

# THE LANDSCAPE OF EMPIRICAL RISK FOR NONCONVEX LOSSES

BY SONG MEI<sup>1</sup>, YU BAI<sup>2</sup> AND ANDREA MONTANARI<sup>3</sup>

*Stanford University*

Most high-dimensional estimation methods propose to minimize a cost function (empirical risk) that is a sum of losses associated to each data point (each example). In this paper, we focus on the case of nonconvex losses. Classical empirical process theory implies uniform convergence of the empirical (or sample) risk to the population risk. While under additional assumptions, uniform convergence implies consistency of the resulting M-estimator, it does not ensure that the latter can be computed efficiently.

In order to capture the complexity of computing M-estimators, we study the landscape of the empirical risk, namely its stationary points and their properties. We establish uniform convergence of the gradient and Hessian of the empirical risk to their population counterparts, as soon as the number of samples becomes larger than the number of unknown parameters (modulo logarithmic factors). Consequently, good properties of the population risk can be carried to the empirical risk, and we are able to establish one-to-one correspondence of their stationary points. We demonstrate that in several problems such as nonconvex binary classification, robust regression and Gaussian mixture model, this result implies a complete characterization of the landscape of the empirical risk, and of the convergence properties of descent algorithms.

We extend our analysis to the very high-dimensional setting in which the number of parameters exceeds the number of samples, and provides a characterization of the empirical risk landscape under a nearly information-theoretically minimal condition. Namely, if the number of samples exceeds the sparsity of the parameters vector (modulo logarithmic factors), then a suitable uniform convergence result holds. We apply this result to nonconvex binary classification and robust regression in very high-dimension.

**1. Introduction.** M-estimation is arguably the most popular approach to high-dimensional estimation. Given data-points  $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n\}$ ,  $\mathbf{z}_i \in \mathbb{R}^d$ , we estimate a parameter vector  $\boldsymbol{\theta} \in \mathbb{R}^p$  via

$$(1.1) \quad \hat{\boldsymbol{\theta}}_n = \arg \min_{\boldsymbol{\theta} \in \Theta_{n,p}} \widehat{R}_n(\boldsymbol{\theta}),$$

---

Received January 2017; revised August 2017.

<sup>1</sup>Supported by Office of Technology Licensing Stanford Graduate Fellowship.

<sup>2</sup>Supported in part by John Duchi's NSF award CAREER-1553086.

<sup>3</sup>Supported by NSF Grant CCF-1319979.

*MSC2010 subject classifications.* Primary 62F10, 62J02; secondary 62H30.

*Key words and phrases.* Nonconvex optimization, empirical risk minimization, landscape, uniform convergence.

$$(1.2) \quad \widehat{R}_n(\boldsymbol{\theta}) \equiv \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\theta}; \mathbf{z}_i).$$

Here,  $\ell : \mathbb{R}^p \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a loss function, and  $\Theta_{n,p}$  is a constraint set. Prominent examples of this general framework include maximum likelihood (ML) estimation [14] and empirical risk minimization [44].

Once the objective (1.1) is formed, it remains to define a computationally efficient scheme to approximate it. Gradient descent is the most frequently applied idea. Assuming for the moment  $\Theta_{n,p} = \mathbb{R}^p$ , this takes the form

$$(1.3) \quad \widehat{\boldsymbol{\theta}}_n(k+1) = \widehat{\boldsymbol{\theta}}_n(k) - h_k \nabla \widehat{R}_n(\widehat{\boldsymbol{\theta}}_n(k)).$$

While a large number of variants and refinements have been developed over the years (projected gradient, accelerated gradient [32], stochastic gradient [39], distributed gradient [42] and so on), these share many of the strengths and weaknesses of the elementary iteration (1.3).

If gradient descent is adopted, the only freedom is in the choice of the loss function  $\ell(\cdot; \cdot)$ . Convexity has been a major guiding principle in this respect. If the function  $\ell(\cdot; \mathbf{z}) : \mathbb{R}^p \rightarrow \mathbb{R}$  is convex, then the empirical risk  $\widehat{R}_n(\cdot)$  is convex as well, and hence gradient descent is globally convergent to an M-estimator (the latter is unique under strict convexity). Also, strong convexity of  $\widehat{R}_n(\cdot)$  can be used to prove optimal statistical guarantees for the M-estimator  $\widehat{\boldsymbol{\theta}}_n$ . This line of thought can be traced back as far as Fisher's argument for the asymptotic efficiency of maximum likelihood estimators [14, 15], and originated many beautiful contributions. In recent years, a flourishing line of research addresses the very high-dimensional regime  $p \gg n$ , by leveraging on suitable restricted strong convexity assumptions [5, 7, 8, 31].

Despite these successes, many problems of practical interest call for nonconvex loss functions. Let us briefly mention a few examples of nonconvex M-estimators that are often preferred by practitioners to their convex counterparts. We will revisit these examples in Section 4.

In binary linear classification, we are given  $n$  pairs  $\mathbf{z}_1 = (y_1, \mathbf{x}_1), \dots, \mathbf{z}_n = (y_n, \mathbf{x}_n)$  with  $y_i \in \{0, 1\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ , and would like to learn a model of the form  $\mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}) = \sigma(\langle \boldsymbol{\theta}_0, \mathbf{x} \rangle)$  with  $\boldsymbol{\theta}_0 \in \mathbb{R}^d$  a parameter vector and  $\sigma : \mathbb{R} \rightarrow [0, 1]$  a threshold function. The nonlinear square loss  $\ell(\boldsymbol{\theta}; y, \mathbf{x}) = (y - \sigma(\langle \boldsymbol{\theta}, \mathbf{x} \rangle))^2$  is commonly used in practice

$$(1.4) \quad \widehat{R}_n(\boldsymbol{\theta}) \equiv \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle))^2.$$

Several empirical studies [9, 34, 45] demonstrate superior robustness and classification accuracy of nonconvex losses in contrast to convex losses (e.g., hinge or logistic loss). The same loss function is commonly used in neural-network models [22].

A similar scenario arises in robust regression. In this case, we are given  $n$  pairs  $\mathbf{z}_1 = (y_1, \mathbf{x}_1), \dots, \mathbf{z}_n = (y_n, \mathbf{x}_n)$  with  $y_i \in \mathbb{R}, \mathbf{x}_i \in \mathbb{R}^d$ , and we assume the linear model  $y_i = \langle \boldsymbol{\theta}_0, \mathbf{x}_i \rangle + \varepsilon_i$ , where the noise terms  $\varepsilon_i$  are i.i.d. with mean zero. Since Huber’s seminal work [18], M-estimators are the method of choice for this problem:

$$(1.5) \quad \widehat{R}_n(\boldsymbol{\theta}) \equiv \frac{1}{n} \sum_{i=1}^n \rho(y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle).$$

Robustness naturally suggests to investigate the use of a nonconvex function  $\rho : \mathbb{R} \rightarrow \mathbb{R}$ , either bounded or increasing slowly at infinity.

Finally, missing data problems famously lead to nonconvex optimization formulations. Consider for instance a mixture of Gaussian problems in which we are given data points  $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^d, \mathbf{z}_i \sim_{\text{i.i.d.}} \sum_{a=1}^k p_a \mathbf{N}(\boldsymbol{\theta}_a, \mathbf{I}_{d \times d})$  (for the sake of simplicity we assume identity covariance and known proportions). The maximum-likelihood problem requires to minimize<sup>4</sup>

$$(1.6) \quad \widehat{R}_n(\boldsymbol{\theta}) \equiv -\frac{1}{n} \sum_{i=1}^n \log \left( \sum_{a=1}^k p_a \phi_d(\mathbf{z}_i - \boldsymbol{\theta}_a) \right),$$

with respect to the cluster centers  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k) \in \mathbb{R}^{d \times k}$ . Other examples include low-rank matrix completion [19], phase retrieval [41], tensor estimation problems [30] and so on.

M-estimation with nonconvex loss functions  $\ell(\cdot; \mathbf{z}) : \mathbb{R}^p \rightarrow \mathbb{R}$  is far less understood than in the convex case. Empirical process theory guarantees uniform convergence of the sample risk  $\widehat{R}_n(\cdot)$  to the population risk  $R(\boldsymbol{\theta}) \equiv \mathbb{E}[\widehat{R}_n(\boldsymbol{\theta})]$  [6]. However, this does not provide a computationally practical scheme, since gradient descent can get stuck in stationary points that are not global minimizers.

In this paper, we present several general results on nonconvex M-estimation and apply them to develop new analysis in each of the three problems mentioned above. We next overview our main results and the paper’s organization, referring to Section 2 for a discussion of related work.

*Uniform convergence of gradient and Hessian.* We prove that, under technical conditions on the loss function  $\ell(\cdot; \cdot)$ ,  $\sup_{\boldsymbol{\theta}} \|\nabla \widehat{R}_n(\boldsymbol{\theta}) - \nabla R(\boldsymbol{\theta})\|_2 \lesssim \sqrt{p(\log n)/n}$  and  $\sup_{\boldsymbol{\theta}} \|\nabla^2 \widehat{R}_n(\boldsymbol{\theta}) - \nabla^2 R(\boldsymbol{\theta})\|_{\text{op}} \lesssim \sqrt{p(\log n)/n}$  (we use  $\lesssim$  to hide constant factors). We refer to Section 3.1 for formal statements.

These results complement the classical analysis that implies uniform convergence of the risk itself, but allow us to control the behavior of stationary points. Note that they guarantee uniform convergence of the gradient and Hessian provided  $n, p \rightarrow \infty$  with  $p(\log p)/n \rightarrow 0$ . Apart from logarithmic factors, this is the optimal condition.

---

<sup>4</sup>Here and below,  $\phi_d(\mathbf{x}) \equiv \exp\{-\|\mathbf{x}\|_2^2/2\}/(2\pi)^{d/2}$  denotes the  $d$ -dimensional standard Gaussian density.

(In this paper, we will refer to the asymptotics  $n, p \rightarrow \infty$  with  $n$  roughly of the same order as  $p$  as *high-dimensional regime*,<sup>5</sup> to contrast it with the low-dimensional analysis for  $n \gg p$ . We will refer to the asymptotics  $n \ll p$  under sparsity assumptions as a *very high-dimensional regime*.)

*Topology of the empirical risk.* As an immediate consequence of the previous result, the structure of the empirical risk function  $\theta \mapsto \widehat{R}_n(\theta)$  is in many cases surprisingly simple. Recall that a Morse function is a twice differentiable function whose stationary points are non-degenerate (i.e., have an invertible Hessian). In particular, stationary points are isolated, and have a well defined index. Assume that the population risk  $R(\theta)$  is *strongly Morse* [i.e., at any stationary point  $\theta$ , all the eigenvalues of the Hessian are bounded away from zero  $|\lambda_i(\nabla^2 R(\theta))| \geq \delta$ ]. Then, for  $n \gtrsim p \log p$ , the stationary points of the empirical risk  $\widehat{R}_n(\theta)$  are in one-to-one correspondence with those of the population risk and have the same index (minima correspond to minima, saddles to saddles, and so on). Weaker conditions ensure this correspondence for local minima alone.

