

# OPTIMAL COMPUTATIONAL AND STATISTICAL RATES OF CONVERGENCE FOR SPARSE NONCONVEX LEARNING PROBLEMS

BY ZHAORAN WANG<sup>\*</sup>, HAN LIU<sup>\*</sup> AND TONG ZHANG<sup>†</sup>

*Princeton University<sup>\*</sup> and Rutgers University<sup>†</sup>*

We provide theoretical analysis of the statistical and computational properties of **penalized  $M$ -estimators** that can be formulated as the solution to a possibly nonconvex optimization problem. Many important estimators fall in this category, including least squares regression with nonconvex regularization, generalized linear models with nonconvex regularization, and sparse elliptical random design regression. For these problems, it is intractable to calculate the global solution due to the nonconvex formulation. In this paper, we propose an **approximate regularization path following method** for solving a variety of learning problems with nonconvex objective functions. Under a unified analytic framework, we simultaneously provide explicit statistical and computational rates of convergence for any local solution attained by the algorithm. Computationally, our algorithm attains a global geometric rate of convergence for calculating the full regularization path, which is optimal among all first-order algorithms. Unlike most existing methods that only attain geometric rates of convergence for one single regularization parameter, our algorithm calculates the full regularization path with the same iteration complexity. In particular, we provide a refined iteration complexity bound to sharply characterize the performance of each stage along the regularization path. Statistically, we provide sharp sample complexity analysis for all the approximate local solutions along the regularization path. In particular, our analysis improves upon existing results by providing a more refined sample complexity bound as well as an exact support recovery result for the final estimator. These results show that the final estimator attains an oracle statistical property due to the usage of nonconvex penalty.

**1. Introduction.** This paper considers the statistical and computational properties of a family of penalized  $M$ -estimators that can be formulated as

$$(1.1) \quad \hat{\beta}_\lambda \in \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \mathcal{L}(\beta) + \mathcal{P}_\lambda(\beta) \right\},$$

---

*AMS 2000 subject classifications:* Primary 62F30, 90C26; secondary 62J12, 90C52

*Keywords and phrases:* nonconvex regularized  $M$ -estimation, path following method, geometric computational rate, optimal statistical rate

where  $\mathcal{L}(\boldsymbol{\beta})$  is a loss function, while  $\mathcal{P}_\lambda(\boldsymbol{\beta})$  is a penalty function with regularization parameter  $\lambda$ . A familiar example is the Lasso estimator (Tibshirani, 1996), in which  $\mathcal{L}(\boldsymbol{\beta}) = \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2/(2n)$  and  $\mathcal{P}_\lambda(\boldsymbol{\beta}) = \lambda\|\boldsymbol{\beta}\|_1$ . Here  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times d}$  is the design matrix,  $\mathbf{y} = (y_1, \dots, y_n)^T \in \mathbb{R}^n$  is the response vector,  $\|\cdot\|_2$  is the Euclidean norm, and  $\|\boldsymbol{\beta}\|_1 = \sum_{j=1}^d |\beta_j|$  is the  $\ell_1$  norm of  $\boldsymbol{\beta}$ . In general, we prefer the settings where both the loss function  $\mathcal{L}(\boldsymbol{\beta})$  and the penalty term  $\mathcal{P}_\lambda(\boldsymbol{\beta})$  in (1.1) are convex, since convexity makes both statistical and computational analysis convenient.

Significant progress has been made on understanding convex penalized  $M$ -estimators (van de Geer, 2000, 2008; Rothman et al., 2008; Wainwright, 2009; Bickel et al., 2009; Zhang, 2009; Koltchinskii, 2009b; Raskutti et al., 2011; Negahban et al., 2012). Meanwhile, penalized  $M$ -estimators with nonconvex loss or penalty functions have recently attracted much interest because of their more attractive statistical properties. For example, unlike the  $\ell_1$  penalty, which induces significant estimation bias for parameters with large absolute values (Zhang and Huang, 2008), nonconvex penalties such as the smoothly clipped absolute deviation (SCAD) penalty (Fan and Li, 2001) and minimax concave penalty (MCP) (Zhang, 2010a) can eliminate this estimation bias and attain more refined statistical rates of convergence. As another example of penalized  $M$ -estimators with nonconvex loss functions, we consider a semiparametric variant of the penalized least squares regression. Recall that a penalized least squares regression estimator can be formulated as

$$\begin{aligned} \hat{\boldsymbol{\beta}}_\lambda &\in \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^d} \left\{ \frac{1}{2n} \|\mathbf{X}\boldsymbol{\beta} - \mathbf{y}\|_2^2 + \mathcal{P}_\lambda(\boldsymbol{\beta}) \right\} \\ &= \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^d} \left\{ \frac{1}{2} (1, -\boldsymbol{\beta}^T) \hat{\mathbf{S}} (1, -\boldsymbol{\beta}^T)^T + \mathcal{P}_\lambda(\boldsymbol{\beta}) \right\}, \end{aligned}$$

where  $\hat{\mathbf{S}} = (\mathbf{y}, \mathbf{X})^T (\mathbf{y}, \mathbf{X}) / n$  is the sample covariance matrix of a random vector  $(Y, \mathbf{X}^T)^T \in \mathbb{R}^{d+1}$ . When the design matrix  $\mathbf{X}$  contains heavy-tail data, we may resort to elliptical random design regression, which is a semiparametric extension of Gaussian random design regression. In detail, we replace the sample covariance matrix  $\hat{\mathbf{S}}$  with a possibly indefinite covariance matrix estimator  $\hat{\mathbf{K}}$  (to be defined in §2.2), which is more robust within the elliptical family. Since  $\hat{\mathbf{K}}$  does not guarantee to be positive semidefinite, the loss function  $\mathcal{L}(\boldsymbol{\beta}) = (1, -\boldsymbol{\beta}^T) \hat{\mathbf{K}} (1, -\boldsymbol{\beta}^T)^T / 2$  could be nonconvex.

Though the global solutions of these nonconvex  $M$ -estimators enjoy nice statistical properties, it is in general computationally intractable to obtain the global solutions. Instead, a more realistic approach is to directly leverage standard optimization procedures to obtain a local solution  $\hat{\boldsymbol{\beta}}_\lambda$  that satisfies

the first-order Karush-Kuhn-Tucker (KKT) condition

$$(1.2) \quad \mathbf{0} \in \partial \left\{ \mathcal{L}(\hat{\beta}_\lambda) + \mathcal{P}_\lambda(\hat{\beta}_\lambda) \right\},$$

where  $\partial(\cdot)$  denotes the subgradient operator.

In the context of least squares regression with nonconvex penalties, several numerical procedures have been proposed to find the local solutions, including local quadratic approximation (LQA) (Fan and Li, 2001), minorize-maximize (MM) algorithm (Hunter and Li, 2005), local linear approximation (LLA) (Zou and Li, 2008), concave convex procedure (CCCP) (Kim et al., 2008), and coordinate descent (Breheny and Huang, 2011; Mazumder et al., 2011). The theoretical properties of the local solutions obtained by these numerical procedures are in general unestablished. Only recently Zhang and Zhang (2012) showed that the gradient descent method initialized at a Lasso solution attains a unique local solution that has the same statistical properties as the global solution; Fan et al. (2014) proved that the LLA algorithm initialized with a Lasso solution attains a local solution with oracle statistical properties. The same conclusion was also obtained by Zhang (2010b); Zhang et al. (2013), where the LLA algorithm was referred to as multi-stage convex relaxation. In recent work, Wang et al. (2013) proposed a calibrated concave-convex procedure (CCCP) along with a high-dimensional BIC criterion that can achieve the oracle estimator. However, these works mainly focused on statistical recovery results, while the corresponding computational complexity results remain unclear. Also, they didn't consider nonconvex loss functions. In addition, their analysis relies on the assumption that all the computation (e.g., solving an optimization problem) can be carried out exactly, which is unrealistic in practice, since practical computational procedures can only attain finite numerical precision in finite time. Moreover, our method only requires the weakest possible minimum signal strength to attain the oracle estimator (Zhang and Zhang, 2012), while the procedures in Wang et al. (2013); Fan et al. (2014) rely on a stronger signal strength which is suboptimal. See §6 for a more detailed discussion.

In this paper, we propose an approximate regularization path following method for solving a general family of penalized  $M$ -estimators with possibly nonconvex loss or penalty functions. Our algorithm leverages the fast local convergence in the proximity of sparse solutions, which is also observed by Nesterov (2013); Wright et al. (2009); Agarwal et al. (2012); Xiao and Zhang (2013). More specifically, we consider a decreasing sequence of regularization parameters  $\{\lambda_t\}_{t=0}^N$ , where  $\lambda_0$  corresponds to an all-zero solution, and  $\lambda_N = \lambda_{\text{tgt}}$  is the target regularization parameter that ensures the obtained estimator to achieve the optimal statistical rate of convergence. For each  $\lambda_t$ , we construct

a sequence of local quadratic approximations of the loss function  $\mathcal{L}(\beta)$ , and utilize a variant of Nesterov's proximal-gradient method (Nesterov, 2013), which iterates over the updating step

$$(1.3) \quad \beta_t^{k+1} \leftarrow \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \mathcal{L}(\beta_t^k) + \nabla \mathcal{L}(\beta_t^k)^T (\beta - \beta_t^k) + \frac{L_t^k}{2} \|\beta - \beta_t^k\|_2^2 + \mathcal{P}_{\lambda_t}(\beta) \right\},$$

where  $k = 1, 2, \dots$ . Here  $\beta_t^k$  and  $L_t^k$  correspond to the  $k$ -th iteration of the proximal-gradient method for regularization parameter  $\lambda_t$ . Here  $L_t^k$  is chosen by an adaptive line-search method, which will be specified in §3.2. Let  $\hat{\beta}_{\lambda_t}$  be an exact local solution satisfying (1.2) with  $\lambda = \lambda_t$ . As illustrated in Figure 1, for each  $\lambda_t$ , our algorithm calculates an approximation  $\hat{\beta}_t$  of the exact local solution  $\hat{\beta}_{\lambda_t}$  up to certain optimization precision. Such approximate local solution  $\hat{\beta}_t$  guarantees to be sparse, and therefore falls into the fast convergence region corresponding to  $\lambda_{t+1}$ . Consequently, the resulting procedure achieves a geometric rate of convergence within each path following stage, and therefore attains a global geometric rate of convergence for calculating the entire regularization path. Moreover, we establish the nonasymptotic statistical rates of convergence and oracle properties for all the approximate and exact local solutions along the full regularization path.

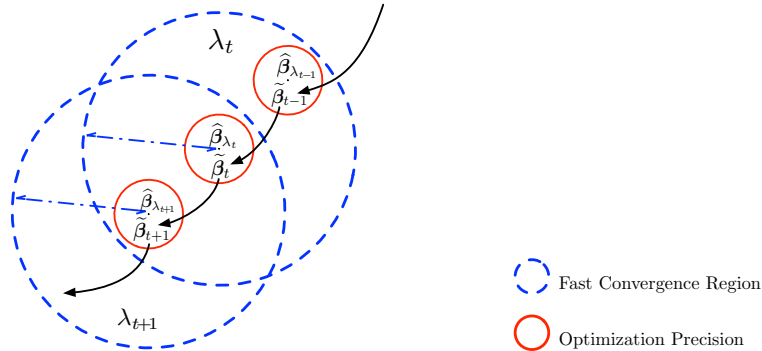


FIG 1. For regularization parameter  $\lambda_t$ ,  $\hat{\beta}_{\lambda_t}$  is an exact local solution satisfying (1.2) with  $\lambda = \lambda_t$ . Within the  $t$ -th path following stage, our algorithm achieves an approximate local solution  $\hat{\beta}_t$ , which approximates the exact local solution  $\hat{\beta}_{\lambda_t}$  up to certain optimization precision. Our approximate path following algorithm ensures that  $\hat{\beta}_t$  is sparse, and therefore falls into the fast convergence region corresponding to regularization parameter  $\lambda_{t+1}$ .

The idea of path following has been well-studied for sparse recovery problems (Efron et al., 2004; Hastie et al., 2005; Park and Hastie, 2007; Zhao and Yu, 2007; Rosset and Zhu, 2007; Friedman et al., 2010; Mazumder et al., 2011; Breheny and Huang, 2011; Xiao and Zhang, 2013; Mairal and Yu, 2012). Compared with these previous works, we consider a broader family

of nonconvex  $M$ -estimators, including nonconvex penalty functions, such as SCAD and MCP, as well as nonconvex loss functions, such as semiparametric elliptical design loss. Moreover, we provide sharp computational and statistical analysis for all the approximate and exact local solutions attained by the proposed approximate path following method along the regularization path.

The contributions of this paper are two folds:

- Computationally, we propose an optimization algorithm that ensures a global geometric rate of convergence for nonconvex sparse learning problems. In detail, recall that  $N$  is the total number of path following stages. Within the  $N$ -th path following stage, we denote by  $\epsilon_{\text{opt}}$  the desired optimization precision of the approximate local solution  $\tilde{\beta}_N$ . We need no more than a logarithmic number of the proximal-gradient update iterations defined in (1.3) to calculate the entire path:

$$\text{Total \# of proximal-gradient iterations} \leq C \log \left( \frac{1}{\epsilon_{\text{opt}}} \right),$$

where  $C > 0$  is a constant. This global geometric rate of convergence is optimal among all first-order methods, because it attains the lower bound for first-order methods on strongly convex and smooth objective function (Nesterov, 2004, Theorem 2.1.12), which is a subclass of the possibly nonconvex objective functions considered in this paper.

- Statistically, we prove that along the full regularization path, all the approximate local solutions obtained by our algorithm enjoy desirable statistical rates of convergence for estimating the true parameter vector  $\beta^*$ . In detail, let  $s^*$  be the number of nonzero entries of  $\beta^*$ , the approximate local solution  $\tilde{\beta}_t$ 's satisfy

$$(1.4) \quad \|\tilde{\beta}_t - \beta^*\|_2 \leq C \lambda_t \sqrt{s^*}, \quad \text{for } t = 1, \dots, N$$

with high probability. In particular, within the  $N$ -th path following stage, we have  $\lambda_N = \lambda_{\text{tgt}} = C' \sqrt{\log d/n}$ . Here  $C$  and  $C'$  are positive constants that do not dependent on  $d$  and  $n$ . In the  $d \gg n$  regime, the final approximate local solution  $\tilde{\beta}_N$  achieves the optimal statistical rate of convergence. Furthermore, we prove that, within the  $t$ -th path following stage, the iterative solution sequence  $\{\beta_t^k\}_{k=0}^\infty$  produced by (1.3) converges towards a unique exact local solution  $\beta_{\lambda_t}$ , which enjoys a more refined oracle statistical property. More specifically, let  $s_1^*$  be the number of “large” nonzero coefficients of  $\beta^*$  and  $s_2^* = s^* - s_1^*$  be the number of “small” nonzero coefficients (detailed definitions of  $s_1^*$

and  $s_2^*$  are provided in Theorem 4.8), we have

$$(1.5) \quad \|\hat{\beta}_{\lambda_t} - \beta^*\|_2 \leq C\sqrt{\frac{s_1^*}{n}} + C'\sqrt{s_2^*}\lambda_t, \quad \text{for } t = 1, \dots, N$$

with high probability. In particular, for the final stage we have  $\lambda_N = \lambda_{\text{tgt}} = C''\sqrt{\log d/n}$ . Here  $C$ ,  $C'$  and  $C''$  are positive constants. Note that the oracle statistical property in (1.5) is significantly sharper than the rate of convergence in (1.4), e.g., when  $s^* = s_1^*$  and  $t = N$ , the right-hand side of (1.4) is of the order of  $\sqrt{s^* \log d/n}$ , while the right-hand side of (1.5) is of the order of  $\sqrt{s^*/n}$ . Moreover, we prove that when the absolute values of the nonzero coefficients of  $\beta^*$  are larger than  $C'''\sqrt{\log d/n}$ ,  $\hat{\beta}_{\lambda_t}$  exactly recovers the support of  $\beta^*$ , i.e.,

$$\text{supp}(\hat{\beta}_{\lambda_t}) = \text{supp}(\beta^*).$$

In summary, our joint analysis of the statistical and computational properties provides a theoretical characterization of the entire regularization path.

In independent work, Loh and Wainwright (2013) discussed similar problems. In detail, they provided sufficient conditions under which local optima have desired theoretical properties, and verified that the approximate local solution attained by the composite gradient descent method satisfies these conditions. Our work differs from theirs in three aspects:

- (i) Our statistical recovery result in (1.4) covers all the approximate local solutions along the entire regularization path. They provided a similar statistical result, but only for the target regularization parameter, i.e.,  $\lambda_N = \lambda_{\text{tgt}}$  in (1.4).
- (ii) As results of independent interest, we prove the oracle statistical properties of the exact local solutions along the regularization path, including the refined statistical rates of convergence in (1.5) and the guarantee of exact support recovery, while they didn't provide such results. Since the statistical result in (1.4) is also achievable using convex regularization, e.g., the  $\ell_1$  penalty, these oracle properties are essential for justifying the benefits of using nonconvex penalty functions.
- (iii) Our analysis technique is different from theirs. In detail, our statistical analysis is embedded in the analysis of the optimization procedure. In particular, we provide fine-grained analysis of the sparsity pattern of all the intermediate solutions obtained from the proximal-gradient iterations. In contrast, they provided characterizations of local solutions under a global restricted strongly convex/smoothness condition.

The rest of this paper is organized as follows. First we briefly introduce some useful notation. In §2 we introduce  $M$ -estimators with possibly noncon-

vex loss and penalty functions. In §3 we present the proposed approximate regularization path following method. In §4 we present the main theoretical results on the computational efficiency and statistical accuracy of the proposed procedure. In §5 we prove the theoretical results in §4. In §6 we provide a detailed comparison between our method and the existing nonconvex procedures. Numerical results are presented in §7.

**Notation:** For  $q \in [1, +\infty)$ , the  $\ell_q$  norm of  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)^T \in \mathbb{R}^d$  is denoted by  $\|\boldsymbol{\beta}\|_q = (\sum_{j=1}^d |\beta_j|^q)^{1/q}$ . Specifically, we define  $\|\boldsymbol{\beta}\|_\infty = \max_{1 \leq j \leq d} \{|\beta_j|\}$  and  $\|\boldsymbol{\beta}\|_0 = \text{card}\{\text{supp}(\boldsymbol{\beta})\}$ , where  $\text{supp}(\boldsymbol{\beta}) = \{j : \beta_j \neq 0\}$  and  $\text{card}\{\cdot\}$  is the cardinality of a set. Correspondingly, we denote the  $\ell_q$  ball  $\{\boldsymbol{\beta} : \|\boldsymbol{\beta}\|_q \leq R\}$  by  $B_q(R)$ . For a set  $S$ , we denote its cardinality by  $|S|$  and its complement by  $\bar{S}$ . For  $S, \bar{S} \subseteq \{1, \dots, d\}$ , we define  $\boldsymbol{\beta}_S \in \mathbb{R}^d$  and  $\boldsymbol{\beta}_{\bar{S}} \in \mathbb{R}^d$  as  $(\boldsymbol{\beta}_S)_j = \mathbf{1}(j \in S) \cdot \beta_j$  and  $(\boldsymbol{\beta}_{\bar{S}})_j = \mathbf{1}(j \notin S) \cdot \beta_j$  for  $j = 1, \dots, d$ , where  $\mathbf{1}(\cdot)$  is the indicator function. We denote all-zero matrices by  $\mathbf{0}$ . For notational simplicity, we use generic absolute constants  $C, C', \dots$ , whose values may change from line to line.

Throughout, we denote the exact and approximate local solutions by  $\hat{\boldsymbol{\beta}}$  and  $\tilde{\boldsymbol{\beta}}$  respectively. We index  $\tilde{\boldsymbol{\beta}}$  with the corresponding regulation parameter  $\lambda$ , e.g.,  $\tilde{\boldsymbol{\beta}}_\lambda$ . For the proposed path following method, we use subscript  $t$  to index the path following stages, e.g. the approximate local solution obtained within the  $t$ -th stage is denoted by  $\tilde{\boldsymbol{\beta}}_t$ . Within the  $t$ -th stage, we index the proximal-gradient iterations with superscript  $k$ , e.g.,  $\boldsymbol{\beta}_t^k$ .

**2. Some Nonconvex Sparse Learning Problems.** Many theoretical results on penalized  $M$ -estimators rely on the condition that the loss and penalty functions are convex, since convexity makes both computational and statistical analysis convenient. However, the statistical performance of the estimator obtained from these convex formulations could be suboptimal in some settings. In the following, we introduce several nonconvex sparse learning problems as motivating examples.

**2.1. Nonconvex Penalty.** Throughout this paper, we consider decomposable penalty functions

$$\mathcal{P}_\lambda(\boldsymbol{\beta}) = \sum_{j=1}^d p_\lambda(\beta_j),$$

e.g., the  $\ell_1$  penalty  $\lambda\|\boldsymbol{\beta}\|_1 = \sum_{i=1}^d \lambda|\beta_j|$ . When the minimum of  $|\beta_j^*| > 0$  is not close to zero, the  $\ell_1$  penalty introduces large bias in parameter estimation. To remedy this effect, Fan and Li (2001) proposed the SCAD penalty

$$(2.1) \quad p_\lambda(\beta_j) = \lambda \int_0^{|\beta_j|} \left\{ \mathbf{1}(z \leq \lambda) + \frac{(a\lambda - z)_+}{(a-1)\lambda} \mathbf{1}(z > \lambda) \right\} dz, \quad a > 2,$$

and Zhang (2010a) proposed the MCP penalty

$$(2.2) \quad p_\lambda(\beta_j) = \lambda \int_0^{|\beta_j|} \left(1 - \frac{z}{\lambda b}\right)_+ dz, \quad b > 0.$$

See Zhang and Zhang (2012) for a detailed survey. These nonconvex penalty functions are illustrated in Figure 2(a). In fact, these nonconvex penalties can be formulated as the sum of the  $\ell_1$  penalty and a concave part

$$(2.3) \quad p_\lambda(\beta_j) = \lambda|\beta_j| + q_\lambda(\beta_j).$$

The concave components  $q_\lambda(\beta_j)$  of SCAD and MCP are illustrated in Figure 2(b), while the corresponding derivatives  $q'_\lambda(\beta_j)$  are illustrated in Figure 2(c). See §A.1 of the supplementary material (Wang et al., 2014b) for the detailed analytical forms of  $p_\lambda(\beta_j)$  and  $q_\lambda(\beta_j)$  for SCAD and MCP.

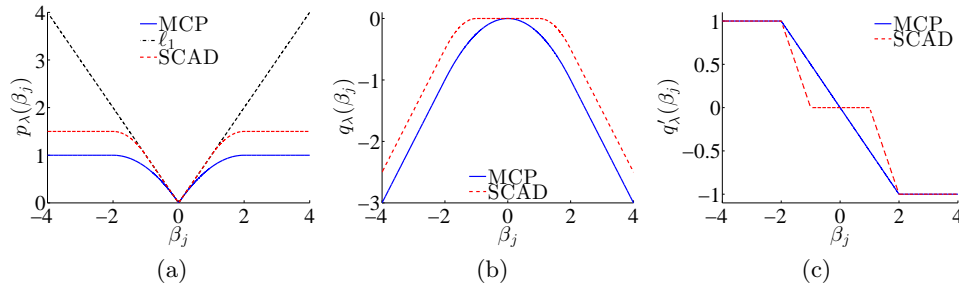


FIG 2. An illustration of nonconvex penalties: (a) Plots of  $p_\lambda(\beta_j)$  for MCP,  $\ell_1$ , and SCAD; (b) Plots of  $q_\lambda(\beta_j)$  for MCP and SCAD; (c) Plots of  $q'_\lambda(\beta_j)$  for MCP and SCAD. Here  $p_\lambda(\beta_j)$  is the penalty function evaluated at the  $j$ -th dimension of  $\beta$ ,  $q_\lambda(\beta_j)$  is the concave component of  $p_\lambda(\beta_j)$ , and  $q'_\lambda(\beta_j)$  is the derivative of  $q_\lambda(\beta_j)$ . Here we set  $a = 2.1$  for SCAD,  $b = 2$  for MCP, and  $\lambda = 1$ .

In fact, our method and theory are not limited to SCAD and MCP. More generally, we only rely on the following regularity conditions on the concave component  $q_\lambda(\beta_j)$ :

### Regularity Conditions on Nonconvex Penalty

- (a)  $q'_\lambda(\beta_j)$  is monotone and Lipschitz continuous, i.e., for  $\beta'_j > \beta_j$ , there exist two constants  $\zeta_- \geq 0$  and  $\zeta_+ \geq 0$  such that

$$-\zeta_- \leq \frac{q'_\lambda(\beta'_j) - q'_\lambda(\beta_j)}{\beta'_j - \beta_j} \leq -\zeta_+ \leq 0;$$

- (b)  $q_\lambda(\beta_j)$  is symmetric, i.e.,  $q_\lambda(-\beta_j) = q_\lambda(\beta_j)$  for any  $\beta_j$ ;

- (c)  $q_\lambda(\beta_j)$  and  $q'_\lambda(\beta_j)$  pass through the origin, i.e.,  $q_\lambda(0) = q'_\lambda(0) = 0$ ;

- (d)  $q'_\lambda(\beta_j)$  is bounded, i.e.,  $|q'_\lambda(\beta_j)| \leq \lambda$  for any  $\beta_j$ ;

- (e)  $q'_\lambda(\beta_j)$  has bounded difference with respect to  $\lambda$ :  $|q'_{\lambda_1}(\beta_j) - q'_{\lambda_2}(\beta_j)| \leq$



$|\lambda_1 - \lambda_2|$  for any  $\beta_j$ .

In regularity condition (a),  $\zeta_-$  and  $\zeta_+$  are two parameters that control the concavity of  $q_\lambda(\beta_j)$ . Note that the second order derivative of a function characterizes its convexity/concavity. Taking  $\beta'_j \rightarrow \beta_j$  in regularity condition (a), we have  $q''_\lambda(\beta_j) \in [-\zeta_-, -\zeta_+]$  (ignoring those  $\beta_j$ 's where  $q''_\lambda(\beta_j)$  doesn't exist), which suggests larger  $\zeta_-$  and  $\zeta_+$  allow  $q_\lambda(\beta_j)$  to be more concave. For SCAD we have  $\zeta_- = 1/(a-1)$  and  $\zeta_+ = 0$ , while for MCP we have  $\zeta_- = 1/b$  and  $\zeta_+ = 0$ . In Figure 2(b) and Figure 2(c), we can verify that regularity conditions (a)-(d) hold for MCP and SCAD. In addition, we illustrate regularity condition (e) for MCP and SCAD in §A.2 of the supplementary material (Wang et al., 2014b).

From (2.3) we have  $\mathcal{P}_\lambda(\boldsymbol{\beta}) = \sum_{j=1}^d p_\lambda(\beta_j) = \lambda \|\boldsymbol{\beta}\|_1 + \sum_{j=1}^d q_\lambda(\beta_j)$ . For notational simplicity, we define

$$(2.4) \quad \mathcal{Q}_\lambda(\boldsymbol{\beta}) = \sum_{j=1}^d q_\lambda(\beta_j) = \mathcal{P}_\lambda(\boldsymbol{\beta}) - \lambda \|\boldsymbol{\beta}\|_1.$$

Hence  $\mathcal{Q}_\lambda(\boldsymbol{\beta})$  denotes the decomposable concave component of the nonconvex penalty  $\mathcal{P}_\lambda(\boldsymbol{\beta})$ .

**2.2. Nonconvex Loss Function.** In this paper, we focus on an example of nonconvex loss function named semiparametric elliptical design regression. More specifically, we have  $n$  pairs of observations  $\mathbf{z}_1 = (y_1, \mathbf{x}_1^T)^T, \dots, \mathbf{z}_n = (y_n, \mathbf{x}_n^T)^T$  of a random vector  $\mathbf{Z} = (Y, \mathbf{X}^T)^T \in \mathbb{R}^{d+1}$  that follows a  $(d+1)$ -dimensional elliptical distribution. (See §A.3 of the supplementary material (Wang et al., 2014b) for a detailed introduction to elliptical distribution.) Then we can verify that  $(Y|\mathbf{X} = \mathbf{x})$  follows a univariate elliptical distribution. If we assume that  $\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}^*$ , the population version of the semiparametric elliptical design regression estimator can be defined as

$$(2.5) \quad \begin{aligned} \check{\boldsymbol{\beta}} &= \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^d} \left\{ \frac{1}{2} \mathbb{E}_{\mathbf{X}, Y} \left( (Y - \mathbf{X}^T \boldsymbol{\beta})^2 \right) + \mathcal{P}_\lambda(\boldsymbol{\beta}) \right\} \\ &= \operatorname{argmin}_{\boldsymbol{\beta} \in \mathbb{R}^d} \left\{ \frac{1}{2} (1, -\boldsymbol{\beta}^T) \boldsymbol{\Sigma}_{\mathbf{Z}} (1, -\boldsymbol{\beta}^T)^T + \mathcal{P}_\lambda(\boldsymbol{\beta}) \right\}. \end{aligned}$$

The above procedure is not practically implementable, since the population covariance matrix  $\boldsymbol{\Sigma}_{\mathbf{Z}}$  in (2.5) is unknown. In practice, we need to estimate the population covariance matrix  $\boldsymbol{\Sigma}_{\mathbf{Z}}$ . For this purpose, we propose a rank-based covariance matrix estimator  $\hat{\mathbf{K}}_{\mathbf{Z}}$ , which is calculated by a two-step procedure described in §A.4 of the supplementary material (Wang et al., 2014b). Since  $\hat{\mathbf{K}}_{\mathbf{Z}}$  is not necessarily positive semidefinite, the loss function

in semiparametric elliptical design regression, i.e.,

$$(2.6) \quad \mathcal{L}(\beta) = \frac{1}{2} (1, -\beta^T) \widehat{\mathbf{K}}_{\mathbf{Z}} (1, -\beta^T)^T,$$

is possibly nonconvex.

**3. Approximate Regularization Path Following Method.** Before we go into details, we first present the high level idea of approximate regularization path following. We then introduce the basic building block of our path following method — a proximal-gradient method tailored to nonconvex problems.

*3.1. Approximate Regularization Path Following.* Fast local geometric convergence in the proximity of sparse solutions has been observed by many authors (Wright et al., 2009; Blumensath and Davies, 2009; Agarwal et al., 2012; Xiao and Zhang, 2013). We exploit such fast local convergence under an approximate path framework to achieve fast global convergence.

**Initialization:** In (1.1), when the regularization parameter  $\lambda$  is sufficiently large, the solution to sparse learning problems is an all-zero vector. Recall that any exact local solution  $\widehat{\beta}_\lambda$  satisfies the first-order optimality condition,  $\mathbf{0} \in \partial\{\mathcal{L}(\widehat{\beta}_\lambda) + \mathcal{P}_\lambda(\widehat{\beta}_\lambda)\}$ . Since the nonconvex penalty  $\mathcal{P}_\lambda(\beta)$  can be formulated as  $\mathcal{P}_\lambda(\beta) = \mathcal{Q}_\lambda(\beta) + \lambda\|\beta\|_1$ , where  $\mathcal{Q}_\lambda(\beta)$  is defined in (2.4), the first-order optimality condition implies there should exist some subgradient  $\xi \in \partial\|\widehat{\beta}_\lambda\|_1$  such that

$$(3.1) \quad \mathbf{0} = \nabla\mathcal{L}(\widehat{\beta}_\lambda) + \nabla\mathcal{Q}_\lambda(\widehat{\beta}_\lambda) + \lambda\xi.$$

Let  $\lambda$  be chosen such that  $\widehat{\beta}_\lambda = \mathbf{0}$ . Then regularity condition (c) implies  $\nabla\mathcal{Q}_\lambda(\mathbf{0}) = \mathbf{0}$ . Meanwhile, since  $\xi \in \partial\|\mathbf{0}\|_1$ , we have  $\|\xi\|_\infty \leq 1$ , which implies  $\|\nabla\mathcal{L}(\mathbf{0})\|_\infty \leq \lambda$  in (3.1). Hence,  $\lambda_0 = \|\nabla\mathcal{L}(\mathbf{0})\|_\infty$  is the smallest regularization parameter such that any exact local solution  $\widehat{\beta}_\lambda$  to the minimization problem (1.1) is all-zero. We choose this  $\lambda_0$  to be the initial parameter of our regularization path.

**Approximate Path Following:** Let  $\lambda_{\text{tgt}} \in (0, \lambda_0)$  be the target regularization parameter in (1.1). In practice, we may choose  $\lambda_{\text{tgt}}$  by cross-validation or the high-dimensional BIC criterion proposed by Wang et al. (2013). We consider a decreasing sequence of regularization parameters  $\{\lambda_t\}_{t=0}^N$ , where

$$(3.2) \quad \lambda_t = \eta^t \lambda_0 \quad (t = 0, \dots, N), \quad \lambda_N = \lambda_{\text{tgt}}, \quad \text{and} \quad \eta \in [0.9, 1).$$

Here  $\eta$  is an absolute constant that doesn't scale with sample size  $n$  and dimension  $d$ . In §4 and §5 we will prove that,  $\eta \in [0.9, 1)$  ensures the global geometric rate of convergence. Consequently, since we have  $\lambda_{\text{tgt}} = \lambda_0 \eta^N$  by

(3.2), the number of path following stages is

$$(3.3) \quad N = \frac{\log(\lambda_0/\lambda_{\text{tgt}})}{\log(\eta^{-1})}.$$

Without loss of generality, we assume that  $\eta$  is properly chosen such that  $N$  is an integer. We will show in §4 that,  $\lambda_{\text{tgt}}$  scales with sample size  $n$  and dimension  $d$ . Since  $\eta$  is a constant, the number of stages  $N$  also scales with  $n$  and  $d$ . Within the  $t$ -th ( $t = 1, \dots, N$ ) path following stage, we aim to obtain a local solution to the minimization problem  $\min_{\beta} \{\mathcal{L}(\beta) + \mathcal{P}_{\lambda_t}(\beta)\}$ .

As shown in Lines 5-9 of Algorithm 1, within the  $t$ -th ( $t = 1, \dots, N-1$ ) path following stage, we employ a variant of proximal-gradient method (Algorithm 3) to obtain an approximate local solution  $\tilde{\beta}_t$  for regularization parameter  $\lambda_t = \eta^t \lambda_0$ . To ensure that each path following stage enjoys a fast geometric rate of convergence, we propose an approximation path following strategy. More specifically, we use the approximate local solution  $\tilde{\beta}_{t-1}$  obtained within the  $(t-1)$ -th path following stage to initialize the  $t$ -th stage (Line 8 and Line 12 of Algorithm 1). Recall that we need to adaptively search for the best  $L_t^k$  ( $k = 0, 1, \dots$ ) in (1.3). To achieve computational efficiency, within the  $(t-1)$ -th path following stage, we store the chosen  $L_{t-1}^k$  at the last proximal-gradient iteration of the  $(t-1)$ -th stage as  $L_{t-1}$ . Within the  $t$ -th stage we initialize the search for  $L_t^0$  with  $L_{t-1}$  (Line 8 and Line 12 of Algorithm 1), which will be explained in §3.2.

**Configuration of Optimization Precision:** We set the optimization precision  $\epsilon_t$  for the  $t$ -th ( $t = 1, \dots, N-1$ ) stage to be  $\lambda_t/4$  (Line 7 of Algorithm 1). Within the  $N$ -th path following stage where  $\lambda_N = \lambda_{\text{tgt}}$  (Line 10), we solve up to high optimization precision  $\epsilon_{\text{opt}} \ll \lambda_{\text{tgt}}/4$  (Line 11). The intuition behind this configuration of optimization precision is explained as follows:

- For  $t = 1, \dots, N-1$ , recall the exact local solution  $\hat{\beta}_{\lambda_t}$  is an estimator of the true parameter vector  $\beta^*$  corresponding to the regularization parameter  $\lambda_t$ . According to high-dimensional statistical theory, the statistical error of  $\hat{\beta}_{\lambda_t}$  should be upper bounded by  $C\lambda_t\sqrt{s^*}$  with high probability, where  $s^* = \|\beta^*\|_0$ . In Lemma 5.1 we will prove that, if the optimization error of the approximate local solution  $\tilde{\beta}_t$  is at most  $\lambda_t/4$ , then  $\tilde{\beta}_t$  lies within a ball of radius  $C'\lambda_t\sqrt{s^*}$  centered at  $\beta^*$  with high probability. That is to say, the approximate local solution  $\tilde{\beta}_t$  has the same order of statistical error as the exact solution  $\hat{\beta}_{\lambda_t}$ , and therefore enjoys desired statistical recovery properties. In particular, in Theorem 5.5 we will prove that,  $\tilde{\beta}_t$  is guaranteed to be sparse, and thus falls into the fast convergence region of the next path following stage.
- However, for  $t = N$ , we need to solve up to high optimization precision

---

**Algorithm 1** The approximate path following method, which solves for a decreasing sequence of regularization parameters  $\{\lambda_t\}_{t=0}^N$ . Within the  $t$ -th path following stage, we employ the proximal-gradient method illustrated in Algorithm 3 to achieve an approximate local solution  $\tilde{\beta}_t$  for  $\lambda_t$ . This approximate local solution is then used to initialize the  $(t+1)$ -th stage.

---

```

1:  $\{\tilde{\beta}_t\}_{t=1}^N \leftarrow \text{Approximate-Path-Following}(\lambda_{\text{tgt}}, \epsilon_{\text{opt}})$ 
2: input:  $\lambda_{\text{tgt}} > 0, \epsilon_{\text{opt}} > 0$  {Here we set  $\epsilon_{\text{opt}} \ll \lambda_{\text{tgt}}/4$ .}
3: parameters:  $\eta \in [0.9, 1), R > 0, L_{\min} > 0, \lambda_0 = \|\nabla \mathcal{L}(\mathbf{0})\|_\infty$ 
   {For logistic loss, we set  $R \in (0, +\infty)$ ; For other loss functions, we set  $R = +\infty$ .}
   {In practice, we set  $L_{\min}$  to be a sufficiently small value, e.g.,  $10^{-6}$ .}
4: initialize:  $\tilde{\beta}_0 \leftarrow \mathbf{0}, L_0 \leftarrow L_{\min}, N \leftarrow \log(\lambda_0/\lambda_{\text{tgt}})/\log(\eta^{-1})$ 
5: for  $t = 1, \dots, N-1$  do
6:    $\lambda_t \leftarrow \eta^t \lambda_0$ 
7:    $\epsilon_t \leftarrow \lambda_t/4$ 
8:    $\{\tilde{\beta}_t, L_t\} \leftarrow \text{Proximal-Gradient}(\lambda_t, \epsilon_t, \tilde{\beta}_{t-1}, L_{t-1}, R)$  as in Algorithm 3
9: end for
10:  $\lambda_N \leftarrow \lambda_{\text{tgt}}$ 
11:  $\epsilon_N \leftarrow \epsilon_{\text{opt}}$ 
12:  $\{\tilde{\beta}_N, L_N\} \leftarrow \text{Proximal-Gradient}(\lambda_N, \epsilon_N, \tilde{\beta}_{N-1}, L_{N-1}, R)$ 
13: return  $\{\tilde{\beta}_t\}_{t=1}^N$ 

```

---

$\epsilon_{\text{opt}} \ll \lambda_{\text{tgt}}/4$ . This is because, even though  $\tilde{\beta}_t$  and  $\hat{\beta}_{\lambda_t}$  both have statistical error of the order  $\lambda_t \sqrt{s^*}$ , in certain regimes (to be specified in Theorem 4.8), the exact local solution  $\hat{\beta}_{\lambda_t}$  can achieve an improved recovery performance (as shown in (1.5)) due to the usage of nonconvex penalties. Therefore, within the final stage we need to obtain an approximate solution  $\tilde{\beta}_N$  as close to the exact local solution  $\hat{\beta}_{\lambda_{\text{tgt}}}$  as possible, so that  $\tilde{\beta}_N$  has a sharper statistical rate of convergence.

In Algorithm 1,  $R > 0$  (Line 3) is a parameter that determines the radius of the constraint used in the proximal-gradient method (Line 8 and Line 12). For least squares loss and semiparametric elliptical design loss, we don't need any constraint. Therefore, we set  $R = +\infty$ . However, for logistic loss we need to impose an  $\ell_2$  constraint of radius  $R \in (0, +\infty)$ . Here  $L_{\min}$  is a parameter used in the proximal-gradient method (Line 3 of Algorithm 3), which is often set to be a sufficiently small value in practice, e.g.,  $L_{\min} = 10^{-6}$ . We will explain with details in §3.2.

**3.2. Proximal-Gradient Method for Nonconvex Problems.** Before we introduce our proximal-gradient method which is tailored to nonconvex problems, we first give a brief introduction to Nesterov's proximal-gradient method

(Nesterov, 2013), which solves the following convex optimization problem

$$(3.4) \quad \text{minimize } \phi_\lambda(\beta), \quad \text{where } \phi_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{P}_\lambda(\beta), \quad \beta \in \Omega.$$

Here  $\mathcal{L}(\beta)$  is convex and differentiable,  $\mathcal{P}_\lambda(\beta)$  is convex but possibly nonsmooth, and  $\Omega$  is a closed convex set.

Recall that  $\beta_t^k$  corresponds to the  $k$ -th iteration of the proximal-gradient method within the  $t$ -th path following stage. Nesterov's proximal-gradient method updates  $\beta_t^k$  to be the minimizer of the following local quadratic approximation of  $\phi_{\lambda_t}(\beta)$  at  $\beta_t^{k-1}$

$$(3.5) \quad \psi_{L_t^k, \lambda_t}(\beta; \beta_t^{k-1}) = \mathcal{L}(\beta_t^{k-1}) + \nabla \mathcal{L}(\beta_t^{k-1})^T (\beta - \beta_t^{k-1}) \\ + \frac{L_t^k}{2} \|\beta - \beta_t^{k-1}\|_2^2 + \mathcal{P}_{\lambda_t}(\beta),$$

where  $L_t^k > 0$  is chosen by line-search.

Nesterov's proximal-gradient method requires that both  $\mathcal{L}(\beta)$  and  $\mathcal{P}_\lambda(\beta)$  in (3.4) are convex. However, in the optimization problem (1.1) considered in this paper,  $\mathcal{L}(\beta)$  and  $\mathcal{P}_\lambda(\beta)$  may be no longer convex. In this case, directly plugging  $\mathcal{L}(\beta)$  and  $\mathcal{P}_\lambda(\beta)$  into Nesterov's proximal-gradient might lead to the phenomenon of bad local optima under a path following scheme, as observed by She (2009, 2012). To extend the proximal-gradient method to nonconvex settings, we adopt an alternative formulation of the objective function.

