulty. In particular, any finite component of $\hat{\beta}_{\mathbf{w}1}$ is asymptotically normal. The result in theorem 2 is also equivalent to the asymptotic normality of the linear combination $\mathbf{B}^T \hat{\beta}_{\mathbf{w}1}$ that was stated in Fan and Peng (2004), where \mathbf{B} is a $q \times s$ matrix, for any given finite number q .

Theorem 2 relates to the results of Portnoy (1985) in a classical setting (corresponding to $p = s$) where he established asymptotic normality of M -estimators when the dimensionality is not higher than $o(n^{2/3})$.

3.2. Covariance estimation

The asymptotic normality (10) allows us to do statistical inference for non-vanishing components. This requires an estimate of the asymptotic covariance matrix of $\hat{\beta}_{\mathbf{w}1}$. Let $\hat{\varepsilon} = \mathbf{Y} - \mathbf{S}^T \hat{\beta}_{\mathbf{w}1}$ be the residual and $\hat{\varepsilon}_i$ be its i th component. A simple substitution estimator of $\sigma_{\mathbf{w}}^2$ is

$$\hat{\sigma}_{\mathbf{w}}^2 = n \sum_{k,l=1}^K w_k w_l \sum_{i=1}^n \psi_k(\hat{\varepsilon}_i) \psi_l(\hat{\varepsilon}_i) / \left\{ \sum_{k=1}^K w_k \sum_{i=1}^n \partial \psi_k(\hat{\varepsilon}_i) \right\}^2.$$

See also the remark before equation (7). Consequently, by result (10), the asymptotic variance–covariance matrix of $\hat{\beta}_{\mathbf{w}1}$ is given by

$$\hat{\sigma}_{\mathbf{w}}^2 (\mathbf{S}^T \mathbf{S})^{-1}. \quad (12)$$

Another possible estimator of the variance and covariance matrix is to apply the standard sandwich formula. In Section 5, through simulation studies, we show that this formula has good properties for p both smaller and larger than n (see Tables 3 and 4 and comments at the end of Section 5.1).

3.3. Choice of weights

Only the factor $\sigma_{\mathbf{w}}^2$ in equation (11) depends on the choice of \mathbf{w} and it is invariant to the scaling of \mathbf{w} . Thus, the optimal choice of weights for maximizing efficiency of the estimator $\hat{\beta}_{\mathbf{w}1}$ is

$$\mathbf{w}_{\text{opt}} = \arg \min_{\mathbf{w}} (\mathbf{w}^T \mathbf{M} \mathbf{w}) \quad \text{subject to } \mathbf{a}^T \mathbf{w} = 1, \mathbf{w} \geq 0, \quad (13)$$

where \mathbf{M} and \mathbf{a} are defined in Section 2 by using an initial estimator, independent of the weighting scheme \mathbf{w} .

Remark 5. The quadratic optimization problem (13) does not have a closed form solution but can easily be solved numerically for a moderate K . The above gain in efficiency, over least squares, could be better understood from the likelihood point of view. Let $f(t)$ denote the unknown error density. The most efficient loss function is the unknown log-likelihood function, $-\log\{f(t)\}$. But, since we have no knowledge of it, the set \mathcal{F}_K , consisting of convex combinations of $\{\rho_k(\cdot)\}_{k=1}^K$ given in equation (3), could be viewed as a collection of basis functions used to approximate it. The broader the set \mathcal{F}_K is, the better it can approximate the log-likelihood function and the more efficient the estimator $\hat{\beta}^a$ in problem (8) becomes. Therefore, we refer to $\rho_{\mathbf{w}}$ as the quasi-likelihood function.

3.4. One-step penalized estimator

The restriction of $\mathbf{w} \geq 0$ guarantees the convexity of $\rho_{\mathbf{w}}$ so that problem (5) becomes a convex optimization problem. However, this restriction may cause substantial loss of efficiency in estimating $\hat{\beta}_{\mathbf{w}1}$ (Table 1). We propose a one-step penalized estimator to overcome this drawback while avoiding non-convex optimization. Let $\hat{\beta}$ be the estimator based on the convex combination of loss functions (5) and $\hat{\beta}_1$ be its non-vanishing components. The one-step estimator is defined as

$$\begin{aligned}\hat{\beta}_{\mathbf{w}1}^{\text{os}} &= \hat{\beta}_1 - \Omega_{n,\mathbf{w}}(\hat{\beta}_1)^{-1} \Phi_{n,\mathbf{w}}(\hat{\beta}_1), \\ \hat{\beta}_{\mathbf{w}2}^{\text{os}} &= \mathbf{0},\end{aligned}\tag{14}$$

where

$$\Phi_{n,\mathbf{w}}(\hat{\beta}_1) = \sum_{i=1}^n \psi_{\mathbf{w}}(Y_i - \mathbf{S}_i^T \hat{\beta}_1) \mathbf{S}_i,$$

Table 1. Asymptotic relative efficiency compared with maximum likelihood estimation

Method [†]		Results for the following distributions $f(\varepsilon)$:						
		DE	t_4	$\mathcal{N}(0, 1)$	$\Gamma(3, 1)$	$\mathcal{B}(3, 5)$	MN_s	MN_1
EWCQR	L_1	1.00	0.80	0.63	0.29	0.41	0.61	0.35
	L_2	0.50	0.35	1.00	0.13	0.68	0.05	0.14
	$L_1 - L_2^+$	1.00	0.85	1.00	0.34	0.68	0.61	0.63
	$L_1 - L_2^-$	1.00	0.85	1.00	0.44	0.80	0.61	0.63
	$K=3$	0.84	0.94	0.86	0.43	0.59	0.76	0.44
	$K=5$	0.83	0.97	0.89	0.47	0.65	0.78	0.50
	$K=9$	0.82	0.97	0.92	0.49	0.68	0.77	0.52
	$K=19$	0.82	0.97	0.94	0.50	0.69	0.75	0.54
WCQR ⁺	$K=29$	0.83	0.97	0.95	0.51	0.71	0.76	0.54
	$K=3$	0.95	0.94	0.87	0.51	0.61	0.76	0.60
	$K=5$	0.96	0.97	0.91	0.59	0.70	0.78	0.69
	$K=9$	0.97	0.98	0.95	0.69	0.78	0.79	0.77
	$K=19$	0.98	0.99	0.98	0.80	0.86	0.80	0.83
WCQR	$K=29$	0.99	0.99	0.99	0.85	0.90	0.80	0.84
	$K=3$	0.95	0.94	0.87	0.51	0.61	0.76	0.61
	$K=5$	0.96	0.97	0.91	0.60	0.72	0.78	0.76
	$K=9$	0.98	0.98	0.95	0.70	0.80	0.79	0.88
	$K=19$	0.99	0.99	0.98	0.81	0.88	0.92	0.95
	$K=29$	0.99	0.99	0.99	0.86	0.92	0.93	0.97

[†]EWCQR, equally weighted composite quantile regression; WCQR, weighted composite quantile regression.

$$\Omega_{n,\mathbf{w}}(\hat{\beta}_1) = \sum_{i=1}^n \partial\psi_{\mathbf{w}}(Y_i - \mathbf{S}_i^T \hat{\beta}_1) \mathbf{S}_i \mathbf{S}_i^T.$$

Theorem 3. Under conditions 1–5 in Appendix A, if $\|\hat{\beta}_1 - \beta_1^*\| = O_p\{\sqrt{(s/n)}\}$, then the one-step estimator $\hat{\beta}_{\mathbf{w}}^{\text{os}}$ (14) enjoys asymptotic normality:

$$\mathbf{b}^T (\mathbf{S}^T \mathbf{S})^{1/2} (\hat{\beta}_{\mathbf{w}1}^{\text{os}} - \beta_1^*) \xrightarrow{\mathcal{D}} \mathcal{N}(0, \sigma_{\mathbf{w}}^2), \quad (15)$$

provided that $s = o(n^{1/3})$, $\partial\psi(\cdot)$ is Lipschitz continuous and

$$\lambda_{\max} \left(\sum_{i=1}^n \|\mathbf{S}_i\|_i \mathbf{S}_i \mathbf{S}_i^T \right) = O(n\sqrt{s}),$$

where $\lambda_{\max}(\cdot)$ denote the maximum eigenvalue of a matrix and $\sigma_{\mathbf{w}}^2$ is defined as in theorem 2.

The one-step estimator (14) overcomes the convexity restriction and is always well defined, whereas estimator (5) is not uniquely defined when convexity of $\rho_{\mathbf{w}}$ is ruined. If we remove the constraint of $w_k \geq 0$ ($k = 1, \dots, K$), the optimal weight vector in expression (13) is equal to

$$\begin{aligned} \mathbf{w}_{\text{opt}} &= \mathbf{M}^{-1} \mathbf{a}, \\ \sigma_{\mathbf{w}_{\text{opt}}}^2 &= (\mathbf{a}^T \mathbf{M}^{-1} \mathbf{a})^{-1}. \end{aligned}$$

This can be significantly smaller than the optimal variance obtained with convexity constraint, especially for multimodal distributions (see Table 1).

The above discussion prompts a further improvement of the penalized adaptive composite quasi-likelihood in Section 2. Use expression (8) to compute the new residuals and new matrix \mathbf{M} and vector \mathbf{a} . Compute the optimal unconstrained weight $\mathbf{w}_{\text{opt}} = \mathbf{M}^{-1} \mathbf{a}$ and the one-step estimator (14).

4. Examples

In this section, we discuss two specific examples of penalized quasi-likelihood regression. The methods proposed are complementary, in the sense that the first is computationally easy but loses some general flexibility whereas the second is computationally intensive but efficient in a broader class of error distributions.

4.1. Penalized composite L_1 – L_2 -regression

First, we consider the combination of L_1 and L_2 loss functions, i.e. $\rho_1(t) = |t - b_0|$ and $\rho_2(t) = t^2$. The nuisance parameter b_0 is the median of the error distribution. Let $\hat{\beta}_{\mathbf{w}}^{L_1-L_2}$ denote the corresponding penalized estimator as the solution to the minimization problem

$$\arg \min_{\beta, b_0} \left\{ w_1 \sum_{i=1}^n |Y_i - b_0 - \mathbf{X}_i^T \beta| + w_2 \sum_{i=1}^n (Y_i - \mathbf{X}_i^T \beta)^2 + n \sum_{j=1}^p \gamma_{\lambda}(|\beta_j^{(0)}|) |\beta_j| \right\}. \quad (16)$$

If the error distribution is symmetric, then $b_0 = 0$ and the minimization problem (16) can be recast as a penalized weighted least squares regression

$$\arg \min_{\beta} \left\{ \sum_{i=1}^n \left(\frac{w_1}{|Y_i - \mathbf{X}_i^T \hat{\beta}^{(0)}|} + w_2 \right) (Y_i - \mathbf{X}_i^T \beta)^2 + n \sum_{j=1}^p \gamma_{\lambda}(|\beta_j^{(0)}|) |\beta_j| \right\}$$

which can be efficiently solved by pathwise co-ordinate optimization (Friedman *et al.*, 2008) or least angle regression (Efron *et al.*, 2004).