*Very high-dimensional regime.* We then extend the above picture to the case in which the number of parameters  $p$  exceeds the number of samples  $n$ , under the assumption that the true parameter vector  $\theta_0$  is  $s_0$ -sparse. This setting is relevant to a large number of applications, ranging from genomics [35] to signal processing [12]. In order to promote sparse estimates, we study the following  $\ell_1$ -regularized nonconvex problem (cf. Section 3.3):

$$(1.7) \quad \begin{aligned} & \text{minimize} && \widehat{R}_n(\theta) + \lambda_n \|\theta\|_1, \\ & \text{subject to} && \|\theta\|_2 \leq r. \end{aligned}$$

We introduce a *generalized gradient linearity* condition on the loss function  $\ell(\cdot, \cdot)$  and prove that—under this condition—the above problem has a unique local minimum for  $n \gtrsim s_0 \log p$ . Again this is a nearly optimal scaling since no consistent estimation is possible when  $n \lesssim s_0$ .

*Applications.* Given a particular M-estimation problem with a suitable statistical model, we combine the above results with an analysis of the population risk  $R(\theta)$  to derive precise characterizations of the empirical risk. In Section 4, we demonstrate that this program can be carried out by studying the three problems outlined below:

1. *Binary linear classification.* We prove that, for<sup>6</sup>  $n \gtrsim d \log d$ , the empirical risk has a unique local minimum, that is also the global minimum. Further, gradient descent converges exponentially to this minimizer:  $\|\hat{\theta}_n(k) - \hat{\theta}_n\|_2 \leq C \|\hat{\theta}(0) - \hat{\theta}_n\|_2 (1 - h/C)^k$ , and enjoys nearly optimal estimation error guarantees:  $\|\hat{\theta}_n -$

<sup>5</sup>The specific asymptotics  $n, p \rightarrow \infty$  with  $n/p$  converging to a constant is also known as “Kolmogorov asymptotics” [40].

<sup>6</sup>Recall that, in this case, the number of parameters  $p$  is equal to the ambient dimension  $d$ .

$\theta_0\|_2 \leq C\sqrt{(d \log n)/n}$ . If the true parameter  $\theta_0$  is  $s_0$ -sparse, for  $n \gtrsim s_0 \log d$ , the  $\ell_1$ -regularized empirical risk has a unique local minimum that is also the global minimum. The minimizer enjoys nearly optimal estimation error guarantees:  $\|\hat{\theta}_n - \theta_0\|_2 \leq C\sqrt{(s_0 \log n)/n}$ .

2. *Robust regression.* We establish similar results for the robust regression model, under technical assumptions on the loss function  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  and on the distribution of the noise  $\varepsilon_i$ . Namely, we prove that the empirical risk has a unique local minimum, that can be found efficiently via gradient descent, provided  $n \gtrsim d \log d$ . If the true parameter  $\theta_0$  is  $s_0$ -sparse, for  $n \gtrsim s_0 \log d$ , the  $\ell_1$ -regularized empirical risk has a unique local minimum.

3. *Mixture of Gaussians.* We consider the special case of two Gaussians with equal proportions, that is,  $k = 2$  with  $p_1 = p_2 = 1/2$ . We prove that, for  $n \gtrsim d \log d$ , the empirical risk has two global minima that are related by exchange of the two Gaussian components  $(\hat{\theta}_1, \hat{\theta}_2)$  and  $(\hat{\theta}_2, \hat{\theta}_1)$ , connected via saddle points. The trust region algorithm converges to one of these two minima when initialized at random. Also the two minima are within nearly optimal statistical errors from the true centers.

1.1. *Notations.* We use normal font for scalars (e.g.,  $a, b, c, \dots$ ) and boldface for vectors  $(\mathbf{x}, \mathbf{w}, \dots)$ . We will typically reserve capital letters for random variables (and capital bold for random vectors). Given  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^m$ , their standard scalar product is denoted by  $\langle \mathbf{u}, \mathbf{v} \rangle \equiv \sum_{i=1}^m u_i v_i$ . The  $\ell_p$  norm of a vector is—as usual—indicated by  $\|\mathbf{x}\|_p$ . The  $m \times m$  identity matrix is denoted by  $\mathbf{I}_{m \times m}$ .

Given a matrix  $\mathbf{M} \in \mathbb{R}^{m \times m}$ , we denote by  $\lambda_i(\mathbf{M})$ ,  $i \in \{1, \dots, m\}$  its eigenvalues in decreasing order, and by  $\|\mathbf{M}\|_{\text{op}} = \max\{\lambda_1(\mathbf{M}), -\lambda_m(\mathbf{M})\}$  its operator norm. Finally, we shall occasionally consider third-order tensors  $\mathbf{T} \in \mathbb{R}^{m \times m \times m}$ . In this case, the operator (or injective) norm is defined as  $\|\mathbf{T}\|_{\text{op}} = \max\{|\langle \mathbf{T}, \mathbf{x}^{\otimes 3} \rangle| : \|\mathbf{x}\|_2 = 1\}$ , where  $\langle \mathbf{T}, \mathbf{x}^{\otimes 3} \rangle = \sum_{i,j,k} T_{ijk} x_i x_j x_k$ .

We let  $\mathbf{B}_q^d(\mathbf{a}, \rho) \equiv \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x} - \mathbf{a}\|_q \leq \rho\}$  be the  $\ell_q$  ball in  $\mathbb{R}^d$  with center  $\mathbf{a}$  and radius  $\rho$ . We will often omit the dimension superscript  $d$  when clear from the context, the subscript  $q$  when  $q = 2$  and the center  $\mathbf{a}$  when  $\mathbf{a} = \mathbf{0}$ . In particular,  $\mathbf{B}(\rho)$  is the Euclidean ball of radius  $\rho$ . For any set  $D \subset \mathbb{R}^d$ , we let  $\partial D$  be the boundary of the set.

We will generally use upper case letters for random variables and lower case for deterministic values (unless the latter are matrices).

**2. Related literature.** While developing a theory on nonconvex M-estimators is an outstanding challenge, several important facts are by now well understood thanks to a stream of beautiful works. We will provide a necessarily incomplete summary in the next paragraphs.

*Uniform convergence of the empirical risk.* Let  $R(\theta) = \mathbb{E} \hat{R}_n(\theta)$  denote the population risk. Under mild conditions on the loss function  $\ell$  and on the sample size,

it is known that with high probability

$$(2.1) \quad \sup_{\theta \in \Theta_{n,p}} |\widehat{R}_n(\theta) - R(\theta)| \leq \varepsilon_n,$$

for some small  $\varepsilon_n \rightarrow 0$  [6, 43]. This immediately implies guarantees for the M-estimator  $\widehat{\theta}$  in  $\ell$ -loss (or prediction error). Under additional conditions on the population risk  $R(\theta)$ , bounds in estimation error can be derived as well.

For general nonconvex losses, uniform convergence results of the form (2.1) do not preclude the existence of multiple local minima of the sample risk  $\widehat{R}_n(\theta)$ . Hence, this theory does not provide—by itself—computationally practical methods to compute  $\widehat{\theta}$ .

*Algorithmic convergence to the “statistical neighborhood.”* In general, gradient descent and other local optimization procedures are expected to converge to local minima of the empirical risk  $\widehat{R}_n(\theta)$ . In several cases, it is proved that every local minimizer  $\widehat{\theta}^{\text{loc}}$  is “statistically good.” More precisely, the estimation error (e.g., the  $\ell_2$  error  $\|\widehat{\theta}^{\text{loc}} - \theta_0\|_2$ ) is within a constant from the minimax rate for the problem at hand. Also, gradient descent converges to such a neighborhood of the true  $\theta_0$  within a small number of iterations. Results of this type have been proved, among others, for linear regression with noisy covariates [24], generalized linear models with nonconvex regularizers [25], robust regression [26] and sparse regression [47].

While these results are very important, they are not completely satisfactory. For instance, one natural question is whether the statistical error might be improved by finding a better local minimum. If, for instance, the estimation error could be improved by a factor 2 by finding a better local minimum, it would be worth in many applications to restart gradient descent at multiple initializations. Also, since convergence to a fixed point is not guaranteed, these approaches come without a clear stopping criterion. Finally, these proofs make use of the restricted strong convexity (RSC) assumption introduced [25, 31], but do not provide any general tool to establish this condition. In contrast, we prove uniform convergence results that can be used to ensure a condition similar to RSC.

To the best of our knowledge, the only proof of unique local minimum of the regularized empirical risk is obtained in a recent paper by Po-Ling Loh [23]. This work assumes the linear regression model  $y_i = \langle \theta, \mathbf{x}_i \rangle + \varepsilon_i$ , and establishes uniqueness for penalized regression with a certain class of bounded regularizers. This result is comparable to our Theorem 8 (see Section 4.4) which uses  $\ell_1$  regularization instead. Note that, in [23], the sample size is required to scale quadratically in the sparsity:  $n \gtrsim s_0^2$ . Our proof technique is substantially different from the one of [23], and we only require  $n \gtrsim s_0 \log d$ .

*Hybrid optimization methods.* It is often difficult to ensure global convergence to a minimizer of the sample risk  $\widehat{R}_n(\cdot)$  or even to a statistical neighborhood of the true parameters. Several papers develop two-stage procedures to overcome this

problem. The first stage constructs a smart initialization  $\hat{\theta}(0)$  that is within a certain large neighborhood of the true parameters. Spectral methods are often used to implement this step. In the second stage, the estimate is refined by gradient descent (or another local procedure) initialized at  $\hat{\theta}(0)$ . This general approach was studied in a number of problems including matrix completion [19], phase retrieval [10] and tensor decomposition [3].

In some cases, the local optimization stage is only proved to converge to a statistical neighborhood of  $\theta_0$ , and hence this style of analysis shares the shortcomings emphasized in the previous paragraph. In others, it is proven to converge to a single point. Further, in practice, the smart initialization is often not needed, and descent algorithms converge from random initialization as well. Finally, as mentioned above, these analyses are typically carried on in a case-by-case manner.

**3. Uniform convergence results.** In this section, we develop our key tools that are uniform convergence results on the gradient and Hessian of the empirical risk. We also establish some of the direct implications of our results. Throughout, the data consists of the i.i.d. random variables  $\{\mathbf{Z}_1, \dots, \mathbf{Z}_n\}$ . We will use  $\{\mathbf{z}_1, \dots, \mathbf{z}_n\}$  if we want to refer to the corresponding realization. The empirical risk is defined by equation (1.2) and the corresponding population risk is  $R(\theta) = \mathbb{E}\hat{R}_n(\theta) = \mathbb{E}\ell(\theta; \mathbf{Z})$ . The true parameter vector  $\theta_0$  satisfies the condition  $\nabla R(\theta_0) = \mathbb{E}[\nabla\ell(\theta_0; \mathbf{Z})] = \mathbf{0}$ .