Recall that the nonconvex penalty can be decomposed as  $\mathcal{P}_\lambda(\beta) = \lambda \|\beta\|_1 + \mathcal{Q}_\lambda(\beta)$ , where  $\mathcal{Q}_\lambda(\beta)$  is defined in (2.4). For notational simplicity, we denote  $\mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$  by  $\tilde{\mathcal{L}}_\lambda(\beta)$ . Therefore, the objective function  $\phi_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{P}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta) + \lambda \|\beta\|_1$  can be reformulated as

$$(3.6) \quad \phi_\lambda(\beta) = \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \|\beta\|_1,$$

where we can view  $\tilde{\mathcal{L}}_\lambda(\beta)$  as a surrogate loss function and  $\lambda \|\beta\|_1$  as a new penalty function. This reformulation ensures the convexity of the new penalty function. Moreover, in Lemma 5.1 we will prove that, the surrogate loss function  $\tilde{\mathcal{L}}_\lambda(\beta)$  is actually strongly convex on a sparse set. Correspondingly, we modify Nesterov's proximal-gradient method to minimize the local quadratic approximation defined as

$$(3.7) \quad \psi_{L_t^k, \lambda_t}(\beta; \beta_t^{k-1}) = \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta - \beta_t^{k-1}) \\ + \frac{L_t^k}{2} \|\beta - \beta_t^{k-1}\|_2^2 + \lambda_t \|\beta\|_1.$$

Note that, unlike (3.5), we use a quadratic approximation to the surrogate loss function  $\tilde{\mathcal{L}}_{\lambda_t}(\beta)$  in (3.7), instead of the original loss function  $\mathcal{L}(\beta)$ . At the  $k$ -th iteration of the proximal-gradient method, we update  $\beta_t^k$  to be the

minimizer of the quadratic approximation defined in (3.7), i.e.,

$$(3.8) \quad \beta_t^k \leftarrow \operatorname{argmin}_{\beta \in \Omega} \left\{ \psi_{L_t^k, \lambda_t}(\beta; \beta_t^{k-1}) \right\}.$$

Now we specify the constraint set  $\Omega$  in (3.8). For  $\mathcal{L}(\beta)$  being least squares or semiparametric elliptical design loss, we set  $\Omega = \mathbb{R}^d$ . For logistic loss, we set  $\Omega = B_2(R)$  with  $R \in (0, +\infty)$ , where  $B_2(R)$  is a centered  $\ell_2$  ball of radius  $R$ . In Lemma 5.1 we will show that, in the setting of logistic loss, the boundedness of  $\|\beta_t^k\|_2$ 's is essential for establishing the strong convexity of the surrogate loss function  $\tilde{\mathcal{L}}_{\lambda_t}(\beta)$  along the full regularization path. To unify the notation, we consider  $\Omega = B_2(R)$  throughout — when the constraint set  $\Omega = \mathbb{R}^d$ , we set  $R = +\infty$ . Correspondingly, we denote (3.8) by

$$(3.9) \quad \beta_t^k \leftarrow \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; R).$$

In the sequel, we provide the closed-form expression of update scheme (3.9):

#### Update Scheme of Proximal-Gradient Method for Nonconvex Problems

- For  $\Omega = \mathbb{R}^d$ , i.e.,  $R = +\infty$ ,  $\mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty)$  is a soft-thresholding operator taking the form of

$$(3.10) \quad \left( \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty) \right)_j = \begin{cases} 0 & \text{if } |\bar{\beta}_j| \leq \lambda_t / L_t^k, \\ \operatorname{sign}(\bar{\beta}_j) (|\bar{\beta}_j| - \lambda_t / L_t^k) & \text{if } |\bar{\beta}_j| > \lambda_t / L_t^k, \end{cases}$$

for  $j = 1, \dots, d$ , where

$$(3.11) \quad \begin{aligned} \bar{\beta} &= \beta_t^{k-1} - \frac{1}{L_t^k} \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) \\ &= \beta_t^{k-1} - \frac{1}{L_t^k} \left( \nabla \mathcal{L}(\beta_t^{k-1}) + \nabla \mathcal{Q}_{\lambda_t}(\beta_t^{k-1}) \right), \end{aligned}$$

and  $\bar{\beta}_j$  is the  $j$ -th dimension of  $\bar{\beta}$ .

- For  $\Omega = B_2(R)$  with  $R \in (0, +\infty)$ ,  $\mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; R)$  can be obtained by projecting  $\mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty)$  defined in (3.10) onto  $B_2(R)$ , i.e.,

$$(3.12) \quad \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; R) = \begin{cases} \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty) & \text{if } \|\mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty)\|_2 < R, \\ \frac{R \cdot \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty)}{\|\mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty)\|_2} & \text{if } \|\mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; +\infty)\|_2 \geq R. \end{cases}$$

See §B.2 of the supplementary material (Wang et al., 2014b) for a detailed derivation. In §B.1 of the supplementary material (Wang et al., 2014b), we

provide the specific forms of  $\nabla \mathcal{L}(\beta)$  and  $\nabla \mathcal{Q}_{\lambda_t}(\beta)$  in (3.11) for the nonconvex problems discussed in §2.

**Line-Search Method:** Before we present the proposed proximal-gradient method in detail, we briefly introduce a line-search algorithm, which adaptively searches for the best quadratic coefficient  $L_t^k$  of the local quadratic approximation (3.7). As shown in Lines 4-7 of Algorithm 2, the main idea of line-search is to iteratively increase  $L_t^k$  by a factor of two and compute the corresponding  $\beta_t^k$ , until the local approximation  $\psi_{L_t^k, \lambda_t}(\beta_t^k; \beta_t^{k-1})$  becomes a tight upper bound of the objective function  $\phi_{\lambda_t}(\beta_t^k)$ . We will theoretically characterize the computational complexity of this line-search algorithm in Remark 4.6, and specify the range of  $L_t^k$  in Theorem 5.5.

---

**Algorithm 2** The line-search method used to search for the best  $L_t^k$  and compute the corresponding  $\beta_t^k$ . Here  $\phi_{\lambda_t}(\beta)$  is the objective function defined in (3.4), and  $\psi_{L_t^k, \lambda_t}(\beta; \beta_t^{k-1})$  is the local quadratic approximation of  $\phi_{\lambda_t}(\beta)$  defined in (3.7).

---

```

1:  $\{\beta_t^k, L_t^k\} \leftarrow \text{Line-Search}(\lambda_t, \beta_t^{k-1}, L_{\text{init}}, R)$ 
2: input:  $\lambda_t > 0, \beta_t^{k-1} \in \mathbb{R}^d, L_{\text{init}} > 0, R > 0$ 
3: initialize:  $L_t^k \leftarrow L_{\text{init}}$ 
4: repeat
5:    $\beta_t^k \leftarrow \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; R)$  as defined in (3.9)
6:   if  $\phi_{\lambda_t}(\beta_t^k) > \psi_{L_t^k, \lambda_t}(\beta_t^k; \beta_t^{k-1})$  then  $L_t^k \leftarrow 2L_t^k$ 
7: until  $\phi_{\lambda_t}(\beta_t^k) \leq \psi_{L_t^k, \lambda_t}(\beta_t^k; \beta_t^{k-1})$ 
8: return  $\{\beta_t^k, L_t^k\}$ 

```

---

**Stopping Criterion:** In the following, we introduce the stopping criterion of our proximal-gradient method. In other words, we specify the optimality conditions that should be satisfied by the approximate solution  $\tilde{\beta}_t$  attained by our proximal-gradient method.

It is known that any exact local solution  $\hat{\beta}_\lambda$  to the optimization problem

$$\text{minimize } \phi_\lambda(\beta), \quad \text{where } \phi_\lambda(\beta) = \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \|\beta\|_1, \quad \beta \in \Omega$$

satisfies the optimality condition, i.e, there exists some  $\xi \in \partial \|\hat{\beta}_\lambda\|_1$  such that

$$(3.13) \quad (\hat{\beta}_\lambda - \beta)^T (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_\lambda) + \lambda \xi) \leq 0, \quad \text{for any } \beta \in \Omega.$$

We can understand this optimality condition as follows: Locally at  $\hat{\beta}_\lambda$ , any feasible direction pointed at  $\hat{\beta}_\lambda$ , i.e.,  $(\hat{\beta}_\lambda - \beta)$  where  $\beta \in \Omega$ , leads to a decrease in the objective function value  $\phi_\lambda(\beta)$ , because as shown in (3.13), such direction forms an obtuse angle with the (sub)gradient vector of  $\phi_\lambda(\beta)$  evaluated at  $\hat{\beta}_\lambda$ . If  $\hat{\beta}_\lambda$  lies in the interior of  $\Omega$ , e.g.,  $\Omega = \mathbb{R}^d$ , then (3.13)

reduces to the well-known first-order KKT condition,<sup>1</sup>

$$(3.14) \quad \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_\lambda) + \lambda \xi = \mathbf{0}, \quad \text{where } \xi \in \partial \|\hat{\beta}_\lambda\|_1.$$

Based on the optimality condition in (3.13), we measure the suboptimality of a  $\beta \in \Omega$  with

$$(3.15) \quad \omega_\lambda(\beta) = \min_{\xi' \in \partial \|\beta\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\beta - \beta')^T}{\|\beta - \beta'\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi') \right\}.$$

To understand this measure of suboptimality, first note that, if  $\beta$  is an exact local solution, then we have  $\omega_\lambda(\beta) \leq 0$  by (3.13). Otherwise, if  $\beta$  is close to some exact local solution, then  $\omega_\lambda(\beta)$  is some small positive value. When  $\beta$  lies in the interior of  $\Omega$ , then (3.15) reduces to a more straightforward

$$(3.16) \quad \omega_\lambda(\beta) = \min_{\xi' \in \partial \|\beta\|_1} \left\{ \|\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi'\|_\infty \right\}.$$

Because for any fixed  $v \in \mathbb{R}^d$ , we have  $(\beta + Cv) \in \Omega$  for  $C > 0$  sufficiently small. Setting  $\beta$  to be this value in (3.15), we have

$$\omega_\lambda(\beta) = \min_{\xi' \in \partial \|\beta\|_1} \max_{v \in \mathbb{R}^d} \left\{ \frac{v^T}{\|v\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi') \right\} = \min_{\xi' \in \partial \|\beta\|_1} \left\{ \|\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi'\|_\infty \right\},$$

where the second equality follows from the duality between  $\ell_1$  and  $\ell_\infty$  norm.

Equipped with the suboptimality measure  $\omega_\lambda(\beta)$  defined in (3.15), now we can define the stopping criterion of our proximal-gradient method to be  $\omega_{\lambda_t}(\beta_t^k) \leq \epsilon_t$ , where  $\epsilon_t > 0$  is the desired optimization precision within the  $t$ -th path following stage (Line 9 of Algorithm 3). Therefore, the proximal-gradient method achieves an approximate local solution  $\tilde{\beta}_t$  with suboptimality  $\epsilon_t$ . Recall that within the  $t$ -th path following stage ( $t = 1, \dots, N-1$ ), we set  $\epsilon_t$  to be  $\lambda_t/4$  (Line 7 of Algorithm 1), while within the  $N$ -th path following stage, we set  $\epsilon_t = \epsilon_{\text{opt}} \ll \lambda_{\text{tgt}}/4$  (Line 11 of Algorithm 1).

**Proposed Proximal-Gradient Method:** We are now ready to present the proposed proximal-gradient method in detail. Recall that, within the  $t$ -th stage of our path following algorithm, we employ the proximal-gradient method to obtain the approximate local solution  $\tilde{\beta}_t$  (Line 8 and Line 12 of Algorithm 1). As shown in Line 8 of Algorithm 3, at the  $k$ -th iteration of our proximal-gradient method, we employ the line-search method (Algorithm 2) to search for the best  $L_t^k$  and calculate the corresponding  $\beta_t^k$ .

At the  $k$ -th iteration of the proximal-gradient method, we set the initial value  $L_{\text{init}}$  of line-search to be  $\max \{L_{\text{min}}, L_t^{k-1}/2\}$  (Line 7 of Algorithm 3),

<sup>1</sup>Because given that  $\hat{\beta}_\lambda$  lies in the interior of  $\Omega$ , we have  $(\hat{\beta}_\lambda + Cv) \in \Omega$  and  $(\hat{\beta}_\lambda - Cv) \in \Omega$  for any fixed  $v \in \mathbb{R}^d$  and  $C > 0$  sufficiently small. Setting  $\beta$  in (3.13) to be these two values, we obtain  $v^T (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_\lambda) + \xi) = 0$ , which implies (3.14) since  $v$  is arbitrarily chosen.



---

**Algorithm 3** The proximal-gradient method for nonconvex problems, which iteratively leverages the line-search method illustrated in Algorithm 2 at each iteration.

---

```

1:  $\{\tilde{\beta}_t, L_t\} \leftarrow \text{Proximal-Gradient}(\lambda_t, \epsilon_t, \beta_t^0, L_t^0, R)$ 
2: input:  $\lambda_t > 0, \epsilon_t > 0, \beta_t^0 \in \mathbb{R}^d, L_t^0 > 0, R > 0$ 
3: parameter:  $L_{\min} > 0$ 
4: initialize:  $k \leftarrow 0$ 
5: repeat
6:    $k \leftarrow k + 1$ 
7:    $L_{\text{init}} \leftarrow \max\{L_{\min}, L_t^{k-1}/2\}$ 
8:    $\beta_t^k, L_t^k \leftarrow \text{Line-Search}(\lambda_t, \beta_t^{k-1}, L_{\text{init}}, R)$  as in Algorithm 2
9: until  $\omega_{\lambda_t}(\beta_t^k) \leq \epsilon_t$  as defined in (3.15)
10:  $\tilde{\beta}_t \leftarrow \beta_t^k$ 
11:  $L_t \leftarrow L_t^k$ 
12: return  $\{\tilde{\beta}_t, L_t\}$ 

```

---

where  $L_{\min} > 0$  is used to prevent  $L_{\text{init}}$  from being too small. In practice,  $L_{\min}$  is often set to be a sufficiently small value, e.g.,  $L_{\min} = 10^{-6}$ . The intuition behind such initialization can be understood as follows: As shown in (3.7),  $L_t^{k-1}$  and  $L_t^k$  are the quadratic coefficients of the local quadratic approximations of the objective function at  $\beta_t^{k-2}$  and  $\beta_t^{k-1}$  respectively. Intuitively speaking,  $\beta_t^{k-2}$  and  $\beta_t^{k-1}$  are close to each other, which implies that  $L_t^{k-1}$  is a good guess for  $L_t^k$ . Hence, we can initialize the line-search method for  $L_t^k$  with a value slightly smaller than  $L_t^{k-1}$ , e.g.,  $L_t^{k-1}/2$ .

When the stopping criterion  $\omega_{\lambda_t}(\beta_t^k) \leq \epsilon_t$  is satisfied (Line 9 of Algorithm 3), the proximal-gradient method stops and outputs the approximate local solution  $\tilde{\beta}_t = \beta_t^k$  (Line 10 of Algorithm 3). We also keep track of  $L_t = L_t^k$  to accelerate the line-search procedure within the next path following stage.

**4. Theoretical Results.** We establish theoretical results on the iteration complexity and statistical performance of our approximate regularization path following method for nonconvex learning problems.

**4.1. Assumptions.** We first list the required assumptions. The first assumption is about the relationship between  $\lambda_{\text{tgt}}$  and  $\|\nabla \mathcal{L}(\beta^*)\|_\infty$ .

**Assumption 4.1.** For least squares loss and logistic loss, we set  $\lambda_{\text{tgt}} = C\sqrt{\log d/n}$ . Meanwhile, for semiparametric elliptical design loss, we set  $\lambda_{\text{tgt}} = C'\|\beta^*\|_1\sqrt{\log d/n}$ . We assume

$$(4.1) \quad \|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda_{\text{tgt}}/8.$$

Assumption 4.1 is a common condition that  $\lambda_{\text{tgt}}$  should be large enough

to dominate the noise. For instance, for least squares loss we have

$$\nabla \mathcal{L}(\beta^*) = \frac{1}{n} \mathbf{X}^T (\mathbf{X} \beta^* - \mathbf{y}),$$

where  $\mathbf{X} \beta^* - \mathbf{y}$  is in fact the noise vector. In Lemma C.1 in §C.1 of the supplementary material (Wang et al., 2014b) we will show that, for least squares loss and logistic loss, we have that  $\|\nabla \mathcal{L}(\beta^*)\|_\infty \leq C \sqrt{\log d/n}$  holds with high probability. Similarly, in Lemma C.2 in §C.1 of the supplementary material (Wang et al., 2014b) we will prove that, for semiparametric elliptical design loss,  $\|\nabla \mathcal{L}(\beta^*)\|_\infty \leq C' \|\beta^*\|_1 \sqrt{\log d/n}$  holds with high probability. Thus, our assumption on  $\lambda_{\text{tgt}}$  and  $\|\nabla \mathcal{L}(\beta^*)\|_\infty$  holds with high probability.

In the sequel, we lay out another assumption on the sparse eigenvalues of  $\nabla^2 \mathcal{L}(\beta)$ , which are defined as follows.

**Definition 4.2** (Sparse Eigenvalues). Let  $s$  be a positive integer. The largest and smallest  $s$ -sparse eigenvalues of the Hessian matrix  $\nabla^2 \mathcal{L}(\beta)$  are

$$\begin{aligned} \rho_+(\nabla^2 \mathcal{L}, s) &= \sup \left\{ \mathbf{v}^T \nabla^2 \mathcal{L}(\beta) \mathbf{v} : \|\mathbf{v}\|_0 \leq s, \|\mathbf{v}\|_2 = 1, \beta \in \mathbb{R}^d \right\}, \\ \rho_-(\nabla^2 \mathcal{L}, s) &= \inf \left\{ \mathbf{v}^T \nabla^2 \mathcal{L}(\beta) \mathbf{v} : \|\mathbf{v}\|_0 \leq s, \|\mathbf{v}\|_2 = 1, \beta \in \mathbb{R}^d \right\}. \end{aligned}$$

For least squares loss and semiparametric elliptical design loss,  $\nabla^2 \mathcal{L}(\beta)$  doesn't depend on  $\beta$ . However, for logistic loss we have

$$(4.2) \quad \nabla^2 \mathcal{L}(\beta) = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \cdot \frac{1}{1 + \exp(-\mathbf{x}_i^T \beta)} \cdot \frac{1}{1 + \exp(\mathbf{x}_i^T \beta)},$$

which depends on  $\beta$ . Note in Definition 4.2, the smallest  $s$ -sparse eigenvalue  $\rho_-(\nabla^2 \mathcal{L}, s)$  is obtained by taking infimum over all  $\beta \in \mathbb{R}^d$ . Consequently, for logistic loss,  $\rho_-(\nabla^2 \mathcal{L}, s)$  is always zero, because in (4.2) we can take  $\beta$  such that  $|\mathbf{x}_i^T \beta| \rightarrow +\infty$  for all nonzero  $\mathbf{x}_i$ 's, which implies that  $\nabla^2 \mathcal{L}(\beta)$  goes to an all-zero matrix. To avoid this degenerate case, for logistic loss we define the sparse eigenvalues by taking infimum/supremum over all  $\beta$  with  $\|\beta\|_2$  bounded instead of over all  $\beta \in \mathbb{R}^d$ .

**Definition 4.3** (Sparse Eigenvalues for Logistic Loss). Let  $s$  be a positive integer. For logistic loss, we define the largest and smallest  $s$ -sparse eigenvalues of  $\nabla^2 \mathcal{L}(\beta)$  to be

$$\begin{aligned} \rho_+(\nabla^2 \mathcal{L}, s, R) &= \sup \left\{ \mathbf{v}^T \nabla^2 \mathcal{L}(\beta) \mathbf{v} : \|\mathbf{v}\|_0 \leq s, \|\mathbf{v}\|_2 = 1, \|\beta\|_2 \leq R \right\}, \\ \rho_-(\nabla^2 \mathcal{L}, s, R) &= \inf \left\{ \mathbf{v}^T \nabla^2 \mathcal{L}(\beta) \mathbf{v} : \|\mathbf{v}\|_0 \leq s, \|\mathbf{v}\|_2 = 1, \|\beta\|_2 \leq R \right\}, \end{aligned}$$

where  $R \in (0, +\infty)$  is an absolute constant such that  $\|\beta^*\|_2 \leq R$ .

In Definition 4.3, we implicitly assume that  $\|\beta^*\|_2$  is upper bounded by some known absolute constant. Although it seems rather restrictive, this

assumption is essential for logistic loss. Otherwise,  $\nabla^2 \mathcal{L}(\beta^*)$  may go to an all-zero matrix when  $\|\beta^*\|_2 \rightarrow +\infty$ . In this case, the curvature of the objective function at  $\beta^*$  is zero, a consistent estimation of  $\beta^*$  is impossible. Although such assumption is necessary for theoretical purposes, we require no prior knowledge about the exact value of  $\|\beta^*\|_2$  in practice, since we can always set  $R$  to be a sufficiently large constant in our algorithm (Line 3 of Algorithm 1). To unify the later analysis for different loss functions, we omit the extra term  $R$  in Definition 4.3 unless its necessary.

Recall that we impose an  $\ell_2$  constraint of radius  $R$  for all the proximal-gradient iterations within each path following stage (Line 8 and Line 12 of Algorithm 1). Therefore, we have  $\|\beta_t^k\|_2 \leq R$  during the whole iterative procedure (for least squares loss and semiparametric elliptical design loss,  $R = +\infty$ ; for logistic loss,  $R \in (0, +\infty)$ ). Now we are ready to present the assumption on the sparse eigenvalues of the Hessian matrix.

**Assumption 4.4.** Let  $s^* = \|\beta^*\|_0$ . We assume:

- There exists an integer  $\tilde{s} > Cs^*$  such that

$$\rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) < +\infty, \quad \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) > 0$$

are two absolute constants. The constant  $C > 0$  is specified in (4.4).

- The concavity parameter  $\zeta_-$  defined in regularity condition (a) satisfies

$$(4.3) \quad \zeta_- \leq C' \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})$$

with constant  $C' < 1$ .

In Assumption 4.4, the constant

$$(4.4) \quad C = 144\kappa^2 + 250\kappa,$$

where  $\kappa$  is a condition number defined as

$$(4.5) \quad \kappa = \frac{\rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_+}{\rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_-}.$$

The constant in (4.4) is rather large for practical purposes. We could expect it to be much smaller if we manage to get smaller constants in the technical proof. However, we mainly focus on providing novel theoretical insights in this paper, without paying too much effort on optimizing constants.

Recall that regularity condition (a) implies  $\zeta_+ \leq \zeta_-$ . Meanwhile, we have  $\rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) \leq \rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})$  by definition. Thus, (4.3) implies

$$(4.6) \quad \zeta_+ \leq C' \rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}),$$

where  $C' < 1$  is the same constant as in (4.3). Therefore, we have  $\kappa \in [1, +\infty)$ . Restrictions (4.3) and (4.6) on the concavity parameters suggest that, the concavity of the concave component  $\mathcal{Q}_\lambda(\beta) = \sum_{j=1}^d q_\lambda(\beta_j)$  of the nonconvex

penalty should not outweigh the convexity of the loss function on a sparse set. It is also worth noting the concavity parameters are independent from the regularization parameter, e.g., for MCP in (2.2),  $b = 1/\zeta_-$  and  $\lambda$  are two independent parameters. Thus, Assumption 4.4 doesn't depend on  $\lambda$  at all.

Assumption 4.4 is closely related to the restricted isometry property (RIP) condition proposed by Candés and Tao (2005). Similar conditions have been studied by Bickel et al. (2009); Raskutti et al. (2010); Negahban et al. (2012); Zhang (2010b); Zhang et al. (2013); Xiao and Zhang (2013). In detail, for least squares loss, the RIP condition assumes there exists an integer  $s$  and some constant  $\delta \in (0, 1)$  such that

$$(4.7) \quad 1 - \delta \leq \rho_-(\nabla^2 \mathcal{L}, s) \leq \rho_+(\nabla^2 \mathcal{L}, s) \leq 1 + \delta.$$

Now we justify Assumption 4.4 for least squares loss with an example.

To show Assumption 4.4 is well defined, we assume the RIP condition in (4.7) holds with  $s = 877s^*$  and  $\delta = 0.01$ . We set the concavity parameters of the nonconvex penalty in (a) to be  $\zeta_+ = 0$  and  $\zeta_- = \rho_-(\nabla^2 \mathcal{L}, s)/20$ , e.g., for MCP defined in (2.2), we take  $b = 1/\zeta_- = 20/\rho_-(\nabla^2 \mathcal{L}, s)$ . In the following, we verify there exists an integer  $\tilde{s} = 438s^*$  that satisfies Assumption 4.4.

First, according to the RIP condition, we have

$$(4.8) \quad \begin{aligned} \rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) &= \rho_+(\nabla^2 \mathcal{L}, 877s^*) = \rho_+(\nabla^2 \mathcal{L}, s) \\ &\leq (1 + \delta) = 1.01 < +\infty, \end{aligned}$$

$$(4.9) \quad \begin{aligned} \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) &= \rho_-(\nabla^2 \mathcal{L}, 877s^*) = \rho_-(\nabla^2 \mathcal{L}, s) \\ &\geq (1 - \delta) = 0.99 > 0. \end{aligned}$$

Second, we calculate the value of  $\tilde{s}$  in detail. Since the condition number  $\kappa$  defined in (4.5) satisfies

$$\begin{aligned} 1 \leq \kappa &= \frac{\rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_+}{\rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_-} = \frac{\rho_+(\nabla^2 \mathcal{L}, s) - \zeta_+}{\rho_-(\nabla^2 \mathcal{L}, s) - \zeta_-} \\ &= \frac{20}{19} \cdot \frac{\rho_+(\nabla^2 \mathcal{L}, s)}{\rho_-(\nabla^2 \mathcal{L}, s)} \leq \frac{20}{19} \cdot \frac{1 + \delta}{1 - \delta} < 1.08. \end{aligned}$$

We now verify that  $\tilde{s}$  satisfies  $\tilde{s} > Cs^*$  in Assumption 4.4, where  $C$  is defined in (4.4). Plugging the range  $1 \leq \kappa < 1.08$  into the definition of  $C$ , we obtain  $C = 144\kappa^2 + 250\kappa < 438$ . Therefore, as long as the RIP condition holds with  $s = 877s^*$  and  $\delta = 0.01$ , we can find an integer  $\tilde{s} = 438s^*$  that satisfies Assumption 4.4, which also implies Assumption 4.4 is a weaker assumption than the RIP condition. For least squares loss, the RIP condition is known to hold for a variety of design matrices with high probability, which implies that Assumption 4.4 also holds with high probability for these designs.

Furthermore, we will justify Assumption 4.4 for  $\mathcal{L}(\beta)$  being semiparametric

elliptical design loss and logistic loss in §C.2 of the supplementary material (Wang et al., 2014b). Also, in the discussion for logistic loss in §C.2, we prove that the assumption of restricted strong convexity/smoothness in Loh and Wainwright (2013) is stronger than our Assumption 4.4.

Hereafter, we use the shorthands

$$(4.10) \quad \rho_+ = \rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}), \quad \rho_- = \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})$$

for notational simplicity.

**4.2. Main Theorems.** We first provide the main results about the computational rate of convergence. We then establish the statistical properties of the local solutions obtained by our approximate path following method.

**4.2.1. Computational Theory.** The next theorem shows that the proposed approximate regularization path following method achieves a global geometric rate of convergence for calculating the entire regularization path, which is the optimal rate among all first-order optimization methods.

Recall that  $\epsilon_{\text{opt}} \ll \lambda_{\text{tgt}}/4$  is the desired optimization precision of the final path following stage (Line 12 of Algorithm 1), and  $N = \log(\lambda_0/\lambda_{\text{tgt}})/\log(\eta^{-1})$  is the total number of approximate path following stages, where  $\eta \in [0.9, 1)$  is an absolute constant. Meanwhile, remind that  $\rho_- = \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) > 0$  is the smallest sparse eigenvalue specified in Assumption 4.4; As defined in regularity condition (a),  $\zeta_- > 0$  is the concavity parameter of the nonconvex penalty, which satisfies (4.3) in Assumption 4.4.

**Theorem 4.5** (Geometric Rate of Convergence). Under Assumption 4.1 and Assumption 4.4, we have the following results:

1. **Geometric Rate of Convergence within the  $t$ -th Stage:** Within the  $t$ -th ( $t = 1, \dots, N$ ) path following stage (Lines 8 and 12 of Algorithm 1), the iterative sequence  $\{\beta_t^k\}_{k=0}^\infty$  produced by the proximal-gradient method (Algorithm 3) converges to a unique local solution  $\hat{\beta}_{\lambda_t}$ .
  - Within the  $t$ -th path following stage ( $t = 1, \dots, N-1$ ), the total number of proximal-gradient iterations (Lines 5-9 of Algorithm 3) is no more than  $C' \log(4C\sqrt{s^*})$ .
  - Within the  $N$ -th stage ( $\lambda_N = \lambda_{\text{tgt}}$ ), the total number of proximal-gradient iterations is no more than  $\max\{0, C' \log(C\lambda_{\text{tgt}}\sqrt{s^*}/\epsilon_{\text{opt}})\}$ .

Here  $s^* = \|\beta^*\|_0$  and

$$(4.11) \quad C = 2\sqrt{21} \cdot \sqrt{\kappa}(1 + \kappa), \quad C' = 2 \Big/ \log\left(\frac{1}{1 - 1/(8\kappa)}\right),$$

where  $\kappa \in [1, +\infty)$  is the condition number defined in (4.5).

2. **Geometric Rate of Convergence over the Full Path:** To compute the entire path, we need no more than

$$(4.12) \quad \underbrace{(N-1)C' \log(4C\sqrt{s^*})}_{1, \dots, (N-1)\text{-th Stages}} + \underbrace{C' \log\left(\frac{C\lambda_{\text{tgt}}\sqrt{s^*}}{\epsilon_{\text{opt}}}\right)}_{N\text{-th Stage}}$$

proximal-gradient iterations, where  $C, C'$  are specified in (4.11).

3. **Geometric Rate of Convergence of Objective Function Value:**

Let  $\tilde{\beta}_t$  be the approximate local solution obtained within the  $t$ -th stage.

- For  $t = 0, \dots, N-1$ , the value of the objective function decays exponentially towards the value at the final exact local solution  $\hat{\beta}_{\lambda_{\text{tgt}}}$ , i.e.,

$$(4.13) \quad \phi_{\lambda_{\text{tgt}}}(\tilde{\beta}_t) - \phi_{\lambda_{\text{tgt}}}(\hat{\beta}_{\lambda_{\text{tgt}}}) \leq C\lambda_0^2 s^* \cdot \eta^{2(t+1)},$$

where  $C = 105/(\rho_- - \zeta_-)$ .

- For  $t = N$ , we have

$$(4.14) \quad \phi_{\lambda_{\text{tgt}}}(\tilde{\beta}_N) - \phi_{\lambda_{\text{tgt}}}(\hat{\beta}_{\lambda_{\text{tgt}}}) \leq (C'\lambda_{\text{tgt}}s^*) \cdot \epsilon_{\text{opt}},$$

where  $C' = 21/(\rho_- - \zeta_-)$ .

PROOF. See the next section for a detailed proof.  $\square$

Result 1 suggests that, within each path following stage, the proximal-gradient algorithm attains a geometric rate of convergence. More specifically, within the  $t$ -th ( $t = 1, \dots, N$ ) stage (Line 8 and Line 12 of Algorithm 1), we only need a logarithmic number of proximal-gradient update iterations (Lines 5-9 of Algorithm 3) to compute an approximate local solution  $\tilde{\beta}_t$ . Furthermore, within the  $t$ -th path following stage, the iterative sequence  $\{\beta_t^k\}_{k=0}^\infty$  produced by Algorithm 3 converges towards a unique local solution  $\hat{\beta}_{\lambda_t}$ . In Theorem 4.8, we will show that  $\hat{\beta}_{\lambda_t}$  enjoys a more refined statistical rate of convergence due to the usage of nonconvex penalty.

Result 2 suggests that our approximate path following method attains a global geometric rate of convergence. From the perspective of high-dimensional statistics, the total number of stages  $N$  scales with dimension  $d$  and sample size  $n$ , because  $N = \log(\lambda_0/\lambda_{\text{tgt}})/\log(\eta^{-1})$ , where  $\eta$  is an absolute constant. From the perspective of optimization, given dimension  $d$  and sample size  $n$ , when the optimization precision  $\epsilon_{\text{opt}}$  is sufficiently small such that in (4.12) the second term dominates its first term, then the total iteration complexity is  $C \log(1/\epsilon_{\text{opt}})$ . In other words, we only need to conduct a logarithmic number of proximal-gradient iterations to compute the full regularization path.

Recall that we measure the suboptimality of an approximate solution with  $\omega_\lambda(\beta)$  defined in (3.15), which doesn't directly reflect the suboptimality of the objective function value. Hence we provide result 3 to characterize the decay of the objective gap  $\phi_{\lambda_{\text{tgt}}}(\tilde{\beta}_t) - \phi_{\lambda_{\text{tgt}}}(\hat{\beta}_{\lambda_{\text{tgt}}})$ . In detail, (4.13) illustrates the exponential decay of the objective gap along the regularization path, i.e.,  $t = 1, \dots, N-1$ , while (4.14) suggests that, the final objective function value evaluated at  $\tilde{\beta}_N$  is sufficiently close to the value at the exact local solution  $\hat{\beta}_{\lambda_{\text{tgt}}}$ , as long as the optimization precision  $\epsilon_{\text{opt}}$  is sufficiently small.

Recall the largest sparse eigenvalue  $\rho_+ = \rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) > 0$  is specified in Assumption 4.4; As defined in regularity condition (a),  $\zeta_+ > 0$  is the concavity parameter of the nonconvex penalty, which satisfies (4.6) in Assumption 4.4;  $L_{\min}$  is a parameter of Algorithm 3 (Line 3).

**Remark 4.6.** Nesterov (2013) proved that the total number of line-search steps (Lines 4-7 of Algorithm 2) within the  $k$ -th proximal-gradient iteration (Line 8 of Algorithm 3) is no more than

$$2(k+1) + \max \left\{ 0, \frac{\log(\rho_+ - \zeta_+) - \log L_{\min}}{\log 2} \right\}.$$

Piecing the above results together, we conclude that, the total number of line-search iterations (Lines 4-7 of Algorithm 2) required to compute the full regularization path is of the same order as (4.12).

**4.2.2. Statistical Theory.** We present two types of statistical results. Recall that  $\tilde{\beta}_t$  is the approximate local solution obtained within the  $t$ -th path following stage, while  $\hat{\beta}_{\lambda_t}$  is the corresponding exact local solution that satisfies the exact optimality condition in (3.13). In Theorem 4.7, we will provide a statistical characterization of all the approximate local solutions  $\{\tilde{\beta}_t\}_{t=1}^N$  attained along the full regularization path. Remind in Theorem 4.5 we prove that within the  $t$ -th stage, the iterative sequence  $\{\beta_t^k\}_{k=0}^\infty$  produced by the proximal-gradient method converges towards a unique exact local solution  $\hat{\beta}_{\lambda_t}$ . In Theorem 4.8, we will provide more refined statistical properties of these exact local solutions  $\{\hat{\beta}_{\lambda_t}\}_{t=1}^N$  along the full regularization path. Since  $\hat{\beta}_{\lambda_N} = \hat{\beta}_{\lambda_{\text{tgt}}}$ , this result justifies the statistical property of the final estimator.

**Theorem 4.7** (Statistical Rates of Convergence of Approximate Local Solutions). Recall that  $\tilde{\beta}_t$  is the approximate local solution obtained within the  $t$ -th path following stage (Line 8 and Line 12 of Algorithm 1). Under Assumption 4.1 and Assumption 4.4, we have

$$(4.15) \quad \|\tilde{\beta}_t - \beta^*\|_2 \leq C\lambda_t \sqrt{s^*}, \quad \text{for } t = 1, \dots, N,$$

where  $s^* = \|\beta^*\|_0$  and  $C = (21/8)/(\rho_- - \zeta_-)$ .

PROOF. See the next section for a detailed proof.  $\square$

Theorem 4.7 provides statistical rates of convergence of all the approximate local solutions attained by our algorithm along the regularization path. Recall that in Assumption 4.1, we set  $\lambda_{\text{tgt}} = C\sqrt{\log d/n}$  for least squares and logistic loss, and  $\lambda_{\text{tgt}} = C'\|\beta^*\|_1\sqrt{\log d/n}$  for semiparametric elliptical design loss. For least squares and logistic loss, taking  $t = N$  in Theorem 4.7, we have

$$\|\tilde{\beta}_N - \beta^*\|_2 \leq \frac{21/8}{\rho_- - \zeta_-} \lambda_{\text{tgt}} \sqrt{s^*} = \frac{21/8 \cdot C}{\rho_- - \zeta_-} \sqrt{\frac{s^* \log d}{n}}.$$

Hence, the final approximate local solution  $\tilde{\beta}_N$  attains the minimax rate of convergence for parameter estimation. Similarly, for semiparametric elliptical design loss, we have

$$\|\tilde{\beta}_N - \beta^*\|_2 \leq \frac{21/8 \cdot C'}{\rho_- - \zeta_-} \|\beta^*\|_1 \sqrt{\frac{s^* \log d}{n}},$$

which suggests that the rate of convergence of the final approximate local solution is also optimal in the regime where  $\|\beta^*\|_1$  is upper bounded by a constant. Moreover, since  $\eta$  is an absolute constant, for  $\tilde{\beta}_{N-K}$  with  $K$  being a positive integer constant, Theorem 4.7 gives

$$\|\tilde{\beta}_{N-K} - \beta^*\|_2 \leq \frac{21/8}{\rho_- - \zeta_-} \lambda_{N-K} \sqrt{s^*} \leq \frac{21/8 \cdot \eta^{-K}}{\rho_- - \zeta_-} \lambda_{\text{tgt}} \sqrt{s^*},$$

which suggests that, the approximate local solution  $\tilde{\beta}_{N-K}$  enjoys the same rate of convergence as the final approximate local solution  $\tilde{\beta}_N$ , but with a larger constant  $C = (21/8) \cdot \eta^{-K}/(\rho_- - \zeta_-) > (21/8)/(\rho_- - \zeta_-)$ .

In independent work, Theorem 1 and Corollaries 1-3 of [Loh and Wainwright \(2013\)](#) showed the approximate local solution  $\tilde{\beta}$  attained by their optimization procedure satisfies  $\|\tilde{\beta} - \beta^*\|_2 \leq C\lambda_{\text{tgt}}\sqrt{s^*}$ . A comparison between Theorem 4.7 and their result suggests that, our approximate local solution  $\tilde{\beta}_N$  obtained within the final path following stage has the same statistical rate of convergence as the approximate local solution attained by their procedure. Meanwhile, Theorem 4.7 provides additional statistical characterizations for the other regularization parameters along the regularization path, i.e.,  $\lambda_1, \dots, \lambda_{N-1}$ .

In the next theorem, we provide a refined statistical rate of convergence. Recall within the  $t$ -th path following stage, the iterative sequence  $\{\beta_t^k\}_{k=0}^\infty$  produced by the proximal-gradient method converges towards a unique exact local solution  $\hat{\beta}_{\lambda_t}$ . The next theorem states that  $\hat{\beta}_{\lambda_t}$  benefits from nonconvex regularization and possesses an improved statistical rate of convergence.

**Theorem 4.8** (Refined Statistical Rates of Convergence of Exact Local So-



lutions). For the regularization parameter  $\lambda_t$ , we assume that the nonconvex penalty  $\mathcal{P}_{\lambda_t}(\boldsymbol{\beta}) = \sum_{j=1}^d p_{\lambda_t}(\beta_j)$  satisfies

$$(4.16) \quad p'_{\lambda_t}(\beta_j) = 0, \quad \text{for } |\beta_j| \geq \nu_t,$$

for some  $\nu_t > 0$ . Let  $S_1^* \cup S_2^* = S^* = \text{supp}(\boldsymbol{\beta}^*)$  with  $|S_1^*| = s_1^*$ ,  $|S_2^*| = s_2^*$  and  $|S^*| = s^* = s_1^* + s_2^*$ . For  $j \in S_1^* \subseteq S^*$ , we assume  $|\beta_j^*| \geq \nu_t$ , while for  $j \in S_2^* \subseteq S^*$ , we assume  $|\beta_j^*| < \nu_t$ . Under Assumption 4.1 and Assumption 4.4, we have

$$(4.17) \quad \|\hat{\boldsymbol{\beta}}_{\lambda_t} - \boldsymbol{\beta}^*\|_2 \leq \underbrace{C \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{S_1^*}}_{S_1^* : \text{Large } |\beta_j|'s} + \underbrace{C' \lambda_t \sqrt{s_2^*}}_{S_2^* : \text{Small } |\beta_j|'s}, \quad \text{for } t = 1, \dots, N,$$

where  $C = 1/(\rho_- - \zeta_-)$  and  $C' = 3/(\rho_- - \zeta_-)$ .

