

# Regularized M-estimators with nonconvexity: Statistical and algorithmic theory for local optima

**Authored by:**

Martin J. Wainwright  
Po-Ling Loh

## **Abstract**

We establish theoretical results concerning all local optima of various regularized M-estimators, where both loss and penalty functions are allowed to be nonconvex. Our results show that as long as the loss function satisfies restricted strong convexity and the penalty function satisfies suitable regularity conditions, any local optimum of the composite objective function lies within statistical precision of the true parameter vector. Our theory covers a broad class of nonconvex objective functions, including corrected versions of the Lasso for errors-in-variables linear models; regression in generalized linear models using nonconvex regularizers such as SCAD and MCP; and graph and inverse covariance matrix estimation. On the optimization side, we show that a simple adaptation of composite gradient descent may be used to compute a global optimum up to the statistical precision  $\epsilon$  in  $\log(1/\epsilon)$  iterations, which is the fastest possible rate of any first-order method. We provide a variety of simulations to illustrate the sharpness of our theoretical predictions.

## **1 Paper Body**

We establish theoretical results concerning local optima of regularized M-estimators, where both loss and penalty functions are allowed to be nonconvex. Our results show that as long as the loss satisfies restricted strong convexity and the penalty satisfies suitable regularity conditions, any local optimum of the composite objective lies within statistical precision of the true parameter vector. Our theory covers a broad class of nonconvex objective functions, including corrected versions of the Lasso for errors-in-variables linear models and regression in generalized linear models using nonconvex regularizers such as SCAD and MCP. On the optimization side, we show that a simple adaptation of composite gradient descent may be used to compute a global optimum up to the statistical precision  $\epsilon$  in  $\log(1/\epsilon)$  iterations, the fastest possible rate for any first-order method. We provide simulations to illustrate the sharpness of our theoretical predictions.

## Introduction

Optimization of nonconvex functions is known to be computationally intractable in general [11, 12]. Unlike convex functions, nonconvex functions may possess local optima that are not global optima, and standard iterative methods such as gradient descent and coordinate descent are only guaranteed to converge to local optima. Although statistical results regarding nonconvex M-estimation often only provide guarantees about the accuracy of global optima, it is observed empirically that the local optima obtained by various estimation algorithms seem to be well-behaved. In this paper, we study the question of whether it is possible to certify “good” behavior, in both a statistical and computational sense, for various nonconvex M-estimators. On the statistical level, we provide an abstract result, applicable to a broad class of (potentially nonconvex) M-estimators, which bounds the distance between any local optimum and the unique minimum of the population risk. Although local optima of nonconvex objectives may not coincide with global optima, our theory shows that any local optimum is essentially as good as a global optimum from a statistical perspective. The class of M-estimators covered by our theory includes the modified Lasso as a special case, but our results are much stronger than those implied by previous work [6]. In addition to nonconvex loss functions, our theory also applies to nonconvex regularizers, shedding new light on a long line of recent work involving the nonconvex SCAD and MCP regularizers [3, 2, 13, 14]. Various methods have been proposed for optimizing convex loss functions with nonconvex penalties [3, 4, 15], but these methods are only guaranteed to generate local optima of the composite objective, which have not been proven to be well-behaved. In contrast, our work provides a set of regularity conditions under which all local optima are guaranteed to lie within a small ball of the population-level minimum, ensuring that standard methods such as projected and composite gradient descent [10] are sufficient for obtaining estimators that lie within statistical error of the

truth. In fact, we establish that under suitable conditions, a modified form of composite gradient descent only requires  $\log(1/\text{stat})$  iterations to obtain a solution that is accurate up to the statistical precision  $\text{stat}$ . Notation. For functions  $f(n)$  and  $g(n)$ , we write  $f(n) \sim g(n)$  to mean that  $f(n) \leq cg(n)$  for some universal constant  $c \in (0, \infty)$ , and similarly,  $f(n) \asymp g(n)$  when  $f(n) \leq c_0 g(n)$  for some universal constant  $c_0 \in (0, \infty)$ . We write  $f(n) \approx g(n)$  when  $f(n) \sim g(n)$  and  $f(n) \asymp g(n)$  hold simultaneously. For a function  $h : \mathbb{R}^p \rightarrow \mathbb{R}$ , we write  $\nabla h$  to denote a gradient or subgradient, if it exists. Finally, for  $q, r \geq 0$ , let  $B_q(r)$  denote the ‘ $q$ ’-ball of radius  $r$  centered around 0.

## Problem formulation

In this section, we develop some general theory for regularized M-estimators. We first establish notation, then discuss assumptions for nonconvex regularizers and losses studied in our paper.