We consider two regimes, a *high-dimensional regime* in which the number of parameters  $p$  is allowed to diverge roughly in proportion with the number of samples  $n$ , and a *very high-dimensional regime* in which the true parameters' vector  $\theta_0$  is sparse and the number of parameters  $p$  can be much larger than  $n$ . We treat these two cases separately because the theory is simpler and more general in the first regime.

**3.1. High-dimensional regime.** In order to avoid technical complications, we will limit optimization to a bounded set, that is, we will let  $\Theta_{n,p} = \mathbf{B}^p(r) \equiv \{\theta \in \mathbb{R}^p, \|\theta\|_2 \leq r\}$  to be the Euclidean ball in  $p$  dimensions.

We begin by stating our assumptions. Assumptions 1 and 2 below quantify the amount of statistical noise in the gradient and Hessian of the loss function.

**ASSUMPTION 1 (Gradient statistical noise).** The gradient of the loss is  $\tau^2$ -sub-Gaussian. Namely, for any  $\lambda \in \mathbb{R}^p$  and  $\theta \in \mathbf{B}^p(r)$ ,

$$(3.1) \quad \mathbb{E}\{\exp(\langle \lambda, \nabla\ell(\theta; \mathbf{Z}) - \mathbb{E}[\nabla\ell(\theta; \mathbf{Z})] \rangle)\} \leq \exp\left(\frac{\tau^2 \|\lambda\|_2^2}{2}\right).$$

**ASSUMPTION 2 (Hessian statistical noise).** The Hessian of the loss, evaluated on a unit vector, is  $\tau^2$ -sub-exponential. Namely, for any  $\lambda \in \mathbf{B}^p(1)$  and  $\theta \in \mathbf{B}^p(r)$ ,

$$(3.2) \quad \mathcal{Z}_{\lambda,\theta} \equiv \langle \lambda, \nabla^2\ell(\theta; \mathbf{Z})\lambda \rangle,$$

$$(3.3) \quad \mathbb{E} \left\{ \exp \left( \frac{1}{\tau^2} |Z_{\lambda, \theta} - \mathbb{E} Z_{\lambda, \theta}| \right) \right\} \leq 2.$$

Our third assumption requires the Hessian of the loss to be a Lipschitz function of the vector of parameters  $\theta$ .

ASSUMPTION 3 (Hessian regularity). The Hessian of the population risk is bounded at one point. Namely, there exists  $\theta_* \in \mathbf{B}^p(r)$  and  $H$  such that  $\|\nabla^2 R(\theta_*)\|_{\text{op}} \leq H$ .

Further, the Hessian of the loss function is Lipschitz continuous with integrable Lipschitz constant. Namely, there exists  $J_*$  (potentially diverging polynomially in  $p$ ) such that

$$(3.4) \quad J(\mathbf{z}) \equiv \sup_{\theta_1 \neq \theta_2 \in \mathbf{B}^p(r)} \frac{\|\nabla^2 \ell(\theta_1; \mathbf{z}) - \nabla^2 \ell(\theta_2; \mathbf{z})\|_{\text{op}}}{\|\theta_1 - \theta_2\|_2},$$

$$(3.5) \quad \mathbb{E}\{J(\mathbf{Z})\} \leq J_*.$$

Further, there exists a constant  $c_h$  such that  $H \leq \tau^2 p^{c_h}$ ,  $J_* \leq \tau^3 p^{c_h}$ .

REMARK 1. The constant  $J_*$  serves as a third derivative control of the loss function, and controls the discretization error in proving the uniform convergence of the Hessian. The sample size will depend on  $H$  and  $J_*$  logarithmically, which is why we assume  $H$  and  $J_*$  to grow at most polynomially in dimension  $p$ .

REMARK 2. Note that  $\nabla \ell$  has the same units<sup>7</sup> as  $1/r$ , and  $\nabla^2 \ell$  has the same units as  $1/r^2$ . Thus,  $\tau$  has the same units as  $1/r$ ,  $H$  has the same units as  $\tau^2$  and  $J_*$  has the same units as  $\tau^3$ . This is the reason why we bound  $H$  and  $J_*$  in the form as in Assumption 3. In this way,  $(r \cdot \tau)$  and  $c_h$  are dimensionless.

Discrete loss functions (e.g., the 0–1 loss) are common within the statistical learning literature, but do not satisfy the above assumption because the gradient and Hessian are not defined everywhere. Note however that these can be well approximated by differentiable losses, with little—if any—practical difference.

We are now in position to state our uniform convergence result.

THEOREM 1. *Under Assumptions 1, 2 and 3 stated above, there exists a universal constant  $C_0$ , such that letting  $C = C_0 \cdot (c_h \vee \log(r\tau/\delta) \vee 1)$ , the following hold:*

---

<sup>7</sup>By this, we mean that the two quantities behave in the same way under a rescaling of the parameters  $\theta$ .



(a) *The sample gradient converges uniformly to the population gradient in Euclidean norm. Namely, if  $n \geq Cp \log p$ , we have*

$$(3.6) \quad \mathbb{P} \left( \sup_{\boldsymbol{\theta} \in \mathbf{B}^p(r)} \|\nabla \widehat{R}_n(\boldsymbol{\theta}) - \nabla R(\boldsymbol{\theta})\|_2 \leq \tau \sqrt{\frac{Cp \log n}{n}} \right) \geq 1 - \delta.$$

(b) *The sample Hessian converges uniformly to the population Hessian in operator norm. Namely, if  $n \geq Cp \log p$ , we have*

$$(3.7) \quad \mathbb{P} \left( \sup_{\boldsymbol{\theta} \in \mathbf{B}^p(r)} \|\nabla^2 \widehat{R}_n(\boldsymbol{\theta}) - \nabla^2 R(\boldsymbol{\theta})\|_{\text{op}} \leq \tau^2 \sqrt{\frac{Cp \log n}{n}} \right) \geq 1 - \delta.$$

3.2. *Topology of the empirical risk.* Theorem 1 immediately implies that the structure of stationary points of the sample risk  $\widehat{R}_n(\cdot)$  must reflect that of the population risk. In order to formalize this intuition, we will discuss its implications for a class of functions that we will call *strongly Morse function*. We will then consider a broader set of functions known as *strict saddle*.

3.2.1. *Strongly Morse functions.* Given a differentiable function  $F : \mathbf{B}^d(r) \rightarrow \mathbb{R}$ , we say that  $\mathbf{x}$  in the interior of the ball  $\mathbf{B}^d(r)$  is critical (or stationary) if  $\nabla F(\mathbf{x}) = 0$ .

Recall that a twice differentiable function  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is said to be a *Morse function* if all its critical points are nondegenerate, that is, have an invertible Hessian. In other words,  $\nabla F(\mathbf{x}) = 0$  implies  $\lambda_i(\nabla^2 F(\mathbf{x})) \neq 0$  for all  $i \in \{1, \dots, d\}$ . Morse functions behave well under differentiable reparametrizations, and hence play a central role in differential topology: we refer to [17] for a readable introduction to this area, and to [13, 28, 29] for additional background. The Supplementary Material [27] contains a brief introduction to the most important notions we use in the proofs.

One key feature of a Morse function  $F$  is that, for any  $\mathbf{x}_0 \in \mathbb{R}^d$ , there exists a neighborhood  $\mathbf{B}(\mathbf{x}_0, \varepsilon)$  of  $\mathbf{x}_0$  such that, within  $\mathbf{B}(\mathbf{x}_0, \varepsilon)$ ,  $F(\mathbf{x})$  is qualitatively well described by its second order Taylor expansion at  $\mathbf{x}_0$ . In particular, if  $\mathbf{x}_0$  is a critical point of  $F$ , then there exists a small neighborhood of  $\mathbf{x}_0$  that does not contain any other critical point:<sup>8</sup> all critical points are isolated.

As a consequence, a morse function can only have a finite number of critical points in a compact domain  $K \subseteq \mathbb{R}^d$ . If this was not the case, that is, if the set of critical points  $S \subseteq K$  was infinite, it would have an accumulation point  $\mathbf{x}_* \in K$ .

---

<sup>8</sup>Since  $F$  is twice differentiable, by Taylor expansion there exists  $\delta(\varepsilon)$  [with  $\delta(\varepsilon) \rightarrow 0$  as  $\varepsilon \rightarrow 0$ ] such that  $\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{x}_0) - \nabla^2 F(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0)\|_2 \leq \|\mathbf{x} - \mathbf{x}_0\|_2 \delta(\varepsilon)$  for all  $\mathbf{x} \in \mathbf{B}(\mathbf{x}_0, \varepsilon)$ . For  $\mathbf{x}_0$  a critical point, assume by contradiction that  $\mathbf{x}$  is another critical point in  $\mathbf{B}(\mathbf{x}_0, \varepsilon)$ . Then we have, for  $\mathbf{H}_0 = \nabla^2 F(\mathbf{x}_0)$ ,  $\mathbf{v} = \mathbf{x} - \mathbf{x}_0$ ,  $\|\mathbf{H}_0 \mathbf{v}\|_2 \leq \|\mathbf{v}\|_2 \delta(\varepsilon)$ . But  $\|\mathbf{H}_0 \mathbf{v}\|_2^2 = \langle \mathbf{v}, \mathbf{H}_0^2 \mathbf{v} \rangle \geq \min_{i \leq d} |\lambda_i(\mathbf{H}_0)|^2 \cdot \|\mathbf{v}\|_2^2$ , which gives a contradiction if we choose  $\varepsilon$  so that  $\delta^2(\varepsilon) < \min_{i \leq d} |\lambda_i(\mathbf{H}_0)|^2$ .

By continuity of the gradient,  $\mathbf{x}_*$  would be itself a critical point, and have infinitely many other critical points in any neighborhood, thus leading to a contradiction.

The index of a nondegenerate critical point  $\mathbf{x}_0$  of a twice differentiable function  $F$  is the number of negative eigenvalues of the Hessian  $\nabla^2 F(\mathbf{x}_0)$ : we will denote this integer by  $\text{Ind}_{\mathbf{x}_0}(F)$ . The Morse lemma characterizes completely the behavior of  $F$  in a neighborhood of  $\mathbf{x}_0$ . Namely, if  $\text{Ind}_{\mathbf{x}_0}(F) = k$ , there exists differentiable coordinates<sup>9</sup>  $\varphi_1(\mathbf{x}), \dots, \varphi_d(\mathbf{x})$  defined on  $\mathbf{B}(\mathbf{x}_0, \varepsilon)$  such that  $F(\mathbf{x}) = F(\mathbf{x}_0) + \sum_{i=1}^{d-k} \varphi_i(\mathbf{x})^2 - \sum_{i=d-k+1}^d \varphi_i(\mathbf{x})^2$ . In other words, all critical points with the same index look alike, modulo differentiable changes of coordinates.