PROOF. See the next section for a detailed proof.  $\square$

In Theorem 4.8, the assumption in (4.16) applies to a variety of nonconvex penalty functions. For SCAD in (2.1), we have  $\nu_t = a\lambda_t$ ; While for MCP in (2.2), we have  $\nu_t = b\lambda_t$ . Theorem 4.8 suggests that, for “small” coefficients such that  $|\beta_j| < \nu_t$ , the second part on the right-hand side of (4.17) has the same recovery performance as in Theorem 4.7, while for “large” coefficients such that  $|\beta_j| \geq \nu_t$ , the first part in (4.17) possesses a more refined rate of convergence. To understand this, we consider an example with  $\mathcal{L}(\boldsymbol{\beta})$  being least squares loss. We assume that  $(Y|\mathbf{X} = \mathbf{x}_i)$  follows a sub-Gaussian distribution with mean  $\mathbf{x}_i^T \boldsymbol{\beta}^*$  and variance proxy  $\sigma^2$ . Moreover, we assume that the columns of  $\mathbf{X}$  are normalized in such a way that  $\max_{j \in \{1, \dots, d\}} \{\|\mathbf{X}_j\|_2\} \leq \sqrt{n}$ . Then we have

$$(4.18) \quad \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{S_1^*} \leq C\sigma \sqrt{\frac{s_1^*}{n}}$$

with high probability. Clearly, this  $\sqrt{s_1^*/n}$  rate of convergence on the right-hand side of (4.18) is significantly faster than the usual  $\sqrt{s^* \log d/n}$  rate, since it gets rid of the  $\log d$  term, and  $s_1^* \leq s^*$ . In fact,  $\nu_t$  is the minimum signal strength above which we are able to obtain this refined rate of convergence. In the examples of SCAD and MCP, we have  $\nu_t = C\lambda_t$ . Recall that  $\{\lambda_t\}_{t=0}^N$  is a decreasing sequence. Hence, we are able to achieve this more refined rate of convergence for smaller and smaller signal strength along the regularization path. Moreover, for  $t = N$ , the minimum signal strength  $\nu_N = \lambda_N = \lambda_{\text{tgt}} = C\sqrt{\log d/n}$ . Hence, the required minimum signal strength goes to zero as the sample size increases. Following a similar proof of Lemma C.1 and Lemma C.2 in the supplementary material (Wang et al., 2014b), we can also obtain similar results for logistic loss and semiparametric elliptical design loss. This

refined rate of convergence is sharper than the result in Theorem 4.7, which is also achievable via convex regularization, e.g., the  $\ell_1$  penalty. Therefore, Theorem 4.8 clearly justifies the benefits of using nonconvex regularization. Moreover, in §6 we will show that our requirement on the minimum signal strength to achieve this refined rate of convergence is optimal, and is a weaker requirement than the suboptimal requirements in Wang et al. (2013); Fan et al. (2014).

In addition to the refined rate of convergence for parameter estimation in Theorem 4.8, in the next theorem we prove that the exact local solution  $\hat{\beta}_{\lambda_t}$  also recovers the support of  $\beta^*$ . Before we present the next theorem, we introduce the definition of an oracle estimator, denoted by  $\hat{\beta}_O$ . Recall that  $S^* = \text{supp}(\beta^*)$ . The oracle estimator  $\hat{\beta}_O$  is defined as

$$(4.19) \quad \hat{\beta}_O = \underset{\substack{\text{supp}(\beta) \subseteq S^* \\ \beta \in \Omega}}{\text{argmin}} \mathcal{L}(\beta),$$

where  $\Omega = \mathbb{R}^d$  for least squares loss and semiparametric elliptical design loss, while  $\Omega = B_2(R)$  for logistic loss with  $R \geq \|\beta^*\|_2$ . In the next lemma, we show that  $\hat{\beta}_O$  is the unique global solution to the minimization problem in (4.19) even for nonconvex loss functions, and has nice statistical properties.

**Lemma 4.9.** Under Assumption 4.4, the oracle estimator  $\hat{\beta}_O$  is the unique global minimizer of (4.19). For  $\mathcal{L}(\beta)$  being least squares loss, we assume that  $(Y|\mathbf{X} = \mathbf{x}_i)$  follows a sub-Gaussian distribution with mean  $\mathbf{x}_i^T \beta^*$  and variance proxy  $\sigma^2$ , then the oracle estimator satisfies

$$(4.20) \quad \|\hat{\beta}_O - \beta^*\|_\infty \leq C\sigma\sqrt{2/\rho_-} \cdot \sqrt{\frac{\log s^*}{n}}$$

with high probability for some constant  $C$ .

**PROOF.** See the supplementary material (Wang et al., 2014b) for a detailed proof.  $\square$

Statistical recovery results similar to (4.20) also hold for logistic loss and semiparametric elliptical design loss under different conditions. These results are omitted here for simplicity. Lemma 4.9 suggests that, for a sufficiently large  $n$  and sufficient minimum signal strength, the oracle estimator  $\hat{\beta}_O$  exactly recovers the support of  $\beta^*$ . More specifically, if the minimum signal strength satisfies  $\min_{j \in S^*} |\beta_j^*| \geq 2\nu$  for  $\nu > 0$ , then with high probability

$$\min_{j \in S^*} |(\hat{\beta}_O)_j| \geq \min_{j \in S^*} |\beta_j^*| - \|\hat{\beta}_O - \beta^*\|_\infty \geq 2\nu - \sigma\sqrt{2/\rho_-} \cdot \sqrt{\frac{\log s^*}{n}},$$

which implies  $\min_{j \in S^*} |(\hat{\beta}_O)_j| \geq \nu > 0$  for  $n$  sufficiently large. Meanwhile, recall that  $\text{supp}(\hat{\beta}_O) \subseteq S^*$  by definition. Hence we have  $\text{supp}(\hat{\beta}_O) = S^*$ .

The next theorem states, under the condition of sufficient minimum signal strength,  $\hat{\beta}_{\lambda_t}$  is the oracle estimator, and exactly recovers the support of  $\beta^*$ .

**Theorem 4.10** (Support Recovery). For the regularization parameter  $\lambda_t$ , suppose that the nonconvex penalty  $\mathcal{P}_{\lambda_t}(\beta) = \sum_{j=1}^d p_{\lambda_t}(\beta_j)$  satisfies (4.16) for some  $\nu_t > 0$ . For least squares loss, we assume that  $(Y|\mathbf{X} = \mathbf{x}_i)$  follows a sub-Gaussian distribution with mean  $\mathbf{x}_i^T \beta^*$  and variance proxy  $\sigma^2$ . Under Assumption 4.1 and Assumption 4.4, if the minimum signal strength satisfies  $\min_{j \in S^*} |\beta_j^*| \geq 2\nu_t$ , then for  $n$  sufficiently large,  $\hat{\beta}_{\lambda_t} = \hat{\beta}_O$ , and  $\text{supp}(\hat{\beta}_{\lambda_t}) = \text{supp}(\hat{\beta}_O) = \text{supp}(\beta^*)$  with high probability.

PROOF. See the next section for a detailed proof.  $\square$

Recall the assumption in (4.16) applies to a variety of nonconvex penalties, including SCAD and MCP, for which we have  $\nu_t = C\lambda_t$  with  $C > 0$ . Hence, the minimum signal strength that is required to achieve the oracle estimator and exact support recovery actually shrinks with the decreasing sequence  $\{\lambda_t\}_{t=0}^N$  along the regularization path. For least squares loss, we have  $\nu_N = C\lambda_{\text{tgt}} = C'\sqrt{\log d/n}$  for  $t = N$ . Hence, within the final path following stage, the required minimum signal strength goes to zero as sample size  $n \rightarrow \infty$ . Furthermore, such requirement on the minimum signal strength for achieving the oracle estimator is optimal, i.e., no weaker requirement exists (Zhang and Zhang, 2012). In §6 we will show that, for least squares loss, some of recent works (Fan et al., 2014; Wang et al., 2013) require a stronger minimum signal strength to achieve the oracle estimator in the same setting of least squares regression. Similar results to Theorem 4.10 also hold for other loss functions, but under different conditions. They are omitted here for simplicity.

**5. Proof of Main Results.** In this section we present the proof sketch of the main results. The desired computational and statistical results rely on the strong convexity of the surrogate loss function  $\tilde{\mathcal{L}}_\lambda(\beta)$ , e.g., we need  $\tilde{\mathcal{L}}_\lambda(\beta)$  to be strongly convex to establish the geometric rate of convergence of the proximal-gradient method within each path following stage. However,  $\tilde{\mathcal{L}}_\lambda(\beta)$  is nonconvex in general, since  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$ , where  $\mathcal{L}(\beta)$  is possibly nonconvex and  $\mathcal{Q}_\lambda(\beta)$  is concave. In the following lemma, we prove that  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$  is strongly convex for  $\beta$  on a sparse set. In a similar way, we establish the strong smoothness of  $\tilde{\mathcal{L}}_\lambda(\beta)$  on a sparse set.

Recall  $\rho_- = \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})$  and  $\rho_+ = \rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})$  are the sparse eigenvalues specified in Assumption 4.4. As defined in regularity condition

(a),  $\zeta_-, \zeta_+ > 0$  are the concavity parameters of the nonconvex penalty, which satisfy (4.3) and (4.6).

**Lemma 5.1.** Let  $\beta, \beta' \in \mathbb{R}^d$  be two sparse vectors, which satisfy  $\|(\beta - \beta')_{\bar{S}^*}\|_0 \leq 2\tilde{s}$ , where  $\tilde{s}$  is specified in Assumption 4.4 and  $S^* = \text{supp}(\beta^*)$ . For  $\mathcal{L}(\beta)$  being logistic loss, we further assume  $\|\beta\|_2 \leq R$  and  $\|\beta'\|_2 \leq R$ , where  $R$  is a constant specified in Definition 4.3. Then the surrogate loss function  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$  satisfies the restricted strong convexity

$$\tilde{\mathcal{L}}_\lambda(\beta') \geq \tilde{\mathcal{L}}_\lambda(\beta) + \nabla \tilde{\mathcal{L}}_\lambda(\beta)^T (\beta' - \beta) + \frac{\rho_- - \zeta_-}{2} \|\beta' - \beta\|_2^2,$$

and the restricted strong smoothness

$$\tilde{\mathcal{L}}_\lambda(\beta') \leq \tilde{\mathcal{L}}_\lambda(\beta) + \nabla \tilde{\mathcal{L}}_\lambda(\beta)^T (\beta' - \beta) + \frac{\rho_+ - \zeta_+}{2} \|\beta' - \beta\|_2^2.$$

PROOF. See §D.2 in the supplementary material (Wang et al., 2014b) for a detailed proof.  $\square$

A similar result has been discussed by Negahban et al. (2012). The main difference is that, our constraint set where  $\tilde{\mathcal{L}}_\lambda(\beta)$  is strongly convex/smooth is a sparse subspace, while that of Negahban et al. (2012) is a cone.

Note that in Lemma 5.1, the strong convexity and smoothness of  $\tilde{\mathcal{L}}_\lambda(\beta)$  rely on the sparsity of  $\beta$  and  $\beta'$ . Hence, we need to establish results regarding the sparsity of  $\beta_t^k$  throughout the whole iterative procedure. In the sequel, we provide several important lemmas: Lemma 5.2 and Lemma 5.3 characterize the statistical properties of any sparse  $\beta$ ; Based on such statistical properties, Lemma 5.4 proves that, any proximal-gradient update iteration with a sparse input produces a sparse output. Equipped with these lemmas, we can establish the sparsity of the solution path by mathematical induction in Theorem 5.5.

The next lemma provides a characterization of any sparse  $\beta$  with certain suboptimality.

**Lemma 5.2.** We assume that  $\beta$  satisfies

$$(5.1) \quad \|\beta_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \omega_\lambda(\beta) \leq \lambda/2$$

with  $\lambda \geq \lambda_{\text{tgt}}$ , where  $\omega_\lambda(\beta)$  is the measure of suboptimality defined in (3.15). For logistic loss, we assume  $\|\beta\|_2 \leq R$ , where  $R > 0$  is a constant specified in Definition 4.3. Under Assumption 4.1 and Assumption 4.4,  $\beta$  satisfies

$$\|\beta - \beta^*\|_2 \leq C\lambda\sqrt{s^*}, \quad \text{where } C = \frac{21/8}{\rho_- - \zeta_-}.$$

Meanwhile, the objective function value evaluated at  $\beta$  satisfies

$$\phi_\lambda(\beta) - \phi_\lambda(\beta^*) \leq C' \lambda^2 s^*, \quad \text{where } C' = \frac{21/2}{\rho_- - \zeta_-}.$$

PROOF. See §D.3 of the supplementary material (Wang et al., 2014b) for a detailed proof.  $\square$

Recall that we use the approximate local solution  $\tilde{\beta}_{t-1}$  obtained within the  $(t-1)$ -th path following stage to be the initialization of the  $t$ -th stage (Line 8 of Algorithm 1), i.e.,  $\beta_t^0 = \tilde{\beta}_{t-1}$ . By setting  $\beta = \tilde{\beta}_{t-1} = \beta_t^0$  and  $\lambda = \lambda_t$  in Lemma 5.2, we can see that, if  $\tilde{\beta}_{t-1}$  is sparse and  $(\lambda_t/2)$ -suboptimal, then the initial point  $\beta_t^0$  of the  $t$ -th stage has nice statistical recovery performance. However, it is unclear whether the rest of  $\beta_t^k$ 's ( $k = 1, 2, \dots$ ) within the  $t$ -th stage also have similar recovery performance. To prove this, we first present Lemma 5.3, which shows that under the condition that  $\beta$  is sparse and  $\phi_\lambda(\beta)$  is close to  $\phi_\lambda(\beta^*)$ ,  $\beta$  has desired statistical properties. After Lemma 5.3, we will explain that if  $\beta_t^0$  satisfies this condition, then all the  $\beta_t^k$ 's ( $k = 1, 2, \dots$ ) within the same path following stage also satisfy this condition, and therefore enjoys nice statistical properties.

**Lemma 5.3.** Suppose that, for  $\lambda \geq \lambda_{\text{tgt}}$ ,  $\beta$  satisfies

$$\|\beta_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \phi_\lambda(\beta) - \phi_\lambda(\beta^*) \leq C \lambda^2 s^*, \quad \text{where } C = \frac{21/2}{\rho_- - \zeta_-}.$$

For logistic loss, we further assume  $\|\beta\|_2 \leq R$ , where  $R$  is a constant specified in Definition 4.3. Under Assumption 4.1 and Assumption 4.4, we have

$$\|\beta - \beta^*\|_2 \leq C' \lambda \sqrt{s^*}, \quad \text{where } C' = \frac{15/2}{\rho_- - \zeta_-}.$$

PROOF. See §D.4 of the supplementary material (Wang et al., 2014b) for a detailed proof.  $\square$

Let  $\lambda = \lambda_t$  and  $\beta = \beta_t^k$  in Lemma 5.3. It suggests that within the  $t$ -th path following stage, all  $\beta_t^k$ 's ( $k = 1, 2, \dots$ ) have nice statistical recovery performance under three sufficient conditions: (i) Each  $\beta_t^k$  is sparse; (ii) The objective function value  $\phi_{\lambda_t}(\beta_t^k)$  is close to  $\phi_{\lambda_t}(\beta^*)$ ; (iii) For logistic loss, we further need  $\|\beta_t^k\|_2 \leq R$ . For condition (ii), recall that if we set  $\beta = \beta_t^0$  and  $\lambda = \lambda_t$  in Lemma 5.2, then  $\beta_t^0$  being sparse and  $(\lambda_t/2)$ -suboptimal implies that  $\phi_{\lambda_t}(\beta_t^0)$  is close to  $\phi_{\lambda_t}(\beta^*)$ . Since the proximal-gradient method ensures the monotone decrease of  $\{\phi_{\lambda_t}(\beta_t^k)\}_{k=0}^\infty$  within the  $t$ -th stage (see Lemma D.1 of the supplementary material (Wang et al., 2014b)), condition (ii) also holds. Meanwhile, condition (iii) obviously holds because of the

$\ell_2$  constraint. To establish the statistical recovery performance of all the  $\beta_t^k$ 's within the  $t$ -th stage, we still need to establish the sparsity of  $\beta_t^k$ 's to guarantee condition (i) holds. To prove this, we present Lemma 5.4, which states that if  $\beta$  is sparse, then a proximal-gradient update operation (3.9) on  $\beta$  produces a sparse solution.

**Lemma 5.4.** Suppose that, for  $\lambda \geq \lambda_{\text{tgt}}$ ,  $\beta$  satisfies

$$\|\beta_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \phi_\lambda(\beta) - \phi_\lambda(\beta^*) \leq C\lambda^2 s^*, \quad \text{and} \quad L < 2(\rho_+ - \zeta_+),$$

where  $C = (21/2)/(\rho_- - \zeta_-)$ . For logistic loss, we assume  $\|\beta\|_2 \leq R$ , where  $R$  is specified in Definition 4.3. Under Assumption 4.1 and Assumption 4.4, the proximal-gradient update operation defined in (3.9) produces a sparse solution, i.e.,

$$\|(\mathcal{T}_{L,\lambda}(\beta; R))_{\bar{S}^*}\|_0 \leq \tilde{s}.$$

Here we set  $R = +\infty$  if the domain  $\Omega$  in (3.8) is  $\mathbb{R}^d$ .

PROOF. See §D.5 of the supplementary material (Wang et al., 2014b) for a detailed proof.  $\square$

For  $\beta = \beta_t^{k-1}$ ,  $\lambda = \lambda_t$  and  $L = L_t^k$ , Lemma 5.4 states that, if  $\beta_t^{k-1}$  is sparse and the objective function value  $\phi_{\lambda_t}(\beta_t^{k-1})$  is close to  $\phi_{\lambda_t}(\beta^*)$ , then  $\beta_t^k = \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; R)$  produced by the proximal-gradient update step (3.8) is also sparse. Within the  $t$ -th path following stage, if  $\beta_t^0$  is sparse,  $\omega_{\lambda_t}(\beta_t^0) \leq \lambda_t/2$ , and for logistic loss  $\|\beta_t^0\|_2 \leq R$ , then by Lemma 5.2 we have

$$\phi_{\lambda_t}(\beta_t^0) - \phi_{\lambda_t}(\beta^*) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda_t^2 s^*.$$

Since  $\{\phi_{\lambda_t}(\beta_t^k)\}_{k=0}^\infty$  decreases monotonically, we have

$$\phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\beta^*) \leq \phi_{\lambda_t}(\beta_t^0) - \phi_{\lambda_t}(\beta^*) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda_t^2 s^*, \quad \text{for } k = 1, 2, \dots$$

Assume that we have  $L_t^k \leq 2(\rho_+ - \zeta_+)$  (which will be proved in Theorem 5.5). Applying Lemma 5.4 recursively, we obtain  $\|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}$  ( $k = 1, 2, \dots$ ). Meanwhile, we have  $\|\beta_t^k\|_2 \leq R$  due to the  $\ell_2$  constraint. Then according to Lemma 5.3, all  $\beta_t^k$ 's within the  $t$ -th path following stage have nice recovery performance, i.e.,

$$\|\beta_t^k - \beta^*\|_2 \leq \frac{15/2}{\rho_- - \zeta_-} \lambda_t \sqrt{s^*}, \quad \text{for } k = 1, 2, \dots$$

Furthermore, by Lemma 5.1 the sparsity of  $\beta_t^k$ 's implies the restricted strong convexity and smoothness of  $\tilde{\mathcal{L}}_{\lambda_t}(\beta)$ , which enable us to establish the geo-

metric rate of convergence within the  $t$ -th path following stage. These results are formally presented in Theorem 5.5.

**Theorem 5.5.** Suppose within the  $t$ -th path following stage, the proximal-gradient method in Algorithm 3 is initialized by  $\beta_t^0$  and  $L_t^0$ , which satisfy

$$\|(\beta_t^0)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \omega_{\lambda_t}(\beta_t^0) \leq \lambda_t/2, \quad \text{and} \quad L_t^0 \leq 2(\rho_+ - \zeta_+).$$

For logistic loss we further assume  $\|\beta_t^0\|_2 \leq R$  with  $R$  specified in Definition 4.3. Then we have the following results:

- For  $k = 1, 2, \dots$ , we have

$$(5.2) \quad \|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \|\beta_t^k - \beta^*\|_2 \leq \frac{15/2}{\rho_- - \zeta_-} \lambda_t \sqrt{s^*}, \quad L_t^k \leq 2(\rho_+ - \zeta_+).$$

- The iterative sequence  $\{\beta_t^k\}_{k=0}^\infty$  converges towards a unique exact local solution  $\hat{\beta}_{\lambda_t}$ , which satisfies  $\|(\hat{\beta}_{\lambda_t})_{\bar{S}^*}\|_0 \leq \tilde{s}$  and the exact optimality condition that  $\omega_{\lambda_t}(\beta_t^k) \leq 0$ .
- To achieve an approximate local solution  $\tilde{\beta}_t$  that satisfies  $\omega_{\lambda_t}(\tilde{\beta}_t) \leq \lambda_t/4$ , we need no more than  $C' \log(4C\sqrt{s^*})$  proximal-gradient iterations defined in Lines 5-9 of Algorithm 3. Here

$$(5.3) \quad C = 2\sqrt{21} \cdot \sqrt{\kappa}(1 + \kappa), \quad C' = 2 \left/ \log \left( \frac{1}{1 - 1/(8\kappa)} \right) \right.$$

- To obtain an approximate local solution  $\tilde{\beta}_t$  such that  $\omega_{\lambda_t}(\tilde{\beta}_t) \leq \epsilon_{\text{opt}}$ , we need no more than  $C' \log(C\lambda_t\sqrt{s^*}/\epsilon_{\text{opt}})$  proximal-gradient iterations. Here  $C$  and  $C'$  are defined in (5.3).

PROOF. See §D.6 of the supplementary material (Wang et al., 2014b) for a detailed proof.  $\square$

To prove the geometric rate of convergence and desired statistical recovery results hold within all path following stages, i.e.,  $t = 0, \dots, N$ , we need to verify that the conditions of Theorem 5.5 hold within each stage. We prove by induction. We assume the initialization of  $(t-1)$ -th path following stage satisfies

$$(5.4) \quad \|(\beta_{t-1}^0)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \omega_{\lambda}(\beta_{t-1}^0) \leq \lambda_t/2, \quad \text{and} \quad L_{t-1}^0 \leq 2(\rho_+ - \zeta_+).$$

Applying Theorem 5.5, we obtain

$$\|(\beta_{t-1}^k)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad L_{t-1}^k \leq 2(\rho_+ - \zeta_+), \quad \text{for } k = 1, 2, \dots$$

Consequently, the approximate solution  $\tilde{\beta}_{t-1}$  produced by the  $(t-1)$ -th stage satisfies  $\|(\tilde{\beta}_{t-1})_{\bar{S}^*}\|_0 \leq \tilde{s}$ , while  $L_{t-1}$  satisfies  $L_{t-1} \leq 2(\rho_+ - \zeta_+)$ . Since

we warm start the  $t$ -th path following stage with  $\beta_t^0 = \tilde{\beta}_{t-1}$  and  $L_t^0 = L_{t-1}$  (Line 8 of Algorithm 1), we have

$$(5.5) \quad \|(\beta_t^0)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad L_t^0 \leq 2(\rho_+ - \zeta_+).$$

Moreover, note that the stopping criterion of the proximal-gradient method ensures  $\omega_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) \leq \lambda_{t-1}/4$  (Line 9 of Algorithm 3), which implies  $\omega_{\lambda_t}(\tilde{\beta}_{t-1}) \leq \lambda_t/2$  according to Lemma D.4 of the supplementary material (Wang et al., 2014b). Thus we have

$$(5.6) \quad \omega_{\lambda_t}(\beta_t^0) \leq \lambda_t/2.$$

Therefore, we know that (5.4) implies (5.5) and (5.6). We will verify (5.5) and (5.6) hold for  $t = 0$  in the proof of Theorem 4.5 in the supplementary material (Wang et al., 2014b). By induction, we have that (5.5) and (5.6) hold for  $t = 0, \dots, N$ . As a consequence of Theorem 5.5, all path following stages have geometric rates of convergence along the solution path, which implies the global geometric rate of convergence in Theorem 4.5. See the supplementary material (Wang et al., 2014b) for a detail proof. Meanwhile, all  $\beta_t^k$ 's have desired statistical properties, i.e.,

$$\|\beta_t^k - \beta^*\|_2 \leq \frac{15/2}{\rho_- - \zeta_-} \lambda_t \sqrt{s^*}, \quad \text{for } t = 1, \dots, N \text{ and } k = 0, 1, \dots,$$

which leads to the statistical rates of convergence of the approximate local solutions  $\{\tilde{\beta}_t\}_{t=1}^N$  in Theorem 4.7, the more refined rates of convergence of the exact local solutions  $\{\hat{\beta}_{\lambda_t}\}_{t=1}^N$  in Theorem 4.8, and the support recovery results in Theorem 4.10. See §D.8-§D.10 of the supplementary material (Wang et al., 2014b) for detailed proofs respectively.

**6. Discussion.** Our work is related to recent works on understanding nonconvex regularization in the context of least squares regression. Zhang (2010a) proposed an MC+ procedure for MCP penalized least squares regression. However, the computation of MC+ might be inefficient because there can be exponentially many switching points on its solution path. To remedy this issue, Zhang (2010b); Zhang et al. (2013) proposed the multi-stage convex relaxation method, which iteratively solves

$$(6.1) \quad \hat{\beta}^k \leftarrow \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \mathcal{L}(\beta) + \sum_{j=1}^d p'_\lambda(|\hat{\beta}_j^{k-1}|) |\beta_j| \right\}, \quad k = 1, 2, \dots,$$

where  $p_\lambda(\beta_j)$  is defined in §2, and the initialization  $\hat{\beta}^0$  is set to be the Lasso estimator corresponding to  $\lambda$ . For  $k$  sufficiently large,  $\hat{\beta}^k$  has the same oracle properties as in Theorem 4.8 and Theorem 4.10. However, for each  $k$  we need to solve the minimization problem in (6.1) exactly, which is not realistic in



practice, since practical optimization methods only attain finite numerical precision in finite iterations. In contrast, we provide simultaneous statistical and computational analysis by explicitly taking the numerical precision into account, and establish the global geometric rate of convergence in terms of iteration complexity for calculating the full regularization path.

This multi-stage convex relaxation method was previously referred to as local linear approximation (LLA), and was analyzed on fixed dimensional models by [Zou and Li \(2008\)](#). [Fan et al. \(2014\)](#) recently provided nonasymptotic analysis of LLA, and proved that LLA finds the oracle estimator in two iterations. However, their results rely on that the Lasso initialization satisfies  $\|\hat{\beta}^0 - \beta^*\|_\infty \leq C\lambda$  with high probability, which requires  $\lambda$  to take the value of  $C'\sqrt{s^* \log d/n}$ . Consequently, their requirement on the minimum signal strength is of the order of  $\sqrt{s^* \log d/n}$ , which is suboptimal. In contrast, we only require a minimum signal strength of the order of  $\sqrt{\log d/n}$ , which is optimal ([Zhang and Zhang, 2012](#)). Also, they didn't analyze the iteration complexity for computing each step of LLA, i.e., solving (6.1).

Very recently, [Wang et al. \(2013\)](#) considered a two-step approach similar to the two-step LLA procedure, named the calibrated CCCP. It differs from the two-step LLA in that, its Lasso initialization  $\hat{\beta}^0$  is obtained using the regularization parameter  $\tau\lambda$ , where  $\tau = o(1)$  and  $\lambda = \sqrt{\log d/n}$ . It attains the oracle estimator under the restricted eigenvalue (RE) condition ([Bickel et al., 2009](#)), but requires the minimum signal strength to be larger than  $Cs^*\sqrt{\log d/n}$ . Under a stronger assumption than the RE condition, namely the relaxed sparse Riesz condition, a minimum signal strength of the order of  $\sqrt{\log d/n}/\tau$  is required. Such requirement is still suboptimal, but is close to the optimal scaling of  $\sqrt{\log d/n}$  in our results, since  $\tau$  can take  $1/\log n$ . They proposed a novel high-dimensional BIC criterion, which can be used to choose the best  $\lambda_{\text{tgt}}$  in our procedure. Also, they provided extensions to logistic regression.

The iterative hard thresholding (IHT) algorithm ([Blumensath and Davies, 2009](#)) can also achieve a local solution with desired statistical recovery performance at a global geometric rate of convergence. However, the theoretical results of IHT are not directly comparable with ours because of the usage of different noise models. If we have to cast the theoretical results of IHT into our model, their results are much weaker than ours. In detail, IHT attains an approximate local solution  $\tilde{\beta}$ , which satisfies

$$(6.2) \quad \|\tilde{\beta} - \beta^*\|_2 \leq 6\|\mathbf{e}\|_2$$

with high probability. Here  $\mathbf{e} \in \mathbb{R}^n$  is the noise vector in their setting, which is often considered to be perturbation noise. Note that a proper normalization

gives  $\mathbf{e} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*)/\sqrt{n}$ , where  $\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*$  is considered to be the sub-Gaussian noise with zero mean and variance proxy  $\sigma^2$  in our setting. Then (6.2) gives

$$(6.3) \quad \|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2 \leq 6\|\mathbf{e}\|_2 = 6\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*\|_2/\sqrt{n} \leq 6 \cdot 4\sigma\sqrt{n}/\sqrt{n} = 24\sigma$$

with high probability. Note that the upper bound on the right-hand side of (6.3) doesn't depend on  $s^*$  and  $d$ , and fails to converge to zero as  $n \rightarrow \infty$ . In summary, casting the results of IHT into our setting of sub-Gaussian noise yields a rather weak result. Also, IHT requires prior knowledge on the true sparsity level  $s^*$  to achieve fast global convergence, while our method doesn't.

In addition, the difference between our work and the independent work by [Loh and Wainwright \(2013\)](#) has been discussed in §1 and §4 with details.

**7. Numerical Results.** We provide numerical results illustrating the computational efficiency and statistical accuracy of the proposed method. In detail, first we illustrate the effectiveness of our method on a problem with both nonconvex loss and penalty functions. Then we conduct comparison between our method and existing nonconvex procedures.

In the first experiment, we consider  $\mathcal{L}(\boldsymbol{\beta})$  being semiparametric elliptical random design loss defined in (2.6) and  $\mathcal{P}_\lambda(\boldsymbol{\beta})$  being the MCP penalty defined in (2.2). We test on a synthetic dataset with  $n = 500$  samples and  $d = 2500$  dimensions. See the supplementary material ([Wang et al., 2014b](#)) for the detailed settings.

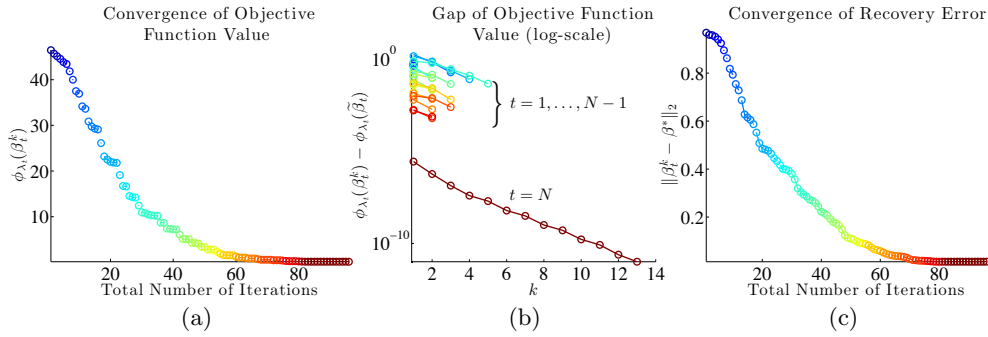


FIG 3. *Semiparametric elliptical design regression with MCP: (a) Plot of the objective function value  $\phi_{\lambda_t}(\beta_t^k)$  along the regularization path; (b) Plot of  $\phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\tilde{\beta}_t)$  (log-scale) within each path following stage; (c) Plot of the recovery error  $\|\beta_t^k - \beta^*\|_2$ . Here we illustrate each path following stage ( $t = 1, \dots, N$ ) with a different color. Note that each point in the figure denotes  $\beta_t^k$ , which corresponds to the  $k$ -th iteration of the proximal-gradient method (Algorithm 3) within the  $t$ -th path following stage.*

As shown in Figure 3(a), the objective function value  $\phi_{\lambda_t}(\beta_t^k)$  is monotone decreasing along the regularization path, as characterized by our theory (see Lemma D.1 of the supplementary material ([Wang et al., 2014b](#))), and

converges eventually.

Figure 3(b) illustrates the geometric rate of convergence within each path following stage. In detail, each line denotes a path following stage. It shows the objective function value gap, i.e.,  $\phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\tilde{\beta}_t)$ , decays exponentially with  $k$  within each stage. Note that

$$\begin{aligned}
 \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\tilde{\beta}_t) &\geq \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\beta_t^k) \\
 (7.1) \quad &\geq \frac{L_t^k}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2 \geq \frac{L_t^k}{2} \frac{\omega_{\lambda_t}^2(\beta_t^k)}{(L_t^k + \rho_+ - \zeta_-)^2}.
 \end{aligned}$$

Here the first inequality is because the objective function  $\phi_{\lambda_t}(\beta_t^k)$  is monotone decreasing, while the second and third inequalities follow from Lemma D.1 and Lemma D.2 of the supplementary material (Wang et al., 2014b) respectively. Therefore,  $\omega_{\lambda_t}(\beta_t^k)$  also decays exponentially within each stage, which implies that we only need a logarithmic number of iterations to attain the desired approximate local solution within each path following stage, as characterized by Theorem 4.5.

Figure 3(b) illustrates the success of the path following scheme in Figure 1: The  $k = 1$  point on each line denotes the initialization of the corresponding path following stage, e.g., the  $t$ -th stage. Recall that such initialization is set to be the approximate local solution  $\tilde{\beta}_{t-1}$  obtained within the  $(t-1)$ -th stage, which falls into the region of optimization precision in Figure 1. Meanwhile, the fast convergence within the  $t$ -th stage suggests that  $\beta_{t-1}$  also falls into the region of fast convergence in Figure 1. Thus, the path following scheme works exactly as we have described in Figure 1 empirically.

Figure 3(c) shows that the  $\ell_2$  recovery error decays towards a small value as the optimization method proceeds, which implies the attained approximate local solution has desired statistical properties, as predicted by Theorem 4.7.

In the second experiment, we compare our method with several existing nonconvex procedures on statistical performance, including LLA (Zou and Li, 2008), the calibrated CCCP (Wang et al., 2013), SparseNet (Mazumder et al., 2011), and the multi-stage convex relaxation method (Zhang, 2010b; Zhang et al., 2013). We consider an example of least squares regression with MCP, where  $n = 200$ ,  $d = 2000$  and  $\|\beta^*\|_0 = 10$ . See the supplementary material (Wang et al., 2014b) for the detailed settings.

We compare the support recovery performance and  $\ell_2$  recovery error of the estimators obtained from these procedures in Table 1, where we use the Lasso estimator and the oracle estimator defined in (4.19) as references. For support recovery, we are interested in the cardinality of true positive sets (TPS) and false positive set (FPS), both of which are defined in Table 1.

Ideally, the cardinality of TPS should be as large as  $\|\beta^*\|_0$  (which is 10 in

TABLE 1

Comparing statistical performance of nonconvex procedures: TPS/FPS denote the true/false positive sets, which are defined as  $\{j \in S^* : \hat{\beta}_j \neq 0\}$  and  $\{j \in \bar{S}^* : \hat{\beta}_j \neq 0\}$  respectively, and  $|\cdot|$  denotes their cardinality. The  $\ell_2$  recovery error is defined as  $\|\hat{\beta} - \beta^*\|_2$ , where  $\hat{\beta}$  is the estimator. Standard deviations are present in the parentheses.

Method	TPS	FPS	$\ell_2$ Error
<b>Approximate Path Following</b>	<b>10</b> (0)	<b>0.180</b> (0.0411)	<b>0.702</b> (0.0278)
SparseNet	10 (0)	0.950 (0.108)	0.848 (0.0230)
Multi-stage Convex Relaxation	10 (0)	2.21 (0.146)	1.28 (0.0753)
LLA	10 (0)	2.98 (0.304)	1.28 (0.0996)
Calibrated CCCP	9.99 (0.01)	3.28 (0.308)	1.40 (0.122)
Lasso	9.98 (0.0141)	31.15 (0.799)	2.63 (0.0460)
Oracle Estimator	10 (0)	0 (0)	0.484 (0.0221)

this example), since a good procedure should exactly identify  $S^* = \text{supp}(\beta^*)$ . Meanwhile, the cardinality of FPS should be close to zero, i.e., few of the coordinates in  $\bar{S}^*$  is wrongly identified as nonzero. Table 1 shows that all nonconvex procedures significantly outperform the Lasso, which produces a less sparse estimator with larger  $\ell_2$  recovery error. In this specific example, our method outperforms the existing nonconvex procedures. Moreover, our method almost recovers  $S^*$  exactly, and achieves a small  $\ell_2$  recovery error that is very close to the  $\ell_2$  error of the oracle estimator, as characterized by Theorem 4.8.

**8. Conclusion.** In this paper, we provided unified theory for penalized  $M$ -estimators with possibly nonconvex loss and penalty functions. These problems are motivated by generalized linear models with nonconvex penalties and semiparametric elliptical design regression, as well as a broad range of other applications. Because it is intractable to compute the global solutions of these problems due to the nonconvex formulation, we need to establish theory that characterizes both the computational and statistical properties of the local solutions obtained by specific algorithms. For this purpose, we proposed an approximate regularization path following method, which serves as a unified framework for solving a variety of high-dimensional sparse learning problems with nonconvexity. Computationally, our method enjoys a fast global geometric rate of convergence for calculating the entire regularization path; Statistically, all the approximate and exact local solutions attained by our method along the regularization path possess sharp statistical rate of convergence in both estimation and support recovery. In particular, we provide sharp theoretical analysis that demonstrates the advantage of using nonconvex penalties. This paper shows that, under suitable conditions, we can efficiently obtain the entire regularization path of a broad class of nonconvex

sparse learning problems.

Our work can be extended in many directions: Our method and theory for least squares loss and logistic loss can be easily extended to other generalized linear models (see §C.2 of the supplementary material (Wang et al., 2014b) for details); For inverse covariance matrix estimation, our work is directly applicable to the Sparse Column Inverse Operator (SCIO) (Liu and Luo, 2012); Meanwhile, it might need more effort than verifying Assumption 4.1 and Assumption 4.4 to adapt the graphical Lasso into our framework, e.g., the optimization algorithm also has to be modified to enforce the positive semidefinite constraint; It is also interesting to consider other loss functions, e.g., quantile regression (Wang et al., 2012), for which Assumption 4.4 may no longer hold.

**Acknowledgement.** We sincerely thank Po-Ling Loh, Martin Wainwright and Yiyuan She for their helpful personal communications. We are grateful to the Editor, Associate Editor and referees for their insightful comments.

Han Liu is supported by NSF Grants III-1116730 and NSF III-1332109, NIH R01MH102339, NIH R01GM083084, and NIH R01HG06841, and FDA HHSF223201000072C. Tong Zhang is supported by the following grants: NSF IIS-1016061, NSF DMS-1007527, and NSF IIS-1250985.

## SUPPLEMENTARY MATERIAL

**Supplementary material for: Optimal Computational and Statistical Rates of Convergence for Sparse Nonconvex Learning Problems** (DOI: [To Be Assigned](#); .pdf). We provide the detailed proof in the supplement (Wang et al., 2014b).

## References.

- AGARWAL, A., NEGAHBAN, S. and WAINWRIGHT, M. J. (2012). Fast global convergence of gradient methods for high-dimensional statistical recovery. *Annals of Statistics* **40** 2452–2482.
- BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *Annals of Statistics* **37** 1705–1732.
- BLUMENSATH, T. and DAVIES, M. E. (2009). Iterative hard thresholding for compressed sensing. *Applied and Computational Harmonic Analysis* **27** 265–274.
- BREHENY, P. and HUANG, J. (2011). Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *Annals of Applied Statistics* **5** 232.
- CANDÉS, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *Annals of Statistics* 2313–2351.
- CANDÉS, E. J. and TAO, T. (2005). Decoding by linear programming. *IEEE Transactions on Information Theory* **51** 4203–4215.
- CATONI, O. (2012). Challenging the empirical mean and empirical variance: a deviation study. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques* **48** 1148–1185.

- EFRON, B., HASTIE, T., JOHNSTONE, I. and TIBSHIRANI, R. (2004). Least angle regression. *Annals of Statistics* **32** 407–499.
- FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* **96** 1348–1360.
- FAN, J., XUE, L. and ZOU, H. (2014). Strong oracle optimality of folded concave penalized estimation. *Annals of Statistics* To appear.
- FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* **33** 1.
- HAN, F. and LIU, H. (2012). Transelliptical component analysis. In *Advances in Neural Information Processing Systems 25*.
- HAN, F. and LIU, H. (2013). Optimal rates of convergence of transelliptical component analysis. *arXiv preprint arXiv:1305.6916* .
- HASTIE, T., ROSSET, S., TIBSHIRANI, R. and ZHU, J. (2005). The entire regularization path for the support vector machine. *Journal of Machine Learning Research* **5** 1391.
- HUNTER, D. R. and LI, R. (2005). Variable selection using mm algorithms. *Annals of Statistics* **33** 1617.
- KIM, Y., CHOI, H. and OH, H.-S. (2008). Smoothly clipped absolute deviation on high dimensions. *Journal of the American Statistical Association* **103** 1665–1673.
- KOLTCHINSKII, V. (2009a). The Dantzig selector and sparsity oracle inequalities. *Bernoulli* **15** 799–828.
- KOLTCHINSKII, V. (2009b). Sparsity in penalized empirical risk minimization. *Ann. Inst. H. Poincaré Probab. Statist* **45** 7–57.
- LIU, H., HAN, F. and ZHANG, C.-H. (2012). Transelliptical graphical models. In *Advances in Neural Information Processing Systems 25*.
- LIU, W. and LUO, X. (2012). High-dimensional sparse precision matrix estimation via sparse column inverse operator. *arXiv preprint arXiv:1203.3896* .
- LOH, P.-L. and WAINWRIGHT, M. J. (2013). Regularized  $M$ -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *arXiv preprint arXiv:1305.2436* .
- MAIRAL, J. and YU, B. (2012). Complexity analysis of the lasso regularization path. *arXiv preprint arXiv:1205.0079* .
- MAZUMDER, R., FRIEDMAN, J. H. and HASTIE, T. (2011). Sparsenet: Coordinate descent with nonconvex penalties. *Journal of the American Statistical Association* **106**.
- NEGAHBAN, S. N., RAVIKUMAR, P., WAINWRIGHT, M. J. and YU, B. (2012). A unified framework for high-dimensional analysis of  $M$ -estimators with decomposable regularizers. *Statistical Science* **27** 538–557.
- NESTEROV, Y. (2004). *Introductory lectures on convex optimization: A basic course*, vol. 87. Springer.
- NESTEROV, Y. (2013). Gradient methods for minimizing composite functions. *Mathematical Programming* **140** 125–161.
- PARK, M. Y. and HASTIE, T. (2007).  $\ell_1$ -regularization path algorithm for generalized linear models. *Journal of the Royal Statistical Society, Series B* **69** 659–677.
- RASKUTTI, G., WAINWRIGHT, M. and YU, B. (2011). Minimax rates of estimation for high-dimensional linear regression over  $\ell_q$ -balls. *IEEE Transactions on Information Theory* **57** 6976–6994.
- RASKUTTI, G., WAINWRIGHT, M. J. and YU, B. (2010). Restricted eigenvalue properties for correlated gaussian designs. *Journal of Machine Learning Research* **99** 2241–2259.
- ROSSET, S. and ZHU, J. (2007). Piecewise linear regularized solution paths. *Annals of Statistics* 1012–1030.
- ROTHMAN, A. J., BICKEL, P. J., LEVINA, E. and ZHU, J. (2008). Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics* **2** 494–515.