If $b_0 \neq 0$, the penalized least squares problem (16) is somewhat different from problem (5) since we have an additional parameter b_0 . Using the same arguments, and treating b_0 as an additional parameter for which we solve in problem (16), we can show that the conclusions of theorems 2 and 3 hold with the asymptotic variance equal to

$$\sigma_{L_1-L_2}^2(\mathbf{w}) = \frac{w_1^2/4 + w_2^2\sigma^2 + w_2w_1B}{\{w_1f(b_0) + w_2\}^2}, \quad (17)$$

where $B = E[\varepsilon\{I(\varepsilon > b_0) - I(\varepsilon < b_0)\}]$ and $f(\cdot)$ is the density of ε . This will hold when b_0 is either known or unknown. Explicit optimization of equation (17) is not trivial and we go through it as follows.

Since $\sigma_{L_1-L_2}^2(\mathbf{w})$ is invariant to the scale of \mathbf{w} , by setting $w_1/w_2 = c\sigma$, we have

$$\sigma_{L_1-L_2}^2(c) = \sigma^2 \frac{c^2/4 + 1 + a_\varepsilon c}{(b_\varepsilon c + 1)^2}, \quad (18)$$

where $a_\varepsilon = B/\sigma$ and $b_\varepsilon = \sigma f(b_0)$. Note that

$$|B| \leq E|\varepsilon|\{I(\varepsilon > b_0) + I(\varepsilon < b_0)\} \leq \sigma.$$

Hence, $|a_\varepsilon| \leq 1$ and $c^2/4 + 1 + a_\varepsilon c = (c/2 + a_\varepsilon)^2 + 1 - a_\varepsilon^2 \geq 0$.

The optimal value of c over $[0, \infty)$ can be easily computed. If $a_\varepsilon b_\varepsilon < 0.5$, then the optimal value is obtained at

$$c_\varepsilon = \frac{2(2b_\varepsilon - a_\varepsilon)_+}{1 - 2a_\varepsilon b_\varepsilon}. \quad (19)$$

In particular, when $2b_\varepsilon \leq a_\varepsilon$, $c_\varepsilon = 0$, and the optimal choice is the least squares estimator. When $a_\varepsilon b_\varepsilon = 0.5$, if $2b_\varepsilon \leq a_\varepsilon$, then the minimizer is $c_\varepsilon = 0$. In all other cases, the minimizer is $c_\varepsilon = \infty$, i.e. we are left to use L_1 -regression alone.

This result shows the limitation of the convex combination, i.e. $c \geq 0$. In many cases, we are left with the least squares or least absolute deviation regression without improving efficiency. The efficiency can be gained and achieved by allowing negative weights via the one-step technique as in Section 3.4. Let

$$g(c) = \frac{c^2/4 + 1 + a_\varepsilon c}{(b_\varepsilon c + 1)^2}.$$

The function $g(c)$ has a pole at $c = -1/b_\varepsilon$ and a unique critical point

$$c_{\text{opt}} = \frac{2(2b_\varepsilon - a_\varepsilon)}{1 - 2a_\varepsilon b_\varepsilon}, \quad (20)$$

provided that $a_\varepsilon b_\varepsilon \neq \frac{1}{2}$. Consequently, the function $g(c)$ cannot have any local maximizer (otherwise, from the local maximizer to the point $c = -1/b_\varepsilon$, there must be a local minimizer, which is also a critical point). Hence, the minimum value is attained at c_{opt} . In other words,

$$\min_{\mathbf{w}} \{\sigma_{L_1-L_2}^2(\mathbf{w})\} = \sigma^2 \min_c \{g(c)\} = d_\varepsilon \sigma^2, \quad (21)$$

where

$$d_\varepsilon = g(c_{\text{opt}}) = \frac{1 - a_\varepsilon^2}{4b_\varepsilon^2 - 4a_\varepsilon b_\varepsilon + 1}. \quad (22)$$

Since the denominator can be written as $(a_\varepsilon - 2b_\varepsilon)^2 + 1 - a_\varepsilon^2$, we have $d_\varepsilon \leq 1$, i.e. it outperforms the least squares estimator, unless $a_\varepsilon = 2b_\varepsilon$. Similarly, it can be shown that

$$d_\varepsilon = \frac{1 - a_\varepsilon^2}{4b_\varepsilon^2\{1 - a_\varepsilon^2 + (2a_\varepsilon - 1/b_\varepsilon)^2/4\}} \leq \frac{1}{4b_\varepsilon^2},$$

i.e. it outperforms the least absolute deviation estimator, unless $a_\varepsilon b_\varepsilon = \frac{1}{2}$.

When the error distribution is symmetric unimodal, $b_\varepsilon \geq 1/\sqrt{12}$, according to chapter 5 of Lehmann (1983). The worst scenario for the L_1 -regression in comparison with the L_2 -regression is the uniform distribution (see chapter 5 of Lehmann (1983)), which has the relative efficiency of merely $\frac{1}{3}$. For such a uniform distribution, $a_\varepsilon = \sqrt{3}/2$ and $b_\varepsilon = 1/\sqrt{12}$, $d_\varepsilon = \frac{3}{4}$, and $c_{\text{opt}} = -2/\sqrt{3}$. Hence, the best L_1 - L_2 -regression is four times better than L_1 -regression alone. More comparisons about the weighted L_1 - L_2 -combination with L_1 - and least squares-regression are given in Table 1.

4.2. Penalized composite quantile regression

Weighted composite quantile regression (WCQR) was first studied by Koenker (1984) in a classical statistical inference setting. Zou and Yuan (2008) used equally weighted composite quantile regression (EWCQR) for penalized model selection with p large but fixed. We show that the efficiency of EWCQR can be substantially improved by properly weighting and extend the work to the case of $p \gg n$. Consider K different quantiles, $0 < \tau_1 < \tau_2 < \dots < \tau_K < 1$. Let

$$\rho_k(t) = \tau_k(t - b_k)_+ + (1 - \tau_k)(t - b_k)_-.$$

The penalized CQR estimator $\hat{\beta}^{\text{CQR}}$ is defined as the solution to the minimization problem

$$\arg \min_{b_1, \dots, b_K, \beta} \left\{ \sum_{k=1}^K w_k \sum_{i=1}^n \rho_k(Y_i - \mathbf{X}_i^T \beta) + n \sum_{j=1}^p \gamma_\lambda(|\beta_j^{(0)}|) |\beta_j| \right\}, \quad (23)$$

where b_k is the estimator of the nuisance parameter $b_k^* = F^{-1}(\tau_k)$, the τ_k th quantile of the error distribution. b_1, \dots, b_K are nuisance parameters and the minimization at expression (23) is done with respect to them also. After some algebra we can confirm that the conclusions of theorems 2 and 3 continue to hold with the asymptotic variance as

$$\sigma_{\text{CQR}}^2(\mathbf{w}) = \sum_{k, k'=1}^K w_k w_{k'} \{ \min(\tau_k, \tau_{k'}) - \tau_k \tau_{k'} \} / \left[\sum_{k=1}^K w_k f\{F^{-1}(\tau_k)\} \right]^2. \quad (24)$$

As shown in Koenker (1984) and Bickel (1973), when $K \rightarrow \infty$, optimally weighted composite quantile regression (OWCQR) is as efficient as the maximum likelihood estimator and always more efficient than EWCQR. Computationally, the minimization problem in equation (23) can be cast as a large-scale linear programming problem by expanding the covariate space with new ancillary variables. Thus, it is computationally intensive to use too many quantiles. In Section 4.3, we can see that usually no more than 10 quantiles are adequate for OWCQR to approach the efficiency of maximum likelihood estimation, whereas determining the optimal value of K in EWCQR seems difficult since the efficiency is not necessarily an increasing function of K (Table 2). Also, some of the weights in \mathbf{w}_{opt} are 0, hence making the OWCQR method computationally less intensive than EWCQR. From our experience in large p and small n situations, this reduction tends to be significant.

Table 2. Optimal weights of convex composite quantile regression with $K = 9$ quantiles

Quantile	Results for the following distributions $f(\varepsilon)$:						
	DE	t_4	$\mathcal{N}(0, 1)$	$\Gamma(3, 1)$	$\mathcal{B}(3, 5)$	MN_s	MN_l
1st	0	0	0.20	0.56	0.39	0.06	0.36
2nd	0	0.12	0.11	0.15	0.10	0.23	0.11
3rd	0	0.14	0.09	0.08	0.11	0.17	0.10
4th	0	0.14	0.08	0.06	0.05	0.10	0.01
5th	1	0.16	0.06	0.05	0	0.14	0
6th	0	0.14	0.08	0.04	0	0	0
7th	0	0.14	0.09	0	0.05	0	0
8th	0	0.12	0.11	0	0.09	0	0
9th	0	0	0.20	0.05	0.20	0.30	0.29

The optimal convex combination of quantile regression uses the weight

$$\mathbf{w}_{\text{opt}}^+ = \arg \min_{\mathbf{w} \geq 0, \mathbf{a}^T \mathbf{w} = 1} (\mathbf{w}^T \mathbf{M} \mathbf{w}), \quad (25)$$

where $\mathbf{a} = (f\{F^{-1}(\tau_1)\}, \dots, f\{F^{-1}(\tau_K)\})^T$ and \mathbf{M} is a $K \times K$ matrix whose (i, j) element is $\min(\tau_i, \tau_j) - \tau_i \tau_j$. The optimal combination of quantile regression, which is obtained by using the one-step procedure, uses the weight

$$\mathbf{w}_{\text{opt}} = \mathbf{M}^{-1} \mathbf{a}. \quad (26)$$

Clearly, both combinations improve the efficiency of EWCQR and the optimal combination is most efficient among the three (see Table 1). When the error distributions are skewed or multimodal, the improvement can be substantial.

4.3. Asymptotic efficiency comparison

In this section, we study the asymptotic efficiency of the proposed estimators under several error distributions. For comparison, we also include L_1 -regression, L_2 -regression and EWCQR. The error distribution ranges from the symmetric to asymmetric distributions: DE, the t -distribution with 4 degrees of freedom (t_4), the normal $\mathcal{N}(0, 1)$, gamma $\Gamma(3, 1)$ and beta $\mathcal{B}(3, 5)$ distributions, a scale mixture of normals distributions (MN_s) $0.1 \mathcal{N}(0, 25) + 0.9 \mathcal{N}(0, 1)$ and a location mixture of normal distributions (MN_l) $0.7 \mathcal{N}(-1, 1) + 0.3 \mathcal{N}(7/3, 1)$. To keep the comparison fair and to satisfy the first assumption of mean 0 error terms, we first centred the error distribution to have mean 0.