Our next definition provides a quantitative version of the notion of Morse functions. We focus on the case in which  $F$  has a bounded domain (a Euclidean ball) because this is the relevant setting for our applications.

**DEFINITION 1.** We say that a twice differentiable function  $F : \mathbf{B}^d(r) \rightarrow \mathbb{R}$  is  $(\varepsilon, \eta)$ -strongly Morse if  $\|\nabla F(\mathbf{x})\|_2 > \varepsilon$  for  $\|\mathbf{x}\|_2 = r$  and, for any  $\mathbf{x} \in \mathbb{R}^d, \|\mathbf{x}\|_2 < r$ , the following holds:

$$(3.8) \quad \|\nabla F(\mathbf{x})\|_2 \leq \varepsilon \quad \Rightarrow \quad \min_{i \in [d]} |\lambda_i(\nabla^2 F(\mathbf{x}))| \geq \eta.$$

The next theorem implies that if the population risk  $R(\cdot)$  is strongly Morse, then the empirical risk retains, with high probability, the same topological structure.

**THEOREM 2.** Under Assumptions 1, 2 and 3, let  $n \geq 4Cp \log n \cdot ((\tau^2/\varepsilon^2) \vee (\tau^4/\eta^2))$ , where  $C = C(\tau^2, \delta, r, c_n)$  is as in the statement of Theorem 1. Then the following happens with probability at least  $1 - \delta$ .

If the population risk  $R : \boldsymbol{\theta} \rightarrow R(\boldsymbol{\theta})$  is  $(\varepsilon, \eta)$ -strongly Morse in  $\mathbf{B}^p(r)$ , then the sample risk  $\widehat{R}_n : \boldsymbol{\theta} \mapsto \widehat{R}_n(\boldsymbol{\theta})$  is  $(\varepsilon/2, \eta/2)$ -strongly Morse in  $\mathbf{B}^p(r)$ . Further there is a one-to-one correspondence between the set of critical points of  $R(\cdot)$ ,  $\mathcal{C} = \{\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(k)}\}$  and the set of critical points of  $\widehat{R}_n(\cdot)$ ,  $\mathcal{C}_n = \{\widehat{\boldsymbol{\theta}}_n^{(1)}, \dots, \widehat{\boldsymbol{\theta}}_n^{(k)}\}$  such that (letting  $\widehat{\boldsymbol{\theta}}_n^{(j)}$  be the point in correspondence with  $\boldsymbol{\theta}^{(j)}$ , for any  $j \in [k]$ )

(a) The index of  $\widehat{\boldsymbol{\theta}}_n^{(j)}$  coincides with the index of  $\boldsymbol{\theta}^{(j)}$ . (In particular, local minima correspond to local minima, and saddles to saddles.)

(b) If we further let  $L = \sup_{\boldsymbol{\theta} \in \mathbf{B}^p(r)} \|\nabla^3 R(\boldsymbol{\theta})\|_{\text{op}}$ , and assume  $n \geq 4Cp \log n / \eta_*^2$  where  $\eta_*^2 = (\varepsilon^2/\tau^2) \wedge (\eta^2/\tau^4) \wedge (\eta^4/(L^2\tau^2))$ , we have, for each  $j \in \{1, \dots, k\}$ ,

$$(3.9) \quad \|\widehat{\boldsymbol{\theta}}_n^{(j)} - \boldsymbol{\theta}^{(j)}\|_2 \leq \frac{2\tau}{\eta} \sqrt{\frac{Cp \log n}{n}}.$$

---

<sup>9</sup>This means that the map  $\mathbf{x} \mapsto (\varphi_1(\mathbf{x}), \dots, \varphi_d(\mathbf{x}))$  is a diffeomorphism.

3.2.2. *Strict saddle functions.* The strong Morse assumption imposes conditions on all the eigenvalues of the Hessian  $\nabla^2 R(\theta)$  at near-critical points, and implies a detailed characterization of the empirical risk. In some applications, only weaker properties can be established for the population risk. These can nevertheless be very useful and Theorem 1 can be used to transfer them to the empirical risk. A useful general notion is the one of *strict saddle* functions, first introduced in [16].

DEFINITION 2. We say that a twice differentiable function  $F : \mathbf{B}^d(r) \rightarrow \mathbb{R}$  is  $(\varepsilon, \eta)$ -*strict saddle* if  $\|\nabla F(\mathbf{x})\|_2 > \varepsilon$  for  $\|\mathbf{x}\|_2 = r$  and, for any  $\mathbf{x} \in \mathbb{R}^d$ ,  $\|\mathbf{x}\|_2 < r$ , the following holds:

$$(3.10) \quad \|\nabla F(\mathbf{x})\|_2 \leq \varepsilon \quad \Rightarrow \quad |\lambda_{\min}(\nabla^2 F(\mathbf{x}))| \geq \eta,$$

where  $\lambda_{\min}(\mathbf{M}) = \min_{i \leq d} \lambda_i(\mathbf{M})$  is the minimum eigenvalue of matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$ .

This definition is completely analogous to the one of strongly Morse functions: the only difference is that we are only imposing a condition on the smallest eigenvalue of the Hessian. In particular, strongly Morse function is a subclass of strict saddle functions.

In strict saddle functions, near critical points are either strongly convex, or have significant negative direction of the Hessian (and hence can be escaped by optimization algorithms). Our definition might seem to impose weaker conditions than the original one in [16], which additionally requires existence of local minima close to convex points. However, Lemma 8 in the Supplementary Material [27] implies that the two definitions are equivalent.

Notice that any local minimum of a strict saddle function is a non-degenerate critical point. Hence, by the same argument in the previous section, local minima are isolated, and there can be finitely many of them in any compact domain. Also, since  $\|\nabla F(\mathbf{x})\|_2 > \varepsilon$  on the boundary, all local minima are in the interior of  $\mathbf{B}^d(r)$ .

THEOREM 3. Under Assumptions 1, 2 and 3, let  $n \geq 4Cp \log n \cdot ((\tau^2/\varepsilon^2) \vee (\tau^4/\eta^2))$ , where  $C = C(\tau^2, \delta, r, c_h)$  is as in the statement of Theorem 1. Then the following happens with probability at least  $1 - \delta$ .

If the population risk  $R : \theta \rightarrow R(\theta)$  is  $(\varepsilon, \eta)$ -strict saddle in  $\mathbf{B}^p(r)$ , then the sample risk  $\widehat{R}_n : \theta \mapsto \widehat{R}_n(\theta)$  is  $(\varepsilon/2, \eta/2)$ -strict saddle in  $\mathbf{B}^p(r)$ . Further, there is a one-to-one correspondence between the set of local minima of  $R(\cdot)$ ,  $\mathcal{C} = \{\theta^{(1)}, \dots, \theta^{(k)}\}$  and the set of local minima of  $\widehat{R}_n(\cdot)$ ,  $\mathcal{C}_n = \{\widehat{\theta}_n^{(1)}, \dots, \widehat{\theta}_n^{(k)}\}$  such that (letting  $\widehat{\theta}_n^{(j)}$  be the local minimum in correspondence with  $\theta^{(j)}$ , for any  $j \in [k]$ ), for each  $j \in \{1, \dots, k\}$ ,

$$(3.11) \quad \|\widehat{\theta}_n^{(j)} - \theta^{(j)}\|_2 \leq \frac{2\tau}{\eta} \sqrt{\frac{Cp \log n}{n}}.$$

3.3. *Very high-dimensional regime.* In the very high-dimensional regime  $n \ll p$ , we will solve the  $\ell_1$ -penalized risk minimization problem

$$(3.12) \quad \begin{aligned} & \text{minimize} && \widehat{R}_n(\boldsymbol{\theta}) + \lambda_n \|\boldsymbol{\theta}\|_1, \\ & \text{subject to} && \|\boldsymbol{\theta}\|_2 \leq r. \end{aligned}$$

We need some additional assumptions. It is fairly straightforward to check them in specific cases; see, for example, Section 4.1. The first assumption is mainly technical, and not overly restrictive: it requires the loss function to have almost surely bounded gradient, in a suitable sense.

ASSUMPTION 4 (Gradient bounds). There exists a constant  $T_*$  such that  $\mathbf{Z}$ -almost surely, for all  $\boldsymbol{\theta} \in \mathbb{B}_2^p(r)$ ,

$$(3.13) \quad \|\nabla \ell(\boldsymbol{\theta}; \mathbf{Z})\|_\infty \leq T_*.$$

Our key structural assumption is stated next. It requires the gradient of the loss function to depend on the parameters only through a linear function of  $\boldsymbol{\theta}$ , possibly dependent on the feature vector  $\mathbf{z}$ . Note that  $\boldsymbol{\theta}_0$  is regarded here as fixed, and hence omitted from the arguments.

ASSUMPTION 5 (Generalized gradient linearity). There exist functions  $g : \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $(t, \mathbf{z}) \mapsto g(t; \mathbf{z})$  and  $\boldsymbol{\psi} : \mathbb{R}^d \rightarrow \mathbb{R}^p$ ,  $\mathbf{z} \mapsto \boldsymbol{\psi}(\mathbf{z})$ , such that

$$(3.14) \quad \langle \nabla \ell(\boldsymbol{\theta}; \mathbf{z}), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle = g(\langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \boldsymbol{\psi}(\mathbf{z}) \rangle; \mathbf{z}).$$

In addition,  $g(t; \mathbf{z})$  is  $L_*$ -Lipschitz to its first argument,  $g(0; \mathbf{z}) = 0$ , and  $\boldsymbol{\psi}(\mathbf{Z})$  is mean-zero and  $\tau^2$ -sub-Gaussian.

As an example, in the case of binary linear classification and robust regression, the data is given as a pair  $\mathbf{z} = (y, \mathbf{x})$ , and there exists a function  $f(t; \mathbf{z})$  such that  $\nabla \ell(\boldsymbol{\theta}; \mathbf{z}) = f(\langle \boldsymbol{\theta} - \boldsymbol{\theta}_0, \mathbf{x} \rangle; \mathbf{z})\mathbf{x}$ . Assumption 5 is satisfied with  $g(t; \mathbf{z}) = tf(t; \mathbf{z})$  provided the latter is Lipschitz as a function of  $t \in \mathbb{R}$ .