- SHE, Y. (2009). Thresholding-based iterative selection procedures for model selection and shrinkage. *Electronic Journal of Statistics* **3** 384–415.
- SHE, Y. (2012). An iterative algorithm for fitting nonconvex penalized generalized linear models with grouped predictors. *Computational Statistics & Data Analysis* **56** 2976–2990.
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B* **58** 267–288.
- VAN DE GEER, S. (2000). *Empirical processes in M-estimation*, vol. 45. Cambridge press.
- VAN DE GEER, S. (2008). High-dimensional generalized linear models and the lasso. *Annals of Statistics* **36** 614–645.
- WAINWRIGHT, M. (2009). Sharp thresholds for high dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming. *IEEE Transactions on Information Theory* **55** 2183–2201.
- WANG, L., KIM, Y. and LI, R. (2013). Calibrating nonconvex penalized regression in ultra-high dimension. *Annals of Statistics* **41** 2505–2536.
- WANG, L., WU, Y. and LI, R. (2012). Quantile regression for analyzing heterogeneity in ultra-high dimension. *Journal of the American Statistical Association* **107** 214–222.
- WANG, Z., LIU, H. and ZHANG, T. (2014a). Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *arXiv preprint arXiv:1306.4960*.
- WANG, Z., LIU, H. and ZHANG, T. (2014b). Supplement to “Optimal computational and statistical rates of convergence for sparse nonconvex learning problems.” DOI: To Be Assigned.
- WRIGHT, S., NOWAK, R. and FIGUEIREDO, M. (2009). Sparse reconstruction by separable approximation. *IEEE Transactions on Signal Processing* **57** 2479–2493.
- XIAO, L. and ZHANG, T. (2013). A proximal-gradient homotopy method for the sparse least-squares problem. *SIAM Journal on Optimization* **23** 1062–1091.
- ZHANG, C.-H. (2010a). Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics* **38** 894–942.
- ZHANG, C.-H. and HUANG, J. (2008). The sparsity and bias of the lasso selection in high-dimensional linear regression. *Annals of Statistics* **36** 1567–1594.
- ZHANG, C.-H. and ZHANG, T. (2012). A general theory of concave regularization for high-dimensional sparse estimation problems. *Statistical Science* **27** 576–593.
- ZHANG, T. (2009). Some sharp performance bounds for least squares regression with  $\ell_1$  regularization. *Annals of Statistics* **37** 2109–2144.
- ZHANG, T. (2010b). Analysis of multi-stage convex relaxation for sparse regularization. *Journal of Machine Learning Research* **11** 1087–1107.
- ZHANG, T. (2013). Multi-stage convex relaxation for feature selection. *Bernoulli* **19** 2277–2293.
- ZHAO, P. and YU, B. (2007). Stagewise lasso. *Journal of Machine Learning Research* **8** 2701–2726.
- ZOU, H. and LI, R. (2008). One-step sparse estimates in nonconcave penalized likelihood models. *Annals of Statistics* **36** 1509.

DEPARTMENT OF OPERATIONS RESEARCH  
AND FINANCIAL ENGINEERING  
PRINCETON UNIVERSITY  
PRINCETON, NEW JERSEY 08544  
USA  
E-MAIL: [zhaoran@princeton.edu](mailto:zhaoran@princeton.edu)  
[hanliu@princeton.edu](mailto:hanliu@princeton.edu)

DEPARTMENT OF STATISTICS  
RUTGERS UNIVERSITY  
PISCATAWAY, NEW JERSEY 08854  
USA  
E-MAIL: [tzhang@stat.rutgers.edu](mailto:tzhang@stat.rutgers.edu)

**SUPPLEMENTARY MATERIAL FOR:  
OPTIMAL COMPUTATIONAL AND STATISTICAL RATES  
OF CONVERGENCE FOR SPARSE NONCONVEX  
LEARNING PROBLEMS**

BY ZHAORAN WANG<sup>\*</sup>, HAN LIU<sup>\*</sup> AND TONG ZHANG<sup>†</sup>

APPENDIX A: NONCONVEX PENALTY AND LOSS FUNCTIONS

We provide detailed descriptions of the nonconvex penalty and loss functions discussed in §2 of Wang et al. (2014a). Specifically, for the nonconvex penalties, i.e., SCAD and MCP, we provide their analytical forms in §A.1, and illustrate regularity condition (e) of Wang et al. (2014a) in §A.2. For the nonconvex loss, i.e., semiparametric elliptical design loss, we provide further details on elliptical distribution in §A.3, and define the two-step elliptical covariance matrix estimation procedure for semiparametric elliptical design regression in §A.4.

**A.1. Analytical Forms of SCAD and MCP.** The SCAD penalty in (2.1) of Wang et al. (2014a) can be written as

$$p_\lambda(\beta_j) = \lambda|\beta_j| \cdot \mathbb{1}(|\beta_j| \leq \lambda) - \frac{\beta_j^2 - 2a\lambda|\beta_j| + \lambda^2}{2(a-1)} \cdot \mathbb{1}(\lambda < |\beta_j| \leq a\lambda) \\ + \frac{(a+1)\lambda^2}{2} \cdot \mathbb{1}(|\beta_j| > a\lambda), \quad a > 2,$$

and the MCP penalty in (2.2) of Wang et al. (2014a) can be written as

$$p_\lambda(\beta_j) = \left( \lambda|\beta_j| - \frac{\beta_j^2}{2b} \right) \cdot \mathbb{1}(|\beta_j| \leq b\lambda) + \frac{b\lambda^2}{2} \cdot \mathbb{1}(|\beta_j| > b\lambda), \quad b > 0.$$

Correspondingly, the specific forms of the concave component  $q_\lambda(\beta_j)$  are

$$q_\lambda(\beta_j) = \begin{cases} \frac{2\lambda|\beta_j| - \beta_j^2 - \lambda^2}{2(a-1)} \cdot \mathbb{1}(\lambda < |\beta_j| \leq a\lambda) \\ \quad + \frac{(a+1)\lambda^2 - 2\lambda|\beta_j|}{2} \cdot \mathbb{1}(|\beta_j| > a\lambda), & \text{SCAD,} \\ -\frac{\beta_j^2}{2b} \cdot \mathbb{1}(|\beta_j| \leq b\lambda) + \left( \frac{b\lambda^2}{2} - \lambda|\beta_j| \right) \cdot \mathbb{1}(|\beta_j| > b\lambda), & \text{MCP.} \end{cases}$$

**A.2. Illustration of Regularity Condition (e) for SCAD and MCP.** We verify regularity condition (e) in Wang et al. (2014a) holds for SCAD



and MCP in Figure 4: For MCP, we illustrate with Figure 4(a); For SCAD, we illustrate with Figure 4(b) (for  $\lambda_2 \geq a\lambda_1$ ) and Figure 4(c).

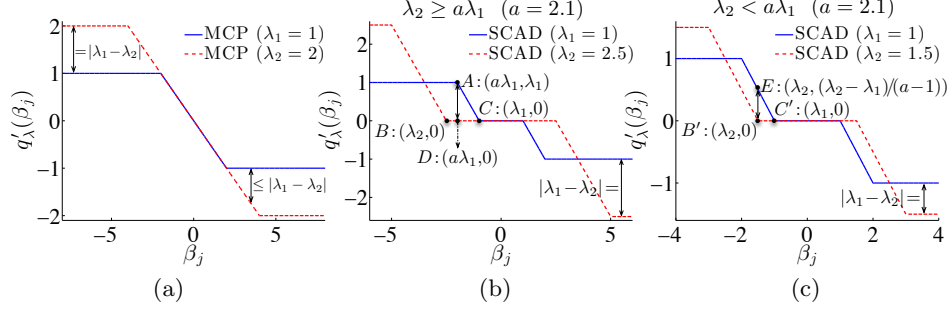


FIG 4. An illustration of regularity condition (e) in Wang et al. (2014a) for MCP and SCAD: (a) Plots of  $q'_{\lambda_1}(\beta_j)$  and  $q'_{\lambda_2}(\beta_j)$  for MCP with  $\lambda_1 = 1$ ,  $\lambda_2 = 2$  and  $b = 2$ ; (b) Plots of  $q'_{\lambda_1}(\beta_j)$  and  $q'_{\lambda_2}(\beta_j)$  for SCAD with  $\lambda_1 = 1$ ,  $\lambda_2 = 2.5$  and  $a = 2.1$ ; (c) Plots of  $q'_{\lambda_1}(\beta_j)$  and  $q'_{\lambda_2}(\beta_j)$  for SCAD with  $\lambda_1 = 1$ ,  $\lambda_2 = 1.5$  and  $a = 2.1$ . Subfigure (a) shows that regularity condition (e) holds for MCP. For SCAD, we consider two cases:  $\lambda_2 \geq a\lambda_1$ , as illustrated in (b);  $\lambda_2 < a\lambda_1$  as illustrated in (c). In the first case,  $|AD| = \lambda_1 \leq (a-1)\lambda_1 \leq |\lambda_1 - \lambda_2|$  since  $a > 2$  and  $\lambda_2 \geq a\lambda_1$ . In the second case,  $|B'E| = (\lambda_2 - \lambda_1)/(a-1) \leq |\lambda_1 - \lambda_2|$ , because the slope of  $EC'$  is  $(-1)/(a-1)$  with  $a > 2$ .

**A.3. Elliptical Distribution.** Before we present the definition of elliptical distribution, we first introduce some notation: If random vectors  $\mathbf{Z}_1$  and  $\mathbf{Z}_2$  have the same distribution, we denote by  $\mathbf{Z}_1 \stackrel{d}{=} \mathbf{Z}_2$ ; The  $d$ -dimensional  $\ell_2$  unit sphere  $\{\mathbf{v} : \|\mathbf{v}\|_2 = 1, \mathbf{v} \in \mathbb{R}^d\}$  is denoted by  $\mathbb{S}^{d-1}$ ; For a matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$ , we define  $\text{diag}(\mathbf{M})$  to be a diagonal matrix with diagonal entries  $[\text{diag}(\mathbf{M})]_{jj} = \mathbf{M}_{jj}$  ( $j = 1, \dots, d$ ).

**Definition A.1** (Elliptical distribution). For  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)^T \in \mathbb{R}^d$  and  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  with  $\text{rank}(\boldsymbol{\Sigma}) = r \leq d$ , a random vector  $\mathbf{W} = (W_1, \dots, W_d)^T$  follows an elliptical distribution, denoted by  $\text{EC}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \Xi)$ , if and only if

$$\mathbf{W} \stackrel{d}{=} \boldsymbol{\mu} + \Xi \mathbf{A} \mathbf{U}.$$

Here  $\mathbf{U}$  is a random vector uniformly distributed on the unit sphere  $\mathbb{S}^{r-1}$ ;  $\Xi \geq 0$  is a scalar random variable independent of  $\mathbf{U}$ ;  $\mathbf{A} \in \mathbb{R}^{d \times r}$  is a deterministic matrix such that  $\mathbf{A} \mathbf{A}^T = \boldsymbol{\Sigma}$ . We call  $\boldsymbol{\Sigma}$  the scatter matrix. The generalized correlation matrix is defined as  $\boldsymbol{\Sigma}^0 = \text{diag}(\boldsymbol{\Sigma})^{-1/2} \cdot \boldsymbol{\Sigma} \cdot \text{diag}(\boldsymbol{\Sigma})^{-1/2}$ . When  $\mathbb{E}(\Xi^2)$  exists,  $\boldsymbol{\Sigma}^0$  is the correlation matrix of  $\mathbf{W}$ .

**Remark A.2.** Note that simultaneously scaling  $\Xi$  and  $\mathbf{U}$  (e.g.,  $\Xi \rightarrow \Xi/C$  and  $\mathbf{U} \rightarrow \mathbf{U}/C$ , where  $C$  is a constant) leads to the same elliptical distribution. To make this model identifiable, we assume  $\mu_j = \mathbb{E}(W_j)$  and  $\boldsymbol{\Sigma}_{jj} = \text{Var}(W_j)$ .

**Remark A.3.** The elliptical distribution family includes a variety of possibly heavy-tail distributions: multivariate Gaussian, multivariate Cauchy, Student's t, logistic, Kotz, symmetric Pearson type-II and type-VII distributions.

**A.4. Rank-based Covariance Matrix Estimation for Semiparametric Elliptical Design Regression.** We provide details of the two-step procedure for estimating the covariance matrix of the elliptically distributed random vector  $\mathbf{Z} = (Y, \mathbf{X}^T)^T \in \mathbb{R}^{d+1}$  discussed in §2.2 of Wang et al. (2014a):

### Elliptical Covariance Matrix Estimation

**S1.** First, we define a rank-based estimator  $\hat{\mathbf{R}}_{\mathbf{Z}}$  of the generalized correlation matrix  $\Sigma_{\mathbf{Z}}^0$  using the Kendall's tau statistic. Let  $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^{d+1}$  with  $\mathbf{z}_i = (z_{i1}, \dots, z_{i(d+1)})^T$  be  $n$  independent observations of  $\mathbf{Z}$ . The Kendall's tau correlation coefficient is defined as

$$\hat{\tau}_{jk}(\mathbf{z}_1, \dots, \mathbf{z}_n) = \begin{cases} \sum_{1 \leq i < i' < n} \frac{2 \operatorname{sign}(z_{ij} - z_{i'j}) \operatorname{sign}(z_{ik} - z_{i'k})}{n(n-1)}, & \text{for } j \neq k, \\ 1, & \text{for } j = k. \end{cases}$$

We define the Kendall's tau correlation matrix estimator as

$$(A.1) \quad \hat{\mathbf{R}}_{\mathbf{Z}} = \left[ (\hat{\mathbf{R}}_{\mathbf{Z}})_{jk} \right] = \left[ \sin \left( \frac{\pi}{2} \hat{\tau}_{jk}(\mathbf{z}_1, \dots, \mathbf{z}_n) \right) \right].$$

Liu et al. (2012); Han and Liu (2012, 2013) showed that  $\hat{\mathbf{R}}_{\mathbf{Z}}$  is a robust estimator of the population generalized correlation matrix  $\Sigma_{\mathbf{Z}}^0$ , and is invariant to different distributions of the generating variable  $\Xi$  within the whole elliptical family.

**S2.** Second, we construct a covariance matrix estimator

$$(A.2) \quad \hat{\mathbf{K}}_{\mathbf{Z}} = \left[ (\hat{\mathbf{K}}_{\mathbf{Z}})_{jk} \right] = \left[ (\hat{\mathbf{R}}_{\mathbf{Z}})_{jk} \cdot \hat{\sigma}_j \hat{\sigma}_k \right],$$

where  $\hat{\sigma}_1, \dots, \hat{\sigma}_{d+1}$  are the estimators of the standard deviations of  $Z_1, \dots, Z_{d+1}$ . We calculate  $\hat{\sigma}_1, \dots, \hat{\sigma}_{d+1}$  using the Catoni's  $M$ -estimator (Catoni, 2012) described in §E. The main advantage of the Catoni's estimator is that, for a fixed level of confidence, it achieves the same deviation behavior as a Gaussian random variable under a weak moment condition.

## APPENDIX B: PROXIMAL-GRADIENT METHOD FOR NONCONVEX PROBLEMS

We provide details of the proximal-gradient method tailored to nonconvex problems. In particular, in §B.1 we provide the optimization update schemes for the specific nonconvex problems discussed in §2 of Wang et al. (2014a). In §B.2 we provide the detailed derivation of the closed-form expression of update scheme (3.9) in Wang et al. (2014a).

**B.1. Optimization Update Schemes for Specific Nonconvex Problems.** To obtain the specific optimization update schemes of the proximal-gradient method for the nonconvex problems discussed in §2 of Wang et al. (2014a), we only need to plug the following specific definitions of  $\nabla \mathcal{L}(\beta)$  and  $\nabla \mathcal{Q}_{\lambda_t}(\beta)$  into (3.11) of Wang et al. (2014a):

- For the (nonconvex) loss functions discussed in §2 of Wang et al. (2014a),

$$\nabla \mathcal{L}(\beta) = \begin{cases} \frac{1}{n} \mathbf{X}^T (\mathbf{X}\beta - \mathbf{y}), & \text{least squares loss,} \\ \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \left( \frac{\exp(\mathbf{x}_i^T \beta)}{1 + \exp(\mathbf{x}_i^T \beta)} - y_i \right), & \text{logistic loss,} \\ \widehat{\mathbf{K}}_{\mathbf{X}} \beta - \widehat{\mathbf{K}}_{\mathbf{X},Y}, & \text{semiparametric elliptical design loss,} \end{cases}$$

where  $\widehat{\mathbf{K}}_{\mathbf{X}} \in \mathbb{R}^{d \times d}$  and  $\widehat{\mathbf{K}}_{\mathbf{X},Y} \in \mathbb{R}^{d \times 1}$  are the submatrices of  $\widehat{\mathbf{K}}_{\mathbf{Z}}$ , which is the semiparametric elliptical covariance matrix estimator defined in (A.2). More specifically,

$$(B.1) \quad \widehat{\mathbf{K}}_{\mathbf{Z}} = \begin{pmatrix} \widehat{\mathbf{K}}_Y & \widehat{\mathbf{K}}_{\mathbf{X},Y}^T \\ \widehat{\mathbf{K}}_{\mathbf{X},Y} & \widehat{\mathbf{K}}_{\mathbf{X}} \end{pmatrix}.$$

- For the nonconvex penalty functions discussed in §2 of Wang et al. (2014a),

$$(\nabla \mathcal{Q}_{\lambda_t}(\beta))_j = \begin{cases} \frac{\lambda_t \text{sign}(\beta_j) - \beta_j}{a-1} \cdot \mathbf{1}(\lambda_t < |\beta_j| \leq a\lambda_t) \\ \quad - \lambda_t \text{sign}(\beta_j) \cdot \mathbf{1}(|\beta_j| > a\lambda_t), & \text{SCAD,} \\ -\frac{\beta_j}{b} \lambda_t \text{sign}(\beta_j) \cdot \mathbf{1}(|\beta_j| \leq b\lambda_t) \\ \quad - \lambda_t \text{sign}(\beta_j) \cdot \mathbf{1}(|\beta_j| > b\lambda_t), & \text{MCP,} \end{cases}$$

where  $a > 2$ ,  $b > 0$ .

**B.2. Derivation of Optimization Update Schemes.** For notational simplicity, we denote  $L_t^k$  by  $L$ ,  $\beta_t^{k-1}$  by  $\beta'$ , and  $\lambda_t$  by  $\lambda$  in the rest of this

section.

**Derivation of (3.10) of Wang et al. (2014a):** If  $\Omega = \mathbb{R}^d$ , then we have

$$\begin{aligned}
 \mathcal{T}_{L,\lambda}(\beta'; +\infty) &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \psi_{L,\lambda}(\beta; \beta') \right\} \\
 &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \tilde{\mathcal{L}}_\lambda(\beta') + \nabla \tilde{\mathcal{L}}_\lambda(\beta')^T (\beta - \beta') + \frac{L}{2} \|\beta - \beta'\|_2^2 + \lambda \|\beta\|_1 \right\} \\
 (B.2) \quad &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \frac{1}{2} \left\| \beta - \underbrace{\left( \beta' - \frac{1}{L} \nabla \tilde{\mathcal{L}}_\lambda(\beta') \right)}_{\bar{\beta}} \right\|_2^2 + \frac{\lambda}{L} \|\beta\|_1 \right\}.
 \end{aligned}$$

It is known that the minimizer of (B.2) can be obtained by soft-thresholding  $\bar{\beta}$  with the threshold of value  $\lambda/L$ , i.e.,

$$(B.3) \quad (\mathcal{T}_{L,\lambda}(\beta'; +\infty))_j = \begin{cases} 0 & \text{if } |\bar{\beta}_j| \leq \lambda/L, \\ \operatorname{sign}(\bar{\beta}_j)(|\bar{\beta}_j| - \lambda/L) & \text{if } |\bar{\beta}_j| > \lambda/L. \end{cases}$$

Therefore we obtain the first update scheme (3.10) in Wang et al. (2014a) for  $\Omega = \mathbb{R}^d$ .

**Derivation of (3.12) of Wang et al. (2014a):** If  $\Omega = B_2(R) = \{\beta : \|\beta\|_2^2 \leq R^2\}$ , by Lagrangian duality we can transform the original optimization problem with constraint into an unconstraint optimization problem. Hence, there exists a Lagrangian multiplier  $\tau \geq 0$  such that

$$\mathcal{T}_{L,\lambda}(\beta'; R) = \operatorname{argmin}_{\beta \in B_2(R)} \left\{ \psi_{L,\lambda}(\beta; \beta') \right\} = \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \psi_{L,\lambda}(\beta; \beta') + \frac{\tau}{2} \|\beta\|_2^2 \right\}.$$

Consequently, based on (B.2) we have

$$\begin{aligned}
 \mathcal{T}_{L,\lambda}(\beta'; R) &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \tilde{\mathcal{L}}_\lambda(\beta') + \nabla \tilde{\mathcal{L}}_\lambda(\beta')^T (\beta - \beta') \right. \\
 &\quad \left. + \frac{L}{2} \|\beta - \beta'\|_2^2 + \lambda \|\beta\|_1 + \frac{\tau}{2} \|\beta\|_2^2 \right\} \\
 &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \frac{L+\tau}{2} \|\beta\|_2^2 - \left( L \cdot \beta' - \nabla \tilde{\mathcal{L}}_\lambda(\beta') \right)^T \beta + \lambda \|\beta\|_1 \right\} \\
 (B.4) \quad &= \operatorname{argmin}_{\beta \in \mathbb{R}^d} \left\{ \frac{1}{2} \left\| \beta - \underbrace{\left( \frac{L}{L+\tau} \beta' - \frac{1}{L+\tau} \nabla \tilde{\mathcal{L}}_\lambda(\beta') \right)}_{\frac{L}{L+\tau} \bar{\beta}} \right\|_2^2 + \frac{\lambda}{L+\tau} \|\beta\|_1 \right\},
 \end{aligned}$$

where  $\bar{\beta} = \beta' - \nabla \tilde{\mathcal{L}}_\lambda(\beta')/L$ . The minimizer of (B.4) can also be obtained by

soft-thresholding, i.e.,

$$(B.5) \quad (\mathcal{T}_{L,\lambda}(\beta'; R))_j = \begin{cases} 0 & \text{if } \frac{L}{L+\tau}|\bar{\beta}_j| \leq \frac{\lambda}{L+\tau}, \\ \text{sign}\left(\frac{L}{L+\tau}\bar{\beta}_j\right)\left(\frac{L}{L+\tau}|\bar{\beta}_j| - \frac{\lambda}{L+\tau}\right) & \text{if } \frac{L}{L+\tau}|\bar{\beta}_j| > \frac{\lambda}{L+\tau}. \end{cases}$$

Comparing (B.5) with (B.3), we have

$$(B.6) \quad \mathcal{T}_{L,\lambda}(\beta'; R) = \frac{L}{L+\tau} \mathcal{T}_{L,\lambda}(\beta'; +\infty).$$

Thus, we can obtain the constraint solution  $\mathcal{T}_{L,\lambda}(\beta'; R)$  by first calculating the unconstraint solution  $\mathcal{T}_{L,\lambda}(\beta'; +\infty)$ , and then rescaling it by a factor of  $L/(L+\tau)$ . Note that here the Lagrangian multiplier  $\tau$  is unknown. We discuss the following two cases:

- If the constraint  $\beta \in B_2(R)$  is inactive, then we have  $\tau = 0$  by complementary slackness, which implies  $\mathcal{T}_{L,\lambda}(\beta'; R) = \mathcal{T}_{L,\lambda}(\beta'; +\infty)$ . Since the constraint is inactive, we have  $\|\mathcal{T}_{L,\lambda}(\beta'; R)\|_2 = \|\mathcal{T}_{L,\lambda}(\beta'; +\infty)\|_2 < R$ .
- If the constraint  $\beta \in B_2(R)$  is active, then we have  $\tau \geq 0$  by complementary slackness. In this case, the minimizer  $\mathcal{T}_{L,\lambda}(\beta'; R)$  lies on the boundary of  $B_2(R)$ . By (B.6) we have

$$\|\mathcal{T}_{L,\lambda}(\beta'; +\infty)\|_2 = \frac{L+\tau}{L} \|\mathcal{T}_{L,\lambda}(\beta'; R)\|_2 = \frac{L+\tau}{L} R \geq R.$$

To obtain  $\mathcal{T}_{L,\lambda}(\beta'; R)$ , we project  $\mathcal{T}_{L,\lambda}(\beta'; +\infty)$  onto  $B_2(R)$ , which can be achieved by setting  $\mathcal{T}_{L,\lambda}(\beta'; R) = R \cdot \mathcal{T}_{L,\lambda}(\beta'; +\infty) / \|\mathcal{T}_{L,\lambda}(\beta'; +\infty)\|_2$ .

Therefore we obtain the second update scheme (3.12) in Wang et al. (2014a) for  $\Omega = B_2(R)$ .

## APPENDIX C: JUSTIFICATION OF ASSUMPTIONS

In §C.1 and §C.2, we prove that Assumption 4.1 and Assumption 4.4 in Wang et al. (2014a) hold with high probability respectively.

**C.1. Justification of Assumption 4.1 in Wang et al. (2014a).** Recall that Assumption 4.1 states that  $\|\nabla \mathcal{L}(\beta^*)\|_\infty$  can be upper bounded by  $\lambda_{\text{tgt}}/8$ . First we provide two lemmas on the upper bound of  $\|\nabla \mathcal{L}(\beta^*)\|_\infty$ .

**Lemma C.1.** For least squares regression with sub-Gaussian noise and logistic regression, we assume that the columns of  $\mathbf{X}$  are normalized in such

a way that  $\max_{j \in \{1, \dots, d\}} \{\|\mathbf{X}_j\|_2\} \leq \sqrt{n}$ . Then we have

$$(C.1) \quad \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty \leq C \sqrt{\frac{\log d}{n}}$$

with probability at least  $1 - d^{-1}$ , where  $C$  is a constant.

PROOF. See Candés and Tao (2007); Zhang and Huang (2008); Zhang (2009); Bickel et al. (2009); Koltchinskii (2009a); Negahban et al. (2012); Wainwright (2009) for a detailed proof.  $\square$

**Lemma C.2.** For semiparametric elliptical design regression, we have

$$(C.2) \quad \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty \leq C \|\boldsymbol{\beta}^*\|_1 \sqrt{\frac{\log d}{n}}$$

with probability at least  $1 - (d+1)^{-5/2} - 2(d+1)^{-3}$ , where  $C$  is a constant.

PROOF. See §E.3 for a detailed proof.  $\square$

Recall in Assumption 4.1 of Wang et al. (2014a) we set  $\lambda_{\text{tgt}} = C \sqrt{\log d/n}$  for least squares and logistic loss, and  $\lambda_{\text{tgt}} = C' \|\boldsymbol{\beta}^*\|_1 \sqrt{\log d/n}$  for semiparametric elliptical design loss. Thus, according to Lemma C.1 and Lemma C.2,  $\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty \leq \lambda_{\text{tgt}}/8$  holds with high probability, i.e., Assumption 4.1 holds with high probability.

**C.2. Justification of Assumption 4.4 in Wang et al. (2014a).** In this section we show that, for semiparametric elliptical design loss and logistic loss, Assumption 4.4 holds with high probability.

**Semiparametric Elliptical Design Loss:** First we provide the following lemma on the largest and smallest sparse eigenvalues of the Hessian matrix  $\nabla^2 \mathcal{L}(\boldsymbol{\beta})$  of semiparametric elliptical design loss.

Let  $n$  be the sample size,  $d$  be the dimension of  $\boldsymbol{\beta}$ , and  $\mathbf{Z} \in \mathbb{R}^{d+1}$  be the elliptically distributed random vector in §2.2 of Wang et al. (2014a). The corresponding covariance matrix estimator  $\widehat{\mathbf{K}}_{\mathbf{Z}} \in \mathbb{R}^{(d+1) \times (d+1)}$  is defined in (A.2), while its submatrix  $\widehat{\mathbf{K}}_{\mathbf{X}} \in \mathbb{R}^{d \times d}$  is defined in (B.1). Hence, the Hessian matrix of semiparametric elliptical design loss is  $\nabla^2 \mathcal{L}(\boldsymbol{\beta}) = \widehat{\mathbf{K}}_{\mathbf{X}}$ . Let  $s$  the sparsity level.

**Lemma C.3.** Under suitable conditions (see Han and Liu (2013) for details), for a sufficiently large  $n$ , there exists an  $s$  such that  $\rho_-(\nabla^2 \mathcal{L}, s) > 0$  and  $\rho_+(\nabla^2 \mathcal{L}, s) < +\infty$  with probability at least  $1 - 4d^{-1} - 6d^{-2}$ . Here  $\rho_+(\nabla^2 \mathcal{L}, s)$  and  $\rho_-(\nabla^2 \mathcal{L}, s)$  are defined in Definition 4.2 of Wang et al. (2014a).

PROOF. See §E.2 for a detailed proof.  $\square$

Equipped with Lemma C.3, we can justify Assumption 4.4. Recall  $s^* = \|\beta^*\|_0$ , where  $\beta^*$  is the true parameter vector. Suppose that Lemma C.3 holds with  $s = Cs^*$ ,  $\rho_+(\nabla^2 \mathcal{L}, s) = C'$  and  $\rho_-(\nabla^2 \mathcal{L}, s) = C''$ , where  $C$  satisfies

$$(C.3) \quad C \geq 2 \left( 144 \cdot \left( \frac{2C'}{C''} \right)^2 + 250 \cdot \left( \frac{2C'}{C''} \right) \right) + 1.$$

Meanwhile, let the concavity parameters of the nonconvex penalty be  $\zeta_+ = 0$  and  $\zeta_- = C''/2$ . In the sequel we verify that, there exists an integer  $\tilde{s} = (C - 1)/2 \cdot s^*$ , where  $C$  satisfies (C.3), that satisfies Assumption 4.4. Note that the condition number  $\kappa$  defined in (4.5) is

$$\begin{aligned} \kappa &= \frac{\rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_+}{\rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_-} = \frac{\rho_+(\nabla^2 \mathcal{L}, Cs^*) - \zeta_+}{\rho_-(\nabla^2 \mathcal{L}, Cs^*) - \zeta_-} = \frac{\rho_+(\nabla^2 \mathcal{L}, s) - \zeta_+}{\rho_-(\nabla^2 \mathcal{L}, s) - \zeta_-} \\ &= \frac{C'}{C'' - C''/2} = \frac{2C'}{C''}. \end{aligned}$$

Since  $\tilde{s} = (C - 1)/2 \cdot s^*$  where  $C$  satisfies (C.3), we have

$$\tilde{s} \geq \left( 144 \cdot \left( \frac{2C'}{C''} \right)^2 + 250 \cdot \left( \frac{2C'}{C''} \right) \right) \cdot s^* = (144\kappa^2 + 250\kappa) \cdot s^*.$$

Thus we find an  $\tilde{s}$  that satisfies the requirements in Assumption 4.4. Therefore, Assumption 4.4 in Wang et al. (2014a) holds with probability at least  $1 - 4d^{-1} - 6d^{-2}$ .

**Logistic Loss:** Remind that, the Hessian matrix of logistic loss is defined in (4.2) of Wang et al. (2014a), while its sparse eigenvalues  $\rho_-(\nabla^2 \mathcal{L}, s, R)$  and  $\rho_+(\nabla^2 \mathcal{L}, s, R)$  are defined in Definition 4.3 of Wang et al. (2014a), where  $R \in (0, +\infty)$  is an absolute constant.

In the sequel, we show that Assumption 4.4 is a weaker assumption than the assumption of restricted strong convexity and smoothness imposed by Loh and Wainwright (2013). Since they proved that their assumption holds with high probability, Assumption 4.4 also holds with high probability.

In fact, it suffices to show if the assumption of restricted strong convexity and smoothness holds, then for a sufficiently large  $n$ , there exists an  $s$  such that  $\rho_-(\nabla^2 \mathcal{L}, s, R) > 0$  and  $\rho_+(\nabla^2 \mathcal{L}, s, R) < +\infty$ . With this claim, we can justify Assumption 4.4 following the same argument as after Lemma C.3. Now we prove the previous claim.

Loh and Wainwright (2013) imposed the following assumption: For  $\beta, \beta' \in$

$\mathbb{R}^d$  such that  $\|\beta\|_2 \leq R$  and  $\|\beta'\|_2 \leq R$ ,  $\mathcal{L}(\beta)$  satisfies

$$(C.4) \quad \mathcal{L}(\beta') - \mathcal{L}(\beta) - \nabla \mathcal{L}(\beta)^T (\beta' - \beta) \leq C \|\beta - \beta'\|_2^2 + C' \cdot \frac{\log d}{n} \|\beta - \beta'\|_1^2,$$

$$(C.5) \quad \mathcal{L}(\beta') - \mathcal{L}(\beta) - \nabla \mathcal{L}(\beta)^T (\beta' - \beta) \geq C'' \|\beta - \beta'\|_2^2 - C''' \cdot \frac{\log d}{n} \|\beta - \beta'\|_1^2.$$

Here all the constants are positive. See equations (28) and (29) of their paper for details.

By Taylor's theorem and the mean value theorem, we have

$$\mathcal{L}(\beta') = \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^T (\beta' - \beta) + \frac{1}{2} (\beta' - \beta)^T \nabla^2 \mathcal{L}(\gamma\beta' + (1-\gamma)\beta) (\beta' - \beta)$$

for some  $\gamma \in [0, 1]$ . Plugging this into the left-hand sides of (C.4) and (C.5), we have

$$(C.6)$$

$$\frac{1}{2} (\beta' - \beta)^T \nabla^2 \mathcal{L}(\gamma\beta' + (1-\gamma)\beta) (\beta' - \beta) \leq C \|\beta' - \beta\|_2^2 + C' \cdot \frac{\log d}{n} \|\beta' - \beta\|_1^2,$$

$$(C.7)$$

$$\frac{1}{2} (\beta' - \beta)^T \nabla^2 \mathcal{L}(\gamma\beta' + (1-\gamma)\beta) (\beta' - \beta) \geq C'' \|\beta' - \beta\|_2^2 - C''' \cdot \frac{\log d}{n} \|\beta' - \beta\|_1^2.$$

Suppose  $\beta$  and  $\beta'$  satisfy  $\|\beta' - \beta\|_0 \leq s$ , which implies  $\|\beta' - \beta\|_1 \leq \sqrt{s} \cdot \|\beta' - \beta\|_2$ . Plugging this upper bound of  $\|\beta' - \beta\|_1$  into the right-hand sides of (C.6) and (C.7), we have

$$(C.8)$$

$$\frac{1}{2} (\beta' - \beta)^T \nabla^2 \mathcal{L}(\gamma\beta' + (1-\gamma)\beta) (\beta' - \beta) \leq \left( C + C' \cdot \frac{s \log d}{n} \right) \cdot \|\beta' - \beta\|_2^2,$$

$$(C.9)$$

$$\frac{1}{2} (\beta' - \beta)^T \nabla^2 \mathcal{L}(\gamma\beta' + (1-\gamma)\beta) (\beta' - \beta) \geq \left( C'' - C''' \cdot \frac{s \log d}{n} \right) \cdot \|\beta' - \beta\|_2^2.$$

In (C.8) and (C.9), taking  $n \geq \max\{2C'''/C'', 2C'/C\} \cdot s \log d/n$ , and dividing  $\|\beta' - \beta\|_2^2$  on both sides, we obtain

$$(C.10) \quad \frac{C''}{2} \leq \frac{1}{2} \cdot \frac{(\beta' - \beta)^T}{\|\beta' - \beta\|_2} \cdot \nabla^2 \mathcal{L}(\gamma\beta' + (1-\gamma)\beta) \cdot \frac{(\beta' - \beta)}{\|\beta' - \beta\|_2} \leq \frac{3C}{2}.$$

Let  $\mathbf{v} = (\beta' - \beta)/\|\beta' - \beta\|_2$ . Obviously,  $\mathbf{v}$  is an arbitrary vector that satisfies  $\|\mathbf{v}\|_2 = 1$  and  $\|\mathbf{v}\|_0 \leq s$ . Taking  $\beta' \rightarrow \beta$ , we have  $C'' \leq \mathbf{v}^T \nabla^2 \mathcal{L}(\beta) \mathbf{v} \leq 3C$  for any  $\beta \leq R$  and any  $\mathbf{v}$  such that  $\|\mathbf{v}\|_2 = 1$  and  $\|\mathbf{v}\|_0 \leq s$ . By Definition 4.3, we have  $\rho_-(\nabla^2 \mathcal{L}, s, R) \geq C' > 0$  and  $\rho_+(\nabla^2 \mathcal{L}, s, R) \leq 3C < +\infty$ .

Therefore, if their assumption holds, then for a sufficiently large  $n$ , there exists some  $s$  such that  $\rho_-(\nabla^2 \mathcal{L}, s, R) > 0$  and  $\rho_+(\nabla^2 \mathcal{L}, s, R) < +\infty$ . Following the same argument as after Lemma C.3, we can show that Assumption



4.4 also holds true. That is to say, Assumption 4.4 of Wang et al. (2014a) is a weaker assumption.

Because their assumption holds with high probability for generalized linear models (not including Poisson model; see Proposition 1 of their paper) with sub-Gaussian design, our Assumption 4.4 holds with high probability in the same setting, including logistic loss as a special case in our paper.

#### APPENDIX D: PROOF OF THEORETICAL RESULTS

To analyze the computational properties of our approximate regularization path following method, we first provide several useful lemmas on the proximal-gradient method that is used within each stage of the path following method.

**D.1. Preliminary Results about Proximal-Gradient Method.** Recall that the objective function can be formulated as  $\phi_{\lambda_t}(\beta) = \tilde{\mathcal{L}}_{\lambda_t}(\beta) + \lambda_t \|\beta\|_1$  where  $\tilde{\mathcal{L}}_{\lambda_t}(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_{\lambda_t}(\beta)$ , while  $\psi_{L_t^k, \lambda_t}(\beta; \beta_t^{k-1})$  is the local quadratic approximation of  $\phi_{\lambda_t}(\beta)$  at  $\beta_t^{k-1}$ , as defined in (3.7) of Wang et al. (2014a). The following lemma, which is adapted from Nesterov (2013), characterizes the decrement of the objective function.

**Lemma D.1.** Under Assumption 4.4, we assume  $\|(\beta_t^{k-1})_{\bar{S}^*}\|_0 \leq \tilde{s}$ , where  $\tilde{s}$  is the positive integer specified in Assumption 4.4. For any  $L_t^k > 0$  and fixed  $\lambda_t \in [\lambda_{\text{tgt}}, \lambda_0]$ , we have

$$\phi_{\lambda_t}(\beta_t^k) \leq \phi_{\lambda_t}(\beta_t^{k-1}) - \frac{L_t^k}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2.$$

Recall that, as defined in (3.15) of Wang et al. (2014a),  $\omega_{\lambda}(\beta)$  characterizes the suboptimality of approximate solutions. The next lemma, which is also adapted from Nesterov (2013), provides an upper bound of  $\omega_{\lambda_t}(\beta_t^k)$  using  $\|\beta_t^k - \beta_t^{k-1}\|_2$ .

**Lemma D.2.** Under the assumptions of Lemma D.1, then we have

$$\omega_{\lambda_t}(\beta_t^k) \leq (L_t^k + \rho_+ - \zeta_-) \|\beta_t^k - \beta_t^{k-1}\|_2,$$

where  $\rho_+ = \rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})$  is the sparse eigenvalue specified in Assumption 4.4; As defined in regularity condition (a),  $\zeta_+ > 0$  is the concavity parameter of the nonconvex penalty, which satisfies (4.6) in Wang et al. (2014a).