Table 1 shows the asymptotic relative efficiency of each estimator compared with maximum likelihood estimation. $L_1-L_2^+$ and L_1-L_2 indicate the optimal convex L_1-L_2 -combination and optimal L_1-L_2 -combination respectively. Whereas L_1 -regression can have higher or lower efficiency than L_2 -regression in different error distributions, $L_1-L_2^+$ and L_1-L_2 -regressions are consistently more efficient than both of them. WCQR⁺ denotes the optimal convex combination of multiple quantile regressions and WCQR represents the optimal combination. In all quantile regressions, quantiles $(1/(K+1), \dots, K/(K+1))$ were used. As shown in Table 1, WCQR⁺ and WCQR always outperform EWCQR and the differences are more significant in the DE distribution and asymmetric distributions such as gamma and beta. In the DE, t_4 - and

$\mathcal{N}(0, 1)$ distributions, nine quantiles are usually adequate for WCQR^+ and WCQR to achieve full efficiency. In $\Gamma(3, 1)$ and $\mathcal{B}(3, 5)$, they need 29 quantiles to achieve efficiency close to that of maximum likelihood estimation whereas the other estimators are significantly inefficient. This difference is most expressed in the multimodal distributions MN_s and MN_l , with WCQR outperforming all. One of the possible problems with EWCQR is that the efficiency does not necessarily increase with K , making the choice of K more difficult. For example, for the DE distribution, the relative efficiency decreases with K . This is understandable, as $K = 1$ is optimal: putting more and odd numbers of quantiles dilutes the weights.

In Table 2 we illustrate both the adaptivity of the proposed composite quasi-maximum-likelihood estimation methodology and computational efficiency of WCQR^+ over EWCQR by showing the positions of zero of the optimal non-negative weight vector $\mathbf{w}_{\text{opt}}^+$. For $K = 9$, only one quantile is needed in the DE case, five and six quantiles are needed for MN_l and MN_s and seven quantiles for the t_4 -, gamma and beta distributions. Only in the normal distribution are all nine quantiles used. Therefore, WCQR^+ can dramatically reduce the computational complexity of EWCQR in large-scale optimization problems where $p \gg n$.

Table 3. Simulation results ($n = 100$ and $p = 12$)

Method		Results for the following distributions $f(\varepsilon)$:					
		DE	t_4	$\mathcal{N}(0, 3)$	$\Gamma(3, 1)$	$\mathcal{B}(3, 5)$	MN_s
L_1	Oracle	0.029†	0.050	0.122	0.082	0.0010	0.043
	Penalized	0.035‡	0.053	0.128	0.097	0.0011	0.051
	(TP, FP)	(3,1.83)	(3,0.8)	(3,0.84)	(3,1)	(3,0.54)	(3,0.93)
	Standard deviation $\times 10^2$	0.646	0.767	0.570	0.950	0.112	0.244
L_2	Oracle	0.047	0.043	0.073	0.064	0.00056	0.083
	Penalized	0.059	0.054	0.106	0.100	0.0011	0.091
	(TP, FP)	(3,0.82)	(3,1.61)	(3,1.89)	(3,1.35)	(3,3.76)	(3,1.47)
	Standard deviation $\times 10^2$	0.779	0.762	0.485	0.869	0.129	0.179
$L_1-L_2^+$	Oracle	0.036	0.043	0.070	0.070	0.00061	0.051
	Penalized	0.037	0.049	0.102	0.099	0.00077	0.058
	(TP, FP)	(3,2.49)	(3,2.39)	(3,1.97)	(3,2.09)	(3,2.42)	(3,2.69)
	Standard deviation $\times 10^1$	0.717	0.702	0.518	0.876	0.095	0.169
L_1-L_2	Oracle	0.036	0.043	0.070	0.063	0.00060	0.051
	Penalized	0.037	0.049	0.102	0.078	0.00063	0.058
	(TP, FP)	(3,2.49)	(3,2.39)	(3,1.97)	(3,2.05)	(3,2.42)	(3,2.69)
	Standard deviation $\times 10^2$	0.717	0.702	0.518	0.846	0.075	0.169
EWCQR	Oracle	0.031	0.046	0.069	0.063	0.00065	0.033
	Penalized	0.042	0.046	0.107	0.074	0.00091	0.040
	(TP, FP)	(3,1.88)	(3,1.57)	(3,2.04)	(3,1.83)	(3,1.88)	(3,1.38)
	Standard deviation $\times 10^2$	0.654	0.562	0.488	0.813	0.087	0.177
WCQR^+	Oracle	0.033	0.047	0.068	0.052	0.00065	0.036
	Penalized	0.039	0.041	0.100	0.054	0.00070	0.037
	(TP, FP)	(3,0.55)	(3,1.47)	(3,0.74)	(3,0.61)	(3,0.98)	(3,0.62)
	Standard deviation $\times 10^1$	0.440	0.612	0.498	0.715	0.071	0.174
WCQR	Oracle	0.033	0.047	0.068	0.048	0.00058	0.028
	Penalized	0.039	0.041	0.100	0.050	0.00062	0.030
	(TP, FP)	(3,0.55)	(3,1.47)	(3,0.74)	(3,0.61)	(3,0.98)	(3,0.62)
	Standard deviation $\times 10^1$	0.440	0.612	0.498	0.650	0.061	0.132

†Median model error of the oracle estimator.

‡Median model error of the penalized estimator.

5. Finite sample study

5.1. Simulated example

In the simulation study, we consider the classical linear model for testing variable selection methods that was used by Fan and Li (2001):

$$y = \mathbf{x}^T \boldsymbol{\beta}^* + \varepsilon, \quad \mathbf{x} \sim N(0, \boldsymbol{\Sigma}_{\mathbf{x}}), \quad (\boldsymbol{\Sigma}_{\mathbf{x}})_{i,j} = 0.5^{|i-j|}.$$

The error vector varies from unimodal to multimodal and heavy- to light-tail distributions in the same way as in Tables 1 and 2, and is centred to have mean 0. The data have $n = 100$ observations. We considered two settings where $p = 12$ and $p = 500$. In both settings, $(\beta_1, \beta_2, \beta_5) = (3, 1.5, 2)$ and the other coefficients are equal to 0. We implemented penalized L_1 - and L_2 -regression, composite L_1 - L_2^+ and L_1 - L_2 , EWCQR, WCQR⁺ and WCQR using quantiles (10%, 20%, ..., 90%). The local linear approximation of the SCAD penalty (6) was used and the

Table 4. Simulation results ($n = 100$ and $p = 500$)

Method		Results for the following distributions $f(\varepsilon)$:					
		DE	t_4	$\mathcal{N}(0, 3)$	$\Gamma(3, 1)$	$\mathcal{B}(3, 5)$	MN_5
Lasso	Oracle	0.039†	0.039	0.035	0.0719	0.062	0.176
	Penalized	1.775‡	1.759	8.687	2.662	1.808	6.497
	(TP,FP)	(3,94.46)	(3,94.26)	(3,96.80)	(3,95.59)	(3,86.88)	(3,96.55)
	Standard deviation $\times 10^2$	3.336	3.257	0.578	3.167	0.989	0.539
L_1	Oracle	0.025	0.031	0.382	0.096	0.0094	0.281
	Penalized	0.035	0.039	1.342	0.131	0.0120	0.514
	(TP,FP)	(3,4.53)	(3,4.47)	(3,5.32)	(3,4.56)	(3,8.10)	(3,4.58)
	Standard deviation $\times 10^2$	0.268	0.274	0.144	0.461	0.215	0.101
L_2	Oracle	0.035	0.043	0.207	0.078	0.0057	0.187
	Penalized	0.093	0.086	1.187	0.175	0.0073	0.764
	(TP,FP)	(3,12.31)	(3,10.64)	(3,11.00)	(3,8.02)	(3,18.75)	(3,16.93)
	Standard deviation $\times 10^1$	0.865	0.828	0.281	0.168	0.396	0.238
L_1 - L_2^+	Oracle	0.193	0.035	0.224	0.080	0.0061	0.195
	Penalized	0.036	0.036	1.160	0.097	0.0077	0.576
	(TP,FP)	(3,17.92)	(3,12.58)	(3,15.87)	(3,15.43)	(3,14.05)	(3,17.92)
	Standard deviation $\times 10^2$	0.226	0.235	0.396	0.144	0.235	0.207
L_1 - L_2	Oracle	0.035	0.035	0.224	0.079	0.0050	0.195
	Penalized	0.036	0.036	1.160	0.095	0.0069	0.576
	(TP,FP)	(3,17.92)	(3,12.58)	(3,15.87)	(3,15.43)	(3,14.05)	(3,17.92)
	Standard deviation $\times 10^2$	0.226	0.235	0.905	0.150	0.190	0.207
EWCQR	Oracle	0.029	0.024	0.252	0.057	0.0064	0.207
	Penalized	0.060	0.070	0.764	0.148	0.0118	0.599
	(TP,FP)	(3,8.71)	(3,8.43)	(3,7.78)	(3,9.59)	(3,9.69)	(3,8.91)
	Standard deviation $\times 10^1$	0.469	0.475	0.153	0.716	0.213	0.139
WCQR ⁺	Oracle	0.028	0.027	0.223	0.050	0.0066	0.204
	Penalized	0.045	0.037	0.595	0.079	0.0076	0.368
	(TP,FP)	(3,3.97)	(3,3.76)	(3,3.93)	(3,3.66)	(3,4.85)	(3,4.05)
	Standard deviation $\times 10^1$	0.244	0.266	0.112	0.273	0.120	0.084
WCQR	Oracle	0.028	0.027	0.223	0.048	0.0048	0.160
	Penalized	0.045	0.037	0.595	0.062	0.0060	0.280
	(TP,FP)	(3,3.97)	(3,3.76)	(3,3.93)	(3,3.66)	(3,4.85)	(3,4.05)
	Standard deviation $\times 10^1$	0.224	0.219	0.112	0.180	0.110	0.060

†Median model error of the oracle estimator.

‡Median model error of the penalized estimator.