## 2.1 Background

Given a collection of  $n$  samples  $Z_{1:n} = \{Z_1, \dots, Z_n\}$ , drawn from

a marginal distribution  $P$  over a space  $Z$ , consider a loss function  $L_n : \mathbb{R}^p \rightarrow \mathbb{R}$ . The value  $L_n(\beta; Z_{1n})$  serves as a measure of the "fit" between a parameter vector  $\beta \in \mathbb{R}^p$  and the observed data. This empirical loss function should be viewed as a surrogate to the population risk function  $L : \mathbb{R}^p \rightarrow \mathbb{R}$ , given by

$L(\beta) := \mathbb{E} Z L_n(\beta; Z_{1n})$ . Our goal is to estimate the parameter vector  $\beta^* := \arg \min_{\beta \in \mathbb{R}^p} L(\beta)$  that minimizes the population risk

risk, assumed to be unique. To this end, we consider a regularized  $M$ -estimator of the form  $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{L_n(\beta; Z_{1n}) + \lambda g(\beta)\}$ ,  $g(\beta) \in \mathbb{R}$

(1)

where  $g : \mathbb{R}^p \rightarrow \mathbb{R}$  is a regularizer, depending on a tuning parameter  $\lambda \geq 0$ , which serves to enforce a certain type of structure on the solution. In all cases, we consider regularizers that are  $p$ -separable across coordinates, and with a slight abuse of notation, we write  $g(\beta) = \sum_{j=1}^p g(\beta_j)$ . Our theory allows for possible nonconvexity in both the loss function  $L_n$  and the regularizer  $g$ . Due to this potential nonconvexity, our  $M$ -estimator also includes a side constraint  $g : \mathbb{R}^p \rightarrow \mathbb{R}_+$ , which we require to be a convex function satisfying the lower bound  $g(\beta) \geq k\|\beta\|_1$ , for all  $\beta \in \mathbb{R}^p$ . Consequently, any feasible point for the optimization problem (1) satisfies the constraint  $k\|\beta\|_1 \leq R$ , and as long as the empirical loss and regularizer are continuous, the Weierstrass extreme value theorem guarantees that a global minimum  $\hat{\beta}$  exists. 2.2

Nonconvex regularizers

We now state and discuss conditions on the regularizer, defined in terms of  $g : \mathbb{R} \rightarrow \mathbb{R}$ . Assumption 1. (i) The function  $g$  satisfies  $g(0) = 0$  and is symmetric around zero (i.e.,  $g(t) = g(-t)$  for all  $t \in \mathbb{R}$ ). (ii) On the nonnegative real line, the function  $g$  is nondecreasing. (iii) For  $t \geq 0$ , the function  $t \mapsto$

$g(t)/t$

is nonincreasing in  $t$ .

(iv) The function  $g$  is differentiable for all  $t \neq 0$  and subdifferentiable at  $t = 0$ , with nonzero subgradients at  $t = 0$  bounded by  $L$ . (v) There exists  $\lambda_0 > 0$  such that  $g_\lambda(t) := g(t) + \lambda t^2$  is convex. Many regularizers that are commonly used in practice satisfy Assumption 1, including the  $\ell_1$ -norm,  $g(t) = k\|t\|_1$ , and the following commonly used nonconvex regularizers: 2

SCAD penalty: This penalty, due to Fan and Li [3], takes the form  $g$  for  $-t \leq 0$ ,  $g(t) = \frac{1}{2}at^2$  for  $0 < t \leq a$ , and  $g(t) = \frac{1}{2}a^2$  for  $t > a$ , for  $a > 0$ . Assumption 1 holds with  $L = 1$  and  $\lambda_0 =$

(2)

$1/a^2$ .

MCP regularizer: This penalty, due to Zhang [13], takes the form  $g(t) = \frac{1}{2}bt^2$  for  $-t \leq 0$ ,  $g(t) = \frac{1}{2}bt^2 - \frac{1}{2}ct^3$  for  $0 < t \leq b/c$ , and  $g(t) = \frac{1}{2}b^2/c$  for  $t > b/c$ , where  $b > 0$  and  $c > 0$ . Assumption 1 holds with  $L = 1$  and  $\lambda_0 =$

(3)

where  $b > 0$  is a fixed parameter. Assumption 1 holds with  $L = 1$  and  $\lambda_0 = 1/b$ . 2.3