#### D.2. Proof of Lemma 5.1 in Wang et al. (2014a).

PROOF. Recall that  $\mathcal{Q}_{\lambda}(\beta)$  is the concave component of the nonconvex penalty  $\mathcal{P}_{\lambda}(\beta)$ , which implies  $-\mathcal{Q}_{\lambda}(\beta)$  is convex. Meanwhile, recall that

$\mathcal{Q}_\lambda(\beta) = \sum_{j=1}^d q_\lambda(\beta_j)$ , where  $q_\lambda(\beta_j)$  satisfies regularity condition (a) in Wang et al. (2014a). Hence we have

$$-\zeta_-(\beta'_j - \beta_j)^2 \leq (q'_\lambda(\beta'_j) - q'_\lambda(\beta_j))(\beta'_j - \beta_j) \leq -\zeta_+(\beta'_j - \beta_j)^2,$$

which implies the convex function  $-\mathcal{Q}_\lambda(\beta)$  satisfies

$$(D.1) \quad \left( \nabla(-\mathcal{Q}_\lambda(\beta')) - \nabla(-\mathcal{Q}_\lambda(\beta)) \right)^T (\beta' - \beta) \leq \zeta_- \|\beta' - \beta\|_2^2,$$

$$(D.2) \quad \left( \nabla(-\mathcal{Q}_\lambda(\beta')) - \nabla(-\mathcal{Q}_\lambda(\beta)) \right)^T (\beta' - \beta) \geq \zeta_+ \|\beta' - \beta\|_2^2.$$

According to Nesterov (2004, Theorem 2.1.5 & Theorem 2.1.9), (D.1) and (D.2) are equivalent definitions of strong smoothness and strong convexity respectively. In other words,  $-\mathcal{Q}_\lambda(\beta)$  satisfies

$$(D.3) \quad -\mathcal{Q}_\lambda(\beta') \leq -\mathcal{Q}_\lambda(\beta) - \nabla \mathcal{Q}(\beta)^T (\beta' - \beta) + \frac{\zeta_-}{2} \|\beta' - \beta\|_2^2,$$

$$(D.4) \quad -\mathcal{Q}_\lambda(\beta') \geq -\mathcal{Q}_\lambda(\beta) - \nabla \mathcal{Q}(\beta)^T (\beta' - \beta) + \frac{\zeta_+}{2} \|\beta' - \beta\|_2^2.$$

For loss function  $\mathcal{L}(\beta)$ , by Taylor's theorem and the mean value theorem, we have

$$(D.5) \quad \begin{aligned} \mathcal{L}(\beta') &= \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^T (\beta' - \beta) \\ &\quad + \frac{1}{2} (\beta' - \beta)^T \nabla^2 \mathcal{L}(\gamma\beta + (1-\gamma)\beta') (\beta' - \beta), \end{aligned}$$

where  $\gamma \in [0, 1]$ . Note that we assume  $\|(\beta' - \beta)_{\bar{s}^*}\|_0 \leq 2\tilde{s}$ , which implies  $\|\beta' - \beta\|_0 \leq s^* + 2\tilde{s}$ . For logistic loss, we assume  $\|\beta\|_2 \leq R$  and  $\|\beta'\|_2 \leq R$ , which implies  $\|\gamma\beta + (1-\gamma)\beta'\|_2 \leq R$  by the convexity of  $\ell_2$  norm. Hence, by Definition 4.2 and Definition 4.3 of Wang et al. (2014a), we have

$$\frac{(\beta' - \beta)^T}{\|\beta' - \beta\|_2} \nabla^2 \mathcal{L}(\gamma\beta + (1-\gamma)\beta') \frac{(\beta' - \beta)}{\|\beta' - \beta\|_2} \in [\rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}), \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})].$$

Plugging this into the right-hand side of (D.5), we have

$$(D.6) \quad \mathcal{L}(\beta') \geq \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^T (\beta' - \beta) + \frac{\rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})}{2} \|\beta' - \beta\|_2^2,$$

$$(D.7) \quad \mathcal{L}(\beta') \leq \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^T (\beta' - \beta) + \frac{\rho_+(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})}{2} \|\beta' - \beta\|_2^2.$$

Recall that  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$ . Subtracting (D.3) from (D.6), and (D.4)

from (D.7), we obtain

$$\begin{aligned}\tilde{\mathcal{L}}_\lambda(\beta') &\geq \tilde{\mathcal{L}}_\lambda(\beta) + \nabla \tilde{\mathcal{L}}_\lambda(\beta)^T(\beta' - \beta) + \frac{\rho_- (\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_-}{2} \|\beta' - \beta\|_2^2 \\ \tilde{\mathcal{L}}_\lambda(\beta') &\leq \tilde{\mathcal{L}}_\lambda(\beta) + \nabla \tilde{\mathcal{L}}_\lambda(\beta)^T(\beta' - \beta) + \frac{\rho_+ (\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) - \zeta_+}{2} \|\beta' - \beta\|_2^2.\end{aligned}$$

Then we conclude the proof.  $\square$

### D.3. Proof of Lemma 5.2 in Wang et al. (2014a).

**PROOF. Statistical Recovery:** Since  $\|\beta_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\|\beta_{\bar{S}^*}^*\|_0 = 0$ , we have  $\|(\beta - \beta^*)_{\bar{S}^*}\| \leq \tilde{s}$ . For logistic loss, we further have  $\|\beta\|_2 \leq R$  and  $\|\beta^*\|_2 \leq R$ . Thus Lemma 5.1 of Wang et al. (2014a) gives

$$(D.8) \quad \tilde{\mathcal{L}}_\lambda(\beta^*) \geq \tilde{\mathcal{L}}_\lambda(\beta) + (\beta^* - \beta)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta) + \frac{\rho_- - \zeta_-}{2} \|\beta^* - \beta\|_2^2,$$

$$(D.9) \quad \tilde{\mathcal{L}}_\lambda(\beta) \geq \tilde{\mathcal{L}}_\lambda(\beta^*) + (\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + \frac{\rho_- - \zeta_-}{2} \|\beta^* - \beta\|_2^2.$$

Adding (D.8) and (D.9) and moving  $(\beta^* - \beta)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta)$  to the left-hand side, we obtain

$$(D.10) \quad (\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta) \geq (\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + (\rho_- - \zeta_-) \|\beta^* - \beta\|_2^2.$$

Let  $\xi \in \partial \|\beta\|_1$  be the subgradient that attains the minimum in

$$\omega_\lambda(\beta) = \min_{\xi' \in \partial \|\beta\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\beta - \beta')^T}{\|\beta - \beta'\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi') \right\}.$$

Then we have

$$(D.11) \quad \omega_\lambda(\beta) = \max_{\beta' \in \Omega} \left\{ \frac{(\beta - \beta')^T}{\|\beta - \beta'\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \right\}.$$

Adding  $\lambda(\beta - \beta^*)^T \xi$  to the both sides of (D.10), we obtain

$$\begin{aligned} &(\beta - \beta^*)^T (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \\ &\geq (\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + (\rho_- - \zeta_-) \|\beta^* - \beta\|_2^2 + \lambda(\beta - \beta^*)^T \xi. \end{aligned}$$

Since  $\beta^* \in \Omega$ , by (D.11) we have

$$(D.12) \quad \frac{(\beta - \beta^*)^T}{\|\beta - \beta^*\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \leq \max_{\beta' \in \Omega} \left\{ \frac{(\beta - \beta')^T}{\|\beta - \beta'\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \right\} = \omega_\lambda(\beta).$$

Recall that we assume  $\omega_\lambda(\beta) \leq \lambda/2$ , we obtain

$$(D.13) \quad (\beta - \beta^*)^T (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \leq \lambda/2 \cdot \|\beta - \beta^*\|_1.$$

Plugging (D.13) into the left-hand side of (D.10), we obtain

$$(D.14) \quad \begin{aligned} & \lambda/2 \cdot \|\beta - \beta^*\|_1 \\ & \geq \underbrace{(\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*)}_{(i)} + (\rho_- - \zeta_-) \|\beta^* - \beta\|_2^2 + \underbrace{\lambda(\beta - \beta^*)^T \xi}_{(ii)}. \end{aligned}$$

Now we provide lower bounds of terms (i) and (ii) in (D.14) respectively.

- **Bounding Term (i) in (D.14):** Recall that  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$ .

We have

$$(D.15) \quad (\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) = \underbrace{(\beta - \beta^*)^T \nabla \mathcal{L}(\beta^*)}_{(i).a} + \underbrace{(\beta - \beta^*)^T \nabla \mathcal{Q}_\lambda(\beta^*)}_{(i).b}.$$

Separating the support of  $\beta - \beta^*$  into  $S^*$  and  $\bar{S}^*$ , we obtain

$$\|\beta - \beta^*\|_1 = \|(\beta - \beta^*)_{\bar{S}^*}\|_1 + \|(\beta - \beta^*)_{S^*}\|_1.$$

Then for term (i).a in (D.15), we have

$$(D.16) \quad \begin{aligned} & (\beta - \beta^*)^T \nabla \mathcal{L}(\beta^*) \\ & \geq -\|\beta - \beta^*\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty \\ & = -\|(\beta - \beta^*)_{\bar{S}^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty - \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty. \end{aligned}$$

For term (i).b in (D.15), we have

$$(D.17) \quad \begin{aligned} & (\beta - \beta^*)^T \nabla \mathcal{Q}_\lambda(\beta^*) \\ & = (\beta - \beta^*)_{S^*}^T (\nabla \mathcal{Q}_\lambda(\beta^*))_{S^*} + (\beta - \beta^*)_{\bar{S}^*}^T (\nabla \mathcal{Q}_\lambda(\beta^*))_{\bar{S}^*}. \end{aligned}$$

Note that  $\mathcal{Q}_\lambda(\beta^*)$  is separable. We have

$$(D.18) \quad \begin{aligned} & (\beta - \beta^*)_{S^*}^T (\nabla \mathcal{Q}_\lambda(\beta^*))_{S^*} = \sum_{j \in S^*} (\beta_j - \beta_j^*) \cdot q'_\lambda(\beta_j^*) \\ & = (\beta - \beta^*)_{S^*}^T \nabla \mathcal{Q}_\lambda(\beta^*), \\ & (\beta - \beta^*)_{\bar{S}^*}^T (\nabla \mathcal{Q}_\lambda(\beta^*))_{\bar{S}^*} = \sum_{j \in \bar{S}^*} (\beta_j - \beta_j^*) \cdot q'_\lambda(\beta_j^*) \\ (D.19) \quad & = \sum_{j \in \bar{S}^*} (\beta_j - \beta_j^*) \cdot q'_\lambda(0) = 0, \end{aligned}$$

where the second equation in (D.19) is because  $\beta_j^* = 0$  for  $j \in \bar{S}^*$ , and the third is by regularity condition (c) that  $q'_\lambda(0) = 0$ . Plugging (D.18) and (D.19) into the right-hand side of (D.17), for term (i).b in (D.15)

we obtain

$$(D.20) \quad \begin{aligned} (\beta - \beta^*)^T \nabla \mathcal{Q}_\lambda(\beta^*) &= (\beta - \beta^*)_{S^*}^T \nabla \mathcal{Q}_\lambda(\beta^*) \\ &\geq -\|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty. \end{aligned}$$

Plugging (D.16) and (D.20) into the right-hand side of (D.15), then for term (i) in (D.14) we obtain

$$(D.21) \quad \begin{aligned} &(\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) \\ &\geq -\|(\beta - \beta^*)_{\bar{S}^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty - \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty \\ &\quad - \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty. \end{aligned}$$

- **Bounding Term (ii) in (D.14):** For term (ii) in (D.14), by separating the support of  $\beta - \beta^*$  into  $S^*$  and  $\bar{S}^*$  we have

$$(D.22) \quad \lambda(\beta - \beta^*)^T \xi = \underbrace{\lambda(\beta - \beta^*)_{S^*}^T \xi_{S^*}}_{(ii).a} + \underbrace{\lambda(\beta - \beta^*)_{\bar{S}^*}^T \xi_{\bar{S}^*}}_{(ii).b}.$$

For term (ii).a in (D.22), since  $\xi \in \partial\|\beta\|_1$ , we have  $\|\xi_{S^*}\|_\infty \leq \|\xi\|_\infty \leq 1$ , which implies

$$(D.23) \quad (\beta - \beta^*)_{S^*}^T \xi_{S^*} \geq -\|\xi_{S^*}\|_\infty \|(\beta - \beta^*)_{S^*}\|_1 \geq -\|(\beta - \beta^*)_{S^*}\|_1.$$

For term (ii).b in (D.22), note that  $\beta_{\bar{S}^*}^* = \mathbf{0}$ . Hence,  $(\beta - \beta^*)_{\bar{S}^*} = \beta_{\bar{S}^*}$ . Recall  $\xi \in \partial\|\beta\|_1$ . For  $\beta_j \neq 0$ , since  $\xi_j = \text{sign}(\beta_j)$ , we have  $\beta_j \xi_j = |\beta_j|$ . For  $\beta_j = 0$ , we have  $\beta_j \xi_j = |\beta_j| = 0$ . Therefore, we obtain

$$(D.24) \quad \begin{aligned} (\beta - \beta^*)_{\bar{S}^*}^T \xi_{\bar{S}^*} &= \beta_{\bar{S}^*}^T \xi_{\bar{S}^*} = \sum_{j \in S^*} \beta_j \xi_j = \sum_{j \in \bar{S}^*} |\beta_j| = \|\beta_{\bar{S}^*}\|_1 \\ &= \|(\beta - \beta^*)_{\bar{S}^*}\|_1. \end{aligned}$$

Plugging (D.23) and (D.24) into the right-hand side of (D.22), we obtain

$$(D.25) \quad \lambda(\beta - \beta^*)^T \xi \geq -\lambda\|(\beta - \beta^*)_{S^*}\|_1 + \lambda\|(\beta - \beta^*)_{\bar{S}^*}\|_1.$$

Plugging (D.21) and (D.25) into the right-hand side of (D.14), we obtain

$$(D.26) \quad \begin{aligned} &\lambda/2 \cdot \|\beta - \beta^*\|_1 \\ &\geq \underbrace{-\|(\beta - \beta^*)_{\bar{S}^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty - \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty - \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty}_{(i) \text{ in (D.14)}} \\ &\quad + (\rho_- - \zeta_-) \|\beta^* - \beta\|_2^2 - \underbrace{\lambda\|(\beta - \beta^*)_{S^*}\|_1 + \lambda\|(\beta - \beta^*)_{\bar{S}^*}\|_1}_{(ii) \text{ in (D.14)}}. \end{aligned}$$

Again, we separate the left-hand side of (D.26) as  $\lambda/2 \cdot \|\beta - \beta^*\|_1 = \lambda/2 \cdot$

$\|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{\bar{S}^*}\|_1 + \lambda/2 \cdot \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1$ . Rearranging the terms, we obtain

$$\begin{aligned}
 (D.27) \quad & (\rho_- - \zeta_-) \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2^2 + \underbrace{(\lambda/2 - \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty)}_{(i)} \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{\bar{S}^*}\|_1 \\
 & \leq (3\lambda/2 + \underbrace{\|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty}_{(ii)} + \underbrace{\|\nabla \mathcal{Q}_\lambda(\boldsymbol{\beta}^*)\|_\infty}_{(iii)}) \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1.
 \end{aligned}$$

For term (ii) in (D.27), by (4.1) in Assumption 4.1 of Wang et al. (2014a) and  $\lambda \geq \lambda_{\text{tgt}}$  we have

$$(D.28) \quad \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty \leq \lambda_{\text{tgt}}/8 \leq \lambda/8.$$

Meanwhile, (D.28) also implies that term (i) in (D.27) is positive. Recall that  $\mathcal{Q}_\lambda(\boldsymbol{\beta}) = \sum_{j=1}^d q_\lambda(\beta_j)$ , where  $q_\lambda(\beta_j)$  satisfies regularity condition (d) in Wang et al. (2014a). Hence for term (iii) in (D.27) we have

$$(D.29) \quad \|\nabla \mathcal{Q}_\lambda(\boldsymbol{\beta}^*)\|_\infty = \max_{1 \leq j \leq d} |q'_\lambda(\beta_j^*)| \leq \lambda.$$

In summary, from (D.27) we obtain

$$\begin{aligned}
 (D.30) \quad & (\rho_- - \zeta_-) \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2^2 \leq (3\lambda/2 + \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_\infty + \|\nabla \mathcal{Q}_\lambda(\boldsymbol{\beta}^*)\|_\infty) \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1 \\
 & \leq (3\lambda/2 + \lambda/8 + \lambda) \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1 \\
 & \leq 21\lambda/8 \cdot \sqrt{s^*} \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_2 \\
 & \leq 21\lambda/8 \cdot \sqrt{s^*} \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2.
 \end{aligned}$$

According to (4.3) of Wang et al. (2014a), we have  $\rho_- - \zeta_- > 0$ . Therefore, (D.30) gives

$$(D.31) \quad \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \leq \frac{21/8}{\rho_- - \zeta_-} \lambda \sqrt{s^*},$$

which implies the first conclusion.

**Results for the Objective Function Value:** Note that on the right-hand side of (D.9), we have  $\rho_- - \zeta_- > 0$ , which gives

$$(D.32) \quad \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*) \geq \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) + (\boldsymbol{\beta}^* - \boldsymbol{\beta})^T \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}).$$

Meanwhile, since  $\boldsymbol{\xi} \in \partial \|\boldsymbol{\beta}\|_1$ , by the convexity of  $\ell_1$  norm we have

$$(D.33) \quad \lambda \|\boldsymbol{\beta}^*\|_1 \geq \lambda \|\boldsymbol{\beta}\|_1 + \lambda (\boldsymbol{\beta}^* - \boldsymbol{\beta})^T \boldsymbol{\xi}.$$

Recall that  $\phi_\lambda(\boldsymbol{\beta}) = \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1$ . Adding (D.32) and (D.33), we obtain

$$(D.34) \quad \phi_\lambda(\boldsymbol{\beta}^*) \geq \phi_\lambda(\boldsymbol{\beta}) + (\boldsymbol{\beta}^* - \boldsymbol{\beta})^T (\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) + \lambda \boldsymbol{\xi}),$$

which implies

$$\phi_\lambda(\boldsymbol{\beta}) - \phi_\lambda(\boldsymbol{\beta}^*) \leq (\boldsymbol{\beta} - \boldsymbol{\beta}^*)^T (\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) + \lambda \boldsymbol{\xi}) \leq \lambda/2 \cdot \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_1.$$

Here the second inequality follows from (D.13), which is a direct consequence of the assumption that  $\omega_\lambda(\beta) \leq \lambda/2$ . Separating the support of  $\beta - \beta^*$  into  $S^*$  and  $\bar{S}^*$ , we obtain

$$\begin{aligned} \phi_\lambda(\beta) - \phi_\lambda(\beta^*) &\leq \lambda/2 \cdot \|\beta - \beta^*\|_1 \\ (D.35) \quad &\leq \lambda/2 \cdot \|(\beta - \beta^*)_{S^*}\|_1 + \lambda/2 \cdot \|(\beta - \beta^*)_{\bar{S}^*}\|_1. \end{aligned}$$

Now we derive an upper bound of  $\|(\beta - \beta^*)_{\bar{S}^*}\|_1$  on the right-hand side of (D.35). On the left-hand side of (D.27), we have  $\rho_- - \zeta_- > 0$ , which gives

$$\begin{aligned} &(\lambda/2 - \|\nabla \mathcal{L}(\beta^*)\|_\infty) \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \\ (D.36) \quad &\leq (3\lambda/2 + \|\nabla \mathcal{L}(\beta^*)\|_\infty + \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty) \|(\beta - \beta^*)_{S^*}\|_1. \end{aligned}$$

Note that in (D.36) we have  $\|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda/8$  by (D.28), and  $\|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty \leq \lambda$  by (D.29). Hence we have

$$(D.37) \quad (\lambda/2 - \lambda/8) \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \leq (3\lambda/2 + \lambda/8 + \lambda) \|(\beta - \beta^*)_{S^*}\|_1,$$

which implies  $\|(\beta - \beta^*)_{\bar{S}^*}\|_1 \leq 7 \|(\beta - \beta^*)_{S^*}\|_1$ . Plugging this into the right-hand side of (D.35), we obtain

$$\begin{aligned} \phi_\lambda(\beta) - \phi_\lambda(\beta^*) &\leq (\lambda/2 + 7\lambda/2) \|(\beta - \beta^*)_{S^*}\|_1 \\ (D.38) \quad &\leq 4\lambda\sqrt{s^*} \|(\beta - \beta^*)_{S^*}\|_2 \leq 4\lambda\sqrt{s^*} \|\beta - \beta^*\|_2. \end{aligned}$$

Plugging the upper bound of  $\|\beta - \beta^*\|_2$  in (D.31) into the right-hand side of (D.38), we obtain

$$\phi_\lambda(\beta) - \phi_\lambda(\beta^*) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*.$$

Hence we reach the second conclusion.  $\square$

#### D.4. Proof of Lemma 5.3 in Wang et al. (2014a).

PROOF. Since  $\|\beta_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\|\beta_{\bar{S}^*}^*\|_0 = 0$ , we have  $\|(\beta - \beta^*)_{\bar{S}^*}\|_0 \leq \tilde{s}$ . For logistic loss, we further have  $\|\beta\|_2 \leq R$  and  $\|\beta^*\|_2 \leq R$ , where  $R$  is specified in Definition 4.3 of Wang et al. (2014a). Therefore, Lemma 5.1 of Wang et al. (2014a) gives

$$(D.39) \quad \tilde{\mathcal{L}}_\lambda(\beta^*) + (\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + \frac{\rho_- - \zeta_-}{2} \|\beta^* - \beta\|_2^2 \leq \tilde{\mathcal{L}}_\lambda(\beta).$$

Recall that  $\phi_\lambda(\beta) = \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \|\beta\|_1$ . Hence, from our assumption that

$$\phi_\lambda(\beta) - \phi_\lambda(\beta^*) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*$$

we obtain

$$(D.40) \quad \tilde{\mathcal{L}}_\lambda(\beta) - \tilde{\mathcal{L}}_\lambda(\beta^*) + \lambda(\|\beta\|_1 - \|\beta^*\|_1) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*.$$

Plugging (D.39) into the left-hand side of (D.40), we have

$$(\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + \frac{\rho_- - \zeta_-}{2} \|\beta^* - \beta\|_2^2 + \lambda(\|\beta\|_1 - \|\beta^*\|_1) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*.$$

Moving  $(\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + \lambda(\|\beta\|_1 - \|\beta^*\|_1)$  to its right-hand side yields

$$(D.41) \quad \begin{aligned} & \frac{\rho_- - \zeta_-}{2} \|\beta^* - \beta\|_2^2 \\ & \leq \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^* - \underbrace{(\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*)}_{(i)} + \underbrace{\lambda(\|\beta^*\|_1 - \|\beta\|_1)}_{(ii)}. \end{aligned}$$

For term (i) in (D.41), following the same way we obtain the lower bound of term (i) in (D.14) (in the proof of Lemma 5.2), we can obtain the same result as in (D.21), which implies

$$(D.42) \quad \begin{aligned} -(\beta - \beta^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) & \leq \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty + \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty \\ & \quad + \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty. \end{aligned}$$

For term (ii) in (D.41), separating the support of  $\beta$  and  $\beta^*$  into  $S^*$  and  $\bar{S}^*$  respectively, we obtain

$$(D.43) \quad \|\beta^*\|_1 - \|\beta\|_1 = \|\beta_{S^*}^*\|_1 + \|\beta_{\bar{S}^*}^*\|_1 - (\|\beta_{S^*}\|_1 + \|\beta_{\bar{S}^*}\|_1).$$

Note that  $\beta_{\bar{S}^*}^* = \mathbf{0}$ , which gives  $\beta_{\bar{S}^*} = \beta_{\bar{S}^*} - \beta_{\bar{S}^*}^* = (\beta - \beta^*)_{\bar{S}^*}$ . Hence, from (D.43) we have

$$(D.44) \quad \begin{aligned} \|\beta^*\|_1 - \|\beta\|_1 & = \|\beta_{S^*}^*\|_1 - \|\beta_{S^*}\|_1 - \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \\ & \leq \|(\beta - \beta^*)_{S^*}\|_1 - \|(\beta - \beta^*)_{\bar{S}^*}\|_1, \end{aligned}$$

where the inequality follows from the triangle inequality. Plugging (D.42) and (D.44) into the right-hand side of (D.41), we obtain

$$(D.45) \quad \begin{aligned} & \frac{\rho_- - \zeta_-}{2} \|\beta^* - \beta\|_2^2 \\ & \leq \underbrace{\|(\beta - \beta^*)_{\bar{S}^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty + \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{L}(\beta^*)\|_\infty + \|(\beta - \beta^*)_{S^*}\|_1 \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty}_{(i) \text{ in (D.41)}} \\ & \quad + \underbrace{\lambda(\|(\beta - \beta^*)_{S^*}\|_1 - \|(\beta - \beta^*)_{\bar{S}^*}\|_1)}_{(ii) \text{ in (D.41)}} + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*. \end{aligned}$$



Rearranging the terms in (D.45), we obtain

$$\begin{aligned}
 (D.46) \quad & \frac{\rho_- - \zeta_-}{2} \|\beta - \beta^*\|_2^2 + \underbrace{(\lambda - \|\nabla \mathcal{L}(\beta^*)\|_\infty) \|(\beta - \beta^*)_{\bar{S}^*}\|_1}_{(i)} \\
 & \leq (\lambda + \underbrace{\|\nabla \mathcal{L}(\beta^*)\|_\infty}_{(ii)} + \underbrace{\|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty}_{(iii)}) \|(\beta - \beta^*)_{S^*}\|_1 + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*.
 \end{aligned}$$

By (4.1) in Assumption 4.1 of Wang et al. (2014a) and  $\lambda \geq \lambda_{\text{tgt}}$ , for term (ii) in (D.46), we have

$$(D.47) \quad \|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda_{\text{tgt}}/8 \leq \lambda/8.$$

Moreover, (D.47) implies that term (i) in (D.46) is positive. For term (iii) in (D.46), since  $\mathcal{Q}_\lambda(\beta) = \sum_{j=1}^d q_\lambda(\beta_j)$ , where  $q_\lambda(\beta_j)$  satisfies regularity condition (d), we have

$$(D.48) \quad \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty \leq \max_{1 \leq j \leq d} |q'_\lambda(\beta_j^*)| \leq \lambda.$$

Therefore, from (D.48) we obtain

$$\begin{aligned}
 & \frac{\rho_- - \zeta_-}{2} \|\beta - \beta^*\|_2^2 \\
 & \leq (\lambda + \|\nabla \mathcal{L}(\beta^*)\|_\infty + \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty) \|(\beta - \beta^*)_{S^*}\|_1 + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^* \\
 & \leq (\lambda + \lambda/8 + \lambda) \|(\beta - \beta^*)_{S^*}\|_1 + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^* \\
 (D.49) \quad & \leq 17/8 \cdot \lambda \|(\beta - \beta^*)_{S^*}\|_1 + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*.
 \end{aligned}$$

To further obtain an upper bound of the right-hand side of (D.49), we discuss two cases regarding the relationship between  $\|(\beta - \beta^*)_{S^*}\|_1$  and  $\lambda s^*$ .

- If  $7/(\rho_- - \zeta_-) \cdot \lambda s^* < \|(\beta - \beta^*)_{S^*}\|_1$ , then we have

$$\frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^* < 3/2 \cdot \lambda \|(\beta - \beta^*)_{S^*}\|_1.$$

Plugging this into the right-hand side of (D.49), we obtain

$$\begin{aligned}
 \frac{\rho_- - \zeta_-}{2} \|\beta - \beta^*\|_2^2 & \leq (17/8 \cdot \lambda + 3/2 \cdot \lambda) \|(\beta - \beta^*)_{S^*}\|_1 \\
 & \leq 29/8 \cdot \lambda \sqrt{s^*} \|(\beta - \beta^*)_{S^*}\|_2 \\
 & \leq 29/8 \cdot \lambda \sqrt{s^*} \|\beta - \beta^*\|_2.
 \end{aligned}$$

Dividing  $\|\beta^* - \beta\|_2$  on both sides, we have

$$(D.50) \quad \|\beta - \beta^*\|_2 \leq \frac{29/4}{\rho_- - \zeta_-} \lambda \sqrt{s^*}.$$

- If  $\|(\beta - \beta^*)_{S^*}\|_1 \leq 7/(\rho_- - \zeta_-) \cdot \lambda s^*$ , then we have

$$17/8 \cdot \lambda \|(\beta - \beta^*)_{S^*}\|_1 < \frac{119/8}{\rho_- - \zeta_-} \lambda^2 s^*.$$

Plugging this into the right-hand side of (D.49), we obtain

$$(D.51) \quad \begin{aligned} \frac{\rho_- - \zeta_-}{2} \|\beta - \beta^*\|_2^2 &\leq \frac{119/8}{\rho_- - \zeta_-} \lambda^2 s^* + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^* \\ &= \frac{203/8}{\rho_- - \zeta_-} \lambda^2 s^*, \end{aligned}$$

which implies

$$(D.52) \quad \|\beta - \beta^*\|_2 \leq \frac{\sqrt{203}/2}{\rho_- - \zeta_-} \lambda \sqrt{s^*}.$$

Combining (D.50) and (D.52), since  $\max\{29/4, \sqrt{203}/2\} \leq 15/2$ , we obtain

$$\|\beta - \beta^*\|_2 < \frac{15/2}{\rho_- - \zeta_-} \lambda \sqrt{s^*}.$$

Hence we conclude the proof.  $\square$

#### D.5. Proof of Lemma 5.4 in Wang et al. (2014a).

PROOF. Recall the proximal-gradient update step defined in (3.8) of Wang et al. (2014a) with  $\Omega = \mathbb{R}^d$ , i.e.,  $R = +\infty$ , takes the form

$$(D.53) \quad (\mathcal{T}_{L,\lambda}(\beta; +\infty))_j = \begin{cases} 0 & \text{if } |\bar{\beta}_j| \leq \lambda/L, \\ \text{sign}(\bar{\beta}_j)(|\bar{\beta}_j| - \lambda/L) & \text{if } |\bar{\beta}_j| > \lambda/L, \end{cases}$$

for  $j = 1, \dots, d$ , where

$$(D.54) \quad \bar{\beta} = \beta - \frac{1}{L} \nabla \tilde{\mathcal{L}}_\lambda(\beta),$$

and  $\bar{\beta}_j$  is the  $j$ -th dimension of  $\bar{\beta}$ . Furthermore, if  $\Omega = B_2(R)$  of radius  $R \in (0, \infty)$ ,  $\mathcal{T}_{L,\lambda}(\beta; R)$  can be obtained by projecting  $\mathcal{T}_{L,\lambda}(\beta; +\infty)$  shown in (D.53) onto  $B_2(R)$ , i.e.,

$$(D.55) \quad \mathcal{T}_{L,\lambda}(\beta; R) = \begin{cases} \mathcal{T}_{L,\lambda}(\beta; +\infty) & \text{if } \|\mathcal{T}_{L,\lambda}(\beta; +\infty)\|_2 < R, \\ \frac{R \cdot \mathcal{T}_{L,\lambda}(\beta; +\infty)}{\|\mathcal{T}_{L,\lambda}(\beta; +\infty)\|_2} & \text{if } \|\mathcal{T}_{L,\lambda}(\beta; +\infty)\|_2 \geq R. \end{cases}$$

Note that  $\mathcal{T}_{L,\lambda}(\beta; +\infty)$  and  $\mathcal{T}_{L,\lambda}(\beta; R)$  have exactly the same sparsity pattern.

Hence we focus on analyzing the sparsity pattern of  $\mathcal{T}_{L,\lambda}(\beta; +\infty)$  in the following.

In fact, update scheme (D.53) defines a soft-thresholding operation on  $\bar{\beta}$  defined in (D.54), with the threshold value  $\lambda/L$ . To show  $\|(\mathcal{T}_{L,\lambda}(\beta; +\infty))_{\bar{S}^*}\|_0 \leq \tilde{s}$ , we need to prove that, for  $j \in \bar{S}^*$ , the number of  $j$ 's such that  $|\bar{\beta}_j| > \lambda/L$  is no more than  $\tilde{s}$ . To achieve this goal, we first reformulate  $\bar{\beta}$  as

$$(D.56) \quad \bar{\beta} = \beta - \frac{1}{L} \nabla \tilde{\mathcal{L}}_\lambda(\beta) = \beta - \frac{1}{L} \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + \frac{1}{L} (\nabla \tilde{\mathcal{L}}_\lambda(\beta^*) - \nabla \tilde{\mathcal{L}}_\lambda(\beta)).$$

Then it suffices to prove there exist integers  $\tilde{s}_1$ ,  $\tilde{s}_2$  and  $\tilde{s}_3$ , which satisfy  $\tilde{s}_1 + \tilde{s}_2 + \tilde{s}_3 \leq \tilde{s}$ , such that

$$(D.57) \quad |\{j \in \bar{S}^* : |\beta_j| \geq 1/4 \cdot \lambda/L\}| \leq \tilde{s}_1,$$

$$(D.58) \quad \left| \left\{ j \in \bar{S}^* : |(\nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| > 1/8 \cdot \lambda/L \right\} \right| \leq \tilde{s}_2,$$

$$(D.59) \quad \left| \left\{ j \in \bar{S}^* : |(\nabla \tilde{\mathcal{L}}_\lambda(\beta)/L - \nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| \geq 5/8 \cdot \lambda/L \right\} \right| \leq \tilde{s}_3.$$

This is because, if (D.57)-(D.59) hold, then there are at most  $\tilde{s}_1 + \tilde{s}_2 + \tilde{s}_3 \leq \tilde{s}$  coordinates  $j \in \bar{S}^*$  such that

$$|\beta_j| + |(\nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| + |(\nabla \tilde{\mathcal{L}}_\lambda(\beta)/L - \nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| > \lambda/L.$$

Since by the triangular inequality (D.56) implies

$$|\bar{\beta}_j| \leq |\beta_j| + |(\nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| + |(\nabla \tilde{\mathcal{L}}_\lambda(\beta)/L - \nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j|,$$

the number of coordinates  $j \in \bar{S}^*$  such that  $|\bar{\beta}_j| > \lambda/L$  is also upper bounded by  $\tilde{s}_1 + \tilde{s}_2 + \tilde{s}_3 \leq \tilde{s}$ . In the following, we will prove (D.58)-(D.59) and specify the corresponding  $\tilde{s}_1$ ,  $\tilde{s}_2$  and  $\tilde{s}_3$ .

**Proof of (D.57):** Note that for  $j \in \bar{S}^*$ , we have  $\beta_j^* = 0$ . Hence we have

$$(D.60) \quad |\{j \in \bar{S}^* : |\beta_j| \geq 1/4 \cdot \lambda/L\}| = |\{j \in \bar{S}^* : |\beta_j - \beta_j^*| \geq 1/4 \cdot \lambda/L\}|.$$

Meanwhile, note that

$$\begin{aligned} & \frac{\lambda}{4L} |\{j \in \bar{S}^* : |\beta_j - \beta_j^*| \geq 1/4 \cdot \lambda/L\}| \\ & \leq \sum_{j \in \bar{S}^*} |\beta_j - \beta_j^*| \cdot \mathbf{1}(|\beta_j - \beta_j^*| \geq 1/4 \cdot \lambda/L) \\ & \leq \sum_{j \in \bar{S}^*} |\beta_j - \beta_j^*| \\ (D.61) \quad & = \|(\beta - \beta^*)_{\bar{S}^*}\|_1. \end{aligned}$$

Plugging (D.61) into the right-hand side of (D.60), we obtain

$$(D.62) \quad |\{j \in \bar{S}^* : |\beta_j| \geq 1/4 \cdot \lambda/L\}| \leq \frac{4L}{\lambda} \|(\beta - \beta^*)_{\bar{S}^*}\|_1.$$

Now we provide an upper bound of  $\|(\beta - \beta^*)_{\bar{S}^*}\|_1$ . Following the same way we derive (D.46) in the proof of Lemma 5.3, we can obtain

$$(D.63) \quad \begin{aligned} & \frac{\rho_- - \zeta_-}{2} \|\beta - \beta^*\|_2^2 + (\lambda - \|\nabla \mathcal{L}(\beta^*)\|_\infty) \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \\ & \leq (\lambda + \|\nabla \mathcal{L}(\beta^*)\|_\infty + \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty) \|(\beta - \beta^*)_{S^*}\|_1 + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*. \end{aligned}$$

According to (4.3) in Wang et al. (2014a), we have  $\rho_- - \zeta_- > 0$ . Hence (D.63) implies

$$(D.64) \quad \begin{aligned} & (\lambda - \|\nabla \mathcal{L}(\beta^*)\|_\infty) \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \\ & \leq (\lambda + \|\nabla \mathcal{L}(\beta^*)\|_\infty + \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty) \|(\beta - \beta^*)_{S^*}\|_1 + \frac{21/2}{\rho_- - \zeta_-} \lambda^2 s^*. \end{aligned}$$

By (4.1) in Assumption 4.1 of Wang et al. (2014a) and  $\lambda \geq \lambda_{\text{tgt}}$ , we have

$$(D.65) \quad \|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda_{\text{tgt}}/8 \leq \lambda/8.$$

Meanwhile, since  $\mathcal{Q}_\lambda(\beta) = \sum_{j=1}^d q_\lambda(\beta_j)$  and  $q_\lambda(\beta_j)$  satisfies regularity condition (d) in Wang et al. (2014a), we have

$$(D.66) \quad \|\nabla \mathcal{Q}_\lambda(\beta^*)\|_\infty = \max_{1 \leq j \leq d} |q'_\lambda(\beta_j^*)| \leq \lambda.$$

Plugging (D.65) and (D.66) into (D.64) and dividing  $\lambda$  on both sides, we obtain

$$(D.67) \quad 7/8 \cdot \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \leq 17/8 \cdot \|(\beta - \beta^*)_{S^*}\|_1 + \frac{21/2}{\rho_- - \zeta_-} \lambda s^*.$$

Now we discuss two cases regarding the relationship between  $\|(\beta - \beta^*)_{S^*}\|_1$  and  $\lambda s^*$ .

- If  $7/(\rho_- - \zeta_-) \cdot \lambda s^* < \|(\beta - \beta^*)_{S^*}\|_1$ , then we have

$$\frac{21/2}{\rho_- - \zeta_-} \lambda s^* \leq 3/2 \cdot \|(\beta - \beta^*)_{S^*}\|_1.$$

Plugging this into the right-hand side of (D.67), we obtain

$$\|(\beta - \beta^*)_{\bar{S}^*}\|_1 \leq 29/7 \cdot \|(\beta - \beta^*)_{S^*}\|_1,$$

which implies

$$(D.68) \quad \begin{aligned} \|(\beta - \beta^*)_{\bar{S}^*}\|_1 & \leq 29/7 \cdot \|(\beta - \beta^*)_{S^*}\|_1 \leq 29/7 \cdot \sqrt{s^*} \|(\beta - \beta^*)_{S^*}\|_2 \\ & \leq 29/7 \cdot \sqrt{s^*} \|\beta - \beta^*\|_2. \end{aligned}$$

Plugging the upper bound of  $\|\beta - \beta^*\|_2$  in Lemma 5.3 of Wang et al. (2014a) into the right-hand side of (D.68), we obtain

$$(D.69) \quad \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \leq 29/7 \cdot \sqrt{s^*} \cdot \frac{15/2}{\rho_- - \zeta_-} \lambda \sqrt{s^*} = \frac{435/14}{\rho_- - \zeta_-} \lambda s^*.$$

- If  $\|(\beta - \beta^*)_{S^*}\|_1 \leq 7/(\rho_- - \zeta_-) \cdot \lambda s^*$ , then plugging this into the right-hand side of (D.67), we obtain

$$(D.70) \quad \|(\beta - \beta^*)_{\bar{S}^*}\|_1 \leq 8/7 \cdot \frac{17/8 \cdot 7 + 21/2}{\rho_- - \zeta_-} \lambda s^* = \frac{29}{\rho_- - \zeta_-} \lambda s^*.$$

Combining (D.69) and (D.70), we obtain

$$\|(\beta - \beta^*)_{\bar{S}^*}\|_1 \leq \frac{\max\{435/14, 29\}}{\rho_- - \zeta_-} \lambda s^* \leq \frac{435/14}{\rho_- - \zeta_-} \lambda s^*.$$

Plugging this into the right-hand side of (D.62), we obtain

$$|\{j \in \bar{S}^* : |\beta_j| \geq 1/4 \cdot \lambda/L\}| \leq \frac{4L}{\lambda} \cdot \frac{435/14}{\rho_- - \zeta_-} \lambda s^* < \frac{125L}{\rho_- - \zeta_-} s^*.$$

Meanwhile, since we assume  $L < 2(\rho_+ - \zeta_+)$ , we have

$$|\{j \in \bar{S}^* : |\beta_j| \geq 1/4 \cdot \lambda/L\}| < 250 \cdot \frac{\rho_+ - \zeta_+}{\rho_- - \zeta_-} \cdot s^* = 250\kappa s^*,$$

where the last equality follows from the definition of the condition number  $\kappa$  in (4.5). Therefore we obtain (D.57) by setting  $\tilde{s}_1 = 250\kappa s^*$ .

**Proof of (D.58):** Recall that  $\nabla \tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$ . Hence we have

$$(D.71) \quad \|(\nabla \tilde{\mathcal{L}}_\lambda(\beta^*))_{\bar{S}^*}\|_\infty \leq \|(\nabla \mathcal{L}(\beta^*))_{\bar{S}^*}\|_\infty + \|(\nabla \mathcal{Q}_\lambda(\beta^*))_{\bar{S}^*}\|_\infty.$$

By (4.1) in Assumption 4.1 of Wang et al. (2014a), we have

$$(D.72) \quad \|(\nabla \mathcal{L}(\beta^*))_{\bar{S}^*}\|_\infty \leq \|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda/8.$$

Recall  $\mathcal{Q}_\lambda(\beta) = \sum_{j=1}^d q_\lambda(\beta_j)$ , where  $q_\lambda(\beta_j)$  satisfies regularity condition (c) that  $q'_\lambda(0) = 0$ . Hence we have

$$(D.73) \quad \|(\nabla \mathcal{Q}_\lambda(\beta^*))_{\bar{S}^*}\|_\infty = \max_{j \in \bar{S}^*} |q'_\lambda(\beta_j^*)| = \max_{j \in \bar{S}^*} |q'_\lambda(0)| = 0,$$

where the second equation follows from the fact that  $\beta_j^* = 0$  for  $j \in \bar{S}^*$ . Plugging (D.73) and (D.72) into the right-hand side of (D.71), we obtain  $\|(\nabla \tilde{\mathcal{L}}_\lambda(\beta^*))_{\bar{S}^*}\|_\infty = \max_{j \in \bar{S}^*} |(\nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| \leq \lambda/8$ . Hence we have

$$|\{j \in \bar{S}^* : |(\nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| > 1/8 \cdot \lambda/L\}| = 0.$$

Therefore, by setting  $\tilde{s}_2 = 0$ , we obtain (D.58).