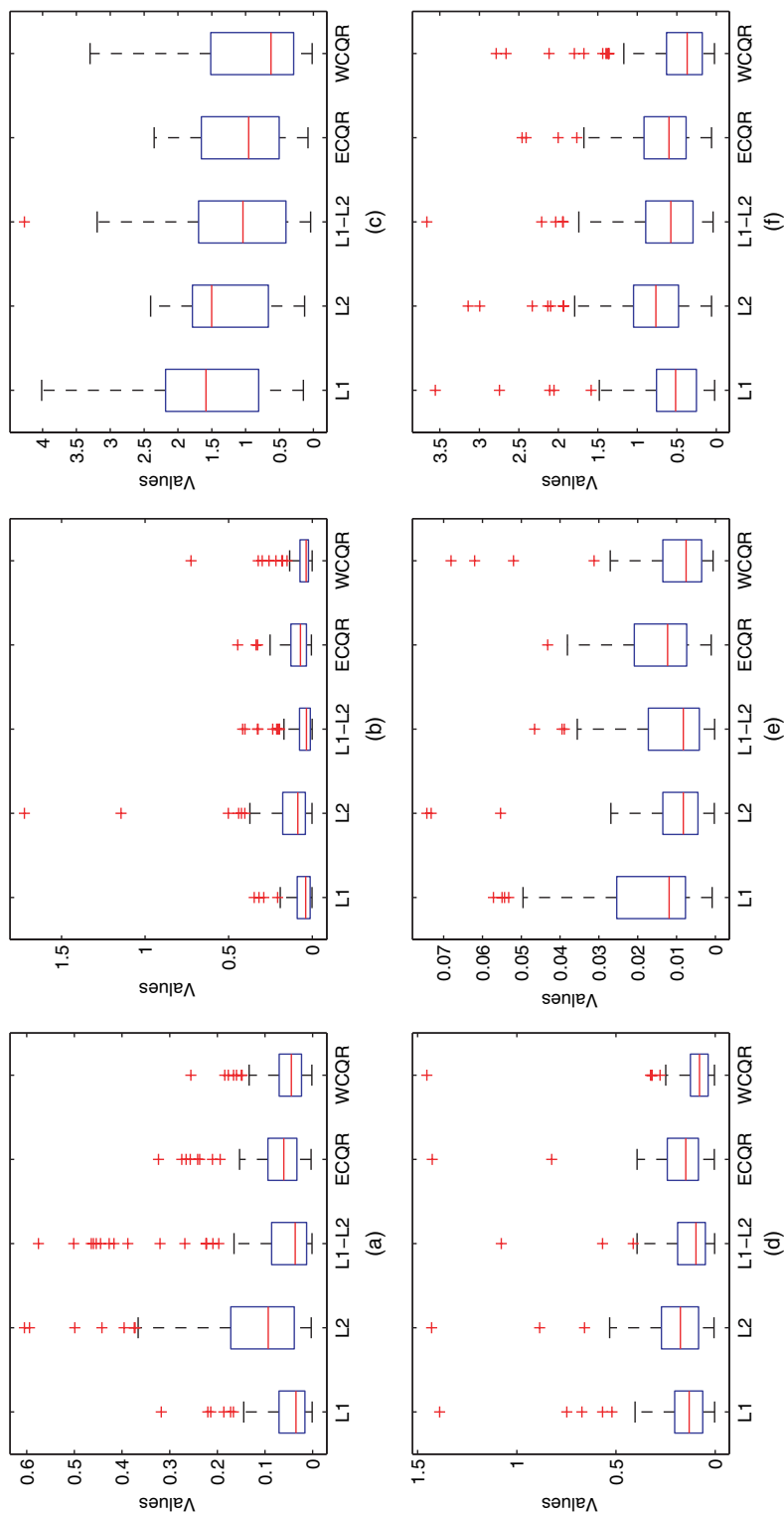


Fig. 1. Boxplots of the median model error of the L_1 , L_2 , L_1-L_2 , ECQR and WCQR methods under various distributional settings with $n = 100$ and $p = 500$: (a) DE distribution; (b) t_4 -distribution; (c) normal distribution; (d) gamma distribution; (e) beta distribution; (f) mixture of normal distributions

tuning parameter in the penalty was selected by using fivefold cross-validation. We compared different methods by

- (a) model error, which is defined as $ME(\hat{\beta}) = (\hat{\beta} - \beta^*)^T E(\mathbf{X}^T \mathbf{X})(\hat{\beta} - \beta^*)$,
- (b) the number of correctly classified non-zero coefficients, i.e. the true positive number TP,
- (c) the number of incorrectly classified zero coefficients, i.e. the false positive number FP, and
- (d) the multiplier $\hat{\sigma}_{\mathbf{w}}$ of the standard error (12).

A total of 100 replications were performed and the median model error, the average of TP and FP, are reported in Table 3. The median model errors of oracle estimators were calculated as the benchmark for comparison.

From the results that are presented in Table 3 and Table 4, we can see that penalized composite L_1 - L_2^+ -regression takes the smaller of the two model errors of L_1 - and L_2 -regression in all distributions except in $\mathcal{B}(3, 5)$ where it outperforms both. As expected, optimal L_1 - L_2 outperforms L_1 - L_2^+ and brings a smaller FP, especially in multimodal and unsymmetric distributions. Also, both L_1 - L_2^+ and L_1 - L_2 perform reasonably well when compared with EWCQR, but with much less computational burden. WCQR⁺ and WCQR in both Table 3 and Table 4 have smaller model errors and smaller FP than EWCQR. Similar conclusions can be drawn from Fig. 1, which compares the boxplots of the model errors of the five methods (WCQR⁺ and L_1 - L_2^+ are not shown) under various distributions in the case $n = 100$ and $p = 500$. For $p \ll n$ in Table 3 we did not include the lasso estimator since it behaves reasonably well in that setting. For $p \gg n$ in Table 4, we included the lasso estimator as a reference. Table 4 shows that the lasso has bigger model errors, more false positive classifications and higher standard errors (usually by a factor of 10) than any other of the five SCAD-based methods discussed.

In addition to the modal error in Tables 3 and 4, we report the multiplier $\hat{\sigma}_{\mathbf{w}}$ of the asymptotic variance (see equation (12)). Being the only part of the standard error that depends on the choice of weights \mathbf{w} and loss functions ρ_k , we explored its behaviour when the dimensionality grows from $p \ll n$ to $p \gg n$. Both Table 3 and Table 4 confirm the stability of the formula throughout the two settings and all five composite quasi-maximum-likelihood estimation methods. Only the lasso estimator, being unable to specify the correct sparsity set when $p \gg n$, inflates $\hat{\sigma}_{\mathbf{w}}$ by one order of magnitude compared with the other composite quasi-maximum-likelihood estimation methods. Note that WCQR⁺ keeps the smallest value of $\hat{\sigma}_{\mathbf{w}}$ and each of L_1 - L_2 , L_1 - L_2^+ , WCQR and WCQR⁺ have smaller standard errors than the classical L_1 -, L_2 - or EWCQR methods.

5.2. Real data example

In this section, we applied the methods proposed to expression quantitative trait locus (EQTL) mapping. Variations in gene expression levels may be related to phenotypic variations such as susceptibility to diseases and response to drugs. Therefore, to understand the genetic basis of gene expression, variation is an important topic in genetics. The availability of genome-wide SNP measurement has made it possible and reasonable to perform the high resolution EQTL mapping on the scale of nucleotides. In our analysis, we conducted the *cis*-EQTL mapping for the gene CCT8. This gene is within the Down syndrome critical region on human chromosome 21, on the minus strand. The overexpression of CCT8 may be associated with Down syndrome phenotypes.

We used the SNP genotype data and gene expression data for the 210 unrelated individuals of the international ‘HapMap’ project, which include 45 Japanese in Tokyo, Japan, 45 Han Chinese in Beijing, China, 60 Utah parents with ancestry from northern and western Europe

(from the collection of the Centre d'Etude du Polymorphisme Humain (CEPH)) and 60 Yoruba parents in Ibadan, Nigeria, and they are available in PLINK format (<http://pngu.mgh.harvard.edu/purcell/plink/>). We included in the analysis more than 2 million SNPs with minor allele frequency greater than 1% and missing data rate less than 5%. The gene expression data were generated by an Illumina Sentrix Human-6 Expression BeadChip and have been normalized (*i*th-quantile normalization across replicates and median normalization across individuals) independently for each population (<ftp://ftp.sanger.ac.uk/pub/genevar/>).

Specifically, we considered the *cis*-candidate region to start 1 megabase upstream of the transcription start site (TSS) of CCT8 and to end 1 megabase downstream of the transcription end site, which includes 1955 SNPs in the Japanese and Chinese, 1978 SNPs in the CEPH and 2146 SNPs in the Yoruba populations. In the following analysis, we grouped Japanese and Chinese together into the Asian population and analysed the three populations Asian, CEPH and Yoruba separately. The additive coding of SNPs (e.g. 0, 1, 2) was adopted and was treated

Table 5. EQTLs for gene CCT8 in the Japanese and Chinese populations ($n = 90$)†

<i>SNP</i>	L_2	$L_1 - L_2^+$ or $L_1 - L_2$	L_1	<i>EWQQR</i>	$WQQR^+$ or $WQQR$	<i>Distance</i> <i>from TSS</i> (kilobases)
rs16981663	−0.11 (0.03)	−0.11 (0.03)	−0.09 (0.04)	−0.10 (0.03)	−0.09 (0.03)	−998
rs16981663‡	0.08 (0.06)	0.08 (0.06)		0.04 (0.06)		−998
rs9981984	−0.12 (0.03)	−0.12 (0.03)	−0.10 (0.04)	−0.09 (0.03)	−0.12 (0.03)	−950
rs7282280				0.05 (0.03)		−231
rs7282280‡				−0.07 (0.05)		−231
rs2245431‡	0.33 (0.10)	0.33 (0.10)	0.36 (0.11)	0.37 (0.09)	0.38 (0.10)	−89
rs2832159	0.21 (0.04)	0.21 (0.04)	0.30 (0.04)	0.20 (0.04)	0.23 (0.04)	13
rs1999321‡	0.11 (0.07)	0.11 (0.07)		0.14 (0.07)		84
rs2832224	0.07 (0.03)	0.07 (0.03)		0.06 (0.03)	0.04 (0.03)	86

†Standard errors of the estimates are reported in parentheses.

‡SNP equal to 2; otherwise SNP is equal to 1.

Table 6. EQTLs for gene CCT8 in the CEPH population ($n = 60$)†

<i>SNP</i>	L_2	$L_1 - L_2^+$ or $L_1 - L_2$	L_1	<i>EWQQR</i>	$WQQR^+$ or $WQQR$	<i>Distance</i> <i>from TSS</i> (kilobases)
rs2831459	0.20 (0.07)	0.20 (0.07)	0.19 (0.08)	0.17 (0.07)	0.18 (0.07)	−999
rs7277536	0.18 (0.09)	0.18 (0.09)	0.09 (0.11)	0.14 (0.09)	0.23 (0.09)	−672
rs7278456‡	0.36 (0.11)	0.36 (0.11)	0.21 (0.13)	0.40 (0.11)	0.35 (0.11)	−663
rs2248610	0.08 (0.04)	0.08 (0.04)	0.09 (0.05)	0.10 (0.05)	0.06 (0.05)	−169
rs965951	0.11 (0.05)	0.11 (0.05)	0.13 (0.06)	0.03 (0.06)	0.12 (0.05)	−13
rs3787662	0.12 (0.06)	0.12 (0.06)	0.08 (0.07)	0.13 (0.06)	0.12 (0.06)	78
rs2832253				0.10 (0.07)		117
rs2832332				0.08 (0.05)		382
rs13046799	−0.16 (0.05)	−0.16 (0.05)	−0.14 (0.06)	−0.14 (0.05)	−0.16 (0.05)	993

†Standard errors of the estimates are reported in parentheses.

‡SNP equal to 2; otherwise SNP is equal to 1.