Nonconvex loss functions and restricted strong convexity

Throughout this paper, we require the loss function  $L_n$  to be differentiable, but we do not require it to be convex. Instead, we impose a weaker condition known as restricted strong convexity (RSC). Such conditions have been discussed in previous literature [9, 1], and involve a lower bound on the remainder in the first-order Taylor expansion of  $L_n$ . In particular, our main statistical result is based on the following RSC condition:  $\frac{1}{2} \log p \leq \frac{1}{2} k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2$ , (4a)  $\frac{1}{2} \log p \leq \frac{1}{2} k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2$ , (4b)  $\frac{1}{2} \log p \leq \frac{1}{2} k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2$  where the  $\lambda_j$ 's are strictly positive constants and the  $\lambda_j$ 's are nonnegative constants. To understand this condition, note that if  $L_n$  were actually strongly convex, then both these RSC inequalities would hold with  $\lambda_1 = \lambda_2 \leq 0$  and  $\lambda_1 = \lambda_2 = 0$ . However, in the high-dimensional setting ( $p \gg n$ ), the empirical loss  $L_n$  can never be strongly convex, but the RSC condition may still hold with strictly positive  $(\lambda_j, \lambda_j)$ . On the other hand, if  $L_n$  is convex (but not strongly convex), the left-hand expression  $q$  in inequality (4) is always  $q$  nonnegative, so inequalities (4a) and (4b) hold  $k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2 \leq \frac{1}{2} \log p$  trivially for  $k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2 \leq \frac{1}{2} \log p$  and  $k^2 \lambda_2^{-1} \|\theta\|_{k^2}^2 \leq \frac{1}{2} \log p$ , respectively. Hence, the RSC inequalities  $n \leq q \leq n$  only enforce a type of strong convexity condition over a cone set of the form  $k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2 \leq \frac{1}{2} \log p$ .

3

### Statistical guarantees and consequences

We now turn to our main statistical guarantees and some consequences for various statistical models. Our theory applies to any vector  $\theta \in \mathbb{R}^p$  that satisfies the first-order necessary conditions to be a local minimum of the program (1):  $\theta + \eta \nabla L_n(\theta) = 0$ ,  $\eta \geq 0$ ,  $\nabla L_n(\theta) = 0$  for all feasible  $\theta \in \mathbb{R}^p$ . (5) When  $\theta$  lies in the interior of the constraint set, condition (5) is the usual zero-subgradient condition. 3.1

### Main statistical results

Our main theorem is deterministic in nature, and specifies conditions on the regularizer, loss function, and parameters, which guarantee that any local optimum  $\theta$  lies close to the target vector  $\theta^* = \arg \min_{\theta} L(\theta)$ . Corresponding probabilistic results will be derived in subsequent sections. For  $\theta \in \mathbb{R}^p$

proofs and more detailed discussion of the results contained in this paper, see the technical report [7]. Theorem 1. Suppose the regularizer  $\lambda$  satisfies Assumption 1,  $L_n$  satisfies the RSC conditions (4) with  $\lambda_1 \leq \lambda_2$ , and  $\theta^*$  is feasible for the objective. Consider any choice of  $\theta$  such that  $\|\theta - \theta^*\|_{k^2}^2 \leq \frac{1}{2} \log p$   $\frac{1}{2} \log p \leq \frac{1}{2} k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2$ ,  $\frac{1}{2} \log p \leq \frac{1}{2} k^2 \lambda_2^{-1} \|\theta\|_{k^2}^2$ , (6)  $L_n \leq 6RL$

3

2

2

2

$16R \max(\lambda_1, \lambda_2) \log p$ . Then any vector  $\theta$  satisfying the first-order necessary conditions (5) satisfies the error bounds  $\|\theta - \theta^*\|_{k^2}^2 \leq \frac{1}{2} \log p$   $\frac{1}{2} \log p \leq \frac{1}{2} k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2$ , and  $k^2 \lambda_2^{-1} \|\theta\|_{k^2}^2 \leq \frac{1}{2} \log p$ , (7)  $4(\lambda_1 \lambda_2)^{-1} \log p$

where  $k = k^2 \lambda_1^{-1} \|\theta\|_{k^2}^2$ . From the bound (7), note that the squared '2'-error grows proportionally with  $k$ , the number of non-zero entries in the target  $\theta^*$

parameter, and with  $\gamma$ . As will be clarified in the following sections, choosing  $\gamma$  proportional to

$$\log p/n$$

and  $R$  proportional to

$$1/\gamma$$

will satisfy the requirements of Theorem 1 w.h.p. for

many statistical models, in which case we have a squared  $\ell_2$ -error that scales as

$$k \log p/n,$$

as expected.

**Remark 1.** It is worthwhile to discuss the quantity  $\gamma$  appearing in the denominator of the bound in Theorem 1. Recall that  $\gamma$  measures the level of curvature of the loss function  $L_n$ , while  $\kappa$  measures the level of nonconvexity of the penalty  $\lambda$ . Intuitively, the two quantities should play opposing roles in our result: Larger values of  $\gamma$  correspond to more severe nonconvexity of the penalty, resulting in worse behavior of the overall objective (1), whereas larger values of  $\gamma$  correspond to more (restricted) curvature of the loss, leading to better behavior. We now develop corollaries for various nonconvex loss functions and regularizers of interest.