**Proof of (D.59):** Consider an arbitrary subset  $S'$  such that

$$(D.74) \quad S' \subseteq \left\{ j : |(\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*))_j| \geq 5/8 \cdot \lambda \right\}.$$

Let  $s' = |S'|$ . In the sequel we provide an upper bound of  $s'$ . Suppose  $\mathbf{v} \in \mathbb{R}^d$  is chosen such that  $v_j = \text{sign}\{(\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*))_j\}$  for  $j \in S'$ , and  $v_j = 0$  for  $j \notin S'$ . Hence we have

$$(D.75) \quad \begin{aligned} \mathbf{v}^T (\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*)) &= \sum_{j \in S'} v_j (\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*))_j \\ &= \sum_{j \in S'} |(\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*))_j| \geq 5/8 \cdot \lambda s'. \end{aligned}$$

Meanwhile, by Cauchy Schwarz inequality we have

$$(D.76) \quad \begin{aligned} \mathbf{v}^T (\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*)) &\leq \|\mathbf{v}\|_2 \|\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*)\|_2 \\ &\leq \sqrt{s'} \|\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*)\|_2, \end{aligned}$$

where the last inequality follows from the fact that  $\|\mathbf{v}\|_2 \leq \sqrt{s'} \|\mathbf{v}\|_\infty = \sqrt{s'}$ , because  $\mathbf{v}$  is chosen such that  $\|\mathbf{v}\|_0 = s'$ . Combining (D.75) and (D.76) gives

$$(D.77) \quad 5/8 \cdot \lambda s' \leq \mathbf{v}^T (\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*)) \leq \sqrt{s'} \|\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*)\|_2.$$

Since  $\|\boldsymbol{\beta}_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\|\boldsymbol{\beta}_{\bar{S}^*}^*\|_0 = 0$ , we have  $\|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{\bar{S}^*}\| \leq \tilde{s}$ . In the setting of logistic loss, we further have  $\|\boldsymbol{\beta}\|_2 \leq R$  and  $\|\boldsymbol{\beta}^*\|_2 \leq R$ , where  $R$  is specified in Definition 4.3 of Wang et al. (2014a). Therefore, Lemma 5.1 in Wang et al. (2014a) implies that  $\tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta})$  is restricted strongly smooth. Hence we have

$$(D.78) \quad \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) \leq \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*) + (\boldsymbol{\beta} - \boldsymbol{\beta}^*)^T \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*) + \frac{\rho_+ - \zeta_+}{2} \|\boldsymbol{\beta}^* - \boldsymbol{\beta}\|_2^2.$$

According to Nesterov (2004, Theorem 2.1.9), the strong smoothness of  $\tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta})$  is equivalent to the Lipschitz continuity of its gradient, i.e.,

$$(D.79) \quad \|\nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}) - \nabla \tilde{\mathcal{L}}_\lambda(\boldsymbol{\beta}^*)\|_2 \leq (\rho_+ - \zeta_+) \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2.$$

Plugging (D.79) into the right-hand side of (D.77), we obtain

$$(D.80) \quad 5/8 \cdot \lambda s' \leq (\rho_+ - \zeta_+) \cdot \sqrt{s'} \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2.$$

Plugging the upper bound of  $\|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2$  in Lemma 5.3 of Wang et al. (2014a) into the right-hand side of (D.80), we obtain

$$(D.81) \quad \begin{aligned} \sqrt{s'} &\leq \frac{8}{5\lambda} \cdot (\rho_+ - \zeta_+) \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2 \\ &\leq \frac{8}{5\lambda} \cdot (\rho_+ - \zeta_+) \cdot \frac{15/2}{\rho_- - \zeta_-} \lambda \sqrt{s^*} = 12\kappa \sqrt{s^*}, \end{aligned}$$

where the last equality follows from the definition of the condition number  $\kappa$

in (4.5) of Wang et al. (2014a). Hence we obtain  $s' \leq 144\kappa^2 s^*$ . Note that  $S'$  is defined as an arbitrary subset of  $\{j : |(\nabla \tilde{\mathcal{L}}_\lambda(\beta) - \nabla \tilde{\mathcal{L}}_\lambda(\beta^*))_j| \geq 5/8 \cdot \lambda\}$  and

$$\begin{aligned} & \left\{ j \in \overline{S^*} : |(\nabla \tilde{\mathcal{L}}_\lambda(\beta) - \nabla \tilde{\mathcal{L}}_\lambda(\beta^*))_j| \geq 5/8 \cdot \lambda \right\} \\ & \subseteq \left\{ j : |(\nabla \tilde{\mathcal{L}}_\lambda(\beta) - \nabla \tilde{\mathcal{L}}_\lambda(\beta^*))_j| \geq 5/8 \cdot \lambda \right\}. \end{aligned}$$

Hence we have

$$\left| \left\{ j \in \overline{S^*} : |(\nabla \tilde{\mathcal{L}}_\lambda(\beta)/L - \nabla \tilde{\mathcal{L}}_\lambda(\beta^*)/L)_j| \geq 5/8 \cdot \lambda/L \right\} \right| \leq 144\kappa^2 s^*.$$

Therefore, by setting  $\tilde{s}_3 = 144\kappa^2 s^*$ , we obtain (D.59).

In summary, we prove that (D.58)-(D.59) hold with  $\tilde{s}_1 = 250\kappa s^*$ ,  $\tilde{s}_2 = 0$  and  $\tilde{s}_2 = 144\kappa^2 s^*$ . In Assumption 4.4, we assume  $\tilde{s} \geq 144\kappa^2 + 250\kappa$ , which implies  $\tilde{s}_1 + \tilde{s}_2 + \tilde{s}_3 \leq \tilde{s}$ . Therefore we have  $\|(\mathcal{T}_{L,\lambda}(\beta; +\infty))_{\overline{S^*}}\|_0 < \tilde{s}$ . Since  $\mathcal{T}_{L,\lambda}(\beta; R)$  has the same sparsity pattern as  $\mathcal{T}_{L,\lambda}(\beta; +\infty)$ , we also have that  $\|(\mathcal{T}_{L,\lambda}(\beta; R))_{\overline{S^*}}\|_0 < \tilde{s}$  for  $R \in (0, +\infty)$ . Hence we conclude the proof.  $\square$

**D.6. Proof of Theorem 5.5 in Wang et al. (2014a).** We first provide a useful lemma. It states that if  $\beta$  is  $\epsilon$ -suboptimal with respect to the regularization parameter  $\lambda$  and sufficiently sparse, then for  $\lambda' \leq \lambda$  the objective function value  $\phi_{\lambda'}(\beta)$  is close to  $\phi_{\lambda'}(\hat{\beta}_{\lambda'})$ . Here  $\hat{\beta}_{\lambda'}$  is the exact local solution corresponding to  $\lambda'$ .

**Lemma D.3.** Let  $\lambda \geq \lambda_{\text{tgt}}$  and  $\lambda' \in [\lambda_{\text{tgt}}, \lambda]$ . Suppose  $\|\beta_{\overline{S^*}}\|_0 \leq \tilde{s}$  and  $\omega_\lambda(\beta) \leq \epsilon$ . Let  $\hat{\beta}_{\lambda'}$  be the exact local solution corresponding to  $\lambda'$ , which satisfies the exact optimality condition in (3.13) of Wang et al. (2014a) and  $\|(\hat{\beta}_{\lambda'})_{\overline{S^*}}\|_0 \leq \tilde{s}$ . For logistic loss, we further assume  $\max\{\|\beta\|_2, \|\hat{\beta}_{\lambda'}\|_2\} \leq R$ , where  $R$  is specified in Definition 4.3 of Wang et al. (2014a). Under Assumption 4.1 and Assumption 4.4 of Wang et al. (2014a), we have

$$\phi_{\lambda'}(\beta) - \phi_{\lambda'}(\hat{\beta}_{\lambda'}) \leq C(\epsilon + 2(\lambda - \lambda')) \cdot (\lambda' + \lambda)s^*, \quad \text{where } C = \frac{21}{\rho_- - \zeta_-}.$$

PROOF. Since  $\|\beta_{\overline{S^*}}\|_0 \leq \tilde{s}$  and  $\|(\hat{\beta}_{\lambda'})_{\overline{S^*}}\|_0 \leq \tilde{s}$ , we have  $\|(\beta - \hat{\beta}_{\lambda'})_{\overline{S^*}}\| \leq 2\tilde{s}$ . In the setting of logistic loss, we further have  $\|\beta\|_2 \leq R$  and  $\|\hat{\beta}_{\lambda'}\|_2 \leq R$ . Therefore, Lemma 5.1 of Wang et al. (2014a) gives

$$\begin{aligned} \tilde{\mathcal{L}}_{\lambda'}(\hat{\beta}_{\lambda'}) & \geq \tilde{\mathcal{L}}_{\lambda'}(\beta) + (\hat{\beta}_{\lambda'} - \beta)^T \nabla \tilde{\mathcal{L}}_{\lambda'}(\beta) + \frac{\rho_- - \zeta_-}{2} \|\hat{\beta}_{\lambda'} - \beta\|_2^2 \\ (D.82) \quad & \geq \tilde{\mathcal{L}}_{\lambda'}(\beta) + (\hat{\beta}_{\lambda'} - \beta)^T \nabla \tilde{\mathcal{L}}_{\lambda'}(\beta), \end{aligned}$$

where the second inequality is because  $\rho_- - \zeta_- > 0$ , which follows from (4.3) in Wang et al. (2014a).

Let  $\xi \in \partial\|\beta\|_1$  be the subgradient that attains the minimum in

$$(D.83) \quad \omega_\lambda(\beta) = \min_{\xi' \in \partial\|\beta\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\beta - \beta')^T}{\|\beta - \beta'\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi') \right\},$$

where  $\Omega = B_2(R)$  in the setting of logistic loss and  $\Omega = \mathbb{R}^d$  in other settings. Since  $\xi$  is a minimizer, we have

$$(D.84) \quad \omega_\lambda(\beta) = \max_{\beta' \in \Omega} \left\{ \frac{(\beta - \beta')^T}{\|\beta - \beta'\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \right\}.$$

By the convexity of  $\ell_1$  norm, we also have

$$(D.85) \quad \lambda' \|\hat{\beta}_{\lambda'}\|_1 \geq \lambda' \|\beta\|_1 + \lambda' \xi^T (\hat{\beta}_{\lambda'} - \beta).$$

Recall that the objective function  $\phi_\lambda(\beta)$  is defined as  $\phi_\lambda(\beta) = \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \|\beta\|_1$ . Adding (D.82) and (D.85), we obtain

$$\phi_{\lambda'}(\hat{\beta}_{\lambda'}) \geq \phi_{\lambda'}(\beta) + (\nabla \tilde{\mathcal{L}}_{\lambda'}(\beta) + \lambda' \xi)^T (\hat{\beta}_{\lambda'} - \beta).$$

Hence we have

$$(D.86) \quad \begin{aligned} & \phi_{\lambda'}(\beta) - \phi_{\lambda'}(\hat{\beta}_{\lambda'}) \\ & \leq (\nabla \tilde{\mathcal{L}}_{\lambda'}(\beta) + \lambda' \xi)^T (\beta - \hat{\beta}_{\lambda'}) \\ & = \left( \overbrace{(\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi)}^{\nabla \tilde{\mathcal{L}}_\lambda(\beta)} + (\nabla \mathcal{Q}_{\lambda'}(\beta) - \nabla \mathcal{Q}_\lambda(\beta)) \right. \\ & \quad \left. + (\lambda' \xi - \lambda \xi) \right)^T (\beta - \hat{\beta}_{\lambda'}) \\ & \leq \underbrace{(\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi)^T (\beta - \hat{\beta}_{\lambda'})}_{(i)} + \underbrace{\|\nabla \mathcal{Q}_{\lambda'}(\beta) - \nabla \mathcal{Q}_\lambda(\beta)\|_\infty}_{(ii)} \underbrace{\|\beta - \hat{\beta}_{\lambda'}\|_1}_{(iv)} \\ & \quad + \underbrace{\|\lambda' \xi - \lambda \xi\|_\infty}_{(iii)} \underbrace{\|\beta - \hat{\beta}_{\lambda'}\|_1}_{(iv)}. \end{aligned}$$

Now we provide upper bounds of terms (i)-(iv) correspondingly.

**Bounding Term (i) in (D.86):** According to (D.84), we have

$$\frac{(\beta - \hat{\beta}_{\lambda'})^T}{\|\beta - \hat{\beta}_{\lambda'}\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \leq \max_{\beta' \in \Omega} \left\{ \frac{(\beta - \beta')^T}{\|\beta - \beta'\|_1} (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi) \right\} = \omega_\lambda(\beta) \leq \epsilon,$$

where the last inequality is our assumption. Therefore we obtain

$$(D.87) \quad (\nabla \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \xi)^T (\beta - \hat{\beta}_{\lambda'}) \leq \epsilon \cdot \|\beta - \hat{\beta}_{\lambda'}\|_1.$$

We will provide an upper bound of  $\|\beta - \hat{\beta}_{\lambda'}\|_1$  when we handle term (iv).



**Bounding Term (ii) in (D.86):** Recall  $\mathcal{Q}_\lambda(\boldsymbol{\beta}) = \sum_{i=1}^d q_\lambda(\beta_j)$ . We have

$$\begin{aligned} \|\nabla \mathcal{Q}_{\lambda'}(\boldsymbol{\beta}) - \nabla \mathcal{Q}_\lambda(\boldsymbol{\beta})\|_\infty &= \max_{1 \leq j \leq d} |q_{\lambda'}(\beta_j) - q_\lambda(\beta_j)| \\ (D.88) \quad &\leq \max_{1 \leq j \leq d} |\lambda' - \lambda| = \lambda - \lambda', \end{aligned}$$

where the inequality follows from regularity condition (e), the last equality is because  $\lambda \geq \lambda'$ .

**Bounding Term (iii) in (D.86):** Since  $\boldsymbol{\xi} \in \partial \|\boldsymbol{\beta}\|_1$ , we have  $\|\boldsymbol{\xi}\|_\infty \leq 1$ . Then we obtain

$$(D.89) \quad \|\lambda' \boldsymbol{\xi} - \lambda \boldsymbol{\xi}\|_\infty = |\lambda' - \lambda| \|\boldsymbol{\xi}\|_\infty \leq |\lambda - \lambda'| = \lambda - \lambda'.$$

**Bounding Term (iv) in (D.86):** Note that

$$(D.90) \quad \|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}_{\lambda'}\|_1 \leq \underbrace{\|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_1}_{(iv).a} + \underbrace{\|\hat{\boldsymbol{\beta}}_{\lambda'} - \boldsymbol{\beta}^*\|_1}_{(iv).b}.$$

For term (iv).a, since  $\boldsymbol{\beta}$  satisfies  $\|\boldsymbol{\beta}_{\bar{S}^*}\|_0 \leq \tilde{s}$ ,  $\omega_\lambda(\boldsymbol{\beta}) \leq \lambda/2$ , and  $\|\boldsymbol{\beta}\|_2 \leq R$  for logistic loss, we have that  $\boldsymbol{\beta}$  satisfies the assumptions of Lemma 5.2 in Wang et al. (2014a). Following the same way we obtain (D.37) in the proof of Lemma 5.2 in Wang et al. (2014a), we can get

$$(\lambda/2 - \lambda/8) \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{\bar{S}^*}\|_1 \leq (3\lambda/2 + \lambda/8 + \lambda) \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1,$$

which implies  $\|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{\bar{S}^*}\|_1 \leq 7 \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1$ . Hence we obtain

$$\begin{aligned} \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_1 &\leq \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{\bar{S}^*}\|_1 + \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1 \leq 8 \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_1 \\ &\leq 8\sqrt{s^*} \|(\boldsymbol{\beta} - \boldsymbol{\beta}^*)_{S^*}\|_2 \\ &\leq 8\sqrt{s^*} \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2. \end{aligned}$$

With the upper bound of  $\|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2$  in Lemma 5.2 of Wang et al. (2014a), we obtain

$$(D.91) \quad \|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_1 \leq \frac{21}{\rho_- - \zeta_-} \lambda s^*.$$

Meanwhile, for term (iv).b, note that we assume  $\hat{\boldsymbol{\beta}}_{\lambda'}$  satisfies  $\|(\hat{\boldsymbol{\beta}}_{\lambda'})_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\|\hat{\boldsymbol{\beta}}_{\lambda'}\|_2 \leq R$  for logistic loss. Since  $\hat{\boldsymbol{\beta}}_{\lambda'}$  is an exact local solution, it satisfies the exact optimality condition  $\omega(\hat{\boldsymbol{\beta}}_{\lambda'}) \leq 0$ , which gives  $\omega(\hat{\boldsymbol{\beta}}_{\lambda'}) < \lambda'/2$ . Hence  $\hat{\boldsymbol{\beta}}_{\lambda'}$  also satisfies the conditions of Lemma 5.2 of Wang et al. (2014a). Similar to (D.91), we have

$$(D.92) \quad \|\hat{\boldsymbol{\beta}}_{\lambda'} - \boldsymbol{\beta}^*\|_1 \leq \frac{21}{\rho_- - \zeta_-} \lambda' s^*.$$

Plugging (D.92) and (D.91) into (D.90), for term (iv) in (D.86), we obtain

$$(D.93) \quad \|\beta - \hat{\beta}_{\lambda'}\|_1 \leq \frac{21}{\rho_- - \zeta_-}(\lambda' + \lambda)s^*.$$

Plugging (D.87)-(D.89) and (D.93) into the right-hand side of (D.86), we obtain

$$\begin{aligned} & \phi_{\lambda'}(\beta) - \phi_{\lambda'}(\hat{\beta}_{\lambda'}) \\ & \leq \underbrace{\epsilon \cdot \frac{21}{\rho_- - \zeta_-}(\lambda' + \lambda)s^*}_{(i) \text{ in (D.86)}} + \underbrace{\left( \frac{(\lambda - \lambda')}{\rho_- - \zeta_-} + \frac{(\lambda - \lambda')}{\rho_- - \zeta_-} \right)}_{(ii) \text{ in (D.86)} \quad (iii) \text{ in (D.86)}} \cdot \underbrace{\frac{21}{\rho_- - \zeta_-}(\lambda' + \lambda)s^*}_{(iv) \text{ in (D.86)}} \\ & \leq \frac{21}{\rho_- - \zeta_-}(\epsilon + 2(\lambda - \lambda')) \cdot (\lambda' + \lambda)s^*, \end{aligned}$$

where the upper bound of term (i) in (D.86) is obtained by plugging (D.93) into the right-hand side of (D.87). Hence we conclude the proof.  $\square$

Now we are ready to prove Theorem 5.5 of Wang et al. (2014a).

**PROOF. Sparsity of  $\{\beta_t^k\}_{k=0}^\infty$  within the  $t$ -th Stage:** In the following, we provide results concerning the sparsity of the sequence  $\{\beta_t^k\}_{k=0}^\infty$  within the  $t$ -th path following stage. In the following we prove this by induction. Note that the initialization satisfies

$$(D.94) \quad \|(\beta_t^0)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \omega_{\lambda_t}(\beta_t^0) \leq \lambda_t/2, \quad \text{and} \quad L_t^0 \leq 2(\rho_+ - \zeta_+).$$

By Lemma 5.2 of Wang et al. (2014a) we have

$$(D.95) \quad \phi_{\lambda_t}(\beta_t^0) - \phi_{\lambda_t}(\beta^*) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda_t^2 s^*,$$

Suppose that, at the  $(k-1)$ -th iteration of the proximal-gradient method (Lines 5-9 of Algorithm 3 in Wang et al. (2014a)), we have

$$(D.96) \quad \|(\beta_t^{k-1})_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad L_t^{k-1} \leq 2(\rho_+ - \zeta_+), \quad \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\beta^*) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda_t^2 s^*,$$

Then by Lemma 5.4 in Wang et al. (2014a), we have that  $\beta_t^k = \mathcal{T}_{L_t^k, \lambda_t}(\beta_t^{k-1}; R)$  satisfies

$$(D.97) \quad \|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}.$$

Note that, in the setting of logistic loss, we always have  $\|\beta_t^k\|_2 \leq R$  for  $k = 0, 1, \dots$  because of the  $\ell_2$  constraint  $\Omega = B_2(R)$ . Since  $\|(\beta_t^{k-1})_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}$  imply  $\|(\beta_t^{k-1} - \beta_t^k)_{\bar{S}^*}\| \leq 2\tilde{s}$ , from Lemma 5.1 of Wang

et al. (2014a) we have

$$(D.98) \quad \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^k) \geq \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta_t^k - \beta_t^{k-1}) \\ + \frac{\rho_- - \zeta_-}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2,$$

$$(D.99) \quad \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^k) \leq \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta_t^k - \beta_t^{k-1}) \\ + \frac{\rho_+ - \zeta_+}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2.$$

Now we prove that (D.99) guarantees the line-search method in Algorithm 2 of Wang et al. (2014a) produces  $L_t^k \leq 2(\rho_+ - \zeta_+)$ . We prove by contradiction: We assume that, when the line-search method stops, it outputs  $L_t^k > 2(\rho_+ - \zeta_+)$ . Recall that we double  $L_t^k$  at each line-search iteration (Line 6 of Algorithm 2 in Wang et al. (2014a)). Then at the line-search iteration right before the line-search method stops, we have  $L_t^{k'} = L_t^k/2 > (\rho_+ - \zeta_+)$ . Recall that the objective function  $\phi_\lambda(\beta) = \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \|\beta\|_1$ . Adding  $\lambda_t \|\beta_t^k\|_1$  to the both sides of (D.99), we obtain

$$\begin{aligned} \phi_{\lambda_t}(\beta_t^k) &= \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^k) + \lambda_t \|\beta_t^k\|_1 \\ &\leq \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta_t^k - \beta_t^{k-1}) \\ &\quad + \frac{\rho_+ - \zeta_+}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2 + \lambda_t \|\beta_t^k\|_1 \\ &\leq \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta_t^k - \beta_t^{k-1}) \\ &\quad + \frac{L_t^{k'}}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2 + \lambda_t \|\beta_t^k\|_1 \\ &= \psi_{L_t^{k'}, \lambda_t}(\beta_t^k; \beta_t^{k-1}), \end{aligned}$$

where the last equality follows from (3.7) of Wang et al. (2014a). The stopping criterion of Algorithm 2 in Wang et al. (2014a) implies that the line-search method should have already stopped and give  $(L_t^k)' = L_t^k/2$ , which contradicts our assumption that the line-search method outputs  $L_t^k$ . Therefore we have

$$(D.100) \quad L_t^k \leq 2(\rho_+ - \zeta_+).$$

Moreover, according to (D.98) and (D.99), Lemma D.1 holds, i.e.,

$$(D.101) \quad \phi_{\lambda_t}(\beta_t^k) \leq \phi_{\lambda_t}(\beta_t^{k-1}) - \frac{L_t^k}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2,$$

which implies

$$\begin{aligned}
 \phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\beta^*) &\leq \phi_{\lambda_t}(\beta_t^{k-1}) - \frac{L_t^k}{2} \|\beta_t^k - \beta_t^{k-1}\|_2^2 - \phi_{\lambda_t}(\beta^*) \\
 (D.102) \quad &\leq \frac{21/2}{\rho_- - \zeta_-} \lambda_t^2 s^*.
 \end{aligned}$$

According to (D.97) and (D.100)-(D.102), now we have

$$(D.103) \quad \begin{aligned} &\|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad L_t^k \leq 2(\rho_+ - \zeta_+), \quad \phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\beta^*) \leq \frac{21/2}{\rho_- - \zeta_-} \lambda_t^2 s^*. \end{aligned}$$

Combining (D.94), (D.96) and (D.103), by induction we prove that (D.103) holds for all  $k = 0, 1, \dots$  within the  $t$ -th path following stage. Furthermore, by Lemma 5.3 of Wang et al. (2014a), all  $\beta_t^k$ 's have nice statistical recovery properties, i.e.,

$$\|\beta_t^k - \beta^*\|_2 \leq \frac{15/2}{\rho_- - \zeta_-} \lambda_t \sqrt{s^*}, \quad \text{for } k = 0, 1, \dots$$

**Convergence to Unique Local Solution:** In the following, we prove that, within the  $t$ -th path following stage, the limit point of the sequence  $\{\beta_t^k\}_{k=0}^\infty$  generated by Algorithm 3 in Wang et al. (2014a) is unique and also an exact local solution. Since  $\|(\beta_t^0)_{\bar{S}^*}\| \leq \tilde{s}$ , the restricted strong convexity of  $\tilde{\mathcal{L}}_\lambda(\beta)$  in Lemma 5.1 of Wang et al. (2014a) implies that the sub-level set

$$\{\beta : \phi_{\lambda_t}(\beta) \leq \phi_{\lambda_t}(\beta_t^0), \quad \|(\beta_t^0 - \beta)_{\bar{S}^*}\| \leq 2\tilde{s}\}$$

is bounded. From (D.101) and (D.103) we have

$$\phi_{\lambda_t}(\beta_t^k) \leq \phi_{\lambda_t}(\beta_t^0) \quad \text{and} \quad \|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad \text{for } k = 1, 2, \dots$$

Thus  $\{\beta_t^k\}_{k=0}^\infty$  is bounded, which implies that  $\{\phi_{\lambda_t}(\beta_t^k)\}_{k=0}^\infty$  is also bounded. Meanwhile, (D.101) implies that  $\{\phi_{\lambda_t}(\beta_t^k)\}_{k=0}^\infty$  decreases monotonically. By the Bolzano-Weierstrass theorem, the limit point of  $\{\phi_{\lambda_t}(\beta_t^k)\}_{k=0}^\infty$  is unique, which implies

$$\lim_{k \rightarrow \infty} \left\{ \phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\beta_t^{k-1}) \right\} = 0.$$

Consequently, by (D.101) we have that, for any limit point of  $\{\beta^k\}_{k=0}^\infty$ ,

$$\lim_{k \rightarrow \infty} \left\{ \|\beta_t^k - \beta_t^{k-1}\|_2 \right\} \leq \frac{2}{L_t^k} \cdot \lim_{k \rightarrow \infty} \left\{ \phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\beta_t^{k-1}) \right\} = 0.$$

Moreover, Lemma D.2 implies

$$\lim_{k \rightarrow \infty} \left\{ \omega_{\lambda_t}(\beta_t^k) \right\} \leq \left( L_t^k + (\rho_+ - \zeta_+) \right) \cdot \lim_{k \rightarrow \infty} \left\{ \|\beta_t^k - \beta_t^{k-1}\|_2 \right\} = 0.$$

In other words, the sequence  $\{\beta_t^k\}_{k=0}^\infty$  has a convergent subsequence, which satisfies  $\lim_{k \rightarrow \infty} \{\omega_{\lambda_t}(\beta_t^k)\} \leq 0$ . Furthermore, it implies that this convergent subsequence of  $\{\beta_t^k\}_{k=0}^\infty$  converges towards an exact local solution  $\hat{\beta}_{\lambda_t}$  that satisfies the optimal condition in (3.13) of Wang et al. (2014a). By (D.103) we have  $\|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}$  ( $k = 1, 2, \dots$ ), which implies  $\|(\hat{\beta}_{\lambda_t})_{\bar{S}^*}\|_0 \leq \tilde{s}$ .

Now we prove the uniqueness of this exact local solution by contradiction. Let  $\xi \in \partial\|\hat{\beta}_{\lambda_t}\|_1$  be the subgradient that attains the minimum in

$$(D.104) \quad \omega_{\lambda_t}(\hat{\beta}_{\lambda_t}) = \min_{\xi' \in \partial\|\hat{\beta}_{\lambda_t}\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\hat{\beta}_{\lambda_t} - \beta')^T}{\|\hat{\beta}_{\lambda_t} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \xi' \right) \right\}.$$

Since  $\omega_{\lambda_t}(\hat{\beta}_{\lambda_t}) \leq 0$ , we have

$$(D.105) \quad \max_{\beta' \in \Omega} \left\{ \frac{(\hat{\beta}_{\lambda_t} - \beta')^T}{\|\hat{\beta}_{\lambda_t} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \xi \right) \right\} \leq 0.$$

We assume there exists another local solution  $\hat{\beta}'_{\lambda_t}$ , which is the limit point of another convergent subsequence of  $\{\beta_t^k\}_{k=0}^\infty$ . Since  $\|(\hat{\beta}'_{\lambda_t})_{\bar{S}^*}\|_0 \leq \tilde{s}$ , we have  $\|(\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t})_{\bar{S}^*}\| \leq 2\tilde{s}$ . In the setting of logistic loss, we have  $\|\hat{\beta}'_{\lambda_t}\|_2 \leq R$  and  $\|\hat{\beta}_{\lambda_t}\|_2 \leq R$  by the  $\ell_2$  constraint. Hence Lemma 5.1 of Wang et al. (2014a) implies

$$(D.106) \quad \begin{aligned} \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}'_{\lambda_t}) &\geq \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + (\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t})^T \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) \\ &\quad + \frac{\rho_- - \zeta_-}{2} \|\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t}\|_2^2. \end{aligned}$$

Meanwhile, the convexity of  $\ell_1$  norm implies

$$(D.107) \quad \lambda_t \|\hat{\beta}'_{\lambda_t}\|_1 \geq \lambda_t \|\hat{\beta}_{\lambda_t}\|_1 + \lambda_t (\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t})^T \xi.$$

Recall that the objective function  $\phi_\lambda(\beta) = \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \|\beta\|_1$ . Adding (D.106) and (D.107), we obtain

$$(D.108) \quad \begin{aligned} &\phi_{\lambda_t}(\hat{\beta}'_{\lambda_t}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \\ &\geq \underbrace{\left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \xi \right)^T (\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t})}_{(i)} + \frac{\rho_- - \zeta_-}{2} \|\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t}\|_2^2. \end{aligned}$$

Since (D.105) implies

$$\begin{aligned} \frac{(\hat{\beta}_{\lambda_t} - \hat{\beta}'_{\lambda_t})^T}{\|\hat{\beta}_{\lambda_t} - \hat{\beta}'_{\lambda_t}\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \xi \right) &\leq \max_{\beta' \in \Omega} \left\{ \frac{(\hat{\beta}_{\lambda_t} - \beta')^T}{\|\hat{\beta}_{\lambda_t} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \xi \right) \right\} \\ &\leq 0, \end{aligned}$$

term (i) in (D.108) is nonnegative. Hence we obtain

$$(D.109) \quad \phi_{\lambda_t}(\hat{\beta}'_{\lambda_t}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \geq \frac{\rho_- - \zeta_-}{2} \|\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t}\|_2^2.$$

Recall we already know that the limit point of  $\{\phi_{\lambda_t}(\beta_t^k)\}_{k=0}^\infty$  is unique, which implies  $\phi_{\lambda_t}(\hat{\beta}'_{\lambda_t}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) = 0$ . Then we obtain  $\|\hat{\beta}'_{\lambda_t} - \hat{\beta}_{\lambda_t}\|_2^2 = 0$ , which contradicts our assumption that  $\hat{\beta}'_{\lambda_t} \neq \hat{\beta}_{\lambda_t}$ . In other words, we prove that the sequence  $\{\beta_t^k\}_{k=0}^\infty$  converges to a unique local solution  $\hat{\beta}_{\lambda_t}$ .

**Geometric Convergence Rate of Algorithm 3 in Wang et al. (2014a):**

Now we establish the geometric rate of convergence of Algorithm 3. According to the stopping criterion of Algorithm 2 in Wang et al. (2014a), we have

$$(D.110) \quad \begin{aligned} \phi_{\lambda_t}(\beta_t^k) &\leq \psi_{L_t^k, \lambda_t}(\beta_t^k; \beta_t^{k-1}) \\ &= \min_{\beta} \left\{ \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta - \beta_t^{k-1}) \right. \\ &\quad \left. + \frac{L_t^k}{2} \|\beta - \beta_t^{k-1}\|_2^2 + \lambda_t \|\beta\|_1 \right\} \\ &\leq \min_{\substack{\beta = \alpha \hat{\beta}_{\lambda_t} + (1-\alpha) \beta_t^{k-1} \\ \alpha \in [0,1]}} \left\{ \overbrace{\tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta - \beta_t^{k-1})}^{(i)} \right. \\ &\quad \left. + \frac{L_t^k}{2} \|\beta - \beta_t^{k-1}\|_2^2 + \lambda_t \|\beta\|_1 \right\}. \end{aligned}$$

For term (i), since  $\|(\beta_t^{k-1})_{\bar{S}^*}\|_0 \leq \tilde{s}$ ,  $\|(\hat{\beta}_{\lambda_t})_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\beta = \alpha \hat{\beta}_{\lambda_t} + (1-\alpha) \beta_t^{k-1}$  with  $\alpha \in [0, 1]$ , we obtain  $\|(\beta - \beta_t^{k-1})_{\bar{S}^*}\|_0 \leq 2\tilde{s}$ . For logistic loss, since  $\|\beta_t^{k-1}\|_2 \leq R$  and  $\|\hat{\beta}_{\lambda_t}\|_2 \leq R$ , we have  $\|\beta\|_2 \leq R$ , since the  $\ell_2$  ball  $B_2(R)$  is a convex set. Applying Lemma 5.1 in Wang et al. (2014a), we have

$$(D.111) \quad \begin{aligned} \tilde{\mathcal{L}}_{\lambda_t}(\beta) &\geq \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta - \beta_t^{k-1}) \\ &\quad + \frac{\rho_- - \zeta_-}{2} \|\beta - \beta_t^{k-1}\|_2^2 \\ &\geq \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1})^T (\beta - \beta_t^{k-1}), \end{aligned}$$

where the second inequality follows from (4.3) in Wang et al. (2014a). Plugging (D.111) into (D.110), we obtain

$$(D.112) \quad \phi_{\lambda_t}(\beta_t^k) \leq \min_{\substack{\beta = \alpha \hat{\beta}_{\lambda_t} + (1-\alpha) \beta_t^{k-1} \\ \alpha \in [0,1]}} \left\{ \tilde{\mathcal{L}}_{\lambda_t}(\beta) + \frac{L_t^k}{2} \|\beta - \beta_t^{k-1}\|_2^2 + \lambda_t \|\beta\|_1 \right\}.$$

Since  $\|(\beta_t^{k-1})_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\|(\hat{\beta}_{\lambda_t})_{\bar{S}^*}\|_0 \leq \tilde{s}$  imply  $\|(\hat{\beta}_{\lambda_t} - \beta_t^{k-1})_{\bar{S}^*}\|_0 \leq 2\tilde{s}$ , Lemma 5.1 in Wang et al. (2014a) implies that the strong convexity of  $\tilde{\mathcal{L}}_{\lambda_t}(\beta)$  holds for  $\hat{\beta}_{\lambda_t}$  and  $\beta_t^{k-1}$ . Hence we have

$$\begin{aligned} \tilde{\mathcal{L}}_{\lambda_t}(\beta) &= \tilde{\mathcal{L}}_{\lambda_t}(\alpha\hat{\beta}_{\lambda_t} + (1-\alpha)\beta^{k-1}) \\ (D.113) \quad &\leq \alpha\tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + (1-\alpha)\tilde{\mathcal{L}}_{\lambda_t}(\beta^{k-1}). \end{aligned}$$

Meanwhile, by the convexity of  $\ell_1$  norm we have

$$\begin{aligned} \lambda_t\|\beta\|_1 &= \lambda_t\|\alpha\hat{\beta}_{\lambda_t} + (1-\alpha)\beta^{k-1}\|_1 \\ (D.114) \quad &\leq \alpha\lambda_t\|\hat{\beta}_{\lambda_t}\|_1 + (1-\alpha)\|\beta^{k-1}\|_1. \end{aligned}$$

Plugging (D.113) and (D.114) into the right-hand side of (D.112), we obtain

$$\begin{aligned} \phi_{\lambda_t}(\beta_t^k) &\leq \min_{\alpha \in [0,1]} \left\{ \alpha \left( \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t\|\hat{\beta}_{\lambda_t}\|_1 \right) \right. \\ &\quad \left. + (1-\alpha) \left( \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k-1}) + \lambda_t\|\beta_t^{k-1}\|_1 \right) \right. \\ &\quad \left. + \frac{L_t^k}{2} \|\alpha\hat{\beta}_{\lambda_t} + (1-\alpha)\beta_t^{k-1} - \beta_t^{k-1}\|_2^2 \right\} \\ &= \min_{\alpha \in [0,1]} \left\{ \alpha\phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) + (1-\alpha)\phi_{\lambda_t}(\beta_t^{k-1}) \right. \\ &\quad \left. + \frac{L_t^k}{2} \|\alpha\hat{\beta}_{\lambda_t} + (1-\alpha)\beta_t^{k-1} - \beta_t^{k-1}\|_2^2 \right\} \\ (D.115) \quad &\leq \min_{\alpha \in [0,1]} \left\{ \phi_{\lambda_t}(\beta_t^{k-1}) - \alpha \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right) \right. \\ &\quad \left. + \frac{\alpha^2 L_t^k}{2} \underbrace{\|\beta_t^{k-1} - \hat{\beta}_{\lambda_t}\|_2^2}_{(i)} \right\}. \end{aligned}$$

For term (i), similar to (D.109), applying the exact optimality condition of  $\hat{\beta}_{\lambda_t}$  and the restricted strong convexity of  $\tilde{\mathcal{L}}_{\lambda_t}(\beta)$ , we obtain

$$\phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \geq \frac{\rho_- - \zeta_-}{2} \|\beta_t^{k-1} - \hat{\beta}_{\lambda_t}\|_2^2.$$

Plugging this into the right-hand side of (D.115), we obtain

$$\begin{aligned} \phi_{\lambda_t}(\beta_t^k) &\leq \min_{\alpha \in [0,1]} \left\{ \phi_{\lambda_t}(\beta_t^{k-1}) - \alpha \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right) \right. \\ (D.116) \quad &\quad \left. + \frac{\alpha^2 L_t^k}{2} \cdot \frac{2}{\rho_- - \zeta_-} \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right) \right\}. \end{aligned}$$

The right-hand side of (D.116) attains its minimum if  $\alpha = (\rho_- - \zeta_-)/(2L_t^k)$ .

Plugging this value of  $\alpha$  into (D.116), we obtain

$$\phi_{\lambda_t}(\beta_t^k) \leq \phi_{\lambda_t}(\beta_t^{k-1}) - \frac{\rho_- - \zeta_-}{4L_t^k} \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right),$$

which implies

$$\begin{aligned} & \phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \\ & \leq \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right) - \frac{\rho_- - \zeta_-}{4L_t^k} \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right) \\ (D.117) \quad & = \left( 1 - \frac{\rho_- - \zeta_-}{4L_t^k} \right) \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right). \end{aligned}$$

Recall that in (D.103) we have  $L_t^k \leq 2(\rho_+ - \zeta_+)$  ( $k = 0, 1, \dots$ ). Plugging in this into the right-hand side of (D.117), we obtain

$$\begin{aligned} \phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) & \leq \left( 1 - \frac{1}{8} \cdot \underbrace{\frac{\rho_- - \zeta_-}{\rho_+ - \zeta_+}}_{1/\kappa} \right) \left( \phi_{\lambda_t}(\beta_t^{k-1}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right) \\ & = \left( 1 - \frac{1}{8\kappa} \right)^2 \left( \phi_{\lambda_t}(\beta_t^{k-2}) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right) \\ & \quad \vdots \\ (D.118) \quad & = \left( 1 - \frac{1}{8\kappa} \right)^k \left( \phi_{\lambda_t}(\beta_t^0) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \right), \end{aligned}$$

where  $\kappa$  is the condition number defined in (4.5) of Wang et al. (2014a). Now we can characterize the total number of proximal-gradient steps required to obtain an approximate solution  $\tilde{\beta}_t = \beta_t^{k+1}$  that satisfies

$$(D.119) \quad \omega_{\lambda_t}(\tilde{\beta}_t) \leq \lambda_t/4 \quad (t = 1, \dots, N-1), \quad \text{or} \quad \omega_{\lambda_t}(\tilde{\beta}) \leq \epsilon_{\text{opt}} \quad (t = N).$$

From Lemma D.2, we have

$$\begin{aligned} \omega_{\lambda_t}(\beta_t^{k+1}) & \leq \left( L_t^{k+1} + (\rho_+ - \zeta_+) \right) \|\beta_t^{k+1} - \beta_t^k\|_2 \\ (D.120) \quad & = L_t^{k+1} \left( 1 + \frac{\rho_+ - \zeta_+}{L_t^{k+1}} \right) \|\beta_t^{k+1} - \beta_t^k\|_2. \end{aligned}$$

Note that the stopping criterion of the line-search method (Line 7 of Algorithm 2 in Wang et al. (2014a)) implies  $L_t^{k+1} \geq \rho_- - \zeta_-$ . Otherwise, we assume that  $L_t^{k+1} < \rho_- - \zeta_-$ . Since  $\|(\beta_t^{k+1})_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}$



imply  $\|(\beta_t^k - \beta_t^{k+1})_{\bar{S}^*}\|_0 \leq 2\tilde{s}$ , by Lemma 5.1 in Wang et al. (2014a) we have

$$\begin{aligned}
& \psi_{L_t^{k+1}, \lambda_t}(\beta_t^{k+1}; \beta_t^k) \\
&= \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^k) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^k)^T (\beta_t^{k+1} - \beta_t^k) + \frac{L_t^{k+1}}{2} \|\beta_t^{k+1} - \beta_t^k\|_2^2 \\
&\quad + \lambda_t \|\beta_t^{k+1}\|_1 \\
&< \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^k) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^k)^T (\beta_t^{k+1} - \beta_t^k) + \frac{\rho_- - \zeta_-}{2} \|\beta_t^{k+1} - \beta_t^k\|_2^2 \\
&\quad + \lambda_t \|\beta_t^{k+1}\|_1 \\
&\leq \tilde{\mathcal{L}}_{\lambda_t}(\beta_t^{k+1}) + \lambda_t \|\beta_t^{k+1}\|_1 \\
&= \phi_{\lambda_t}(\beta_t^{k+1}).
\end{aligned}$$

Here the first equality is from the definition in (3.7) of Wang et al. (2014a), the first inequality is from our assumption that  $L_t^{k+1} < \rho_- - \zeta_-$ , the second inequality follows from the restricted strong convexity by Lemma 5.1 in Wang et al. (2014a). However, this contradicts the stopping criterion  $\phi_{\lambda_t}(\beta_t^{k+1}) \leq \psi_{L_t^{k+1}, \lambda_t}(\beta_t^{k+1}; \beta_t^k)$ . Therefore we have proved  $L_t^{k+1} \geq \rho_- - \zeta_-$ . From (D.120) we have

$$\begin{aligned}
\omega_{\lambda_t}(\beta_t^{k+1}) &\leq L_t^{k+1} \left( 1 + \frac{\rho_+ - \zeta_+}{\rho_- - \zeta_-} \right) \|\beta_t^{k+1} - \beta_t^k\|_2 \\
(D.121) \quad &= L_t^{k+1} (1 + \kappa) \|\beta_t^{k+1} - \beta_t^k\|_2.
\end{aligned}$$

Moreover, by Lemma D.1 we have

$$\frac{L_t^{k+1}}{2} \|\beta_t^{k+1} - \beta_t^k\|_2^2 \leq \phi_{\lambda}(\beta_t^k) - \phi_{\lambda}(\beta_t^{k+1}).$$

Plugging this into the right-hand side of (D.121), we obtain

$$\begin{aligned}
\omega_{\lambda_t}(\beta_t^{k+1}) &\leq (1 + \kappa) L_t^{k+1} \|\beta_t^{k+1} - \beta_t^k\|_2 \\
&\leq (1 + \kappa) \sqrt{2L_t^{k+1} (\phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\beta_t^{k+1}))}.
\end{aligned}$$

According to (D.101), the sequence  $\{\phi_{\lambda_t}(\beta_t^k)\}_{k=0}^\infty$  decreases monotonically. Therefore, we have  $\phi_{\lambda_t}(\beta_t^{k+1}) \geq \phi_{\lambda_t}(\hat{\beta}_{\lambda_t})$ , which implies

$$(D.122) \quad \omega_{\lambda_t}(\beta_t^{k+1}) \leq (1 + \kappa) \sqrt{2L_t^{k+1} (\phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}))}.$$

Now we provide an upper bound of the right-hand side of (D.122). Recall that in (D.103) we have  $L_t^k \leq 2(\rho_+ - \zeta_+)$  ( $k = 0, 1, \dots$ ), and in (D.118) we have  $\phi_{\lambda_t}(\beta_t^k) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \leq (1 - 1/(8\kappa))^k (\phi_{\lambda_t}(\beta_t^0) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}))$ . Note that we assume  $\|(\beta_t^0)_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $\omega_{\lambda_t}(\beta_t^0) \leq \lambda_t/2$ . In Lemma D.3, we set

$\lambda' = \lambda = \lambda_t$  and  $\epsilon = \lambda_t/2$ , then we have

$$\phi_{\lambda_t}(\beta_t^0) - \phi_{\lambda_t}(\hat{\beta}_{\lambda_t}) \leq \frac{21}{\rho_- - \zeta_-} \lambda_t^2 s^*.$$

Plugging these into the right-hand side of (D.122), we obtain

$$\begin{aligned} \omega_{\lambda_t}(\beta_t^{k+1}) &\leq (1 + \kappa) \sqrt{4(\rho_+ - \zeta_+) \cdot \left(1 - \frac{1}{8\kappa}\right)^k \frac{21}{\rho_- - \zeta_-} \lambda_t^2 s^*} \\ &= (1 + \kappa) \sqrt{84\kappa \left(1 - \frac{1}{8\kappa}\right)^k} \cdot \lambda_t \sqrt{s^*}. \end{aligned}$$

Therefore, for  $t = 1, \dots, N-1$ , to ensure that  $\beta_t^{k+1}$  satisfies  $\omega_{\lambda_t}(\beta_t^{k+1}) \leq \lambda_t/4$ , it suffices to make  $k$  satisfy

$$(1 + \kappa) \sqrt{84\kappa \left(1 - \frac{1}{8\kappa}\right)^k} \cdot \lambda_t \sqrt{s^*} \leq \lambda_t/4,$$

which implies

$$k \geq 2 \log(8\sqrt{21} \cdot \sqrt{\kappa}(1 + \kappa) \cdot \sqrt{s^*}) / \log\left(1 - \frac{1}{8\kappa}\right).$$

Similarly, for  $t = N$ , to ensure that  $\beta_t^{k+1}$  satisfies  $\omega_{\lambda_t}(\beta_t^{k+1}) \leq \epsilon_{\text{opt}}$ ,  $k$  should satisfy

$$k \geq 2 \log(2\sqrt{21} \cdot \sqrt{\kappa}(1 + \kappa) \cdot \sqrt{s^*} \lambda_t / \epsilon_{\text{opt}}) / \log\left(1 - \frac{1}{8\kappa}\right).$$

Therefore we conclude the proof of Theorem 5.5.  $\square$

**D.7. Proof of Theorem 4.5 in Wang et al. (2014a).** First we present a useful lemma. It ensures that the approximate solution  $\tilde{\beta}_{t-1}$ , which is obtained from the  $(t-1)$ -th path following stage, is  $(\lambda_t/2)$ -suboptimal with respect to regularization parameter  $\lambda_t$ , i.e.,  $\omega_{\lambda_t}(\tilde{\beta}_{t-1}) \leq \lambda_t/2$ .