Table 7. EQTLs of gene CCT8 in the Yoruba population ($n = 60$)†

<i>SNP</i>	L_2	$L_1 - L_2^+$ or $L_1 - L_2$	L_1	<i>EWQQR</i>	<i>WCQR</i> ⁺ or <i>WCQR</i>	<i>Distance</i> <i>from TSS</i> (kilobases)
rs9982023‡			0.12 (0.05)	0.14 (0.04)		−531
rs1236427				0.15 (0.04)		−444
rs2831972	−0.22 (0.06)	−0.22 (0.06)	−0.16 (0.07)	−0.30 (0.05)	−0.30 (0.06)	−360
rs2091966‡	−0.21 (0.11)	−0.21 (0.11)	−0.57 (0.16)	−0.39 (0.13)	−0.20 (0.11)	−358
rs2832010	−0.04 (0.03)	−0.04 (0.03)	−0.18 (0.08)	−0.32 (0.05)	−0.07 (0.03)	−336
rs2832024			0.14 (0.09)	0.26 (0.06)		−332
rs2205413	−0.08 (0.04)	−0.08 (0.04)	−0.15 (0.05)	−0.16 (0.04)	−0.04 (0.03)	−330
rs2205413‡				−0.29 (0.05)		−330
rs2832042‡	0.14 (0.04)	0.14 (0.04)	0.23 (0.05)	0.23 (0.04)	0.13 (0.04)	−330
rs2832053‡				−0.12 (0.13)		−315
rs2832053			0.09 (0.04)	0.06 (0.02)		−315
rs8130766	−0.01 (0.03)	−0.01 (0.03)	−0.14 (0.05)	−0.10 (0.03)	−0.04 (0.03)	−296
rs16983288‡	−0.13 (0.07)	−0.13 (0.07)	−0.28 (0.08)	−0.28 (0.05)	−0.15 (0.06)	−288
rs16983303	−0.06 (0.03)	−0.06 (0.03)	−0.10 (0.02)	−0.15 (0.03)	−0.09 (0.03)	−283
rs8134601‡	0.18 (0.11)	0.18 (0.11)	0.15 (0.12)	0.16 (0.07)	0.19 (0.10)	−266
rs8134601	−0.16 (0.12)	−0.16 (0.12)	0.08 (0.15)	0.25 (0.11)	−0.17 (0.11)	−266
rs7276141‡			−0.06 (0.12)	0.15 (0.11)		−264
rs7281691	0.23 (0.10)	0.23 (0.10)	−0.03 (0.13)	−0.18 (0.10)	0.26 (0.09)	−263
rs7281691‡	−0.14 (0.09)	−0.14 (0.09)	−0.05 (0.13)	−0.23 (0.09)	−0.12 (0.09)	−263
rs1006903‡	−0.01 (0.05)	−0.01 (0.05)	0.13 (0.06)	0.07 (0.04)	0.01 (0.05)	−246
rs7277685			0.07 (0.05)	0.06 (0.03)		−240
rs9982426	0.02 (0.03)	0.02 (0.03)	0.12 (0.04)	0.18 (0.05)		−238
rs2832115				−0.08 (0.05)		−225
rs11910981	−0.09 (0.03)	−0.09 (0.03)	−0.15 (0.03)	−0.19 (0.03)	−0.08 (0.03)	−160
rs2243503				0.07 (0.06)		−133
rs2243552			0.10 (0.03)	0.03 (0.05)		−128
rs2247809	0.01 (0.06)	0.01 (0.06)	0.18 (0.07)	0.26 (0.06)	0.01 (0.05)	−116
rs878797‡	0.11 (0.06)	0.11 (0.06)	0.26 (0.07)	0.23 (0.05)	0.05 (0.06)	−55
rs6516887	0.07 (0.04)	0.07 (0.04)	−0.05 (0.07)	−0.09 (0.04)	0.07 (0.04)	−44
rs8128844			−0.10 (0.06)	−0.17 (0.05)	0.02 (0.04)	−24
rs965951‡	0.10 (0.10)	0.10 (0.10)	0.28 (0.11)	0.26 (0.08)	0.13 (0.09)	−13
rs2070610			0.18 (0.05)	0.17 (0.04)		−0
rs2832159	0.06 (0.06)	0.06 (0.06)	−0.04 (0.07)	−0.20 (0.06)	0.11 (0.05)	13
rs2832178‡	−0.16 (0.06)	−0.16 (0.06)	−0.16 (0.08)	−0.20 (0.06)	−0.24 (0.06)	34
rs2832186			−0.06 (0.07)	0.12 (0.05)		38
rs2832190‡	−0.41 (0.06)	−0.41 (0.06)	−0.25 (0.11)	−0.20 (0.08)	−0.49 (0.06)	42
rs2832190	−0.22 (0.05)	−0.22 (0.05)	−0.16 (0.05)	−0.26 (0.06)	−0.28 (0.05)	42
rs7275293			0.13 (0.12)	0.32 (0.09)		54
rs16983792	−0.11 (0.04)	−0.11 (0.04)	−0.10 (0.05)	−0.18 (0.04)	−0.14 (0.04)	82
rs2251381‡			−0.11 (0.08)	−0.15 (0.05)		85
rs2251517‡	−0.25 (0.05)	−0.25 (0.05)	−0.26 (0.07)	−0.27 (0.05)	−0.28 (0.05)	86
rs2251517	−0.11 (0.04)	−0.11 (0.04)	−0.19 (0.05)	−0.23 (0.03)	−0.15 (0.04)	86
rs2832225				−0.07 (0.03)		87
rs7283854			0.10 (0.04)	0.13 (0.02)		443

†Standard errors of the estimates are reported in parentheses.

‡SNP equal to 2; otherwise SNP is equal to 1.

as categorical variables instead of continuous variables to allow non-additive effects, i.e. two dummy variables will be created for categories 1 and 2 respectively. The category 0 represents the major, normal, population. The missing SNP measurements were imputed as 0s. The response variable is the gene expression level of gene CCT8, measured by microarray.

In the first step, the analysis-of-variance F -statistic was computed for each SNP independently

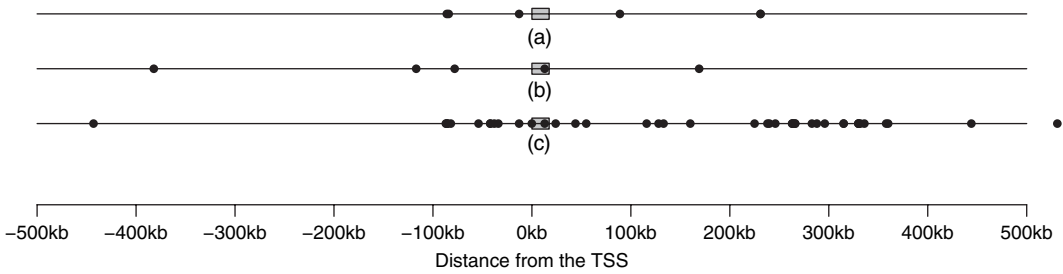


Fig. 2. Chromosome locations of identified EQTLs of gene CCT8 (■, CCT8's coding region): EQTLs selected by any of the five methods are shown: (a) Japanese and Chinese; (b) CEPH; (c) Yoruba

and a version of the independent screening method of Fan and Lv (2008) was implemented. This method is particularly computationally efficient in ultrahigh dimensional problems and here we retained the top 100 SNPs with the largest F -statistics. In the second step, we applied to the screened data the penalized L_2 - and L_1 -regression, $L_1-L_2^+$, L_1-L_2 , EWCQR, WCQR⁺ and WCQR with local linear approximation of the SCAD penalty. All the four composite quantile regressions used quantiles at (10%, ..., 90%). The lasso was used as the initial estimator and the tuning parameter in both the lasso and the SCAD penalty was chosen by fivefold cross-validation. In all three populations, the L_1-L_2 - and $L_1-L_2^+$ -regressions reduced to L_2 -regression. This is not unexpected owing to the gene expression normalization procedure. In addition, WCQR reduced to WCQR⁺. The selected SNPs, their coefficients and distances from the TSS are summarized in Tables 5–7.

In the Asian population (Table 5), the five methods are reasonably consistent in not only variables selection but also coefficient estimation (in terms of signs and order of magnitude). WCQR uses the weights (0.19, 0.11, 0.02, 0, 0.12, 0.09, 0.18, 0.19, 0.10). There are four SNPs which were chosen by all five methods. Two of them, rs2832159 and rs2245431, up-regulate gene expression, whereas rs9981984 and rs16981663 down-regulate gene expression. EWCQR selects the largest set of SNPs, whereas L_1 -regression selects the smallest set. In the CEPH population (Table 6), all five methods consistently selected the same seven SNPs with only EWCQR choosing two additional SNPs. WCQR uses the weight (0.19, 0.21, 0, 0.04, 0.03, 0.07, 0.1, 0.21, 0.15). The coefficient estimations were also highly consistent. Deutsch *et al.* (2005) performed a similar *cis*-EQTL mapping for gene CCT8 using the same CEPH data as here. They considered a 100-kilobase region surrounding the gene, which contains 41 SNPs. Using analysis of variance with correction for multiple tests, they identified four EQTLs (rs965951, rs2832159, rs8133819 and rs2832160), among which rs965951 had the smallest p -value. Our analysis verified rs965951 to be an EQTL but did not find the other SNPs to be associated with the gene expression of CCT8. In other words, conditioning on the presence of SNP rs965951 the other three make little additional contributions. The analysis of the Yoruba population yields a large number of EQTLs (Table 7). EWCQR again selects the largest set of 44 EQTLs. L_1 -regression selects 38 EQTLs. L_2 -regression and WCQR both select 27 SNPs, 26 of which are the same. WCQR uses the weight (0.1, 0, 0.17, 0.16, 0.11, 0.3, 0, 0, 0.16). The coefficients that are estimated by the different methods are mostly consistent (in terms of signs and order of magnitude), except that the coefficient estimates for rs8134601, rs7281691, rs6516887 and rs2832159 by EWCQR and L_1 -regression have different signs from those of L_2 -regression and WCQR. The EQTLs are almost all within 500 kilobases upstream of the TSS or 500 kilobases downstream of the transcription end site (Fig. 2) and mostly from 100 kilobases upstream of the TSS to 350 kilobases downstream of the transcription end site.

6. Discussion

In this paper, a robust and efficient penalized quasi-likelihood approach is introduced for model selection with non-polynomial dimensionality. It is shown that such an adaptive learning technique has a strong oracle property. As specific examples, two complementary methods of penalized composite L_1 – L_2 -regression and OWCQR are introduced and they are shown to have good efficiency and model selection consistency in ultrahigh dimensional space. Numerical studies show that our method is adaptive to unknown error distributions and outperforms the lasso (Tibshirani, 1996) and EWCQR (Zou and Yuan, 2008).

The penalized composite quasi-likelihood method can also be used in sure independence screening (Fan and Lv, 2008; Fan and Song, 2010) or the iterated version (Fan *et al.*, 2009), resulting in robust variable screening and selection. In this case, the marginal regression coefficients or contributions will be ranked and thresholded (Fan and Lv, 2008; Fan and Song, 2010). It can also be applied to the aggregation problems of classification (Bickel *et al.*, 2009) where the usual L_2 risk function could be replaced with a composite quasi-likelihood function. The idea can also be used to choose the loss functions in machine learning. For example, one can adaptively combine the hinge loss function in the support vector machine, the exponential loss in AdaBoost and the logistic loss function in logistic regression to yield a more efficient classifier.

Acknowledgements

This research was partially supported by National Science Foundation grants DMS-0704337 and DMS-0714554 and National Institutes of Health grant R01-GM072611. The bulk of the work was conducted while Weiwei Wang was a postdoctoral fellow at Princeton University.

Appendix A: Regularity conditions

Let D_k be the set of discontinuity points of $\psi_k(t)$, which is a subgradient of ρ_k . Assume that the distribution of error terms F_ε is sufficiently smooth that $F_\varepsilon(\cup_{k=1}^K D_k) = 0$. Additional regularity conditions on ψ_k are needed, as in Bai *et al.* (1992).