### 3.2 Corrected linear regression

We begin by considering the case of high-dimensional linear regression with systematically corrupted observations. Recall that in the framework of ordinary linear regression, we have the model  $y_i = h^T x_i + \epsilon_i$ ,  $\epsilon_i \sim \mathcal{N}(0, \sigma^2)$ ,  $x_i \in \mathbb{R}^p$ ,  $i = 1, \dots, n$ ,

$$(8)$$

$$y_i = h^T x_i + \epsilon_i$$

where  $h \in \mathbb{R}^p$  is the unknown parameter vector and  $\{(x_i, y_i)\}_{i=1}^n$  are observations. Following Loh and Wainwright [6], assume we instead observe pairs  $\{(z_i, y_i)\}_{i=1}^n$ , where the  $z_i$ 's are systematically corrupted versions of the corresponding  $x_i$ 's. Some examples include the following: (a) Additive noise: Observe  $z_i = x_i + w_i$ , where  $w_i \perp x_i$ ,  $E[w_i] = 0$ , and  $\text{cov}[w_i] = \Sigma_w$ . (b) Missing data: For  $\epsilon_j \in [0, 1]$ , observe  $z_i = \epsilon_j x_i + (1 - \epsilon_j) 0$  such that for each component  $j$ , we independently observe  $z_{ij} = x_{ij}$  with probability  $1 - \epsilon_j$ , and  $z_{ij} = 0$  with probability  $\epsilon_j$ . We use the population and empirical loss functions  $L(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h^T z_i)^2$  and

$$L_n(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h^T z_i)^2$$

and

$$L_n(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h^T z_i)^2$$

$$(9)$$

$\hat{h}$  where  $(\hat{h}, \hat{b})$  are estimators for  $(h, \sigma^2)$  depending on  $\{(z_i, y_i)\}_{i=1}^n$ . Then  $\hat{h} = \arg \min_{h \in \mathbb{R}^p} L(h)$ . From the formulation (1), the corrected linear regression estimator is given by

$$\hat{h} = \arg \min_{h \in \mathbb{R}^p} \left( \frac{1}{n} \sum_{i=1}^n (y_i - h^T z_i)^2 + \lambda \|h\|_2^2 \right). \quad (10)$$

We now state a corollary in the case of additive noise (model (a)), where we take  $T = Z^T Z$  and

and

$$\begin{aligned} \mathbf{b} &= \\ \mathbf{Z}^T \mathbf{y} - \mathbf{n} \\ (11) \end{aligned}$$

$\mathbf{b}$  in equation (11) is always negative-definite, so the empirical loss function When  $p \leq n$ , the matrix  $\mathbf{b}$  are applicable to missing data (model  $L_n$  previously defined (9) is nonconvex. Other choices of  $\mathbf{b}$  (b)), and also lead to nonconvex programs [6]. 4

Corollary 1. Suppose we have i.i.d. observations  $\{(\mathbf{z}_i, y_i)\}_{i=1}^n$  from a corrupted linear model with sub-Gaussian additive noise. Suppose  $(\mathbf{z}, R)$  are chosen such that  $\mathbf{z}^T \mathbf{z}$  is feasible and  $r \log p \leq c_0 \leq c_1 \log p$ . Then given a sample size  $n \geq C \max\{R^2, k\} \log p$ , any local optimum  $\mathbf{z}^*$  of the nonconvex program (10) satisfies the estimation error bounds  $\|\mathbf{z}^* - \mathbf{z}^*\|_2 \leq k^{-1/2} \epsilon$ , and  $k^2 \epsilon^2 \leq \mathbf{z}^{*T} \mathbf{z}^* \leq k^2 \epsilon^2$  with probability at least  $1 - c_1 \exp(-c_2 \log p)$ , where  $k^2 \leq k_0 = k$ . Remark 2. When  $\mathbf{z}^T \mathbf{z} = k^2 k_1$  and  $g(\mathbf{z}) = k^2 k_1$ , taking  $\epsilon = \log n / p$  and  $R = b_0 k$  for some constant  $b_0 \leq k^2$  yields the required scaling  $n \asymp k \log p$ . Hence, the bounds in Corollary 1 agree with bounds in Theorem 1 of Loh and Wainwright [6]. Note, however, that the latter results are stated only for a global minimum  $\mathbf{z}^*$  of the program (10), whereas Corollary 1 is a much stronger e Theorem 2 of our earlier paper [6] provides an indirect result holding for any local minimum  $\mathbf{z}^*$ . route for establishing similar bounds on  $k^2 \epsilon^2 \leq \mathbf{z}^{*T} \mathbf{z}^* \leq k^2 \epsilon^2$ , since the projected gradient descent algorithm may become stuck in local minima. In contrast, our argument here does not rely on an algorithmic proof and applies to a more general class of (possibly nonconvex) penalties. Corollary 1 also has important consequences in the case where pairs  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  from the linear model (8) are observed without corruption and  $\mathbf{z}^T \mathbf{z}$  is nonconvex. Then the empirical loss  $L_n$  is equivalent to the least-squares loss, modulo a constant factor. Much existing work [3, 14] only establishes statistical consistency of global minima and then provides specialized algorithms for obtaining specific local optima that are provably close to global optima. In contrast, our results demonstrate that any optimization algorithm converging to a local optimum suffices. 3.3