**Lemma D.4.** Let  $\tilde{\beta}_{t-1}$  ( $t = 1, \dots, N$ ) be the approximate solution obtained from the  $(t-1)$ -th path following stage (Line 8 of Algorithm 1 in Wang et al. (2014a)). If  $\omega_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) \leq \lambda_{t-1}/4$ . Under Assumption 4.1 and Assumption 4.4 in Wang et al. (2014a), we have

$$\omega_{\lambda_t}(\tilde{\beta}_{t-1}) \leq \lambda_t/2,$$

where  $\lambda_t = \eta \lambda_{t-1}$  with  $\eta \in [0.9, 1)$ .

**PROOF.** Consider the regularization parameter  $\lambda_{t-1}$ . Let  $\xi \in \partial \|\tilde{\beta}_{t-1}\|_1$

be the subgradient that attains the minimum in

(D.123)

$$\omega_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) = \min_{\xi' \in \partial \|\tilde{\beta}_{t-1}\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) + \lambda_{t-1} \xi' \right) \right\},$$

which implies

$$(D.124) \quad \omega_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) = \max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) + \lambda_{t-1} \xi \right) \right\}.$$

Now we consider regularization parameter  $\lambda_t$ . We have

$$(D.125) \quad \begin{aligned} \omega_{\lambda_t}(\tilde{\beta}_{t-1}) &= \min_{\xi' \in \partial \|\tilde{\beta}_{t-1}\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\tilde{\beta}_{t-1}) + \lambda_t \xi' \right) \right\} \\ &\leq \max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\tilde{\beta}_{t-1}) + \lambda_t \xi \right) \right\}, \end{aligned}$$

where  $\xi$  is defined as the minimizer of (D.123). Recall that  $\nabla \tilde{\mathcal{L}}_{\lambda_t}(\tilde{\beta}_{t-1}) = \nabla \mathcal{L}(\tilde{\beta}_{t-1}) + \nabla \mathcal{Q}_{\lambda_t}(\tilde{\beta}_{t-1})$ . We have

$$\begin{aligned} \nabla \tilde{\mathcal{L}}_{\lambda_t}(\tilde{\beta}_{t-1}) + \lambda_t \xi &= \left( \nabla \mathcal{L}(\tilde{\beta}_{t-1}) + \nabla \mathcal{Q}_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) + \lambda_t \xi \right) + (\lambda_{t-1} \xi - \lambda_t \xi) \\ &\quad + \left( \nabla \mathcal{Q}_{\lambda_t}(\tilde{\beta}_{t-1}) - \nabla \mathcal{Q}_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) \right). \end{aligned}$$

Plugging this into the right-hand side of (D.125), we obtain

(D.126)

$$\begin{aligned} \omega_{\lambda_t}(\tilde{\beta}_{t-1}) &\leq \underbrace{\max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) + \lambda_{t-1} \xi \right) \right\}}_{(i)} \\ &\quad + \underbrace{\max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} (\lambda_{t-1} \xi - \lambda_t \xi) \right\}}_{(ii)} \\ &\quad + \underbrace{\max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} \left( \nabla \mathcal{Q}_{\lambda_t}(\tilde{\beta}_{t-1}) - \nabla \mathcal{Q}_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) \right) \right\}}_{(iii)}. \end{aligned}$$

According to (D.124), term (i) in (D.126) is equal to  $\omega_{\lambda_{t-1}}(\tilde{\beta}_{t-1})$ , which is upper bounded by  $\lambda_{t-1}/4$  by our assumption. For term (ii) in (D.126), we

have

$$\begin{aligned} \max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} (\lambda_{t-1} \xi - \lambda_t \xi) \right\} &\leq \max_{\beta' \in \mathbb{R}^d} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} (\lambda_{t-1} \xi - \lambda_t \xi) \right\} \\ &= \|\lambda_{t-1} \xi - \lambda_t \xi\|_\infty \\ &\leq \lambda_{t-1} - \lambda_t, \end{aligned}$$

where first inequality is due to the duality between  $\ell_1$  and  $\ell_\infty$  norm, while the second inequality is due to the fact that  $\lambda_{t-1} > \lambda_t$  and  $\|\xi\|_\infty \leq 1$ , which follows from  $\xi \in \partial \|\tilde{\beta}_{t-1}\|_1$ . Similarly, for term (iii) we have

$$\begin{aligned} \max_{\beta' \in \Omega} \left\{ \frac{(\tilde{\beta}_{t-1} - \beta')^T}{\|\tilde{\beta}_{t-1} - \beta'\|_1} (\nabla \mathcal{Q}_{\lambda_t}(\tilde{\beta}_{t-1}) - \nabla \mathcal{Q}_{\lambda_{t-1}}(\tilde{\beta}_{t-1})) \right\} \\ \leq \|\nabla \mathcal{Q}_{\lambda_t}(\tilde{\beta}_{t-1}) - \nabla \mathcal{Q}_{\lambda_{t-1}}(\tilde{\beta}_{t-1})\|_\infty \\ = \max_{1 \leq j \leq d} |q'_{\lambda_t}((\tilde{\beta}_{t-1})_j) - q'_{\lambda_{t-1}}((\tilde{\beta}_{t-1})_j)| \\ \leq \lambda_{t-1} - \lambda_t, \end{aligned}$$

where the second inequality follows from regularity condition (e) in Wang et al. (2014a). Hence, from (D.126) we obtain

$$\begin{aligned} \omega_{\lambda_t}(\tilde{\beta}_{t-1}) &\leq \overbrace{\lambda_{t-1}/4}^{(i) \text{ in (D.126)}} + \overbrace{\lambda_{t-1} - \lambda_t}^{(ii) \text{ in (D.126)}} + \overbrace{\lambda_{t-1} - \lambda_t}^{(iii) \text{ in (D.126)}} \\ &\leq (1/(4\eta) + 1/\eta - 1 + 1/\eta - 1)\lambda_t \leq \lambda_t/2, \end{aligned}$$

where the last inequality is obtained by plugging in  $\eta \in [0.9, 1)$ . Hence we conclude the proof.  $\square$

Now we are ready to prove Theorem 4.5 in Wang et al. (2014a).

**PROOF. Geometric Rate of Convergence within Each Stage:** The stopping criterion of Algorithm 3 (Line 9) in Wang et al. (2014a) implies

$$\omega_{\lambda_{t-1}}(\tilde{\beta}_{t-1}) \leq \lambda_{t-1}/4, \quad \text{for } t = 1, \dots, N.$$

By Lemma D.4 we have

$$(D.127) \quad \omega_{\lambda_t}(\tilde{\beta}_{t-1}) \leq \lambda_t/2, \quad \text{for } t = 1, \dots, N.$$

Recall we initialize the  $t$ -th stage with  $\tilde{\beta}_{t-1} = \beta_t^0$  and  $L_{t-1} = L_t^0$  (Line 8 of Algorithm 1). By Theorem 5.5 in Wang et al. (2014a), as long as  $\|(\tilde{\beta}_{t-1})_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $L_{(t-1)} \leq 2(\rho_+ - \zeta_+)$ , we have

$$\|(\beta_t^k)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad L_t^k \leq 2(\rho_+ - \zeta_+), \quad \text{for } k = 1, 2, \dots,$$

which implies  $\|(\tilde{\beta}_t)_{\bar{S}^*}\|_0 \leq \tilde{s}$  and  $L_t \leq 2(\rho_+ - \zeta_+)$ . Remind that we initialize the entire path following procedure with  $\tilde{\beta}_0 = \mathbf{0}$  and  $L_0 = L_{\min} \leq 2(\rho_+ - \zeta_+)$  (Line 4 of Algorithm 1 in Wang et al. (2014a)). By induction we obtain

$$\|(\tilde{\beta}_t)_{\bar{S}^*}\|_0 \leq \tilde{s}, \quad L_t \leq 2(\rho_+ - \zeta_+), \quad \text{for } t = 1, \dots, N.$$

By setting  $\lambda = \lambda_t$  and  $\tilde{\beta} = \tilde{\beta}_t$  ( $t = 1, \dots, N$ ) in Theorem 5.5 of Wang et al. (2014a), we obtain that, within the  $t$ -th stage ( $t = 1, \dots, N - 1$ ), the total number of proximal-gradient iterations is no more than

$$2 \log(8\sqrt{21} \cdot \sqrt{\kappa}(1 + \kappa) \cdot \sqrt{s^*}) / \log\left(\frac{1}{1 - 1/(8\kappa)}\right),$$

while within the  $N$ -th stage, the total number of proximal-gradient steps is no more than

$$2 \log(2\sqrt{21} \cdot \sqrt{\kappa}(1 + \kappa) \cdot \sqrt{s^*} \lambda_{\text{tgt}} / \epsilon_{\text{opt}}) / \log\left(\frac{1}{1 - 1/(8\kappa)}\right).$$

Hence we obtain the first conclusion.

**Geometric Rate of Convergence over the Full Path:** Now we prove the second statement about the total number of proximal-gradient steps along the entire solution path. The total number of path following stages is

$$N = \log(\lambda_{\text{tgt}} / \lambda_0) / \log \eta.$$

Together with the first result, we have that the total number of proximal-gradient steps is no more than

$$(N - 1)C' \log(4C\sqrt{s^*}) + C' \log(C\sqrt{s^*} \lambda_{\text{tgt}} / \epsilon_{\text{opt}}).$$

where

$$C = 2\sqrt{21} \cdot \sqrt{\kappa}(1 + \kappa), \quad C' = 2 / \log\left(\frac{1}{1 - 1/(8\kappa)}\right).$$

**Geometric Rate of Convergence of the Objective Function Values:**

Now we prove the third statement concerning the objective function value. For  $t = 1, \dots, N - 1$ , by (D.127) we have  $\omega_{\lambda_{t+1}}(\tilde{\beta}_t) \leq \lambda_{t+1}/2$ . Setting  $\lambda' = \lambda_{\text{tgt}}$ ,  $\lambda = \lambda_{t+1}$ ,  $\beta = \tilde{\beta}_t$  and  $\epsilon = \lambda_{t+1}/2$  in Lemma D.3, we obtain

$$\phi_{\lambda_{\text{tgt}}}(\tilde{\beta}_t) - \phi_{\lambda_{\text{tgt}}}(\hat{\beta}_{\lambda_{\text{tgt}}}) \leq \frac{21}{\rho_- - \zeta_-} (\lambda_{t+1}/2 + 2(\lambda_{t+1} - \lambda_{\text{tgt}})) \cdot (\lambda_{\text{tgt}} + \lambda_{t+1}) s^*.$$

Since  $\lambda_{\text{tgt}} \leq \lambda_{t+1}$ , we have

$$\phi_{\lambda_{\text{tgt}}}(\tilde{\beta}_t) - \phi_{\lambda_{\text{tgt}}}(\hat{\beta}_{\lambda_{\text{tgt}}}) \leq \frac{21}{\rho_- - \zeta_-} (\lambda_{t+1}/2 + 2\lambda_{t+1}) \cdot 2\lambda_{t+1} s^* = \frac{105 \cdot \lambda_{t+1}^2 s^*}{\rho_- - \zeta_-}.$$

Since  $\lambda_{t+1} = \eta^{t+1}\lambda_0$ , we obtain

$$\phi_{\lambda_{\text{tgt}}}(\tilde{\beta}_t) - \phi_{\lambda_{\text{tgt}}}(\hat{\beta}_{\lambda_{\text{tgt}}}) \leq \eta^{2(t+1)} \frac{105 \cdot \lambda_0^2 s^*}{\rho_- - \zeta_-}, \quad \text{for } t = 1, \dots, N-1.$$

Similarly, for  $t = N$ , we have  $\omega_{\lambda_{\text{tgt}}}(\tilde{\beta}_N) \leq \epsilon_{\text{opt}}$ . By setting  $\lambda = \lambda' = \lambda_{\text{tgt}}$  and  $\epsilon = \epsilon_{\text{opt}}$  in Lemma D.3, we have

$$\phi_{\lambda_{\text{tgt}}}(\tilde{\beta}_t) - \phi_{\lambda_{\text{tgt}}}(\hat{\beta}_{\lambda_{\text{tgt}}}) \leq \frac{21 \cdot \lambda_{\text{tgt}} s^*}{\rho_- - \zeta_-} \epsilon_{\text{opt}}.$$

Therefore we conclude the proof of Theorem 4.5 in Wang et al. (2014a).  $\square$

#### D.8. Proof of Theorem 4.7 in Wang et al. (2014a).

PROOF. Recall  $\tilde{\beta}_t$  is the approximate local solution obtained from the  $t$ -th path following stage (Lines 8 and 12 of Algorithm 1 in Wang et al. (2014a)). Therefore, it satisfies the stopping criterion of the proximal-gradient method (Line 9 of Algorithm 3 in Wang et al. (2014a)), i.e., for  $t = 1, \dots, N-1$  we have  $\omega_{\lambda_t}(\tilde{\beta}_t) \leq \lambda_t/4 < \lambda_t/2$ , while for  $t = N$  we have  $\omega_{\lambda_t}(\tilde{\beta}_t) \leq \epsilon_{\text{opt}} \ll \lambda_{\text{tgt}}/4 < \lambda_t/2$ . Meanwhile, by (5.2) in Theorem 5.5 of Wang et al. (2014a),  $\tilde{\beta}_t$  satisfies  $\|(\tilde{\beta}_t)_{\bar{S}^*}\|_0 \leq \tilde{s}$ . For logistic loss, we further have  $\|\tilde{\beta}_t\|_2 \leq R$  due to the  $\ell_2$  constraint. Therefore Lemma 5.2 in Wang et al. (2014a) gives

$$\|\tilde{\beta}_t - \beta^*\|_2 \leq \frac{21/8}{\rho_- - \zeta_-} \lambda_t \sqrt{s^*}, \quad \text{for } t = 1, \dots, N,$$

which concludes the proof.  $\square$

#### D.9. Proof of Theorem 4.8 in Wang et al. (2014a).

PROOF. We denote the subgradients by  $\xi^* \in \partial\|\beta^*\|_1$  and  $\hat{\xi} \in \partial\|\hat{\beta}_{\lambda_t}\|_1$ . In particular, we set  $\hat{\xi}$  to be the subgradient that attains the minimum in

$$\omega_{\lambda_t}(\hat{\beta}_{\lambda_t}) = \min_{\xi' \in \partial\|\hat{\beta}_{\lambda_t}\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\hat{\beta}_{\lambda_t} - \beta')^T}{\|\hat{\beta}_{\lambda_t} - \beta'\|_1} \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \xi' \right) \right\}.$$

Recall that  $\hat{\beta}_{\lambda_t}$  satisfies the exact optimality condition that  $\omega_{\lambda_t}(\hat{\beta}_{\lambda_t}) \leq 0$ , hence we have

$$(D.128) \quad \max_{\beta' \in \Omega} \left\{ (\hat{\beta}_{\lambda_t} - \beta')^T \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \hat{\xi} \right) \right\} \leq 0.$$

Theorem 5.5 of Wang et al. (2014a) gives  $\|(\hat{\beta}_{\lambda_t})_{\bar{S}^*}\|_0 \leq \tilde{s}$ . Since

$$\|(\hat{\beta}_{\lambda_t} - \beta^*)_{\bar{S}^*}\|_0 \leq \tilde{s},$$

according to Lemma 5.1 of Wang et al. (2014a) the restricted convexity holds for  $\tilde{\mathcal{L}}_{\lambda_t}(\beta)$  at  $\beta_t$  and  $\beta^*$ , i.e.,

$$(D.129) \quad \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) \geq \tilde{\mathcal{L}}_{\lambda_t}(\beta^*) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\beta^*)^T (\hat{\beta}_{\lambda_t} - \beta^*) + \frac{\rho_- - \zeta_-}{2} \|\hat{\beta}_{\lambda_t} - \beta^*\|_2^2,$$

$$(D.130) \quad \tilde{\mathcal{L}}_{\lambda_t}(\beta^*) \geq \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t})^T (\beta^* - \hat{\beta}_{\lambda_t}) + \frac{\rho_- - \zeta_-}{2} \|\beta^* - \hat{\beta}_{\lambda_t}\|_2^2.$$

Meanwhile, by the convexity of  $\ell_1$  norm, we have

$$(D.131) \quad \lambda_t \|\hat{\beta}_{\lambda_t}\|_1 \geq \lambda_t \|\beta^*\|_1 + \lambda_t (\hat{\beta}_{\lambda_t} - \beta^*)^T \xi^*,$$

$$(D.132) \quad \lambda_t \|\beta^*\|_1 \geq \lambda_t \|\hat{\beta}_{\lambda_t}\|_1 + \lambda_t (\beta^* - \hat{\beta}_{\lambda_t})^T \hat{\xi}.$$

Recall that  $\tilde{\mathcal{L}}_{\lambda}(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_{\lambda}(\beta)$ . Adding (D.129)-(D.132), we obtain

$$(D.133) \quad 0 \geq \underbrace{\left( \nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^* \right)^T (\hat{\beta}_{\lambda_t} - \beta^*)}_{(i)} \\ + \underbrace{\left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \hat{\xi} \right)^T (\beta^* - \hat{\beta}_{\lambda_t}) + (\rho_- - \zeta_-) \|\hat{\beta}_{\lambda_t} - \beta^*\|_2^2}_{(ii)}.$$

According to (D.128) we have

$$\left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \hat{\xi} \right)^T (\hat{\beta}_{\lambda_t} - \beta^*) \leq \max_{\beta' \in \Omega} \left\{ (\hat{\beta}_{\lambda_t} - \beta')^T \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda \hat{\xi} \right) \right\} \leq 0,$$

which implies that term (ii) in (D.133) is nonnegative. Moving term (i) in (D.133) to its left-hand side, we obtain

(D.134)

$$\begin{aligned} & (\rho_- - \zeta_-) \|\hat{\beta}_{\lambda_t} - \beta^*\|_2^2 \\ & \leq \left( \nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^* \right)^T (\hat{\beta}_{\lambda_t} - \beta^*) \\ & \leq \min_{\xi^* \in \partial \|\beta^*\|_1} \left\{ \sum_{j=1}^d \left( \left| \left( \nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^* \right)_j \right| \cdot \left| (\beta^* - \hat{\beta}_{\lambda_t})_j \right| \right) \right\}. \end{aligned}$$

In the sequel, we decompose the summation on the right-hand side of (D.134) into three parts:  $j \in \bar{S}^*$ ,  $j \in S_1^*$  and  $j \in S_2^*$ , where  $S_1^* = \{j : |\beta_j| \geq \nu_t\}$  and  $S_2^* = \{j : |\beta_j| < \nu_t\}$ . Here  $\nu_t > 0$  is defined in (4.16) of Wang et al. (2014a).

- For  $j \in \bar{S}^*$ , by regularity condition (c) in Wang et al. (2014a), we have

$$\left( \nabla \mathcal{Q}_{\lambda_t}(\beta^*) \right)_j = q'_{\lambda_t}(\beta_j^*) = q'_{\lambda_t}(0) = 0, \quad \text{for } j \in \bar{S}^*.$$

By (4.1) in Assumption 4.1 of Wang et al. (2014a), we have

$$\begin{aligned} \max_{j \in \overline{S}^*} \left| (\nabla \mathcal{L}(\beta^*))_j \right| &\leq \max_{1 \leq j \leq d} \left| (\nabla \mathcal{L}(\beta^*))_j \right| = \|\nabla \mathcal{L}(\beta^*)\|_\infty \\ &\leq \lambda_{\text{tgt}}/8 \leq \lambda_t/8 < \lambda_t. \end{aligned}$$

Hence we have

$$\max_{j \in \overline{S}^*} \left| (\nabla \mathcal{L}(\beta^*) + \mathcal{Q}_{\lambda_t}(\beta^*))_j \right| \leq \lambda_t.$$

Meanwhile, since  $\xi^* \in \partial \|\beta^*\|_1$ , we have  $\lambda_t \xi_j^* \in [-\lambda_t, \lambda_t]$ . Therefore, for any  $j \in \overline{S}^*$ , we can always find a  $\xi_j^*$  such that

$$\left| (\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*))_j + \lambda_t \xi_j^* \right| = 0,$$

which implies

$$\min_{\xi^* \in \partial \|\beta^*\|_1} \left\{ \left| (\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^*)_j \right| \right\} = 0, \quad \text{for } j \in \overline{S}^*.$$

Thus we obtain

(D.135)

$$\min_{\xi^* \in \partial \|\beta^*\|_1} \left\{ \sum_{j \in \overline{S}^*} \left| (\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^*)_j \right| \cdot \left| (\beta^* - \hat{\beta}_{\lambda_t})_j \right| \right\} = 0.$$

- For  $j \in S_1^* \subseteq S^*$ , we have  $|\beta_j^*| \geq \nu_t$ . Recall that  $\mathcal{P}_\lambda(\beta) = \mathcal{Q}_\lambda(\beta) + \lambda \|\beta\|_1$ . By our assumption on  $\mathcal{P}_{\lambda_t}(\beta)$  in (4.16) of Wang et al. (2014a), we have

$$(\nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^*)_j = p'_{\lambda_t}(\beta_j^*) = 0, \quad \text{for } j \in S_1^*,$$

which implies

(D.136)

$$\begin{aligned} &\min_{\xi^* \in \partial \|\beta^*\|_1} \left\{ \sum_{j \in S_1^*} \left( \left| (\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^*)_j \right| \cdot \left| (\beta^* - \hat{\beta}_{\lambda_t})_j \right| \right) \right\} \\ &= \sum_{j \in S_1^*} \left| (\nabla \mathcal{L}(\beta^*))_j \right| \cdot \left| (\beta^* - \hat{\beta}_{\lambda_t})_j \right| \\ &\leq \|(\nabla \mathcal{L}(\beta^*))_{S_1^*}\|_2 \cdot \|\beta^* - \hat{\beta}_{\lambda_t}\|_2. \end{aligned}$$

- For  $j \in S_2^* \subseteq S^*$ , we have  $|\beta_j^*| < \nu_t$ . According to (4.1) in Assumption 4.1 of Wang et al. (2014a), we have

$$\max_{j \in S_2^*} \left| (\nabla \mathcal{L}(\beta^*))_j \right| \leq \max_{1 \leq j \leq d} \left| (\nabla \mathcal{L}(\beta^*))_j \right| = \|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda_t/8 \leq \lambda_t/8.$$



Meanwhile we have

$$\max_{j \in S_2^*} \left| (\nabla \mathcal{Q}_{\lambda_t}(\beta^*))_j \right| = \max_{j \in S_2^*} |q'_{\lambda_t}(\beta_j^*)| \leq \max_{1 \leq j \leq d} |q'_{\lambda_t}(\beta_j^*)| \leq \lambda_t,$$

where the last inequality follows from regularity condition (d) in Wang et al. (2014a). Also, since  $\xi^* \in \partial \|\beta^*\|_1$ , we have  $|\xi_j^*| \leq 1$ . Therefore we obtain that, for  $j \in S_2^*$ ,

$$\begin{aligned} & \left| (\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^*)_j \right| \\ & \leq \max_{j \in S_2^*} \left| (\nabla \mathcal{L}(\beta^*))_j \right| + \max_{j \in S_2^*} \left| (\nabla \mathcal{Q}_{\lambda_t}(\beta^*))_j \right| + \lambda_t \leq 3\lambda_t. \end{aligned}$$

which implies

$$\begin{aligned} & \min_{\xi^* \in \partial \|\beta^*\|_1} \left\{ \sum_{j \in S_2^*} \left| (\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{Q}_{\lambda_t}(\beta^*) + \lambda_t \xi^*)_j \right| \cdot \left| (\beta^* - \hat{\beta}_{\lambda_t})_j \right| \right\} \\ & \leq 3\lambda_t \sum_{j \in S_2^*} \left| (\beta^* - \hat{\beta}_{\lambda_t})_j \right| \\ & = 3\lambda_t \|(\beta^* - \hat{\beta}_{\lambda_t})_{\overline{S_2^*}}\|_1 \\ & \leq 3\lambda_t \sqrt{s^*} \|(\beta^* - \hat{\beta}_{\lambda_t})_{\overline{S_2^*}}\|_2 \\ & \leq 3\lambda_t \sqrt{s_2^*} \|\beta^* - \hat{\beta}_{\lambda_t}\|_2. \end{aligned}$$

Plugging (D.135)-(D.137) into the right-hand side of (D.134), we obtain

$$\|\hat{\beta}_{\lambda_t} - \beta^*\|_2 \leq \frac{1}{\rho_- - \zeta_-} \left( \|(\nabla \mathcal{L}(\beta^*))_{S_1^*}\|_2 + 3\lambda_t \sqrt{s_2^*} \right),$$

which concludes the proof of Theorem 4.8 in Wang et al. (2014a).  $\square$

**D.10. Proof for Lemma 4.9 and Theorem 4.10.** First, we prove Lemma 4.9 in Wang et al. (2014a), which states that the oracle estimator  $\hat{\beta}_O$  is uniquely defined and has nice statistical recovery property.

PROOF. To prove that the global minimizer of (4.19) in Wang et al. (2014a) is unique even for nonconvex loss functions, in the following we show that  $\mathcal{L}(\beta)$  is strongly convex on the sparse set  $\{\beta : \text{supp}(\beta) \subseteq S^*\}$ . We assume that  $\beta$  and  $\beta'$  satisfy  $\text{supp}(\beta) \subseteq S^*$  and  $\text{supp}(\beta') \subseteq S^*$ . By Taylor's theorem and the mean value theorem, we have

$$\begin{aligned} & \text{(D.138)} \\ & \mathcal{L}(\beta') = \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^T (\beta' - \beta) + \frac{1}{2} (\beta' - \beta)^T \nabla^2 \mathcal{L}(\gamma \beta' + (1-\gamma)\beta) (\beta' - \beta), \end{aligned}$$

where  $\gamma \in [0, 1]$ . Note that we have  $\|\beta' - \beta\|_0 = s^* < s^* + 2\tilde{s}$ . By Definition 4.2 and Definition 4.3 in Wang et al. (2014a), we have

$$\frac{(\beta' - \beta)^T}{\|\beta' - \beta\|_2} \nabla^2 \mathcal{L}(\gamma\beta + (1 - \gamma)\beta') \frac{(\beta' - \beta)}{\|\beta' - \beta\|_2} \geq \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}).$$

Plugging this into the right-hand side of (D.138), we obtain

$$(D.139) \quad \mathcal{L}(\beta') \geq \mathcal{L}(\beta) + \nabla \mathcal{L}(\beta)^T (\beta' - \beta) + \frac{\rho_-}{2} \|\beta' - \beta\|_2^2,$$

where  $\rho_- = \rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s})$  is a positive constant according to Assumption 4.4. Note that (D.139) holds for any  $\beta$  and  $\beta'$  such that  $\text{supp}(\beta) \subseteq S^*$  and  $\text{supp}(\beta') \subseteq S^*$ . Therefore,  $\mathcal{L}(\beta)$  is strongly convex on this sparse set, which implies the minimizer of (4.19) in Wang et al. (2014a) is unique.

Now we prove the statistical recovery property of the oracle estimator  $\hat{\beta}_O$  in the setting where  $\mathcal{L}(\beta)$  is least squares loss. Let  $\hat{\beta}'_O, \beta^{*'} \in \mathbb{R}^{s^*}$  be the restrictions of  $\hat{\beta}_O, \beta^* \in \mathbb{R}^d$  to  $S^*$  respectively, and  $\mathbf{X}_{S^*} \in \mathbb{R}^{n \times s^*}$  be a new matrix containing the columns of  $\mathbf{X}$ , i.e.,  $\mathbf{X}_j$ , that satisfy  $j \in S^*$ . Since  $\hat{\beta}'_O$  is the solution to the ordinary least squares problem

$$\hat{\beta}'_O = \underset{\beta' \in \mathbb{R}^{s^*}}{\text{argmin}} \frac{1}{2n} \|\mathbf{X}_{S^*} \beta' - \mathbf{y}\|_2^2,$$

it has the closed-form expression of

$$\hat{\beta}'_O = (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \mathbf{y}.$$

Here we still need to prove that  $\mathbf{X}_{S^*}^T \mathbf{X}_{S^*} \in \mathbb{R}^{s^* \times s^*}$  is invertible. Note that the smallest eigenvalue of  $\mathbf{X}_{S^*}^T \mathbf{X}_{S^*}$  is defined as

$$\Lambda_{\min}(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*}) = \inf \left\{ \mathbf{v}^T \mathbf{X}_{S^*}^T \mathbf{X}_{S^*} \mathbf{v} : \|\mathbf{v}\|_2 = 1, \mathbf{v} \in \mathbb{R}^{s^*} \right\},$$

which satisfies

$$\begin{aligned} \Lambda_{\min}(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*}) &= \inf \left\{ \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} : \|\mathbf{v}\|_2 = 1, \mathbf{v} \in \mathbb{R}^d, \text{supp}(\mathbf{v}) = S^* \right\} \\ &\geq \inf \left\{ \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} : \|\mathbf{v}\|_2 = 1, \mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_0 \leq s^* \right\} \\ &\geq \inf \left\{ \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} : \|\mathbf{v}\|_2 = 1, \mathbf{v} \in \mathbb{R}^d, \|\mathbf{v}\|_0 \leq s^* + 2\tilde{s} \right\} \\ (D.140) \quad &= n\rho_-(\nabla^2 \mathcal{L}, s^* + 2\tilde{s}) \\ &> 0. \end{aligned}$$

Here the first and second inequality are due to  $\{\mathbf{v} : \text{supp}(\mathbf{v}) = S^*\} \subseteq \{\mathbf{v} : \|\mathbf{v}\|_0 \leq s^*\} \subseteq \{\mathbf{v} : \|\mathbf{v}\|_0 \leq s^* + 2\tilde{s}\}$ , while the second equality follows from Definition 4.2 in Wang et al. (2014a), since for least squares loss  $\nabla^2 \mathcal{L}(\beta) = \mathbf{X}^T \mathbf{X}/n$ , and the last inequality follows from Assumption 4.4 in Wang et al. (2014a). Therefore the smallest eigenvalue of  $\mathbf{X}_{S^*}^T \mathbf{X}_{S^*}$  is positive,

which implies that  $\mathbf{X}_{S^*}^T \mathbf{X}_{S^*}$  is invertible.

By our assumption on  $(Y|\mathbf{X} = \mathbf{x}_i)$ , we have  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon} = \mathbf{X}_{S^*}\boldsymbol{\beta}^{*'} + \boldsymbol{\epsilon}$ , where  $\boldsymbol{\epsilon} \in \mathbb{R}^n$  is a zero mean sub-Gaussian random vector with independent entries and variance proxy  $\sigma^2$ . Therefore, we have

$$\begin{aligned} \hat{\boldsymbol{\beta}}'_O - \boldsymbol{\beta}^{*'} &= (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \mathbf{y} - \boldsymbol{\beta}^{*'} = (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T (\mathbf{X}\boldsymbol{\beta}^* + \boldsymbol{\epsilon}) - \boldsymbol{\beta}^{*'} \\ &= (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon}. \end{aligned}$$

Now we provide an upper bound of  $\|\hat{\boldsymbol{\beta}}'_O - \boldsymbol{\beta}^{*'}\|_\infty$ . Note that the  $j$ -th entry of  $(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon} \in \mathbb{R}^{s^*}$  could be denoted as  $\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon}$ . Here  $\mathbf{e}_j \in \mathbb{R}^{s^*}$  denotes a vector that is all-zero except an “1” in its  $j$ -th coordinate. Hence, for any  $j$ ,  $\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon}$  is sub-Gaussian with variance proxy  $\|\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T\|_2^2 \sigma^2$ . Therefore we have

$$\mathbb{P}(|\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon}| > t) \leq 2 \exp\left(-t^2 / (\|\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T\|_2^2 \sigma^2)\right),$$

which implies

$$\begin{aligned} &\mathbb{P}\left(\max_{j \in \{1, \dots, s^*\}} |\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon}| > t\right) \\ &\leq 2s^* \exp\left(-t^2 / \left(\max_{j \in \{1, \dots, s^*\}} \|\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T\|_2^2 \sigma^2\right)\right). \end{aligned}$$

Taking  $t = C \max_{j \in \{1, \dots, s^*\}} \|\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T\|_2 \sigma \cdot \sqrt{2 \log s^*}$  with  $C > 0$ , we have that

$$\begin{aligned} \|\hat{\boldsymbol{\beta}}'_O - \boldsymbol{\beta}^{*'}\|_\infty &= \|(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon}\|_\infty \\ &= \max_{j \in \{1, \dots, s^*\}} |\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \boldsymbol{\epsilon}| \\ (D.141) \quad &\leq C \max_{j \in \{1, \dots, s^*\}} \|\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T\|_2 \sigma \cdot \sqrt{2 \log s^*} \end{aligned}$$

holds with probability at least  $1 - 2 \exp(-C^2)/s^*$ . In other words, there exists a constant  $C > 0$  sufficiently large such that (D.141) holds with high probability. Note that, for any  $j \in \{1, \dots, d\}$

$$\begin{aligned} \|\mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T\|_2^2 &= \mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \mathbf{X}_{S^*} (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{e}_j^T \\ &= \mathbf{e}_j(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{e}_j^T \\ &\leq \Lambda_{\max}((\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1}) \\ &= 1/\Lambda_{\min}(\mathbf{X}_{S^*}^T \mathbf{X}_{S^*}) \\ &\leq 1/(n\rho_-), \end{aligned}$$

where the last inequality follows from (D.140). Plugging this into (D.141),

we obtain

$$\|\widehat{\beta}'_O - \beta^{*'}\|_\infty \leq C\sigma\sqrt{2/\rho_-} \cdot \sqrt{\frac{\log s^*}{n}}.$$

We remind that  $\widehat{\beta}'_O$  and  $\beta^{*'}$  are the restrictions of  $\widehat{\beta}_O$  and  $\beta^*$  to  $S^*$ , and  $\text{supp}(\widehat{\beta}_O) \subseteq S^*$ . Therefore we obtain

$$\|\widehat{\beta}_O - \beta^*\|_\infty \leq C\sigma\sqrt{2/\rho_-} \cdot \sqrt{\frac{\log s^*}{n}},$$

which concludes the proof.  $\square$

Now we prove Theorem 4.10 in Wang et al. (2014a).

PROOF. Let  $\widehat{\xi} \in \partial\|\widehat{\beta}_{\lambda_t}\|_1$ . We set  $\widehat{\xi}$  as the subgradient which attains the minimum in

$$\omega_{\lambda_t}(\widehat{\beta}_{\lambda_t}) = \min_{\xi' \in \partial\|\widehat{\beta}_{\lambda_t}\|_1} \max_{\beta' \in \Omega} \left\{ \frac{(\widehat{\beta}_{\lambda_t} - \beta')^T}{\|\widehat{\beta}_{\lambda_t} - \beta'\|_1} \left( \nabla \widetilde{\mathcal{L}}_{\lambda_t}(\widehat{\beta}_{\lambda_t}) + \lambda_t \xi' \right) \right\}.$$

Since  $\widehat{\beta}_{\lambda_t}$  satisfies the exact optimality condition that  $\omega_{\lambda_t}(\widehat{\beta}_{\lambda_t}) \leq 0$ , we have

$$(D.142) \quad \max_{\beta' \in \Omega} \left\{ (\widehat{\beta}_{\lambda_t} - \beta')^T \left( \nabla \widetilde{\mathcal{L}}_{\lambda_t}(\widehat{\beta}_{\lambda_t}) + \lambda_t \widehat{\xi} \right) \right\} \leq 0.$$

Now we prove that there exists some  $\xi_O \in \partial\|\widehat{\beta}_O\|_1$ , such that  $\widehat{\beta}_O$  satisfies the same exact optimality condition

$$(D.143) \quad \max_{\beta' \in \Omega} \left\{ (\widehat{\beta}_O - \beta')^T \left( \nabla \widetilde{\mathcal{L}}_{\lambda_t}(\widehat{\beta}_O) + \lambda_t \xi_O \right) \right\} \leq 0.$$

Recall that  $\widetilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$ . In (D.143), we have

$$(D.144) \quad \begin{aligned} & (\widehat{\beta}_O - \beta')^T \left( \nabla \widetilde{\mathcal{L}}_{\lambda_t}(\widehat{\beta}_O) + \lambda_t \xi_O \right) \\ &= \underbrace{\sum_{j \in S^*} (\widehat{\beta}_O - \beta')_j \left( \nabla \widetilde{\mathcal{L}}_{\lambda_t}(\widehat{\beta}_O) + \lambda_t \xi_O \right)_j}_{(i)} + \underbrace{\sum_{j \in S^*} (\widehat{\beta}_O - \beta')_j \left( \nabla \widetilde{\mathcal{L}}_{\lambda_t}(\widehat{\beta}_O) + \lambda_t \xi_O \right)_j}_{(ii)}. \end{aligned}$$

For term (i) in (D.144), according to Lemma 4.9 of Wang et al. (2014a) we have, for  $n$  sufficiently large,

$$(D.145) \quad |(\widehat{\beta}_O)_j| \geq |\beta_j^*| - \|\widehat{\beta}_O - \beta^*\|_\infty \geq 2\nu_t - \sigma\sqrt{2/\rho_-} \cdot \sqrt{\frac{\log s^*}{n}} \geq \nu_t.$$

Recall that  $\mathcal{P}_\lambda(\beta) = \mathcal{Q}_\lambda(\beta) + \lambda\|\beta\|_1$ . Hence we have

$$(D.146) \quad \left( \nabla \mathcal{Q}_{\lambda_t}(\hat{\beta}_O) + \lambda_t \xi_O \right)_j = \left( \nabla \mathcal{P}_{\lambda_t}(\hat{\beta}_O) \right)_j = p'_{\lambda_t} \left( (\hat{\beta}_O)_j \right) = 0,$$

where the last equality is from (4.16) of Wang et al. (2014a). Then we have

$$(D.147) \quad \begin{aligned} & \sum_{j \in S^*} (\hat{\beta}_O - \beta')_j \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_O) + \lambda_t \xi_O \right)_j \\ &= \sum_{j \in S^*} (\hat{\beta}_O - \beta')_j \left( \nabla \mathcal{L}(\hat{\beta}_O) + \nabla \mathcal{Q}_{\lambda_t}(\hat{\beta}_O) + \lambda_t \xi_O \right)_j \\ &= \sum_{j \in S^*} (\hat{\beta}_O - \beta')_j \left( \nabla \mathcal{L}(\hat{\beta}_O) \right)_j, \end{aligned}$$

where the first equality follows from  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{Q}_\lambda(\beta)$ , and the second follows from (D.146). We remind that  $\hat{\beta}_O$  is the global solution to the oracle minimization problem in (4.19) of Wang et al. (2014a). Hence  $\hat{\beta}_O$  satisfies the exact optimality condition of (4.19)

$$\max_{\beta' \in \Omega} \left\{ \sum_{j \in S^*} (\hat{\beta}_O - \beta')_j \left( \nabla \mathcal{L}(\hat{\beta}_O) \right)_j \right\} \leq 0.$$

Thus, taking maximum over  $\beta' \in \Omega$  on both sides of (D.147), we have that, the maximum of term (i) over  $\beta' \in \Omega$  is upper bounded by zero.