THEOREM 4. Under Assumptions 2, 3, 4 and 5 stated above, there exists a constant  $C_1$  that depends on  $(r, \tau^2, c_h, \delta)$ , and a universal constant  $C_0$  such that letting  $C_2 = C_0 \cdot (c_h \vee \log(r\tau/\delta) \vee 1)$ , the following hold:

(a) The sample directional gradient converges uniformly to the population directional gradient, along the direction  $(\boldsymbol{\theta} - \boldsymbol{\theta}_0)$ . Namely, we have

$$(3.15) \quad \begin{aligned} & \mathbb{P} \left( \sup_{\boldsymbol{\theta} \in \mathbb{B}_2^p(r) \setminus \{\mathbf{0}\}} \frac{|\langle \nabla \widehat{R}_n(\boldsymbol{\theta}) - \nabla R(\boldsymbol{\theta}), \boldsymbol{\theta} - \boldsymbol{\theta}_0 \rangle|}{\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_1} \leq (T_* + L_*\tau) \sqrt{\frac{C_1 \log(np)}{n}} \right) \\ & \geq 1 - \delta. \end{aligned}$$

(b) *The sample restricted Hessian converges uniformly to the population restricted Hessian in the set  $B_2^p(r) \cap B_0^p(s_0)$  for any  $s_0 \leq p$ . Namely, as  $n \geq C_2 s_0 \log(np)$  we have*

$$\begin{aligned} & \mathbb{P}\left(\sup_{\theta \in B_2^p(r) \cap B_0^p(s_0), \mathbf{v} \in B_2^p(1) \cap B_0^p(s_0)} | \langle \mathbf{v}, (\nabla^2 \widehat{R}_n(\theta) - \nabla^2 R(\theta)) \mathbf{v} \rangle | \right) \\ & \leq \tau^2 \sqrt{\frac{C_2 s_0 \log(np)}{n}} \geq 1 - \delta. \end{aligned}$$

**4. Applications.**

4.1. *Binary linear classification: High-dimensional regime.* As mentioned in the Introduction, in this case we are given  $n$  pairs  $\mathbf{z}_1 = (y_1, \mathbf{x}_1), \dots, \mathbf{z}_n = (y_n, \mathbf{x}_n)$  with  $y_i \in \{0, 1\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ , whereby  $\mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}) = \sigma(\langle \theta_0, \mathbf{x} \rangle)$  (hence  $p = d$  in this case). We estimate  $\theta_0$  by minimizing the nonlinear square loss (1.4), which we copy here for the reader’s convenience:

$$\begin{aligned} (4.1) \quad & \text{minimize} \quad \widehat{R}_n(\theta) \equiv \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\langle \theta, \mathbf{x}_i \rangle))^2, \\ & \text{subject to} \quad \|\theta\|_2 \leq r. \end{aligned}$$

This can be regarded as a smooth version of the 0–1 loss.

We collect below the technical assumptions on this model.

ASSUMPTION 6 (Binary linear classification). (a) The activation  $z \mapsto \sigma(z)$  is three times differentiable with  $\sigma'(z) > 0$  for all  $z$ , and has bounded first, second and third derivatives. Namely, for some constant  $L_\sigma > 0$ :

$$(4.2) \quad \max\{\|\sigma'\|_\infty, \|\sigma''\|_\infty, \|\sigma'''\|_\infty\} \leq L_\sigma.$$

(b) The feature vector  $\mathbf{X}$  has zero mean and is  $\tau^2$ -sub-Gaussian, that is,  $\mathbb{E}[e^{\langle \lambda, \mathbf{X} \rangle}] \leq e^{\frac{\tau^2 \|\lambda\|_2^2}{2}}$  for all  $\lambda \in \mathbb{R}^d$ .

(c) The feature vector  $\mathbf{X}$  spans all directions in  $\mathbb{R}^d$ , that is,  $\mathbb{E}[\mathbf{X}\mathbf{X}^T] \geq \underline{\gamma} \tau^2 \mathbf{I}_{d \times d}$  for some  $0 < \underline{\gamma} < 1$ .

Assumption 6(a) is satisfied by many classical activation functions, a prominent example being the logistic (or sigmoid) function  $\sigma_L(z) = (1 + e^{-z})^{-1}$ .

Our main results on binary linear classification are summarized in the theorem below.

**THEOREM 5.** *Under Assumption 6, further assume  $\|\theta_0\|_2 \leq r/3$ . There exist positive constants  $C_1, C_2$  and  $h_{\max}$  depending on parameters  $(L_\sigma, r, \tau^2, \underline{\gamma}, \delta)$  and the activation function  $\sigma(\cdot)$ , but independent of  $n$  and  $d$ , such that, if  $n \geq \overline{C}_1 d \log d$ , the following hold with probability at least  $1 - \delta$ :*

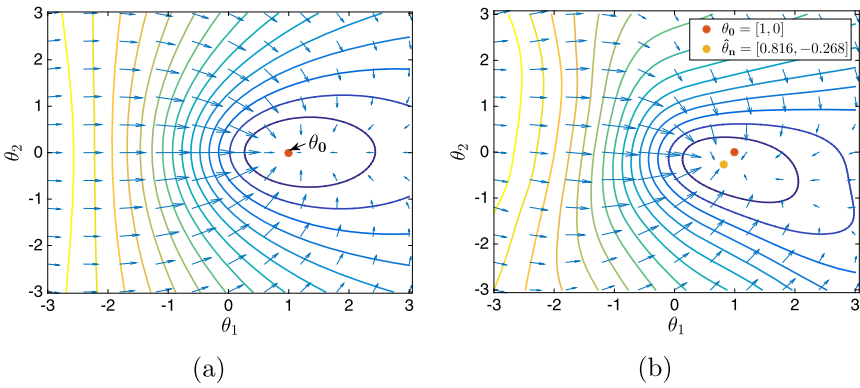


FIG. 1. Binary linear classification: (a) Population risk for  $d = 2$ . (b) A realization of the empirical risk for  $d = 2$  and  $n/d = 20$ .

(a) The empirical risk function  $\theta \mapsto \widehat{R}_n(\theta)$  has a unique local minimizer in  $\mathbf{B}^d(\mathbf{0}, r)$ , that is, the global minimizer  $\widehat{\theta}_n$ .

(b) Gradient descent with fixed step size  $h_k = h \leq h_{\max}$  converges exponentially fast to the global minimizer, for any initialization  $\theta_s \in \mathbf{B}^d(\theta_0, 2r/3)$ :  $\|\widehat{\theta}_n(k) - \widehat{\theta}_n\|_2 \leq C_1 \|\theta_s - \widehat{\theta}_n\|_2 (1 - h/C_1)^k$ .

(c) We have  $\|\widehat{\theta}_n - \theta_0\|_2 \leq C_2 \sqrt{(d \log n)/n}$ .

The proof of this theorem can be found in Section E.1 in the Supplementary Material [27], and is based on the following two-step strategy. First, we study the population risk  $R(\theta)$ , and establish its qualitative properties using analysis. In particular, our results imply that  $R(\theta)$  is strongly Morse in the domain  $\mathbf{B}^d(\mathbf{0}, r)$  (but we prove that an even stronger characterization). Second, we use our uniform convergence result (Theorem 1) to prove that the same properties carry over to the sample risk  $\widehat{R}_n(\theta)$ . Figure 1 presents a small numerical example that illustrates how the qualitative features of the population risk apply to the empirical risk as well.

A few remarks are in order. First of all, the convergence rate of gradient descent [at point (b)] is independent of the dimension  $d$  and number of samples  $n$ . In other words,  $O(\log(1/\varepsilon))$  iterations are sufficient to converge within distance  $\varepsilon$  from the global minimizer. Classical theory of empirical risk minimization only concerns the statistical properties of the optimum, but does not provide efficient algorithms.

Next, note that our condition on the sample size  $n$  is nearly optimal. Indeed, it is information-theoretically impossible to estimate  $\theta_0$  from less than  $n < d$  binary samples. Finally, the convergence rate at point (c) also nearly matches the optimal (parametric) rate  $\sqrt{d/n}$ .

4.2. Binary linear classification: Very high-dimensional regime. As in the previous section, we are given  $n$  pairs  $\mathbf{z}_1 = (y_1, \mathbf{x}_1), \dots, \mathbf{z}_n = (y_n, \mathbf{x}_n)$  with  $y_i \in$

$\{0, 1\}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ , and  $\mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}) = \sigma(\langle \boldsymbol{\theta}_0, \mathbf{x} \rangle)$ . However,  $\boldsymbol{\theta}_0$  is assumed to be sparse, and the number of samples  $n$  is allowed to be much smaller than the ambient dimension  $d = p$ . We adopt again the nonlinear square loss (1.4), but now use a  $\ell_2$ -constrained  $\ell_1$ -regularized risk minimization, as per equation (3.12), which we rewrite here explicitly for the reader’s ease:

$$(4.3) \quad \begin{aligned} &\text{minimize} && \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle))^2 + \lambda_n \|\boldsymbol{\theta}\|_1, \\ &\text{subject to} && \|\boldsymbol{\theta}\|_2 \leq r. \end{aligned}$$

The very high-dimensional regime  $d \gg n$  is of interest in many contexts. In machine learning, the number of parameters  $p$  can increase when a large number of additional features are added to the model (for instance, nonlinear functions of an original set of features). In signal processing,  $\boldsymbol{\theta}_0$  represents an unknown signal, of which we measure noisy random linear projections  $\langle \mathbf{x}_i, \boldsymbol{\theta}_0 \rangle$ ,  $i \in [n]$ , quantized to *one single bit*. This scenario is relevant to group testing [4] and analog-to-digital conversion [20, 21], and has been studied under the name of “one-bit compressed sensing;” see [36] and references therein.

In the very high-dimensional regime, we need additional assumptions on the distribution of  $\mathbf{X}$  as well as the activation function  $\sigma$ .

ASSUMPTION 7 (Fast-decaying activation). The activation function  $\sigma$  satisfy  $\sup_{t \in \mathbb{R}} \{|\sigma'(t)t|, |\sigma''(t)t|\} \leq C_\sigma$  for some absolute constant  $C_\sigma$ .

ASSUMPTION 8 (Continuous and bounded features). The feature vector  $\mathbf{X}$  has a density  $p(\cdot)$  in  $\mathbb{R}^d$ , that is,  $\mathbb{P}(\mathbf{X} \in A) = \int_A p(\mathbf{x}) \, d\mathbf{x}$  for all Borel sets  $A \subseteq \mathbb{R}^d$ . In addition, the feature vector is bounded:  $\|\mathbf{X}\|_\infty \leq M\tau$ , and  $|\langle \mathbf{X}, \boldsymbol{\theta}_0 / \|\boldsymbol{\theta}_0\|_2 \rangle| \leq M\tau$  almost surely, with  $\boldsymbol{\theta}_0$  the ground truth parameter. Here,  $M$  is a dimensionless constant greater than 1.

REMARK 3. Assumption 7 holds popular examples of activation functions, such as the logistic  $\sigma_L(z) = (1 + e^{-z})^{-1}$  or probit  $\sigma_P(z) = \Phi(z)$ .