### Generalized linear models

Moving beyond linear regression, we now consider the case where observations are drawn from a generalized linear model (GLM). Recall that a GLM is characterized by the conditional distribution

$y_i | \mathbf{x}_i \sim P(y_i = \eta(\mathbf{x}_i^T \boldsymbol{\beta}) + \epsilon_i)$  where  $\eta$  is a scale parameter and  $\epsilon$  is the cumulant function. By standard properties of exponential families [8, 5], we have  $\eta(\mathbf{x}_i^T \boldsymbol{\beta}) = E[y_i | \mathbf{x}_i, \boldsymbol{\beta}]$ . In our analysis, we assume there exists  $\eta_0 > 0$  such that  $\eta(t) \geq \eta_0$  for all  $t \in \mathbb{R}$ . This boundedness assumption holds in various settings, including linear regression, logistic regression, and multinomial regression. The bound is required to establish both statistical consistency results in the present section and fast global convergence guarantees for our optimization algorithms in Section 4. We will assume that  $\boldsymbol{\beta}$  is sparse and optimize the penalized maximum likelihood program  $\min_{\boldsymbol{\beta}} \sum_{i=1}^n g(\boldsymbol{\beta}^T \mathbf{x}_i - y_i) + \lambda \|\boldsymbol{\beta}\|_1$ .

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n g(\boldsymbol{\beta}^T \mathbf{x}_i - y_i) + \lambda \|\boldsymbol{\beta}\|_1$$

(12)

We then have the following corollary: Corollary 2. Suppose we have i.i.d. observations  $\{(x_i, y_i)\}_{i=1}^n$  from a GLM, where the  $x_i$ 's are sub-Gaussian. Suppose  $(\beta, R)$  are chosen such that  $\beta$  is feasible and  $R \leq C \log p / \epsilon$ . Given a sample size  $n \geq CR \log p$ , any local optimum  $\hat{\beta}$  of the nonconvex program (12) satisfies  $\|\hat{\beta} - \beta\|_2 \leq C \sqrt{k \log p / n}$ , and  $\|\hat{\beta}\|_2 \leq C \sqrt{k \log p}$ , with probability at least  $1 - \exp(-c \log p)$ , where  $k = k_0 = k$ . 5

4

Optimization algorithm

We now describe how a version of composite gradient descent may be applied to efficiently optimize the nonconvex program (1). We focus on a version of the optimization problem with the side function  $o(\beta) = \frac{1}{2} \|\beta\|_2^2$ . We may then write the program (1) as

$$\min_{\beta} L_n(\beta) + \frac{1}{2} \|\beta\|_2^2. \quad (14)$$

The objective function then decomposes nicely into a sum of a differentiable but nonconvex function and a possibly nonsmooth but convex penalty. Applied to the representation (14), the composite gradient descent procedure of Nesterov [10] produces a sequence of iterates  $\{\beta_t\}_{t=0}^T$  via the updates

2 t

$$\beta_t = \beta_{t-1} - \eta \nabla L_n(\beta_{t-1})$$

$$\beta_t = \arg \min_{\beta} L_n(\beta) + \frac{1}{2} \|\beta\|_2^2$$

$$+ \frac{1}{2} \|\beta_t - \beta_{t-1}\|_2^2, \quad \beta_{t+1} = \arg \min_{\beta} (15)$$

$$\beta_t = \arg \min_{\beta} L_n(\beta) + \frac{1}{2} \|\beta\|_2^2 \text{ where}$$

1 ?

is the stepsize. Define the Taylor error around  $\beta_2$  in the direction  $\beta_1 - \beta_2$  by  $T(\beta_1, \beta_2) = L_n(\beta_1) - L_n(\beta_2) - \langle \nabla L_n(\beta_2), \beta_1 - \beta_2 \rangle$ . For all vectors  $\beta_1, \beta_2 \in B_2(3) \cap B_1(R)$ , we require the following form of restricted strong convexity:  $\log p \leq \frac{1}{2} \|\beta_1 - \beta_2\|_2^2 + \frac{1}{2} \|\beta_1 - \beta_2\|_2^2$ , (17a)  $\|\beta_1 - \beta_2\|_2^2 \leq \frac{1}{2} \|\beta_1 - \beta_2\|_2^2 + \frac{1}{2} \|\beta_1 - \beta_2\|_2^2$ , (17b) The conditions (17) are similar but not identical to the earlier RSC conditions (4). The main difference is that we now require the Taylor difference to be bounded below uniformly over  $\beta_1, \beta_2 \in B_2(3) \cap B_1(R)$ , as opposed to for a fixed  $\beta_2 = \beta$ . We also assume an upper bound:  $\log p \leq T(\beta_1, \beta_2) \leq \frac{1}{2} \|\beta_1 - \beta_2\|_2^2 + \frac{1}{2} \|\beta_1 - \beta_2\|_2^2$ , for all  $\beta_1, \beta_2 \in B_2(3) \cap B_1(R)$ , (18) a condition referred to as restricted smoothness in past work [1]. Throughout this section, we assume  $\beta_i \in B_2(3)$  for all  $i$ , where  $\beta$  is the coefficient ensuring the convexity of the function  $g(\beta)$  from equation (13). Furthermore, we define  $\beta = \min\{\beta_1, \beta_2\}$  and  $\beta = \max\{\beta_1, \beta_2, \beta_3\}$ . The following theorem applies to any population loss function  $L$  for which the population minimizer  $\beta$  is  $k$ -sparse and  $k \leq k_0 = k$ , and under the scaling  $n \geq Ck \log p$ , for a constant  $C$  depending on the  $\beta_i$ 's and  $\beta$ 's. We show that the composite gradient updates (15) exhibit a type of globally geometric convergence in terms of the quantity