Condition 1. The function ψ_k satisfies $E\{\psi_k(\varepsilon_1 + c)\} = a_k c + o(|c|)$ as $|c| \rightarrow 0$, for some $a_k > 0$. For sufficiently small $|c|$, $g_{kl}(c) = E[\{\psi_k(\varepsilon_1 + c) - \psi_k(\varepsilon_1)\}\{\psi_l(\varepsilon_1 + c) - \psi_l(\varepsilon_1)\}]$ exists and is continuous at $c = 0$, where $k, l = 1, \dots, K$. The error distribution satisfies the Cramér condition $E\{|\psi_w(\varepsilon_i)|^m\} \leq m!RK^{m-2}$, for some constants R and K .

This condition implies that $E\{\psi_k(\varepsilon_i)\} = 0$, which is an unbiased score function of parameter β . It also implies that $E\{\partial\psi_k(\varepsilon_i)\} = a_k$ exists. The following two conditions are important for establishing sparsity properties of parameter $\hat{\beta}_w$ by controlling the penalty weighting scheme \mathbf{d} and the regularization parameter λ_n .

Condition 2. Assume that $D_n = \max\{d_j : j \in \mathcal{M}_*\} = o(n^{\alpha_1 - \alpha_0/2})$ and $\lambda_n D_n = O(n^{-(1+\alpha_0)/2})$. In addition, $\liminf \min\{d_j : j \in \mathcal{M}_*^c\} > 0$.

The first statement is to ensure that the bias term in theorem 2 is negligible. It is needed to control the bias due to the convex penalty. The second requirement is to make sure that the weights \mathbf{d} in the second part are uniformly large so that the vanishing coefficients are estimated as 0. It can also be regarded as a normalization condition, since the actual weights in the penalty are $\{\lambda_n d_j\}$.

The lasso estimator will not satisfy the first requirement of condition 2 unless λ_n is small and $\alpha_1 > \alpha_0/2$. Nevertheless, under the sparse representation condition (Zhao and Yu, 2006), Fan and Lv (2010) showed that, with probability tending to 1, the lasso estimator is model selection consistent with $\|\hat{\beta}_1 - \beta_1^*\|_\infty = O\{n^{-\gamma} \log(n)\}$, when the minimum signal $\beta_n^* = \min\{|\beta_j^*|, j \in \mathcal{M}_*\} \geq n^{-\gamma} \log(n)$. They also showed that the same result holds for the SCAD-type estimators under weaker conditions. Using one of them as the initial estimator, the weight $d_j = \gamma_\lambda(\hat{\beta}_j^0)/\lambda$ in problem (8) would satisfy condition 2, on a set with probability

tending to 1. This is due to the fact that with $\gamma_\lambda(\cdot)$ given by expression (6), for $j \in \mathcal{M}_*^c$, $d_j = \gamma_\lambda(0)/\lambda = 1$, whereas, for $j \in \mathcal{M}_*$, $d_j \leq \gamma_\lambda(\beta_n^*/2)/\lambda = 0$, as long as $\beta_n^* \gg n^{-\gamma} \log(n) = O(\lambda_n)$. In other words, the results of theorems 1 and 2 are applicable to the penalized estimator (8) with data-driven weights.

Condition 3. The regularization parameter $\lambda_n \gg n^{-1/2 + (\alpha_0 - 2\alpha_1)_+ / 2 + \alpha_2}$, where parameter α_1 is defined in condition 5 and $\alpha_2 \in [0, \frac{1}{2}]$, is a constant, bounded by the restriction in condition 4.

We use the following notation throughout the proof. Let \mathbf{B} be a matrix. Denote by $\lambda_{\min}(\mathbf{B})$ and $\lambda_{\max}(\mathbf{B})$ the minimum and maximum eigenvalue of the matrix \mathbf{B} when it is a square symmetric matrix. Let $\|\mathbf{B}\| = \lambda_{\max}^{1/2}(\mathbf{B}^T \mathbf{B})$ be the operator norm and $\|\mathbf{B}\|_\infty$ the largest absolute value of the elements in \mathbf{B} . As a result, $\|\cdot\|$ is the Euclidean norm when applied to a vector. Define $\|\mathbf{B}\|_{2, \infty} = \max_{\|\mathbf{v}\|_2=1} (\|\mathbf{B}\mathbf{v}\|_\infty)$.

Condition 4. The matrix $\mathbf{S}^T \mathbf{S}$ satisfies $C_1 n \leq \lambda_{\min}(\mathbf{S}^T \mathbf{S}) \leq \lambda_{\max}(\mathbf{S}^T \mathbf{S}) \leq C_2 n$ for some positive constants C_1 and C_2 . There exists $\xi > 0$ such that

$$\sum_{i=1}^n (\|\mathbf{S}_i\|/n^{1/2})^{2+\xi} \rightarrow 0,$$

where \mathbf{S}_i^T is the i th row of \mathbf{S} . Furthermore, assume that the design matrix satisfies $\|\mathbf{X}\|_\infty = O(n^{1/2 - (\alpha_0 - 2\alpha_1)_+ / 2 - \alpha_2})$ and $\max_{j \notin \mathcal{M}_*} (\|\mathbf{X}_j^*\|^2) = O(n)$, where \mathbf{X}_j^* is the j th column of \mathbf{X} .

Condition 5. Assume that

$$\begin{aligned} \sup_{\beta \in \mathcal{B}(\beta_1^*, \beta_n^*)} [\|\mathbf{Q} \text{diag}\{\partial \psi_{\mathbf{w}}(\beta)\} \mathbf{S}\|_{2, \infty}] &= O(n^{1-\alpha_1}), \\ \max_{\beta \in \mathcal{B}(\beta_1^*, \beta_n^*)} (\lambda_{\min}^{-1}[\mathbf{S}^T \text{diag}\{\partial \psi_{\mathbf{w}}(\beta)\} \mathbf{S}]) &= O_P(n^{-1}), \end{aligned}$$

where $\mathcal{B}(\beta_1^*, \beta_n^*)$ is an s -dimensional ball centred at β_1^* with radius β_n^* and $\text{diag}\{\partial \psi_{\mathbf{w}}(\beta)\}$ is the diagonal matrix with i th element equal to $\partial \psi_{\mathbf{w}}(Y_i - \mathbf{S}_i^T \beta)$.

Appendix B: Lemmas

Recall that $\mathbf{X} = (\mathbf{S}, \mathbf{Q})$ and $\mathcal{M}_* = \{1, \dots, s\}$ is the true model.

Lemma 1. Under conditions 2 and 4, the penalized quasi-likelihood $L_n(\beta)$ defined by equation (9) has a unique global minimizer $\hat{\beta} = (\hat{\beta}_1^T, \mathbf{0}^T)^T$, if

$$\sum_{i=1}^n \psi_{\mathbf{w}}(Y_i - \mathbf{X}_i^T \hat{\beta}) \mathbf{S}_i + n \lambda_n \mathbf{d}_{\mathcal{M}_*} \circ \text{sgn}(\hat{\beta}_1) = \mathbf{0}, \quad (27)$$

$$\|\mathbf{z}(\hat{\beta})\|_\infty < n \lambda_n, \quad (28)$$

where

$$\mathbf{z}(\hat{\beta}) = \mathbf{d}_{\mathcal{M}_*^c}^{-1} \circ \sum_{i=1}^n \psi_{\mathbf{w}}(Y_i - \mathbf{X}_i^T \hat{\beta}) \mathbf{Q}_i,$$

$\mathbf{d}_{\mathcal{M}_*}$ and $\mathbf{d}_{\mathcal{M}_*^c}$ stand for the subvectors of \mathbf{d} , consisting of its first s elements and the last $p - s$ elements respectively and sgn and ‘ \circ ’ (the Hadamard product) in condition (27) are taken co-ordinatewise. Conversely, if $\hat{\beta}$ is a global minimizer of $L_n(\beta)$, then condition (27) holds and inequality (28) holds with strict inequality replaced with non-strict inequality.

Proof. Under conditions 2 and 4, $L_n(\beta)$ is strictly convex. Necessary conditions (27) and (28) are direct consequences of the Karush–Kuhn–Tucker conditions of optimality. The sufficient condition follows from similar arguments to those in the proof of theorem 1 in Fan and Lv (2010) and the strict convexity of the function $L(\beta)$.

Lemma 2. Under conditions 1–5 we have that

$$\|\hat{\beta}^0 - \beta^*\|_2 = O_P\{\sqrt{(s/n)} + \lambda_n \|\mathbf{d}_0\|\},$$

where \mathbf{d}_0 is the subvector of \mathbf{d} , consisting of its first s elements.

Proof. Since $\hat{\beta}_2^\circ = \beta_2^* = 0$, we need to consider only the subvector of the first s components. Let us first show the existence of the biased oracle estimator. We can restrict our attention to the s -dimensional subspace $\{\beta \in \mathbb{R}^p : \beta_{\mathcal{M}_0^c} = \mathbf{0}\}$. Our aim is to show that

$$P[\inf_{\|\mathbf{u}\|=1} \{L_n(\beta_1^* + \gamma_n \mathbf{u}, \mathbf{0}) > L_n(\beta^*)\}] \rightarrow 1, \quad (29)$$

for sufficiently large γ_n . Here, there is a minimizer inside the ball $\|\beta_1 - \beta_1^*\| < \gamma_n$, with probability tending to 1. Using the strict convexity of $L_n(\beta)$, this minimizer is the unique global minimizer.

By the Taylor series expansion at $\gamma_n = 0$, we have

$$L_n(\beta_1^* + \gamma_n \mathbf{u}, \mathbf{0}) - L_n(\beta_1^*, \mathbf{0}) = T_1 + T_2,$$

where

$$\begin{aligned} T_1 &= -\gamma_n \sum_{i=1}^n \psi_{\mathbf{w}}(\varepsilon_i) \mathbf{S}_i^T \mathbf{u} + \frac{1}{2} \gamma_n^2 \sum_{i=1}^n \partial \psi_{\mathbf{w}}(\varepsilon_i - \bar{\gamma}_n \mathbf{S}_i^T \mathbf{u}) (\mathbf{S}_i^T \mathbf{u})^2 \\ &= -I_1 + I_2, \\ T_2 &= n \lambda_n \sum_{j=1}^s d_j (|\beta_j^* + \gamma_n u_j| - |\beta_j^*|); \end{aligned}$$

here $\bar{\gamma}_n \in [0, \gamma_n]$. By the Cauchy–Schwarz inequality,

$$|T_2| \leq n \gamma_n \lambda_n \|\mathbf{d}_0\| \|\mathbf{u}\| = n \gamma_n \lambda_n \|\mathbf{d}_0\|.$$

For all $\|\mathbf{u}\| = 1$, we have

$$|I_1| \leq \gamma_n \left\| \sum_{i=1}^n \psi_{\mathbf{w}}(\varepsilon_i) \mathbf{S}_i \right\|$$

and

$$E \left\| \sum_{i=1}^n \psi_{\mathbf{w}}(\varepsilon_i) \mathbf{S}_i \right\| \leq \left[E \{ \psi_{\mathbf{w}}^2(\varepsilon) \} \sum_{i=1}^n \|\mathbf{S}_i\|^2 \right]^{1/2} = [E \{ \psi_{\mathbf{w}}^2(\varepsilon) \} \text{tr}(\mathbf{S}^T \mathbf{S})]^{1/2},$$

which is of order $O\{\sqrt{ns}\}$ by condition 4. Hence, $I_1 = O_p\{\gamma_n \sqrt{ns}\}$ uniformly in \mathbf{u} .