Also note that Assumption 8 requires  $|\langle \mathbf{X}, \boldsymbol{\theta}_0 \rangle| / \|\boldsymbol{\theta}_0\|_2 \leq M\tau$  to hold only when  $\boldsymbol{\theta}_0$  is the fixed ground truth parameter and not uniformly over all  $s_0$ -sparse vectors. For unbounded sub-Gaussian feature vectors, this assumption does not hold directly. However, for any dataset  $\{\mathbf{X}_i\}_{i=1}^n$  with  $\mathbf{X}_i$  independent  $\tau^2$ -sub-Gaussian, with high probability  $\sup_{i \in [n]} \{\|\mathbf{X}_i\|_\infty, |\langle \mathbf{X}_i, \boldsymbol{\theta}_0 / \|\boldsymbol{\theta}_0\|_2 \rangle|\} \leq C\sqrt{\log(nd)}\tau$ . The next theorem can then be supplemented by a truncation argument at level  $M = C\sqrt{\log(nd)}$ , leading to the same conclusions with an additional  $\log(nd)$  factor in the error bound.

In the statement of the following theorem, for convenience, we will also assume  $n \leq d^{100}$ . This is a technical assumption so that we can bound  $\log(nd) \leq$

$101 \log(d)$ . And since we are considering the very high-dimensional regime, it is not meaningful to discuss  $n > d^{100}$ .

**THEOREM 6.** *Under Assumptions 6, 7 and 8, further assume  $\|\theta_0\|_0 \leq s_0$ ,  $\|\theta_0\|_2 \leq r/2$  and  $n \leq d^{100}$ . Then there exist constants  $C_n, C_\lambda, C_s$  and  $\varepsilon_0$  depending on  $(L_\sigma, C_\sigma, r, \tau^2, \underline{\gamma}, \delta)$  and the activation function  $\sigma(\cdot)$ , but independent of  $n, d, s_0$  and  $M$ , such that as  $n \geq C_n s_0 \log d$  and  $\lambda_n \geq C_\lambda M \sqrt{(\log d)/n}$ , the following hold with probability at least  $1 - \delta$ :*

- (a) *Any stationary point of problem (4.3) is in  $B_2^d(\theta_0, C_s((M^2 s_0 \log d)/n + s_0 \lambda_n^2)^{1/2})$ .*
- (b) *As long as  $n$  is large enough such that  $n \geq C_n s_0 \log^2 d$  and  $C_s((M^2 s_0 \log d)/n + s_0 \lambda_n^2)^{1/2} \leq \varepsilon_0$ , the problem has a unique local minimizer  $\hat{\theta}_n$  which is also the global minimizer.*

As in the previous section, our proof makes a crucial use of the sparse uniform convergence result, Theorem 4, together with an analysis of the population risk.

**REMARK 4.** Let us emphasize that Theorem 6 leaves open the existence of a fast algorithm to find the global optimizer  $\hat{\theta}_n$ . However [33], Theorem 3, implies that, by running  $k$  steps of projected gradient descent, we can find an estimate  $\hat{\theta}_n(k)$  which has a subgradient of order  $O(1/k)$ . While we expect this sequence to converge to  $\hat{\theta}_n$ , we defer this question to future work.

Theorem 6 establishes a nearly optimal upper bound on the  $\ell_2$  estimation error  $\|\hat{\theta}_n - \theta_0\|_2$ . Indeed this error is within a logarithmic factor from the error achieved by an oracle estimator that is given the exact support of  $\theta_0$ . For comparison, [36, 37] proves  $\|\hat{\theta}_n^{\text{LP}} - \theta_0\|_2 \lesssim (s_0/n)^{1/4} (\log p/s_0)^{1/4}$  for a linear programming formulation, under the more restrictive assumption of Gaussian feature vectors  $\mathbf{x}_i \sim N(\mathbf{0}, I_{d \times d})$ . This analysis was generalized in [1] to feature vectors with i.i.d. entries, although with the same estimation error bound. The optimal rate  $\|\hat{\theta}_n^{\text{cvx}} - \theta_0\|_2 \lesssim (s_0/n) \log(p/s_0)$  was obtained only recently in [38], again for standard Gaussian feature vectors.

Let us finally emphasize that the estimator defined here uses a bounded loss function and is potentially more robust to outliers than other approaches that use a convex loss (e.g., logistic loss).

**4.3. Robust regression: High-dimensional regime.** In robust regression, we are given data  $\mathbf{z}_1 = (y_1, \mathbf{x}_1), \dots, \mathbf{z}_n = (y_n, \mathbf{x}_n)$  with  $y_i \in \mathbb{R}, \mathbf{x}_i \in \mathbb{R}^d$ , and we assume the linear model  $y_i = \langle \theta_0, \mathbf{x}_i \rangle + \varepsilon_i$ , where the noise terms  $\varepsilon_i$  are i.i.d. with mean zero. Also in this case we have  $p = d$ . We use the loss (1.5), which we copy here



for the reader’s convenience:

$$(4.4) \quad \begin{aligned} &\text{minimize} && \frac{1}{n} \sum_{i=1}^n \rho(y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle), \\ &\text{subject to} && \|\boldsymbol{\theta}\|_2 \leq r. \end{aligned}$$

Classical choices for loss function  $t \mapsto \rho(t)$  are the Huber loss [18] which is convex with  $\rho_{\text{Huber}}(t) = |t| - \text{const.}$  for  $t$  large enough, and Tukey’s bisquare loss, which is bounded and defined as

$$(4.5) \quad \rho_{\text{Tukey}}(t) = \begin{cases} 1 - (1 - (t/t_0)^2)^3 & \text{for } |t| \leq t_0, \\ 1 & \text{for } |t| \geq t_0. \end{cases}$$

It is common to define the associated score function as  $\psi(t) = \rho'(t)$ .

We next formulate our assumptions.

ASSUMPTION 9 (Robust regression). (a) The score function  $z \mapsto \psi(z)$  is twice differentiable and odd in  $z$  with  $\psi(z) \geq 0$  for all  $z \geq 0$ , and has bounded zero, first and second derivatives. Namely, for some constant  $L_\psi > 0$ :

$$(4.6) \quad \max\{\|\psi\|_\infty, \|\psi'\|_\infty, \|\psi''\|_\infty\} \leq L_\psi.$$

(b) The noise  $\varepsilon$  has a symmetric distribution, that is, such that  $\varepsilon$  is distributed as  $-\varepsilon$ . Further, defining  $g(z) \equiv \mathbb{E}_\varepsilon\{\psi(z + \varepsilon)\}$  we have  $g(z) > 0$  for all  $z > 0$ , as well as  $g'(0) > 0$ .

(c) The feature vector  $\mathbf{X}$  has zero mean and is  $\tau^2$ -sub-Gaussian, that is,  $\mathbb{E}[e^{\langle \boldsymbol{\lambda}, \mathbf{X} \rangle}] \leq e^{\frac{\tau^2 \|\boldsymbol{\lambda}\|_2^2}{2}}$  for all  $\boldsymbol{\lambda} \in \mathbb{R}^d$ .

(d) The feature vector  $\mathbf{X}$  spans all directions in  $\mathbb{R}^d$ , that is,  $\mathbb{E}[\mathbf{X}\mathbf{X}^\top] \succeq \underline{\gamma} \tau^2 \mathbf{I}_{d \times d}$  for some  $0 < \underline{\gamma} < 1$ .

Note that the condition  $g(z) \equiv \mathbb{E}_\varepsilon\{\psi(z + \varepsilon)\} > 0$  for all  $z > 0$  and  $g'(0) > 0$  are quite mild, and holds—for instance—if the noise has a density that is strictly positive for all  $\varepsilon$ , and decreasing for  $\varepsilon > 0$ .

THEOREM 7. Under Assumption 9, further assume  $\|\boldsymbol{\theta}_0\|_2 \leq r/3$ . Then there exist positive constants  $C_1, C_2$  and  $h_{\max}$  depending on parameters  $(L_\psi, r, \tau^2, \underline{\gamma}, \delta)$ , the loss function  $\rho(\cdot)$ , and the law of noise  $\mathbb{P}_\varepsilon$  but independent of  $n$  and  $d$ , such that as  $n \geq C_1 d \log d$ , the robust regression estimator satisfies the following with probability at least  $1 - \delta$ :

(a) The empirical risk function  $\mathbf{w} \mapsto \widehat{R}_n(\boldsymbol{\theta})$  has a unique local minimizer in  $\mathbf{B}^d(r)$ , that is, the global minimizer  $\widehat{\boldsymbol{\theta}}_n$ .

(b) *Gradient descent with fixed step size*  $h_k = h \leq h_{\max}$  *converges exponentially fast to the global minimizer, for any initialization*  $\theta_s \in \mathbf{B}^d(\theta_0, 2r/3)$ :

$$\|\hat{\theta}_n(k) - \hat{\theta}_n\|_2 \leq C_1 \|\theta_s - \hat{\theta}_n\|_2 (1 - h/C_1)^k.$$

(c) *We have*  $\|\hat{\theta}_n - \theta_0\|_2 \leq C_2 \sqrt{(d \log n)/n}$ .

The proof follows the same two steps strategy as for the binary classification problem. In particular, we obtain a precise characterization of the population risk, which (in particular) is strongly Morse.

4.4. *Robust regression: Very high-dimensional regime.* As in the previous section, we are given  $n$  pairs  $\mathbf{z}_1 = (y_1, \mathbf{x}_1), \dots, \mathbf{z}_n = (y_n, \mathbf{x}_n)$  with  $y_i \in \mathbb{R}$ ,  $\mathbf{x}_i \in \mathbb{R}^d$ , and we assume the linear model  $y_i = \langle \theta_0, \mathbf{x}_i \rangle + \varepsilon_i$ , where the noise terms  $\varepsilon_i$  are i.i.d. with mean zero. However,  $\theta_0$  is assumed to be sparse, while the number of samples  $n$  is much smaller than the ambient dimension  $d = p$ . We adopt again the loss (1.5), but now use a  $\ell_2$ -constrained  $\ell_1$ -regularized risk minimization, as per equation (3.12), which we rewrite here explicitly for the reader's ease:

$$(4.7) \quad \begin{aligned} & \text{minimize} && \frac{1}{n} \sum_{i=1}^n \rho(y_i - \langle \theta, \mathbf{x}_i \rangle) + \lambda_n \|\theta\|_1, \\ & \text{subject to} && \|\theta\|_2 \leq r. \end{aligned}$$

Like the case of very high-dimensional binary classification, we also need continuous and bounded feature assumptions, that is, Assumption 8 and we need a fast decaying assumption on  $\psi = \rho'$ .

ASSUMPTION 10 (Fast-decaying score function). The score function  $\psi$  satisfies  $\sup_{t \in \mathbb{R}} \{|\psi(t)t|\} \leq C_\psi$  for some absolute constant  $C_\psi$ .