For term (ii) in (D.144), remind that  $(\hat{\beta}_O)_j = 0$  for  $j \in \bar{S}^*$ . According to regularity condition (c) we have

$$\left( \nabla \mathcal{Q}_{\lambda_t}(\hat{\beta}_O) \right)_j = 0.$$

Meanwhile, for any  $j \in \bar{S}^*$  it holds that

$$(D.148) \quad \left| \left( \nabla \mathcal{L}(\hat{\beta}_O) \right)_j \right| \leq \left\| \nabla \mathcal{L}(\hat{\beta}_O) \right\|_\infty \leq \left\| \mathbf{X}^T (\mathbf{y} - \mathbf{X} \hat{\beta}_O) / n \right\|_\infty.$$

On the right-hand side of (D.148), we have

$$(D.149) \quad \begin{aligned} \mathbf{y} - \mathbf{X} \hat{\beta}_O &= \mathbf{y} - \mathbf{X}_{S^*} (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T \mathbf{y} \\ &= \mathbf{X}_{S^*} \beta_{S^*}^* + \epsilon - \mathbf{X}_{S^*} (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T (\mathbf{X}_{S^*} \beta_{S^*}^* + \epsilon) \\ &= (\mathbf{I} - \mathbf{X}_{S^*} (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T) \epsilon. \end{aligned}$$

Note that  $\mathbf{X}_{S^*} (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T$  is a projection matrix, which further implies  $\mathbf{I} - \mathbf{X}_{S^*} (\mathbf{X}_{S^*}^T \mathbf{X}_{S^*})^{-1} \mathbf{X}_{S^*}^T$  is also a projection matrix. We define the right-hand side of (D.149) to be  $\epsilon'$ . Hence we have that  $\epsilon'$  is sub-Gaussian with variance proxy no larger than that of  $\epsilon$ , and

$$\left| \left( \nabla \mathcal{L}(\hat{\beta}_O) \right)_j \right| \leq \left\| \nabla \mathcal{L}(\hat{\beta}_O) \right\|_\infty = \left\| \mathbf{X}^T \epsilon' / n \right\|_\infty.$$

Note that

$$\nabla \mathcal{L}(\beta^*) = \mathbf{X}^T(\mathbf{y} - \mathbf{X}\beta^*)/n = \mathbf{X}^T\epsilon/n.$$

Following the same proof for  $\|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda_{\text{tgt}}/8$ , from (D.148) and (D.149) we obtain

$$\left| \left( \nabla \mathcal{L}(\hat{\beta}_O) \right)_j \right| \leq \|\nabla \mathcal{L}(\hat{\beta}_O)\|_\infty \leq \lambda_{\text{tgt}}/8.$$

Therefore, since  $\xi_O \in \partial \|\hat{\beta}_O\|_1$ , for

$$\left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_O) + \lambda_t \xi_O \right)_j = \left( \nabla \mathcal{L}(\hat{\beta}_O) + \nabla \mathcal{Q}_{\lambda_t}(\hat{\beta}_O) + \lambda_t \xi_O \right)_j = \left( \nabla \mathcal{L}(\hat{\beta}_O) + \lambda_t \xi_O \right)_j,$$

we can set  $(\xi_O)_j = -\left( \nabla \mathcal{L}(\hat{\beta}_O) / \lambda_t \right)_j$ . Then we obtain, for any  $j \in \bar{S}^*$ ,

$$\left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_O) + \lambda_t \xi_O \right)_j = 0,$$

which further implies that term (ii) in (D.144) is zero. In summary, taking maximum over  $\beta' \in \Omega$  on both sides of (D.144), we obtain (D.143).

Now we are ready to prove that  $\hat{\beta}_{\lambda_t} = \hat{\beta}_O$ . Recall that the oracle estimator satisfies  $\text{supp}(\hat{\beta}_O) \subseteq S^*$ . Meanwhile, from Theorem 5.5 of Wang et al. (2014a) we have  $\|(\hat{\beta}_{\lambda_t})_{\bar{S}^*}\|_0 \leq \tilde{s}$ . Hence, we have  $\|(\hat{\beta}_{\lambda_t} - \hat{\beta}_O)_{\bar{S}^*}\|_0 \leq \tilde{s}$ , and Lemma 5.1 of Wang et al. (2014a) yields

$$\begin{aligned} \text{(D.150)} \quad \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) &\geq \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_O) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_O)^T (\hat{\beta}_{\lambda_t} - \hat{\beta}_O) \\ &\quad + \frac{\rho_- - \zeta_-}{2} \|\hat{\beta}_{\lambda_t} - \hat{\beta}_O\|_2^2, \end{aligned}$$

$$\begin{aligned} \text{(D.151)} \quad \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_O) &\geq \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t})^T (\hat{\beta}_O - \hat{\beta}_{\lambda_t}) \\ &\quad + \frac{\rho_- - \zeta_-}{2} \|\hat{\beta}_O - \hat{\beta}_{\lambda_t}\|_2^2. \end{aligned}$$

Meanwhile, by the convexity of  $\ell_1$  norm, we have

$$\text{(D.152)} \quad \lambda_t \|\hat{\beta}_{\lambda_t}\|_1 \geq \lambda_t \|\hat{\beta}_O\|_1 + \lambda_t (\hat{\beta}_{\lambda_t} - \hat{\beta}_O)^T \xi_O,$$

$$\text{(D.153)} \quad \lambda_t \|\hat{\beta}_O\|_1 \geq \lambda_t \|\hat{\beta}_{\lambda_t}\|_1 + \lambda_t (\hat{\beta}_O - \hat{\beta}_{\lambda_t})^T \hat{\xi}.$$

Adding (D.150)-(D.153), we obtain

$$\begin{aligned} 0 &\geq \underbrace{\left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \hat{\xi} \right)^T (\hat{\beta}_O - \hat{\beta}_{\lambda_t})}_{\text{(i)}} + \underbrace{\left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_O) + \lambda_t \xi_O \right)^T (\hat{\beta}_{\lambda_t} - \hat{\beta}_O)}_{\text{(ii)}} \\ &\quad + (\rho_- - \zeta_-) \|\hat{\beta}_{\lambda_t} - \hat{\beta}_O\|_2^2. \end{aligned}$$

According to (D.142), we have

$$\begin{aligned} (\hat{\beta}_{\lambda_t} - \hat{\beta}_O)^T \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \hat{\xi} \right) &\leq \max_{\beta' \in \Omega} \left\{ (\hat{\beta}_{\lambda_t} - \beta')^T \left( \nabla \tilde{\mathcal{L}}_{\lambda_t}(\hat{\beta}_{\lambda_t}) + \lambda_t \hat{\xi} \right) \right\} \\ &\leq 0, \end{aligned}$$

which implies term (i) is nonnegative. Similarly, according to (D.143), term (ii) is also nonnegative. Hence we have

$$(\rho_- - \zeta_-) \|\hat{\beta}_{\lambda_t} - \hat{\beta}_O\|_2^2 \leq 0.$$

By (4.3) of Wang et al. (2014a) we have  $\rho_- - \zeta_- > 0$ , which implies  $\hat{\beta}_{\lambda_t} = \hat{\beta}_O$ . Thus, we conclude that  $\hat{\beta}_{\lambda_t}$  is the oracle estimator  $\hat{\beta}_O$ . Moreover, (D.145) implies that  $\min_{j \in S^*} |(\hat{\beta}_O)_j| > 0$ . Together with  $\text{supp}(\hat{\beta}_O) \subseteq S^*$ , we have  $\text{supp}(\hat{\beta}_{\lambda_t}) = \text{supp}(\hat{\beta}_O) = \text{supp}(\beta^*)$ .  $\square$

## APPENDIX E: THEORETICAL RESULTS ABOUT SEMIPARAMETRIC ELLIPTICAL DESIGN REGRESSION

In this section, we first introduce the Catoni's  $M$ -estimator of standard deviation, then we provide the detailed proofs of some necessary results on semiparametric elliptical design regression.<sup>2</sup>

**E.1. Catoni's  $M$ -Estimator of Standard Deviation.** Catoni (2012) proposed a novel estimator for the mean and standard deviation of heavy-tail distributions. Let  $\mathbf{Z} = (Z_1, \dots, Z_{d+1})$  be the elliptically distributed random vector defined in §2.2 of Wang et al. (2014a). We consider the estimator of the marginal mean  $\mathbb{E}(Z_j)$  ( $j = 1, \dots, d+1$ ). Let  $h : \mathbb{R} \rightarrow \mathbb{R}$  be a continuous strictly increasing function satisfying

$$-\log(1 - x + x^2/2) \leq h(x) \leq \log(1 + x + x^2/2).$$

For instance, we choose  $h(\cdot)$  to be

$$h(x) = \begin{cases} \log(1 + x + x^2/2), & \text{if } x \geq 0, \\ -\log(1 - x + x^2/2), & \text{otherwise.} \end{cases}$$

Let  $\delta \in (0, 1)$  be such that  $n \geq 2 \log(1/\delta)$ . We introduce

$$(E.1) \quad a_\delta = \sqrt{2 \log(1/\delta) \left/ \left( nv + \frac{2nv \log(1/\delta)}{n - 2 \log(1/\delta)} \right) \right.},$$

---

<sup>2</sup>§E.1, Lemma E.1 and Corollary E.2 come from an unpublished internal technical report. We provide them here for completeness.

where  $v$  is an upper bound of  $\text{Var}(Z_j)$  for all  $j$ . Catoni's estimator of  $\mathbb{E}(Z_j)$  is defined as  $\hat{\mu}_j = \hat{\mu}_j(n, \delta)$  such that

$$(E.2) \quad \sum_{i=1}^n h(\alpha_\delta(z_{i,j} - \hat{\mu}_j)) = 0, \quad j = 1, \dots, d+1,$$

where  $z_{i,j}$  is the  $i$ -th ( $i = 1, \dots, n$ ) realizations of  $Z_j$ . As  $h(\cdot)$  is differentiable everywhere, we can solve (E.2) with Newton's method efficiently. Similarly we can estimate  $\mathbb{E}(Z_j^2)$  with  $\hat{m}_j$  defined in a similar way. Then we obtain an estimator of the marginal standard deviation  $\sigma_j$

$$(E.3) \quad \hat{\sigma}_j = \sqrt{\hat{m}_j - \hat{\mu}_j^2}, \quad j = 1, \dots, d+1.$$

**E.2. Proof of Lemma C.3.** To establish results concerning the smallest sparse eigenvalue for  $\hat{\mathbf{K}}_{\mathbf{X}}$ , we need to prove several concentration results. The next lemma and proposition provide the concentration inequality for Catoni's estimator of marginal standard deviation, which is defined in (E.3). We first consider the estimator of variance in the following lemma.

**Lemma E.1.** Let  $\mathbf{X} = (X_1, \dots, X_d)^T$  be a random vector and  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be  $n$  independent realizations of  $\mathbf{X}$  with  $\text{Var}(X_j) = v_j$  and  $\mathbb{E}(X_j^4) \leq M$ , for  $j = 1, \dots, d$ . We assume that

$$\max_{1 \leq j \leq d} \{|\mathbb{E}(X_j)|\} \leq \mu_{\max}, \quad v_{\max} = \max_{1 \leq j \leq d} \{v_j\}.$$

For the estimator  $\hat{v}_j = \hat{m}_j - \hat{\mu}_j^2$  with  $\hat{m}_j$  and  $\hat{\mu}_j$  defined in (E.2), if  $n > 5 \log d$ , we have, with probability at least  $1 - 2d^{-3}$ ,

$$\max_{1 \leq j \leq d} \{|v_j - \hat{v}_j|\} \leq C \sqrt{\frac{\log d}{n}},$$

where  $C$  is a constant.

PROOF. For  $j \in \{1, \dots, d\}$ , we use  $\hat{m}_j$  to estimate  $\mathbb{E}(X_j^2)$ . Catoni (2012) showed that

$$\mathbb{P}\left(|\hat{m}_j - \mathbb{E}(X_j^2)| > t\right) \leq \exp\left(-\frac{nt^2}{M}\right).$$

Taking a union bound, we have

$$\mathbb{P}\left(\max_{1 \leq j \leq d} \{|\hat{m}_j - \mathbb{E}(X_j^2)|\} > t\right) \leq d \exp\left(-\frac{nt^2}{M}\right),$$

or equivalently, with probability at least  $1 - d^{-3}$ ,

$$(E.4) \quad \max_{1 \leq j \leq d} \{|\hat{m}_j - \mathbb{E}(X_j^2)|\} \leq 2\sqrt{M} \sqrt{\frac{\log d}{n}}.$$



Meanwhile, we use  $\hat{\mu}_j$  to estimate  $\mathbb{E}(X_j)$ . By similar arguments as above, we have

$$(E.5) \quad \max_{1 \leq j \leq d} \left\{ |\hat{\mu}_j - \mathbb{E}(X_j)| \right\} \leq 2\sqrt{v_{\max}} \sqrt{\frac{\log d}{n}}$$

with probability at least  $1 - d^{-3}$ .

Note that

$$\max_{1 \leq j \leq d} \left\{ \left| \hat{\mu}_j^2 - (\mathbb{E}(X_j))^2 \right| \right\} \leq \max_{1 \leq j \leq d} \left\{ |\hat{\mu}_j - \mathbb{E}(X_j)| \right\} \cdot \max_{1 \leq j \leq d} \left\{ |\hat{\mu}_j + \mathbb{E}(X_j)| \right\}.$$

Since we assume that  $\max_{1 \leq j \leq d} \{\mathbb{E}(X_j)\} \leq \mu_{\max}$ , we have

$$(E.6) \quad \max_{1 \leq j \leq d} \left\{ \left| \hat{\mu}_j^2 - (\mathbb{E}(X_j))^2 \right| \right\} \leq \left( 4\mu_{\max} + 4\sqrt{v_{\max}} \sqrt{\frac{\log d}{n}} \right) \cdot \sqrt{v_{\max}} \sqrt{\frac{\log d}{n}}$$

with probability at least  $1 - d^{-3}$ . Since  $\log d/n < 1$ , from (E.6) we have,

$$(E.7) \quad \max_{1 \leq j \leq d} \left\{ \left| \hat{\mu}_j^2 - (\mathbb{E}(X_j))^2 \right| \right\} \leq (4\mu_{\max} + 4\sqrt{v_{\max}}) \cdot \sqrt{v_{\max}} \sqrt{\frac{\log d}{n}}.$$

Combining (E.4) and (E.7), we have, with probability at least  $1 - 2d^{-3}$ ,

$$\max_{1 \leq j \leq d} \left\{ |\hat{m}_j - \hat{\mu}_j^2 - \text{Var}(X_j)| \right\} \leq C \sqrt{\frac{\log d}{n}},$$

where  $C = 2\sqrt{M} + (4\mu_{\max} + 4\sqrt{v_{\max}})\sqrt{v_{\max}}$ .  $\square$

We use  $\hat{\sigma}_j = \sqrt{\hat{v}_j}$  to estimate  $\sigma_j = \sqrt{v_j}$ . Using Lemma E.1, we derive a concentration inequality for  $\hat{\sigma}_j$  in the following corollary.

**Corollary E.2.** Let  $\sigma_j = \sqrt{v_j}$  and  $\hat{\sigma}_j = \sqrt{\hat{v}_j}$  for  $j = 1, \dots, d$ . By assuming  $\sigma_j \geq \sigma_{\min} > 0$  for all  $j = 1, \dots, d$ , we have, with probability at least  $1 - 2d^{-3}$ ,

$$\max_{1 \leq j \leq d} \left\{ |\sigma_j - \hat{\sigma}_j| \right\} \leq C \sqrt{\frac{\log d}{n}},$$

where  $C$  is a constant.

PROOF. By Lemma E.1, we have, with probability at least  $1 - 2d^{-3}$ ,

$$\max_{1 \leq j \leq d} \left\{ |v_j - \hat{v}_j| \right\} \leq C \sqrt{\frac{\log d}{n}}.$$

Since  $|v_j - \hat{v}_j| = |\sigma_j - \hat{\sigma}_j| \cdot |\sigma_j + \hat{\sigma}_j|$ , it follows that

$$\max_{1 \leq j \leq d} \left\{ |\sigma_j - \hat{\sigma}_j| \right\} \leq \frac{C}{\min_{1 \leq j \leq d} \left\{ |\sigma_j + \hat{\sigma}_j| \right\}} \sqrt{\frac{\log d}{n}} \leq \frac{C}{\sigma_{\min}} \sqrt{\frac{\log d}{n}}.$$

As we assume that  $\sigma_j > \sigma_{\min}$  for all  $j$ , we conclude the proof.  $\square$

Before we establish the sparse eigenvalue condition for  $\widehat{\mathbf{K}}_{\mathbf{X}}$ , we provide a concentration result of  $\widehat{\mathbf{R}}_{\mathbf{X}}$  in the following lemma.

**Lemma E.3** (Han and Liu (2013)). Let  $\mathbf{x}_1, \dots, \mathbf{x}_n$  be  $n$  realizations of a random vector  $\mathbf{X} \sim \text{EC}_d(0, \boldsymbol{\Sigma}_{\mathbf{X}}, \Xi)$  as in Definition A.1. We assume that the smallest eigenvalue of the generalized correlation matrix  $\boldsymbol{\Sigma}_{\mathbf{X}}^0$  is strictly positive. Under the sign sub-Gaussian condition (see Han and Liu (2013) for more details), the correlation matrix estimator  $\widehat{\mathbf{R}}_{\mathbf{X}}$  defined in (A.1) satisfies that, with probability at least  $1 - 2d^{-1} - d^{-2}$ ,

$$\sup_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{|\mathbf{v}^T (\widehat{\mathbf{R}}_{\mathbf{X}} - \boldsymbol{\Sigma}_{\mathbf{X}}^0) \mathbf{v}|}{\|\mathbf{v}\|_2^2} \right\} \leq C \sqrt{\frac{s \log d}{n}}$$

for  $s \in \{1, \dots, d\}$  and a sufficiently large  $n$ .

We now prove Lemma C.3.

PROOF. Here we denote the diagonal matrix that has  $x_1, \dots, x_d$  on its diagonal by  $\text{diag}(x_1, \dots, x_d)$ . Let  $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_d)$  and  $\widehat{\mathbf{D}} = \text{diag}(\widehat{\sigma}_1, \dots, \widehat{\sigma}_d)$ . First we consider the smallest sparse eigenvalue, which satisfies

$$\begin{aligned} \rho_-(\nabla^2 \mathcal{L}, s) &= \inf_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{\mathbf{v}^T \widehat{\mathbf{K}}_{\mathbf{X}} \mathbf{v}}{\|\mathbf{v}\|_2^2} \right\} \\ &= \inf_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{(\widehat{\mathbf{D}} \mathbf{v})^T \widehat{\mathbf{R}}_{\mathbf{X}} (\widehat{\mathbf{D}} \mathbf{v})}{\|\widehat{\mathbf{D}} \mathbf{v}\|_2^2} \cdot \frac{\|\widehat{\mathbf{D}} \mathbf{v}\|_2^2}{\|\mathbf{v}\|_2^2} \right\} \\ &\geq \inf_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{\mathbf{v}^T \widehat{\mathbf{R}}_{\mathbf{X}} \mathbf{v}}{\|\mathbf{v}\|_2^2} \right\} \cdot \min_{1 \leq j \leq d} \{\widehat{\sigma}_j\}. \end{aligned} \tag{E.8}$$

The first term on the right-hand side of (E.8) is the smallest sparse eigenvalue of  $\widehat{\mathbf{R}}_{\mathbf{X}}$ . Since we have from Lemma E.3 that, with probability at least  $1 - 2d^{-1} - d^{-2}$ ,

$$\sup_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{|\mathbf{v}^T (\widehat{\mathbf{R}}_{\mathbf{X}} - \boldsymbol{\Sigma}_{\mathbf{X}}^0) \mathbf{v}|}{\|\mathbf{v}\|_2^2} \right\} \leq C \sqrt{\frac{s \log d}{n}}.$$

Then for a sufficiently large  $n$ , we have

$$\mathbf{v}^T (\boldsymbol{\Sigma}_{\mathbf{X}}^0 - \widehat{\mathbf{R}}_{\mathbf{X}}) \mathbf{v} \leq C \sqrt{\frac{s \log d}{n}} \leq \frac{1}{2} \Lambda_{\min}(\boldsymbol{\Sigma}_{\mathbf{X}}^0), \quad \text{for } \|\mathbf{v}\|_0 \leq s.$$

Here  $\Lambda_{\min}(\boldsymbol{\Sigma}_{\mathbf{X}}^0)$  denotes the smallest eigenvalue of  $\boldsymbol{\Sigma}_{\mathbf{X}}^0$ , which is strictly

positive by assumption. Then we obtain

$$\frac{1}{2}\Lambda_{\min}(\Sigma_X^0) \leq \mathbf{v}^T \Sigma_X^0 \mathbf{v} - \frac{1}{2}\Lambda_{\min}(\Sigma_X^0) \leq \mathbf{v}^T \widehat{\mathbf{R}}_X \mathbf{v}, \quad \text{for } \|\mathbf{v}\|_0 \leq s.$$

Taking infimum over both sides, we get

$$(E.9) \quad \inf_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{\mathbf{v}^T \widehat{\mathbf{R}}_X \mathbf{v}}{\|\mathbf{v}\|_2^2} \right\} \geq \frac{1}{2}\Lambda_{\min}(\Sigma_X^0) > 0.$$

We now consider  $\min_{1 \leq j \leq d} \{\widehat{\sigma}_j\}$  in (E.8). In Corollary E.2 we prove that, with probability at least  $1 - 2d^{-3}$ ,

$$|\sigma_j - \widehat{\sigma}_j| \leq C' \sqrt{\frac{\log d}{n}}, \quad \text{for } 1 \leq j \leq d,$$

where  $C'$  is a constant. For a sufficiently large  $n$ , we have

$$\widehat{\sigma}_j \geq \frac{1}{2}\sigma_j > 0, \quad \text{for } 1 \leq j \leq d$$

with the same probability. Taking minimum over both sides, we get

$$(E.10) \quad \min_{1 \leq j \leq d} \{\widehat{\sigma}_j\} \geq \frac{1}{2} \min_{1 \leq j \leq d} \{\sigma_j\} > 0$$

with probability at least  $1 - 2d^{-2}$ . Plugging (E.9) and (E.10) into the right-hand side of (E.8), we reach the conclusion that  $\rho_-(\nabla^2 \mathcal{L}, s) > 0$  holds with probability at least  $1 - 2d^{-1} - 3d^{-2}$ .

Now we consider the largest sparse eigenvalue, which satisfies

$$(E.11) \quad \begin{aligned} \rho_+(\nabla^2 \mathcal{L}, s) &= \sup_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{\mathbf{v}^T \widehat{\mathbf{R}}_X \mathbf{v}}{\|\mathbf{v}\|_2^2} \right\} \\ &= \sup_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{(\widehat{\mathbf{D}}\mathbf{v})^T \widehat{\mathbf{R}}_X (\widehat{\mathbf{D}}\mathbf{v})}{\|\widehat{\mathbf{D}}\mathbf{v}\|_2^2} \cdot \frac{\|\widehat{\mathbf{D}}\mathbf{v}\|_2^2}{\|\mathbf{v}\|_2^2} \right\} \\ &\leq \sup_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{\mathbf{v}^T \widehat{\mathbf{R}}_X \mathbf{v}}{\|\mathbf{v}\|_2^2} \right\} \cdot \max_{1 \leq j \leq d} \{\widehat{\sigma}_j\}. \end{aligned}$$

The first term on the right-hand side of (E.11) is the largest sparse eigenvalue of  $\widehat{\mathbf{R}}_X$ . Since we have from Lemma E.3 that, with probability at least  $1 - 2d^{-1} - d^{-2}$ ,

$$\sup_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{|\mathbf{v}^T (\widehat{\mathbf{R}}_X - \Sigma_X^0) \mathbf{v}|}{\|\mathbf{v}\|_2^2} \right\} \leq C \sqrt{\frac{s \log d}{n}}.$$

Then for a sufficiently large  $n$ , we have

$$\mathbf{v}^T (\widehat{\mathbf{R}}_{\mathbf{X}} - \boldsymbol{\Sigma}_{\mathbf{X}}^0) \mathbf{v} \leq C \sqrt{\frac{s \log d}{n}} \leq \frac{1}{2} \Lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{X}}^0), \quad \text{for } \|\mathbf{v}\|_0 \leq s.$$

Here  $\Lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{X}}^0)$  denotes the largest eigenvalue of  $\boldsymbol{\Sigma}_{\mathbf{X}}^0$ . Then we obtain

$$\mathbf{v}^T \widehat{\mathbf{R}}_{\mathbf{X}} \mathbf{v} \leq \mathbf{v}^T \boldsymbol{\Sigma}_{\mathbf{X}}^0 \mathbf{v} + \frac{1}{2} \Lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{X}}^0) \leq \frac{3}{2} \Lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{X}}^0), \quad \text{for } \|\mathbf{v}\|_0 \leq s.$$

Taking supremum over both sides, we get

$$(E.12) \quad \sup_{\|\mathbf{v}\|_0 \leq s} \left\{ \frac{\mathbf{v}^T \widehat{\mathbf{R}}_{\mathbf{X}} \mathbf{v}}{\|\mathbf{v}\|_2^2} \right\} \leq \frac{1}{2} \Lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{X}}^0) < +\infty.$$

We now consider  $\max_{1 \leq j \leq d} \{\widehat{\sigma}_j\}$  in (E.11). In Corollary E.2 we prove that, with probability at least  $1 - 2d^{-3}$ ,

$$|\sigma_j - \widehat{\sigma}_j| \leq C' \sqrt{\frac{\log d}{n}}, \quad \text{for } 1 \leq j \leq d,$$

where  $C'$  is a constant. For a sufficiently large  $n$ , we have

$$\widehat{\sigma}_j \leq \frac{3}{2} \sigma_j < +\infty, \quad \text{for } 1 \leq j \leq d$$

with the same probability. Taking minimum over both sides, we get

$$(E.13) \quad \max_{1 \leq j \leq d} \{\widehat{\sigma}_j\} \leq \frac{3}{2} \max_{1 \leq j \leq d} \{\sigma_j\} < +\infty$$

with probability at least  $1 - 2d^{-2}$ . Plugging (E.12) and (E.13) into the right-hand side of (E.11), we reach the conclusion that  $\rho_+(\nabla^2 \mathcal{L}, s) < +\infty$  holds with probability at least  $1 - 2d^{-1} - 3d^{-2}$ . Thus we conclude the proof.  $\square$

### E.3. Proof of Lemma C.2.

PROOF. For semiparametric elliptical design regression, we have

$$\nabla \mathcal{L}(\boldsymbol{\beta}^*) = \widehat{\mathbf{K}}_{\mathbf{X},Y} - \widehat{\mathbf{K}}_{\mathbf{X}} \boldsymbol{\beta}^* = \widehat{\mathbf{K}}_{\mathbf{X},Y} - \boldsymbol{\Sigma}_{\mathbf{X},Y} + \boldsymbol{\Sigma}_{\mathbf{X},Y} - \widehat{\mathbf{K}}_{\mathbf{X}} \boldsymbol{\beta}^*,$$

where  $\widehat{\mathbf{K}}_{\mathbf{X}} \in \mathbb{R}^{d \times d}$  and  $\widehat{\mathbf{K}}_{\mathbf{X},Y} \in \mathbb{R}^{d \times 1}$  are the submatrices of  $\widehat{\mathbf{K}}_{\mathbf{Z}} \in \mathbb{R}^{(d+1) \times (d+1)}$  defined in (B.1). Since  $\mathbb{E}(Y|\mathbf{X} = \mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}^*$ , we have

$$\boldsymbol{\Sigma}_{\mathbf{X},Y} = \mathbb{E}(\mathbf{X}Y) = \mathbb{E}(\mathbf{X} \mathbf{X}^T \boldsymbol{\beta}^*) = \boldsymbol{\Sigma}_{\mathbf{X}} \boldsymbol{\beta}^*.$$

Hence we have

$$\begin{aligned} \|\nabla \mathcal{L}(\boldsymbol{\beta}^*)\|_{\infty} &= \|\widehat{\mathbf{K}}_{\mathbf{X},Y} - \boldsymbol{\Sigma}_{\mathbf{X},Y} + \boldsymbol{\Sigma}_{\mathbf{X}} \boldsymbol{\beta}^* - \widehat{\mathbf{K}}_{\mathbf{X}} \boldsymbol{\beta}^*\|_{\infty} \\ &\leq \|\widehat{\mathbf{K}}_{\mathbf{X},Y} - \boldsymbol{\Sigma}_{\mathbf{X},Y}\|_{\infty} + \|\boldsymbol{\Sigma}_{\mathbf{X}} \boldsymbol{\beta}^* - \widehat{\mathbf{K}}_{\mathbf{X}} \boldsymbol{\beta}^*\|_{\infty}. \end{aligned}$$

Before we upper bound the two terms on the right-hand side, we establish a concentration inequality for  $\widehat{\mathbf{K}}_{\mathbf{Z}}$ . Let  $\mathbf{D}_{\mathbf{Z}} = \text{diag}(\sigma_1, \dots, \sigma_{d+1})$  and  $\widehat{\mathbf{D}}_{\mathbf{Z}} =$

$\text{diag}(\hat{\sigma}_1, \dots, \hat{\sigma}_{d+1})$ , where  $\sigma_1, \dots, \sigma_{d+1}$  are the marginal standard deviations of  $\mathbf{Z} \in \mathbb{R}^{(d+1)} = (Y, \mathbf{X})^T$  while  $\hat{\sigma}_1, \dots, \hat{\sigma}_{d+1}$  are the corresponding Catoni's estimators defined in (E.3). We have

$$\boldsymbol{\Sigma}_{\mathbf{Z}} = \mathbf{D}_{\mathbf{Z}} \boldsymbol{\Sigma}_{\mathbf{Z}}^0 \mathbf{D}_{\mathbf{Z}}, \quad \hat{\mathbf{K}}_{\mathbf{Z}} = \hat{\mathbf{D}}_{\mathbf{Z}} \hat{\mathbf{R}}_{\mathbf{Z}} \hat{\mathbf{D}}_{\mathbf{Z}},$$

where  $\hat{\mathbf{R}}_{\mathbf{Z}}$  is the rank-based estimator of the generalized correlation matrix  $\boldsymbol{\Sigma}_{\mathbf{Z}}^0$  defined in (A.1). Han and Liu (2012) proved that, with probability at least  $1 - (d+1)^{-5/2}$ ,

$$\|\hat{\mathbf{R}}_{\mathbf{Z}} - \boldsymbol{\Sigma}_{\mathbf{Z}}^0\|_{\max} \leq C \sqrt{\frac{\log(d+1)}{n}},$$

where  $\|\mathbf{M}\|_{\max} = \max_{1 \leq i, j \leq d} \{|M_{i,j}|\}$  for  $\mathbf{M} \in \mathbb{R}^{d \times d}$ . We have

(E.14)

$$\begin{aligned} & \|\hat{\mathbf{D}}_{\mathbf{Z}} \hat{\mathbf{R}}_{\mathbf{Z}} \hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}} \boldsymbol{\Sigma}_{\mathbf{Z}}^0 \mathbf{D}_{\mathbf{Z}}\|_{\max} \\ &= \|\mathbf{D}_{\mathbf{Z}} (\hat{\mathbf{R}}_{\mathbf{Z}} - \boldsymbol{\Sigma}_{\mathbf{Z}}^0) \mathbf{D}_{\mathbf{Z}} + (\hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}}) \hat{\mathbf{R}}_{\mathbf{Z}} \mathbf{D}_{\mathbf{Z}} + \hat{\mathbf{D}}_{\mathbf{Z}} \hat{\mathbf{R}}_{\mathbf{Z}} (\hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}})\|_{\max} \\ &\leq \|\mathbf{D}_{\mathbf{Z}} (\hat{\mathbf{R}}_{\mathbf{Z}} - \boldsymbol{\Sigma}_{\mathbf{Z}}^0) \mathbf{D}_{\mathbf{Z}}\|_{\max} + \|(\hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}}) \hat{\mathbf{R}}_{\mathbf{Z}} \mathbf{D}_{\mathbf{Z}}\|_{\max} \\ &\quad + \|\hat{\mathbf{D}}_{\mathbf{Z}} \hat{\mathbf{R}}_{\mathbf{Z}} (\hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}})\|_{\max} \\ &\leq \|\mathbf{D}_{\mathbf{Z}}\|_{\max}^2 \|\hat{\mathbf{R}}_{\mathbf{Z}} - \boldsymbol{\Sigma}_{\mathbf{Z}}^0\|_{\max}^2 + \|\mathbf{D}_{\mathbf{Z}}\|_{\max} \|\hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}}\|_{\max} \\ &\quad + \|\hat{\mathbf{D}}_{\mathbf{Z}}\|_{\max} \|\hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}}\|_{\max}. \end{aligned}$$

Following similar arguments in Corollary E.2, we have

$$\|\hat{\mathbf{D}}_{\mathbf{Z}} - \mathbf{D}_{\mathbf{Z}}\|_{\max} \leq C \sqrt{\frac{\log(d+1)}{n}}, \quad \|\hat{\mathbf{D}}_{\mathbf{Z}}\|_{\max} \leq \|\mathbf{D}_{\mathbf{Z}}\|_{\max} + C \sqrt{\frac{\log(d+1)}{n}}$$

with probability at least  $1 - 2(d+1)^{-3}$ . We assume that  $\sigma_j$  ( $1 \leq j \leq d+1$ ) is upper bounded, from (E.14) we have, with probability at least  $1 - (d+1)^{-5/2} - 2(d+1)^{-3}$ ,

$$\|\boldsymbol{\Sigma}_{\mathbf{Z}} - \hat{\mathbf{K}}_{\mathbf{Z}}\|_{\max} \leq C \sqrt{\frac{\log(d+1)}{n}},$$

which implies that with the same probability,

$$\|\hat{\mathbf{K}}_{\mathbf{X}, Y} - \boldsymbol{\Sigma}_{\mathbf{X}, Y}\|_{\infty} \leq C \sqrt{\frac{\log(d+1)}{n}},$$

$$\|\boldsymbol{\Sigma}_{\mathbf{X}} \boldsymbol{\beta}^* - \hat{\mathbf{K}}_{\mathbf{X}} \boldsymbol{\beta}^*\|_{\infty} \leq \|\boldsymbol{\beta}^*\|_1 \|\boldsymbol{\Sigma}_{\mathbf{X}} - \hat{\mathbf{K}}_{\mathbf{X}}\|_{\max} \leq C \|\boldsymbol{\beta}^*\|_1 \sqrt{\frac{\log(d+1)}{n}}.$$

Then we reach the conclusion.  $\square$

## APPENDIX F: DETAILED SETTINGS OF NUMERICAL EXPERIMENTS

The detailed settings of the first numerical experiment in §7 of Wang et al. (2014a) are as follows:

- The design matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  contains  $n = 500$  independent realizations of a random vector  $\mathbf{X} \in \mathbb{R}^d$  with  $d = 2500$ , which follows a  $t$ -distribution with 5 degrees of freedom, zero mean and correlation matrix  $\Sigma_{\mathbf{X}}^0$ . We set the correlation matrix  $\Sigma_{\mathbf{X}}^0$  to be  $(\Sigma_{\mathbf{X}}^0)_{i,j} = 0.8^{|i-j|}$  ( $1 \leq i, j \leq d$ ). Meanwhile, in the  $i$ -th data sample the response  $y_i$  follows a univariate  $t$ -distribution with 5 degrees of freedom, mean  $\mathbf{x}_i^T \boldsymbol{\beta}^*$  and variance 0.01. Here  $\mathbf{x}_i^T$  is the  $i$ -th row of the design matrix  $\mathbf{X}$ , and  $\boldsymbol{\beta}^*$  is the true parameter vector specified as follows.
- For the true parameter vector  $\boldsymbol{\beta}^* \in \mathbb{R}^d$ , we set its first 100 coordinates to be independent realizations of a standard univariate Gaussian distribution (zero mean and unit variance), and other coordinates to be zero, i.e., we set  $s^* = |\text{supp}(\boldsymbol{\beta}^*)| = 100$ .
- For the sequence of regularization parameters  $\{\lambda_t\}_{t=0}^N$ , we set  $\lambda_{\text{tgt}} = 0.05$  by cross-validation. Remind  $\lambda_0 = \|\nabla \mathcal{L}(\mathbf{0})\|_{\infty} = \|\widehat{\mathbf{K}}_{\mathbf{X},Y}\|_{\infty}$ , where  $\widehat{\mathbf{K}}_{\mathbf{X},Y} \in \mathbb{R}^d$  is in (B.1). We fix the random seed in MATLAB to be 2. In this setting, we observe  $\lambda_0 = 2.8516$ . We set  $\eta = 0.9015$ , so that the total number of regularization parameters is  $N = \log(\lambda_{\text{tgt}}/\lambda_0)/\log \eta = 39$ .
- For the MCP penalty defined in (2.2) of Wang et al. (2014a), we set  $b = 1.1$ . Meanwhile, we set the optimization precision within the  $N$ -th path following stage to be  $\epsilon_{\text{opt}} = 10^{-6}$ , and  $L_{\min} = 10^{-6}$ .

The detailed settings of the second numerical experiment in §7 are as follows:

- The design matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  contains  $n = 200$  independent realizations of a random vector  $\mathbf{X} \in \mathbb{R}^d$  with  $d = 2000$ , which follows a zero mean Gaussian distribution with covariance matrix  $(\Sigma_{\mathbf{X}})_{i,j} = 0.9 \cdot \mathbb{1}(i \neq j) + \mathbb{1}(i = j)$ . Meanwhile, we set  $\mathbf{y} - \mathbf{X}\boldsymbol{\beta}^*$  to be an  $n$ -dimensional Gaussian random vector with zero mean and covariance matrix  $\mathbf{I}$ ;  $\boldsymbol{\beta}^*$  is set to be zero on its first 1990 dimensions, and takes  $+2$  and  $-2$  with equal probability on its last 10 dimensions.
- For the sequence of regularization parameters  $\{\lambda_t\}_{t=0}^N$ , we set  $\lambda_{\text{tgt}}$  by cross-validation, and  $\lambda_0 = \|\nabla \mathcal{L}(\mathbf{0})\|_{\infty} = \|\mathbf{X}^T \mathbf{y}\|_{\infty}/n$ . Other parameters are set to be the same as in the previous experiment.
- We compare with LLA (Zou and Li, 2008), the calibrated CCCP (Wang et al., 2013), SparseNet (Mazumder et al., 2011), and the multi-stage

convex relaxation method (Zhang, 2010b; Zhang et al., 2013). For SparseNet we use the same sequence of regularization parameters as in our setting. For the other procedures, we employ the `glmnet` package (Friedman et al., 2010) to compute the Lasso problem in (6.1) of Wang et al. (2014a) at each stage.

- For the multi-stage convex relaxation method, we set the maximum number of stages to be 20. For the calibrated CCCP, we set the tuning parameter  $\tau = 1/\log(n)$  as suggested by Wang et al. (2013). To be fair, each method selects its most suitable regularization parameter using cross-validation. We repeat the experiment for 1000 times.

## REFERENCES

- BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *Annals of Statistics* **37** 1705–1732.
- CANDÉS, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *Annals of Statistics* 2313–2351.
- CATONI, O. (2012). Challenging the empirical mean and empirical variance: a deviation study. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques* **48** 1148–1185.
- FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* **33** 1.
- HAN, F. and LIU, H. (2012). Transelliptical component analysis. In *Advances in Neural Information Processing Systems* 25.
- HAN, F. and LIU, H. (2013). Optimal rates of convergence of transelliptical component analysis. *arXiv preprint arXiv:1305.6916*.
- KOLTCHINSKII, V. (2009a). The Dantzig selector and sparsity oracle inequalities. *Bernoulli* **15** 799–828.
- LIU, H., HAN, F. and ZHANG, C.-H. (2012). Transelliptical graphical models. In *Advances in Neural Information Processing Systems* 25.
- LOH, P.-L. and WAINWRIGHT, M. J. (2013). Regularized  $M$ -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *arXiv preprint arXiv:1305.2436*.
- MAZUMDER, R., FRIEDMAN, J. H. and HASTIE, T. (2011). Sparsenet: Coordinate descent with nonconvex penalties. *Journal of the American Statistical Association* **106**.
- NEGAHBAN, S. N., RAVIKUMAR, P., WAINWRIGHT, M. J. and YU, B. (2012). A unified framework for high-dimensional analysis of  $m$ -estimators with decomposable regularizers. *Statistical Science* **27** 538–557.
- NESTEROV, Y. (2004). *Introductory lectures on convex optimization: A basic course*, vol. 87. Springer.
- NESTEROV, Y. (2013). Gradient methods for minimizing composite functions. *Mathematical Programming* **140** 125–161.
- WAINWRIGHT, M. (2009). Sharp thresholds for high dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming. *IEEE Transactions on Information Theory* **55** 2183–2201.
- WANG, L., KIM, Y. and LI, R. (2013). Calibrating nonconvex penalized regression in ultra-high dimension. *Annals of Statistics* **41** 2505–2536.
- WANG, Z., LIU, H. and ZHANG, T. (2014a). Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *arXiv preprint arXiv:1306.4960*.

- ZHANG, T. (2010b). Analysis of multi-stage convex relaxation for sparse regularization. *Journal of Machine Learning Research* **11** 1087-1107.
- ZHANG, C.-H. and HUANG, J. (2008). The sparsity and bias of the lasso selection in high-dimensional linear regression. *Annals of Statistics* **36** 1567-1594.
- ZHANG, T. (2009). Some sharp performance bounds for least squares regression with  $\ell_1$  regularization. *Annals of Statistics* **37** 2109-2144.
- ZHANG, T. (2013). Multi-stage convex relaxation for feature selection. *Bernoulli* **19** 2277-2293.
- ZOU, H. and LI, R. (2008). One-step sparse estimates in nonconcave penalized likelihood models. *Annals of Statistics* **36** 1509.

DEPARTMENT OF OPERATIONS RESEARCH  
AND FINANCIAL ENGINEERING  
PRINCETON UNIVERSITY  
PRINCETON, NEW JERSEY 08544  
USA  
E-MAIL: [zhaoran@princeton.edu](mailto:zhaoran@princeton.edu)  
[hanliu@princeton.edu](mailto:hanliu@princeton.edu)

DEPARTMENT OF STATISTICS  
RUTGERS UNIVERSITY  
PISCATAWAY, NEW JERSEY 08854  
USA  
E-MAIL: [tzhang@stat.rutgers.edu](mailto:tzhang@stat.rutgers.edu)