Under the stated scaling on the sample size, we are guaranteed that  $\hat{\beta} \rightarrow (0, 1)$ . Let

$b = 2 \log \frac{1}{\epsilon} \left( \frac{1}{2} \log 2 + \text{RL} \left( \frac{1}{T} \right) \right) := \frac{1}{2} + 1 + \log \log \frac{1}{\epsilon} + \log \log \frac{1}{\epsilon} \frac{1}{2}$  where  $\text{RL}(\frac{1}{T}) := \text{Ln}(\frac{1}{T}) + \frac{1}{2} \log \frac{1}{T}$ , and define  $\text{stat} := k^2 b \frac{1}{\epsilon} \frac{1}{k^2} \cdot \frac{1}{\epsilon} :=$

$1 + \frac{1}{2} \log \frac{1}{\epsilon} + \frac{1}{2} \log \frac{1}{\epsilon}$

, where  $\frac{1}{2} \log \frac{1}{\epsilon} :=$

(19)

(20)

Theorem 2. Suppose  $\text{Ln}$  satisfies the RSC/RSM conditions (17) and (18), and suppose  $\frac{1}{2}$  satisfies Assumption 1. Suppose  $\hat{\beta}$  is any global minimum of the program (14), with  $\left( \frac{1}{2} \right) \leq \log p \leq 4 \log p \leq R \leq c$ , and  $\frac{1}{2} \leq \max k \leq \text{Ln}(\frac{1}{2}) k^2, \frac{1}{2} \leq n \leq n c^2$

stat Then for any stepsize  $\frac{1}{2} \leq \max\{\frac{1}{3} \leq \frac{1}{2}, \frac{1}{2}\}$  and tolerance  $\frac{1}{2} \leq 1 \leq$ , we have

$4 \leq \frac{1}{2} \leq k \log p \leq 2 \leq b \leq 2 \leq k^2 \leq t \leq \frac{1}{2} k^2 + \frac{1}{2} 128 \leq \text{stat}, \frac{1}{2} \leq T \leq \frac{1}{2} \left( \frac{1}{2} \right) \cdot \frac{1}{2} \leq \frac{1}{2} \leq n$

6

(21)

Remark 3. Note that for the optimal choice of tolerance parameter  $\frac{1}{2}$  stat, the bound in inequality (21) takes the form  $\frac{1}{2}$ , meaning successive iterates are guaranteed to converge to a region  $b$ . Combining Theorems 1 and 2, we have within statistical accuracy of the true global optimum  $\frac{1}{2} \leq \log p \leq t \leq b, \frac{1}{2} \leq T \leq (c \leq \text{stat}) \cdot \max k^2 \leq \frac{1}{2} k^2, k^2 \leq \frac{1}{2} k^2 = O(n)$

5

Simulations

In this section, we report the results of simulations for two versions of the loss function  $\text{Ln}$ , corresponding to linear and logistic regression, and three penalty functions: Lasso, q SCAD, and MCP. In all cases, we chose regularization parameters  $R =$

1.1  $\frac{1}{2}$

$\frac{1}{2} \leq \frac{1}{2} \left( \frac{1}{2} \leq \frac{1}{2} \right)$  and  $\frac{1}{2} =$

$\log p \leq n$ .

Linear regression: In the case of linear regression, we simulated covariates corrupted by additive noise according to the mechanism described in Section 3.2, giving the estimator

$\frac{1}{2} T X^T X y^T Z \leq b \leq \arg \min \frac{1}{2} \leq \frac{1}{2} w \leq \frac{1}{2} \leq \frac{1}{2} + \frac{1}{2} \left( \frac{1}{2} \right) \cdot \left( \frac{1}{2} \right) \leq n \leq n g, \frac{1}{2} \left( \frac{1}{2} \right) \leq R$

2 We generated i.i.d. samples  $x_i \sim N(0, I)$  and  $i \sim N(0, (0.1)^2)$ , and set  $w = (0.2)^2 I$ .

Logistic regression: In the case of logistic regression, we generated i.i.d. samples  $x_i \sim N(0, I)$ .