THEOREM 8. *Under Assumptions 6, 8 and 10, further assume*  $\|\theta_0\|_0 \leq s_0$ ,  $\|\theta_0\|_2 \leq r/2$ , *and*  $n \leq d^{100}$ . *Then there exist constants*  $C_n, C_\lambda, C_s$  *and*  $\varepsilon_0$  *depending on*  $(L_\psi, C_\psi, r, \tau^2, \underline{\gamma}, \delta)$ , *the loss function*  $\rho$ , *and the law of noise*  $\mathbb{P}_\varepsilon$ , *but independent of*  $n, d, s_0$  *and*  $M$ , *such that as*  $n \geq C_n s_0 \log d$  *and*  $\lambda_n \geq C_\lambda M \sqrt{(\log d)/n}$ , *the following hold with probability at least*  $1 - \delta$ :

(a) *Any stationary point of problem (4.7) is in*  $\mathbf{B}_2^d(\theta_0, C_s((M^2 s_0 \log d)/n + s_0 \lambda_n^2)^{1/2})$ .

(b) *As long as*  $n$  *is large enough such that*  $n \geq C_n s_0 \log^2 d$  *and*  $C_s((M^2 s_0 \log d)/n + s_0 \lambda_n^2)^{1/2} \leq \varepsilon_0$ , *the problem has a unique local minimizer*  $\hat{\theta}_n$  *which is also the global minimizer.*

The proof of this theorem is almost the same as the proof of Theorem 6. We will omit the proof to avoid redundancies.

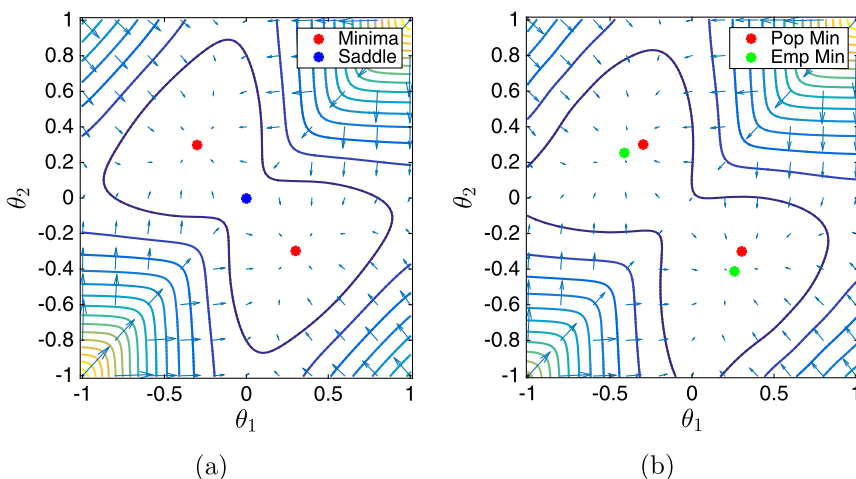


FIG. 2. Gaussian mixture model: (a) Population risk for  $d = 1$ . (b) A realization of the empirical risk for  $d = 1$ , and  $n = 30$ .

4.5. *Gaussian mixture model.* In the applications considered so far, the population risk has a unique stationary point which is also the global minimum. We used our uniform convergence theorems to prove that the empirical risk has the same property, and hence can be optimized efficiently.

In order to illustrate our approach on an example with multiple local minima, we consider clustering within a simple Gaussian mixture model. We are given data points  $\mathbf{z}_1, \dots, \mathbf{z}_n \in \mathbb{R}^d$ , with  $\mathbf{z}_i$  drawn from a mixture of two Gaussians, in equal proportions,  $\mathbf{z}_i \sim (1/2)\mathcal{N}(\boldsymbol{\theta}_{0,1}, \mathbf{I}_{d \times d}) + (1/2)\mathcal{N}(\boldsymbol{\theta}_{0,2}, \mathbf{I}_{d \times d})$ . Define the separation parameter  $D = \|\boldsymbol{\theta}_{0,2} - \boldsymbol{\theta}_{0,1}\|_2/2$ . We want to estimate the centers  $\boldsymbol{\theta}_{0,1}, \boldsymbol{\theta}_{0,2}$  by solving the maximum likelihood problem [here  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \in \mathbb{R}^{2d}$ ]

$$(4.8) \quad \text{minimize} \quad \widehat{R}_n(\boldsymbol{\theta}) \equiv -\frac{1}{n} \sum_{i=1}^n \log \left( \sum_{a=1}^2 \phi_d(\mathbf{z}_i - \boldsymbol{\theta}_a) \right).$$

In this case, the population risk has at least two global minima related by the symmetry under exchange of the two components:  $\boldsymbol{\theta}_+ = (\boldsymbol{\theta}_{0,1}, \boldsymbol{\theta}_{0,2})$  and  $\boldsymbol{\theta}_- = (\boldsymbol{\theta}_{0,2}, \boldsymbol{\theta}_{0,1})$ , as well as a saddle point  $\boldsymbol{\theta}_s = ((\boldsymbol{\theta}_{0,1} + \boldsymbol{\theta}_{0,2})/2, (\boldsymbol{\theta}_{0,1} + \boldsymbol{\theta}_{0,2})/2)$ . This is a common phenomenon: symmetries lead to multiple minima of the risk function. In a recent paper, Xu, Hsu and Maleki [46] prove that these are the only critical points. A related analysis was carried out by Daskalakis, Tzamos, Christos and Zampetakis [11] in order to study the behavior of the EM algorithm. See Figure 2 for an illustration.

**THEOREM 9.** *Let  $\widehat{R}_n(\boldsymbol{\theta})$  be the empirical risk for an equal-proportion mixture of two Gaussians. Then there exist constants  $C_1, C_2$  and  $C_3$  depending on  $(D, \delta)$*

but independent of  $n$  and  $d$ , such that as  $n \geq C_1 d \log d$ , the following holds with probability at least  $1 - \delta$ :

(a) Inside  $\mathbf{B}^{2d}(\boldsymbol{\theta}_s, C_2)$ , the empirical risk has exactly two local minima  $\hat{\boldsymbol{\theta}}_+$ ,  $\hat{\boldsymbol{\theta}}_-$  related by an exchange of the two classes.

(b) For any initialization  $\hat{\boldsymbol{\theta}}_0 \in \mathbf{B}^{2d}(\boldsymbol{\theta}_s, C_2)$ , the trust region algorithm will converge to one of the local minima.

(c) The local minima satisfy

$$(4.9) \quad \|\hat{\boldsymbol{\theta}}_+ - \boldsymbol{\theta}_+\|_2 \leq C_3 \sqrt{\frac{d \log n}{n}}, \quad \|\hat{\boldsymbol{\theta}}_- - \boldsymbol{\theta}_-\|_2 \leq C_3 \sqrt{\frac{d \log n}{n}}.$$

As in previous examples, we obtain a precise characterization of the population risk, building on [46], and then transfer the result to empirical risk using our uniform convergence results. Our analysis implies—in particular—that the population risk is strict saddle.

**5. Numerical experiments.** We carried out extensive numerical experiments in order to verify how accurate is our theory. Sections 5.1 to 5.3 present simulations for the nonconvex binary classification and robust regression models studied in Section 4. Sections 5.4 present illustrations with real data. We will present simulations for the Gaussian mixture model (Section H.1) and binary classification using the Australian credit dataset (Section H.2) in the Supplementary Material [27].

5.1. *Binary linear classification: High-dimensional regime.* Figures 3(a), 3(b), 4(a), 4(b) report our results for the nonconvex binary classification model of Section 4.1.

We consider i.i.d. predictors  $\mathbf{X}_i \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_{d \times d})$ , and generate labels  $Y_i \in \{0, 1\}$  with  $\mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}) = \sigma(\langle \boldsymbol{\theta}_0, \mathbf{x} \rangle)$  where  $\sigma(u) = \sigma_L(u) = (1 + e^{-u})^{-1}$  is the logistic activation. We perform gradient descent [cf. equation (1.3) to minimize the empirical risk (1.4)], with a minor revision in practice: we will project the points back into  $\mathbf{B}^d(r)$  if the iteration points fall out of the ball, with  $r = 3\|\boldsymbol{\theta}_0\|_2$ . The step size is fixed to be  $h = 1$ .

In order to test the hypothesis that the landscape is simple (i.e., it has a unique local minimum), we run projected gradient descent starting from multiple random initializations  $\boldsymbol{\theta}_s \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_{d \times d}/d)$ . If the landscape is simple, we expect the iterates  $\hat{\boldsymbol{\theta}}_n(k)$  to converge to the same global minimizer with no dependence on the initialization. If the landscape is rough, projected gradient descent will converge to different points depending on the initialization. Given a maximum number of iterations  $k_{\max}$ , we define the following quantity, depending on the data  $(\mathbf{Y}, \mathbf{X}) \equiv \{(Y_i, \mathbf{X}_i)\}_{1 \leq i \leq n}$ ,

$$(5.1) \quad S_{\mathbf{Y}, \mathbf{X}} = \sqrt{\text{Tr}(\widehat{\text{Var}}_{\text{init}}(\hat{\boldsymbol{\theta}}_n(k_{\max}) | \mathbf{Y}, \mathbf{X}))},$$

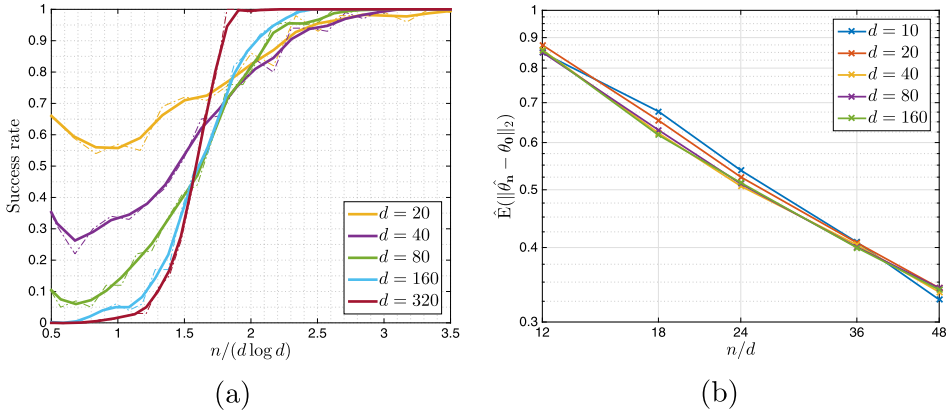


FIG. 3. Binary linear classification, high dimensional: (a) Success rate versus  $n/(d \log d)$  for several ambient dimensions  $d$ , with  $\|\theta_0\|_2 = 3$  (dashed lines are empirical averages, continuous lines are a smoothed version); (b) Estimation error  $\hat{\mathbb{E}}[\|\hat{\theta}_n - \theta_0\|_2]$  versus  $n/d$ , for  $\|\theta_0\|_2 = 1$ .