Finally, we deal with I_2 . Let $H_i(c) = \inf_{|v| \leq c} \{\partial \psi_{\mathbf{w}}(\varepsilon_i - v)\}$. By lemma 3.1 of Portnoy (1984), we have

$$\begin{aligned} I_2 &\geq \gamma_n^2 \sum_{i=1}^n H_i(\gamma_n |\mathbf{S}_i^T \mathbf{u}|) (\mathbf{S}_i^T \mathbf{u})^2 \\ &\geq c \gamma_n^2 n, \end{aligned}$$

for a positive constant c . Combining all the above results, we have with probability tending to 1 that

$$L_n(\beta_1^* + \gamma_n \mathbf{u}, \mathbf{0}) - L_n(\beta_1^*, \mathbf{0}) \geq n \gamma_n [c \gamma_n - O_p\{\sqrt{(s/n)}\} - \lambda_n \|\mathbf{d}_0\|],$$

where the right-hand side is larger than 0 when $\gamma_n = B\{\sqrt{(s/n)} + \lambda_n \|\mathbf{d}_0\|\}$ for a sufficiently large $B > 0$. Since the objective function is strictly convex, there is a unique minimizer $\hat{\beta}_1^\circ$ such that

$$\|\hat{\beta}_1^\circ - \beta_1^*\| = O_p\{\sqrt{(s/n)} + \lambda_n \|\mathbf{d}_0\|\}.$$

Lemma 3. Under the conditions of theorem 2,

$$(\mathbf{b}^T \mathbf{A}_n \mathbf{b})^{-1/2} \sum_{i=1}^n \psi_{\mathbf{w}}(\varepsilon_i) \mathbf{b}^T \mathbf{S}_i \xrightarrow{\mathcal{D}} \mathcal{N}(0, 1) \quad (30)$$

where $\mathbf{A}_n = E\{\psi_{\mathbf{w}}^2(\varepsilon) \mathbf{S}^T \mathbf{S}\}$.

Proof. By condition 1, since \mathbf{S}_i is independent of $\psi_{\mathbf{w}}(\varepsilon_i)$, we have $E\{\psi_{\mathbf{w}}(\varepsilon_i)\mathbf{S}_i\} = 0$, and

$$\text{var}\left\{(\mathbf{b}^T \mathbf{A}_n \mathbf{b})^{-1/2} \sum_{i=1}^n \psi_{\mathbf{w}}(\varepsilon_i) \mathbf{b}^T \mathbf{S}_i\right\} = 1. \quad (31)$$

To complete the proof of lemma 3, we need to check only the Lyapunov condition. By condition 1, $E\{|\psi_{\mathbf{w}}(\varepsilon)|^{2+\xi}\} < \infty$. Furthermore, condition 4 implies that

$$\mathbf{b}^T \mathbf{A}_n \mathbf{b} = E\{\psi_{\mathbf{w}}^2(\varepsilon) \mathbf{b}^T \mathbf{S} \mathbf{S}^T \mathbf{b}\} \geq c_1 n,$$

for a positive constant c_1 . Using these together with the Cauchy–Schwarz inequality, we have

$$\begin{aligned} \sum_{i=1}^n E\{|\mathbf{b}^T \mathbf{A}_n \mathbf{b}|^{-1/2} \psi_{\mathbf{w}}(\varepsilon_i) \mathbf{b}^T \mathbf{S}_i|^{2+\xi}\} &= O(1) \sum_{i=1}^n |n^{-1/2} \mathbf{b}^T \mathbf{S}_i|^{2+\xi} \\ &= O(1) \sum_{i=1}^n |n^{-1/2} \|\mathbf{S}_i\||^{2+\xi}, \end{aligned}$$

which tends to 0 by condition 4. This completes the proof.

The following Bernstein inequality can be found in lemma 2.2.11 of van der Vaart and Wellner (1996).

Lemma 4. Let Y_1, \dots, Y_n be independent random variables with zero mean such that $E(|Y_i|^m) \leq m! M^{m-2} v_i / 2$, for every $m \geq 2$ (and all i) and some constants M and v_i . Then

$$P(|Y_1 + \dots + Y_n| > t) \leq 2 \exp\left\{-\frac{t^2}{2(v + Mt)}\right\},$$

for $v \geq v_1 + \dots + v_n$.

Then the following inequality (32) is a consequence of the previous Bernstein inequality. Let $\{Y_i\}$ satisfy the condition of lemma 4 with $v_i \equiv 1$. For a given sequence $\{a_i\}$, $E(|a_i Y_i|^m) \leq m! |a_i| M^{m-2} a_i^2 / 2$. A direct application of lemma 4 yields

$$P(|a_1 Y_1 + \dots + a_n Y_n| > t) \leq 2 \exp\left[-\frac{t^2}{2\left\{\sum_{i=1}^n a_i^2 + M \max_i(|a_i|t)\right\}}\right]. \quad (32)$$

Appendix C: Proofs of theorems

C.1. Proof of theorem 1

We need to show only that $\hat{\beta}^0$ is the unique minimizer of $L(\beta)$ in \mathbb{R}^p on a set Ω_n which has a probability tending to 1. Since $\hat{\beta}_1^0$ already satisfies condition (27), we only need to check condition (28).

We now define the set Ω_n . Let

$$\xi = (\xi_1, \dots, \xi_p)^T = \sum_{i=1}^n \psi_{\mathbf{w}}(Y_i - \mathbf{X}_i^T \beta^*) \mathbf{X}_i$$

and consider the event $\Omega_n = \{\|\xi_{\mathcal{M}_s^c}\|_{\infty} \leq u_n \sqrt{n}\}$ with u_n being chosen later. Then, by condition 1 and Bernstein's inequality, it follows directly from condition (32) that

$$P(|\xi_j| > t) \leq 2 \exp\left\{-\frac{t^2}{2(\|\mathbf{X}_j^*\|^2 R + tK\|\mathbf{X}_j^*\|_{\infty})}\right\},$$

where \mathbf{X}_j^* is the j th column of \mathbf{X} . Taking $t = u_n \sqrt{n}$, we have

$$P(|\xi_j| > u_n \sqrt{n}) \leq 2 \exp\left\{-\frac{u_n^2}{2(R\|\mathbf{X}_j^*\|^2/n + Ku_n\|\mathbf{X}_j^*\|_{\infty}/\sqrt{n})}\right\} \leq \exp(-cu_n^2), \quad (33)$$

for some positive constant $c > 0$, by condition 4. Thus, by using the union bound, we conclude that

$$P(\Omega_n) \geq 1 - \sum_{j \in \mathcal{M}_s^c} P(|\xi_j| > u_n \sqrt{n}) \geq 1 - 2(p-s) \exp(-cu_n^2).$$

We now check whether condition (27) holds on the set Ω_n . Let $\psi_{\mathbf{w}}(\beta)$ be the n -dimensional vector with i th element $\psi_{\mathbf{w}}(Y_i - \mathbf{X}_i^T \beta)$. Then, by condition 2

$$\begin{aligned} \|\mathbf{z}(\hat{\beta}^0)\|_{\infty} &\leq \|\mathbf{d}_{\mathcal{M}_n^c}^{-1} \circ \xi_{\mathcal{M}_n^c}\|_{\infty} + \|\mathbf{d}_{\mathcal{M}_n^c}^{-1} \circ \mathbf{Q}^T \{\psi_{\mathbf{w}}(\hat{\beta}^0) - \psi_{\mathbf{w}}(\beta^*)\}\|_{\infty} \\ &= O[n^{1/2} u_n + \|\mathbf{Q}^T \text{diag}\{\partial \psi_{\mathbf{w}}(\mathbf{v})\} \mathbf{S}(\hat{\beta}_1^0 - \beta_1^*)\|_{\infty}] \end{aligned} \quad (34)$$

where \mathbf{v} lies between $\hat{\beta}^0$ and β_1^* . By condition 5, the second term in equation (34) is bounded by

$$O(n^{1-\alpha_1}) \|\hat{\beta}_1^0 - \beta_1^*\| = O_P[n^{1-\alpha_1} \{\sqrt{(s/n)} + \lambda_n \|\mathbf{d}_0\|\}],$$

where the equality follows from lemma 2. By the choice of parameters,

$$(n\lambda_n)^{-1} \|\mathbf{z}(\hat{\beta}^0)\|_{\infty} = O\{n^{-1/2} \lambda_n^{-1} (u_n + n^{(\alpha_0-2\alpha_1)/2}) + D_n n^{\alpha_0/2-\alpha_1}\} = o(1),$$

by taking $u_n = n^{(\alpha_0-2\alpha_1)/2+\alpha_2}$. Hence, by lemma 1, $\hat{\beta}^0$ is the unique global minimizer.

C.2. Proof of theorem 2

By theorem 1, $\hat{\beta}_{w1} = \hat{\beta}_1^0$ almost surely. It follows from lemma 2 that

$$\|\hat{\beta}_{w1} - \beta_1^*\| = O_P\{(\lambda_n D_n + 1/\sqrt{n})\sqrt{s}\}.$$

This establishes the first part of the theorem 2.

Let

$$Q_n(\beta_1) = \sum_{i=1}^n \psi_{\mathbf{w}}(Y_i - \mathbf{S}_i^T \beta_1) \mathbf{S}_i.$$

By Taylor's expansion at the point β_1^* , we have

$$Q_n(\hat{\beta}_{w1}) = Q_n(\beta_1^*) + \partial Q_n(\mathbf{v})(\hat{\beta}_{w1} - \beta_1^*),$$

where \mathbf{v} lies between the points $\hat{\beta}_{w1}$ and β_1^* and

$$\partial Q_n(\mathbf{v}) = - \sum_{i=1}^n \partial \psi_{\mathbf{w}}(Y_i - \mathbf{S}_i^T \mathbf{v}) \mathbf{S}_i \mathbf{S}_i^T. \quad (35)$$

By lemma 2, $\|\mathbf{v} - \beta_1^*\| \leq \|\hat{\beta}_{w1} - \beta_1^*\| = o_P(1)$.