Since  $\eta(t) = \log(1 + \exp(t))$ , the program (12) becomes

$\left( \frac{1}{2} \leq n \leq 1 \leq \left\{ \log(1 + \exp(h^T x_i)) \leq y_i h^T x_i \right\} + \frac{1}{2} \left( \frac{1}{2} \right) \cdot \left( \frac{1}{2} \right) \leq b \leq \arg \min n$

$i=1 \leq g, \frac{1}{2} \left( \frac{1}{2} \right) \leq R$  We optimized the programs (22) and (23) using the composite

gradient updates (15). Figure 1 shows the results of corrected linear regression

with Lasso, SCAD, and MCP  $\frac{1}{2}$  regularizers for three different problem sizes  $p$ .

In each case,  $\frac{1}{2}$  is a  $k$ -sparse vector with  $k = b \leq pc$ , where the nonzero entries



were generated from a normal distribution and the vector was then rescaled so  $\|b\|_2 = 1$ . As predicted by Theorem 1, the curves corresponding to the same penalty function stack up nicely when the estimation error  $\|b\|_2$  is plotted against the rescaled sample size  $n/(k \log p)$ , and the  $\ell_2$  error decreases to zero as the number of samples increases, showing that the estimators (22) and (23) are statistically consistent. We chose the parameter  $a = 3.7$  for SCAD and  $b = 3.5$  for MCP. comparing penalties for corrected linear regression 0.5

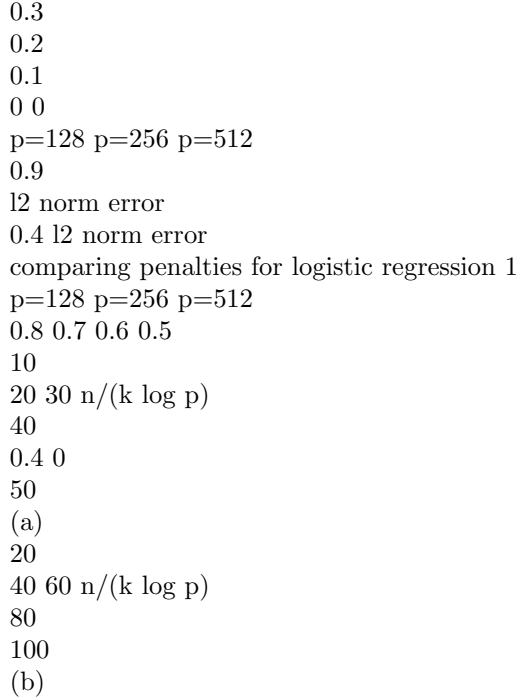


Figure 1. Plots showing statistical consistency of (a) linear and (b) logistic regression with Lasso, SCAD, and MCP. Each point represents an average over 20 trials. The estimation error  $\|b\|_2$  is plotted against the rescaled sample size  $n/(k \log p)$ . Lasso, SCAD, and MCP results are represented by solid, dotted, and dashed lines, respectively.

The simulations in Figure 2 depict the optimization-theoretic conclusions of Theorem 2. Each panel shows two different families of curves, corresponding to statistical error (red) and optimization error (blue).

The vertical axis measures the  $\ell_2$ -error on a log scale, while the horizontal axis tracks the iteration number. The curves were obtained by running composite gradient descent from 10 random starting points. We used  $p = 128$ ,  $k = bpc$ , and  $n = b20k \log pc$ . As predicted by our theory, the optimization error decreases at a linear rate until it falls to the level of statistical error. Panels (b) and (c) provide simulations for two values of the SCAD parameter  $a$ ; the larger choice  $a = 3.7$  corresponds to a higher level of curvature and produces a tighter cluster of local optima. log error plot for corrected linear regression with MCP,  $b = 1.5$  2 opt err stat err 0

log error plot for corrected linear regression with SCAD,  $a = 3.7$  2 opt err  
 stat err 0  
 log error plot for corrected linear regression with SCAD,  $a = 2.5$  1 opt err  
 stat err 0  
 ?2  
 ?2  
 ?1  
 ?8  
 —2) ?  
 ?3  
 t  
 ?4  
 log(—  
 ?6  
 log(—  
 log(—  
 t  
 t  
 ?  
 ?  
 \*  
 ?4  
 \*  
 —2)  
 \*  
 —2)  
 ?2  
 ?6  
 ?4 ?5 ?6  
 ?8  
 ?10 ?12 0  
 200  
 400 600 iteration count  
 800  
 ?10 0  
 1000  
 ?7 200  
 400  
 (a)  
 600 800 iteration count  
 1000  
 ?8 0  
 1200  
 200  
 400  
 (b)

600 800 iteration count

1000

1200

(c)

Figure 2. Plots illustrating linear rates of convergence for corrected linear regression with MCP and

SCAD. Red lines depict statistical error  $\log \|b - \hat{b}\|_2$  and blue lines depict optimization error

$\log \|b - \hat{b}^k\|_2$ . As predicted by Theorem 2, the optimization error decreases linearly up to statistical  $\log \|b - \hat{b}\|_2$  accuracy. Each plot shows the solution trajectory for 10 initializations of composite gradient descent. Panel (a) shows results for MCP; panels (b) and (c) show results for SCAD with different values of  $a$ .