By using condition (28), we have

$$Q_n(\hat{\beta}_{w1}) + n\lambda_n \mathbf{d}_0 \circ \text{sgn}(\hat{\beta}_{w1}) = 0,$$

or equivalently

$$\hat{\beta}_{w1} - \hat{\beta}_1^* = -\partial Q_n(\mathbf{v})^{-1} Q_n(\beta_1^*) - \partial Q_n(\mathbf{v})^{-1} n\lambda_n \mathbf{d}_0 \circ \text{sgn}(\hat{\beta}_{w1}). \quad (36)$$

Note that $\|\mathbf{d}_0 \circ \text{sgn}(\hat{\beta}_{w1})\| = \|\mathbf{d}_0\|$. We have, for any vector \mathbf{u} ,

$$|\mathbf{u}^T \partial Q_n(\mathbf{v})^{-1} \mathbf{d}_0 \circ \text{sgn}(\hat{\beta}_{w1})| \leq \|\partial Q_n(\mathbf{v})^{-1}\| \|\mathbf{u}\| \|\mathbf{d}_0\|.$$

Consequently, for any unit vector \mathbf{b} ,

$$\begin{aligned} \|\mathbf{b}^T (\mathbf{S}^T \mathbf{S})^{1/2} \partial Q_n(\mathbf{v})^{-1} \mathbf{d}_0 \circ \text{sgn}(\hat{\beta}_{w1})\| &\leq \lambda_{\max}^{1/2} (\mathbf{S}^T \mathbf{S}) \lambda_{\min}^{-1} \{\partial Q_n(\mathbf{v})\} D_n \sqrt{s} \\ &= O_P\{D_n \sqrt{(s/n)}\}, \end{aligned}$$

by using conditions 4 and 5. This shows that the second term in equation (35), when multiplied by the vector $\mathbf{b}^T (\mathbf{S}^T \mathbf{S})^{1/2}$, is of order

$$O_P\{\lambda_n D_n \sqrt{(sn)}\} = o_P(1),$$

by condition 2. Therefore, we need to establish the asymptotic normality of the first term in equation (35). This term is identical to the situation that was dealt with by Portnoy (1985). Using his result, the second conclusion of theorem 2 follows. This completes the proof.

C.3. Proof of theorem 3

First, by Taylor series expansion,

$$\Phi_{n,w}(\hat{\beta}_1) = \Phi_{n,w}(\beta_1^*) + \Omega_{n,w}(\bar{\beta}_1)(\hat{\beta}_1 - \beta_1^*), \quad (37)$$

where $\bar{\beta}_1$ lies between β_1^* and $\hat{\beta}_1$. Consequently,

$$\|\bar{\beta}_1 - \hat{\beta}_1\| \leq \|\beta_1^* - \hat{\beta}_1\| = o_P(1).$$

By the definition of the one-step estimator (14) and equation (37), we have

$$\hat{\beta}_{w1}^{\text{os}} - \beta_1^* = \Omega_{n,w}(\hat{\beta}_1)^{-1} \Phi_{n,w}(\beta_1^*) + \mathbf{R}_n, \quad (38)$$

where

$$\mathbf{R}_n = \Omega_{n,w}(\hat{\beta}_1)^{-1} \{\Omega_{n,w}(\hat{\beta}_1) - \Omega_{n,w}(\bar{\beta}_1)\}(\hat{\beta}_1 - \beta_1^*).$$

We first deal with the remainder term. Note that

$$\|\mathbf{R}_n\| \leq \|\Omega_{n,w}(\hat{\beta}_1)^{-1}\| \|\Omega_{n,w}(\hat{\beta}_1) - \Omega_{n,w}(\bar{\beta}_1)\| \|\hat{\beta}_1 - \beta_1^*\| \quad (39)$$

and

$$\Omega_{n,w}(\hat{\beta}_1) - \Omega_{n,w}(\bar{\beta}_1) = \sum_{i=1}^n f_i(\hat{\beta}_1, \bar{\beta}_1) \mathbf{S}_i \mathbf{S}_i^T, \quad (40)$$

where $f_i(\hat{\beta}_1, \bar{\beta}_1) = \partial\psi(Y_i - \mathbf{S}_i^T \hat{\beta}_1) - \partial\psi(Y_i - \mathbf{S}_i^T \bar{\beta}_1)$. By Lipschitz continuity, we have

$$|f_i(\hat{\beta}_1, \bar{\beta}_1)| \leq C \|\mathbf{S}_i\| \|\hat{\beta}_1 - \bar{\beta}_1\|,$$

where C is the Lipschitz coefficient of $\partial\psi_w(\cdot)$. Let \mathbf{I}_s be the identity matrix of order s and

$$b_n = \lambda_{\max} \left(\sum_{i=1}^n \|\mathbf{S}_i\| \mathbf{S}_i \mathbf{S}_i^T \right).$$

By equation (40), we have

$$\Omega_{n,w}(\hat{\beta}_1) - \Omega_{n,w}(\bar{\beta}_1) \leq C \|\hat{\beta}_1 - \bar{\beta}_1\| \sum_{i=1}^n \|\mathbf{S}_i\| \mathbf{S}_i \mathbf{S}_i^T \leq C \|\hat{\beta}_1 - \bar{\beta}_1\| b_n \mathbf{I}_s.$$

Hence, none of the eigenvalues of the matrix is larger than $C \|\hat{\beta}_1 - \bar{\beta}_1\| b_n$. Similarly, by equation (40),

$$\Omega_{n,w}(\hat{\beta}_1) - \Omega_{n,w}(\bar{\beta}_1) \geq -C \|\hat{\beta}_1 - \bar{\beta}_1\| \sum_{i=1}^n \|\mathbf{S}_i\| \mathbf{S}_i \mathbf{S}_i^T \geq -C \|\hat{\beta}_1 - \bar{\beta}_1\| b_n \mathbf{I}_s,$$

and all its eigenvalues should be at least $-C \|\hat{\beta}_1 - \bar{\beta}_1\| b_n$. Consequently,

$$\|\Omega_{n,w}(\hat{\beta}_1) - \Omega_{n,w}(\bar{\beta}_1)\| \leq C \|\hat{\beta}_1 - \bar{\beta}_1\| b_n.$$

By condition 5 and the assumption of $\hat{\beta}_1$, it follows from inequality (39) that

$$\|\mathbf{R}_n\| = O_P(s/n \times b_n/n) = O_P(s^{3/2}/n).$$

Thus, for any unit vector \mathbf{b} ,

$$\mathbf{b}^T (\mathbf{S}^T \mathbf{S})^{1/2} \mathbf{R}_n \leq \lambda_{\max}^{1/2} (\mathbf{S}^T \mathbf{S}) \|\mathbf{R}_n\| = O_P(s^{3/2}/n^{1/2}) = o_P(1).$$

The main term in equation (37) can be handled by using lemma 3 and the same method as Portnoy (1985). This completes the proof.

References

- Bai, Z. D., Rao, C. R. and Wu, Y. (1992) M -estimation of multivariate linear regression parameters under a convex discrepancy function. *Statist. Sin.*, **2**, 237–254.
- Bickel, P. J. (1973) On some analogues to linear combinations of order statistics in the linear model. *Ann. Statist.*, **1**, 597–616.
- Bickel, P. J., Ritov, Y. and Tsybakov, A. B. (2008) Hierarchical selection of variables in sparse high-dimensional regression. *Preprint arXiv:0801.1158v1*.

- Bickel, P. J., Ritov, Y. and Tsybakov, A. B. (2009) Simultaneous analysis of lasso and dantzig selector. *Ann. Statist.*, **37**, 1705–1732.
- Deutsch, S., Lyle, R., Dermitzakis, E. T., Attar, H., Subrahmanyam, L., Gehrig, C., Parand, L., Gagnebin, M., Rougemont, J., Jongeneel, C. V. and Antonarakis, S. E. (2005) Gene expression variation and expression quantitative trait mapping of human chromosome 21 genes. *Hum. Molec. Genet.*, **14**, 3741–3749.
- Efron, B., Hastie, T., Johnstone, I. and Tibshirani, R. (2004) Least angle regression. *Ann. Statist.*, **32**, 407–499.
- Fan, J. and Li, R. (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Statist. Ass.*, **96**, 1348–1360.
- Fan, J. and Lv, J. (2008) Sure independence screening for ultrahigh dimensional feature space (with discussion). *J. R. Statist. Soc. B*, **70**, 849–911.
- Fan, J. and Lv, J. (2010) Properties of non-concave penalized likelihood with NP-dimensionality. Submitted to *Ann. Statist.*
- Fan, J. and Peng, H. (2004) Nonconcave penalized likelihood with a diverging number of parameters. *Ann. Statist.*, **32**, 928–961.
- Fan, J., Samworth, R. and Wu, Y. (2009) Ultrahigh dimensional variable selection: beyond the linear model. *J. Mach. Learn. Res.*, **10**, 1829–1853.
- Fan, J. and Song, R. (2010) Sure independence screening in generalized linear models with NP-dimensionality. *Ann. Statist.*, to be published.
- Frank, I. E. and Friedman, J. H. (1993) A Statistical view of some chemometrics regression tools. *Technometrics*, **35**, 109–148.
- Friedman, J. H., Hastie, T., Hofling, H. and Tibshirani, R. (2008) Pathwise coordinate optimization. *Ann. Appl. Statist.*, **1**, 302–332.
- Huang, J., Horowitz, J. L. and Ma, S. (2008) Asymptotic properties of bridge estimators in sparse high-dimensional regression models. *Ann. Statist.*, **36**, 587–613.
- Huber, P. J. (1964) Robust estimation of location parameter. *Ann. Math. Statist.*, **35**, 73–101.
- Kim, Y., Choi, H. and Oh, H. (2008) Smoothly clipped absolute deviation on high dimensions. *J. Am. Statist. Ass.*, **103**, 1656–1673.
- Koenker, R. (1984) A note on L-estimates for linear models. *Statist. Probab. Lett.*, **2**, 323–325.
- Lehmann, E. L. (1983) *Theory of Point Estimation*, p. 506. New York: Wiley.
- Li, Y. and Zhu, J. (2008) L_1 -norm quantile regression. *J. Computnl Graph. Statist.*, **17**, 163–185.
- Portnoy, S. (1984) Asymptotic behavior of M -estimators of p regression parameters when p^2/n is large; I, Consistency. *Ann. Statist.*, **12**, 1298–1309.
- Portnoy, S. (1985) Asymptotic behavior of M -estimators of p regression parameters when p^2/n is large; II, Normal approximation. *Ann. Statist.*, **13**, 1403–1417.
- Tibshirani, R. (1996) Regression shrinkage and selection via the lasso. *J. R. Statist. Soc. B*, **58**, 267–288.
- van der Vaart, A. W. and Wellner, J. A. (1996) *Weak Convergence and Empirical Processes*, p. 508. New York: Springer.
- Wu, Y. and Liu, Y. (2009) Variable selection in quantile regression. *Statist. Sin.*, **37**, 801–817.
- Xie, H. and Huang, J. (2009) SCAD-penalized regression in high-dimensional partially linear models. *Ann. Statist.*, **37**, 673–696.
- Yuan, M. and Lin, Y. (2007) On the non-negative garrotte estimator. *J. R. Statist. Soc. B*, **69**, 143–161.
- Zhao, P. and Yu, B. (2006) On model selection consistency of Lasso. *J. Mach. Learn. Res.*, **7**, 2541–2567.
- Zou, H. (2006) The adaptive LASSO and its oracle properties. *J. Am. Statist. Ass.*, **101**, 1418–1429.
- Zou, H. and Li, R. (2008) One-step sparse estimates in nonconcave penalized likelihood models. *Ann. Statist.*, **36**, 1509–1533.
- Zou, H. and Yuan, M. (2008) Composite quantile regression and the oracle model selection theory. *Ann. Statist.*, **36**, 1108–1126.
- Zou, H. and Zhang, H. H. (2009) On the adaptive elastic-net with a diverging number of parameters. *Ann. Statist.*, **37**, 1733–1751.