Figure 3 provides analogous results to Figure 2 for logistic regression, using  $p = 64$ ,  $k = b \log p$ , and  $n = b^2 k \log p$ . The plot shows solution trajectories for 20 different initializations of composite gradient descent. Again, the log optimization error decreases at a linear rate up to the level of statistical error, as predicted by Theorem 2. Whereas the convex Lasso penalty yields a unique b SCAD and MCP produce multiple local optima. local/global optimum ?, log error plot for logistic regression with Lasso 1

log error plot for logistic regression with SCAD,  $a = 3.7 \ 0.5$

opt err stat err

0

—2)

?1

?6 500

1000 iteration count

(a)

1500

2000

?

?1.5

t

?

?5

?7 0

?1.5

t

?2

log(—

log(—

t

log(—

?4

?1

\*

—2)

0.5  
 \*  
 \*  
 2  
 — ) ?  
 3  
 opt err stat err  
 0  
 0.5  
 1 2  
 log error plot for logistic regression with MCP,  $b = 3 \ 0.5$   
 opt err stat err  
 0  
 2.5  
 2 2.5  
 3  
 3  
 3.5  
 3.5  
 4 0  
 500  
 1000 iteration count  
 (b)  
 1500  
 2000  
 4 0  
 500  
 1000 iteration count  
 1500  
 2000  
 (c)

Figure 3. Plots showing linear rates of convergence on a log scale for logistic regression. Red lines depict statistical error and blue lines depict optimization error. (a) Lasso penalty. (b) SCAD penalty. (c) MCP. Each plot shows the solution trajectory for 20 initializations of composite gradient descent.

# 6 Discussion

We have analyzed theoretical properties of local optima of regularized M-estimators, where both the loss and penalty function are allowed to be non-convex. Our results are the first to establish that all local optima of such non-convex problems are close to the truth, implying that any optimization method guaranteed to converge to a local optimum will provide statistically consistent solutions. We show that a variant of composite gradient descent may be used to obtain near-global optima in linear time, and verify our theoretical results with simulations. Acknowledgments PL acknowledges support from a Hertz Foundation Fellowship and an NSF Graduate Research Fellowship. MJW and PL were

also partially supported by grants NSF-DMS-0907632 and AFOSR09NL184. The authors thank the anonymous reviewers for helpful feedback. 8

## 2 References

- [1] A. Agarwal, S. Negahban, and M. J. Wainwright. Fast global convergence of gradient methods for high-dimensional statistical recovery. *Annals of Statistics*, 40(5):2452?2482, 2012.
- [2] P. Breheny and J. Huang. Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *Annals of Applied Statistics*, 5(1):232?253, 2011.
- [3] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96:1348?1360, 2001.
- [4] D. R. Hunter and R. Li. Variable selection using MM algorithms. *Annals of Statistics*, 33(4):1617?1642, 2005.
- [5] E.L. Lehmann and G. Casella. *Theory of Point Estimation*. Springer Verlag, 1998.
- [6] P. Loh and M.J. Wainwright. High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. *Annals of Statistics*, 40(3):1637?1664, 2012.
- [7] P. Loh and M.J. Wainwright. Regularized M -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *arXiv e-prints*, May 2013. Available at <http://arxiv.org/abs/1305.2436>.
- [8] P. McCullagh and J. A. Nelder. *Generalized Linear Models (Second Edition)*. London: Chapman & Hall, 1989.
- [9] S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for highdimensional analysis of M -estimators with decomposable regularizers. *Statistical Science*, 27(4):538?557, December 2012. See *arXiv* version for lemma/propositions cited here.
- [10] Y. Nesterov. Gradient methods for minimizing composite objective function. *CORE Discussion Papers 2007076*, Universit Catholique de Louvain, Center for Operations Research and Econometrics (CORE), 2007.
- [11] Y. Nesterov and A. Nemirovskii. *Interior Point Polynomial Algorithms in Convex Programming*. SIAM studies in applied and numerical mathematics. Society for Industrial and Applied Mathematics, 1987.
- [12] S. A. Vavasis. Complexity issues in global optimization: A survey. In *Handbook of Global Optimization*, pages 27?41. Kluwer, 1995.
- [13] C.-H. Zhang. Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics*, 38(2):894?942, 2010.
- [14] C.-H. Zhang and T. Zhang. A general theory of concave regularization for high-dimensional sparse estimation problems. *Statistical Science*, 27(4):576?593, 2012.
- [15] H. Zou and R. Li. One-step sparse estimates in nonconcave penalized likelihood models. *Annals of Statistics*, 36(4):1509?1533, 2008.