where the variance is taken over the random initializations  $\theta_s$ . In words,  $S_{\mathbf{Y}, \mathbf{X}}$  is the spread of the limit points of projected gradient descent, for the instance  $(\mathbf{Y}, \mathbf{X})$ . We then define the empirical success probability as

$$(5.2) \quad \hat{\mathbb{P}}_{\text{succ}} \equiv \hat{\mathbb{P}}(S_{\mathbf{Y}, \mathbf{X}} \leq \varepsilon).$$

In Figure 3(a), we plot our results for the empirical success rate, for several values of  $n, d$ . In this experiment, we take  $\|\theta_0\|_2 = 3$ . For each pair  $(n, d)$ , we generate 100 instances  $(Y_i, \mathbf{X}_i)$  and run projected gradient descent from 10 random initializations. We use  $k_{\max} = 10^4$  iterations and tolerance  $\varepsilon = 10^{-2}$  though results seem to be fairly insensitive to these parameters. For each dimension  $d$ , the success rate goes rapidly from 0 to 1 as the number of samples  $n$  crosses a threshold. We plot the success probability as function of the rescaled number of samples  $n/(d \log d)$ . On this scale, curves for different dimension cross each other, and become steeper as  $d$  increases. This is consistent with Theorem 5. This also suggests a sharp phase transition at  $n_*(d)$  which is roughly of order  $d \log d$ . It is a fascinating open question whether a sharp threshold actually exists.<sup>10</sup>

Figure 3(b) illustrates the behavior of the estimation error  $\|\hat{\theta}_n - \theta_0\|_2$  achieved by gradient descent. In all the following experiments, we will take  $\|\theta_0\|_2 = 1$ . We plot the estimation error (averaged over 100 random instances)  $\hat{\mathbb{E}}[\|\hat{\theta}_n - \theta_0\|_2]$  versus  $n/d$ . Curves for different dimensions collapse, and are consistent with the optimal rate  $\|\hat{\theta}_n - \theta_0\|_2 = \Theta(\sqrt{d/n})$ .

Figure 4(a) shows the convergence of gradient descent for several values of  $n$  and  $d$ , for fixed  $n/d = 20$ . Namely, we plot the distance from the global minimizer

<sup>10</sup>When convergence to a single global minimum fails, we observe that often projected gradient actually convergence to the boundary of  $\mathbb{B}^d(r)$ .

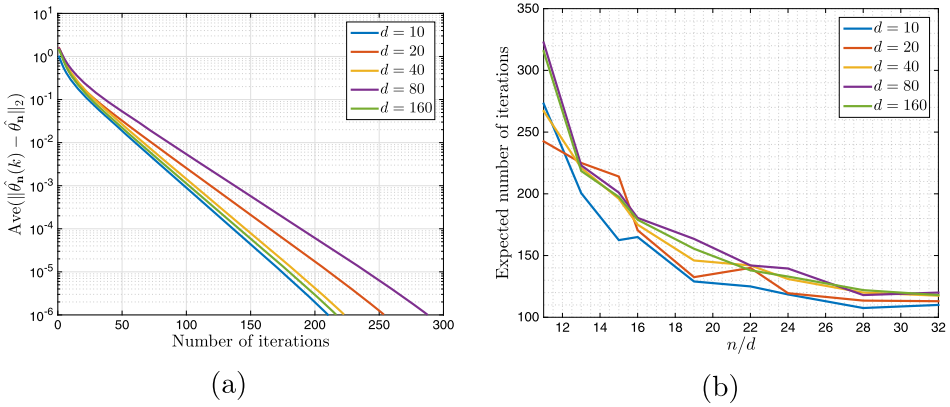


FIG. 4. Binary linear classification, high dimensional: (a) The convergence of the gradient descent algorithm. Here  $\|\theta_0\|_2 = 1, n/d = 20$ . The y-axis is on a log-scale; (b) Minimum number of iterations needed to achieve average distance  $10^{-4}$  from the global optimizer.

as a function of the number of iterations  $k$ , estimated using 100 realizations  $(\mathbf{Y}, \mathbf{X})$ . Since there is a small probability that gradient descent fails to find unique minimizer, we average the distance from the global minimizer over the results between the (0.05, 0.95) quantiles of these 100 instances. Convergence to the global minimizer appears to be exponential as predicted by Theorem 5. Also, convergence is fairly independent of the dimension for fixed  $n/d$ .

Finally, Figure 4(b) shows the number of iterations needed to achieve the  $\varepsilon = 10^{-4}$  optimization error. We run 100 instances, and we plot the expected number of iteration, by averaging the results between the (0.05, 0.95) quantiles of these 100 instances. When  $n/d$  is small, the landscape is not very smooth, and convergence is slower. When  $n/d$  grows, the number of iterations decreases and converges to a constant. This is also predicted by Theorem 5: the landscape of empirical risk will be as smooth as the landscape of population risk, as  $n \geq Cd \log d$ .

5.2. Binary linear classification: Very high-dimensional regime. In Figures 5, 6(a), 6(b), we present our results on nonconvex binary linear classification in the very high-dimensional regime. Data  $(Y_i, \mathbf{X}_i)$  were generated as in the previous section, with  $\theta_0$  a vector  $k$  nonzero entries all of size  $1/\sqrt{k}$ . We use proximal gradient descent to solve problem (3.12) with  $r = 10$ .

In Figure 5, we use random initializations  $\theta_s \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_{d \times d}/d)$ , and plot the empirical standard deviation of the resulting iterates  $\text{std}(\hat{\theta}_n(i)) = \text{Tr}(\widehat{\text{Var}}(\hat{\theta}_n(i)))^{1/2}$ . Note that the variance is taken over the random initializations, for a same realization of the data  $(\mathbf{Y}, \mathbf{X})$ , and hence captures smoothness (or roughness) of the empirical risk landscape. The standard deviation appears to converge exponentially fast to 0, confirming that indeed proximal gradient is converging to the unique local minimizer, as anticipated by Theorem 6.

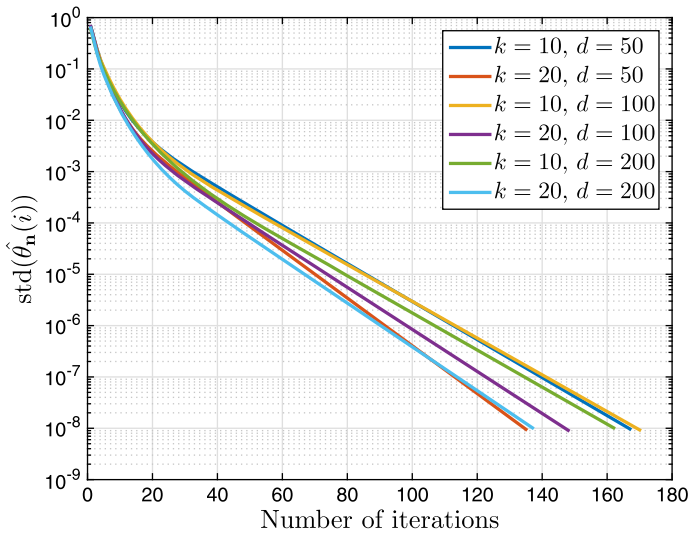


FIG. 5. Binary linear classification, very high dimensional. The standard deviation of each iteration point with respect to random initialization.

In Figure 6(a), we plot the expected distance from the global minimizer  $\hat{\theta}_n$  for each iterates. Proximal gradient appears to converge exponentially fast for  $n \gg k \log^2(d)$ .

5.3. Robust linear regression. In Figures 7, 8(a), 8(b), we present simulations for robust regression. We generated random covariates  $\mathbf{X}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{d \times d})$  and re-

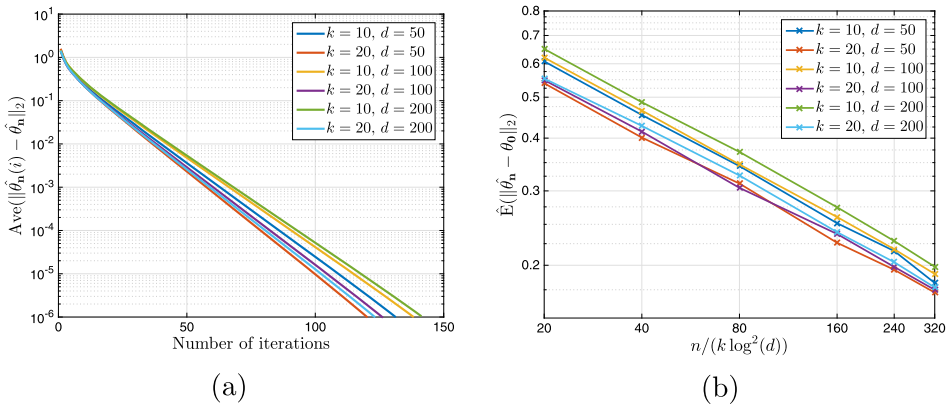


FIG. 6. Binary linear classification, very high-dimensional regime: (a) The convergence of proximal gradient descent. Here,  $\|\theta_0\|_2 = 1$ , and  $n/(k \log^2(d)) = 20$ , and  $\lambda_n = 1/100 \cdot \sqrt{\log^2(d)/n}$ . (b) Convergence of the statistical error.

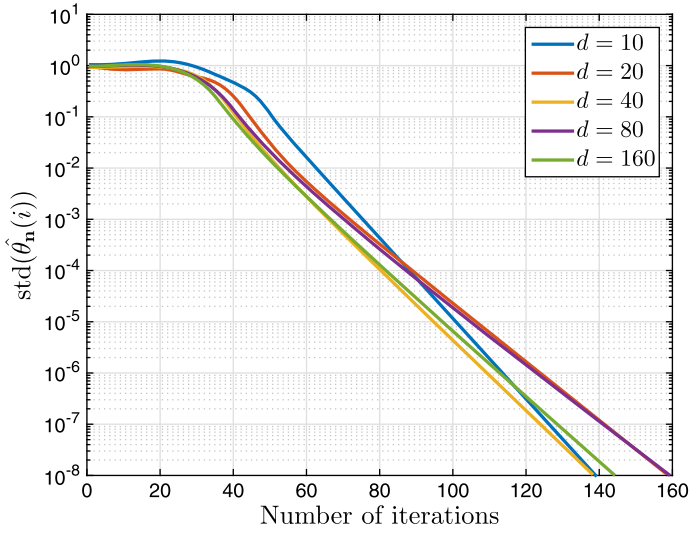
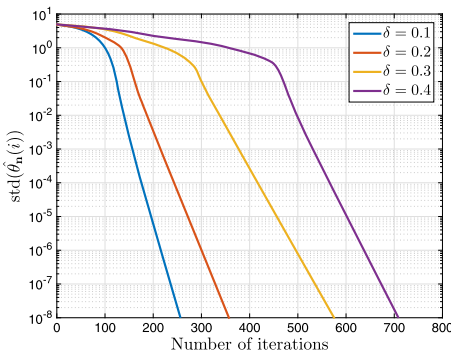


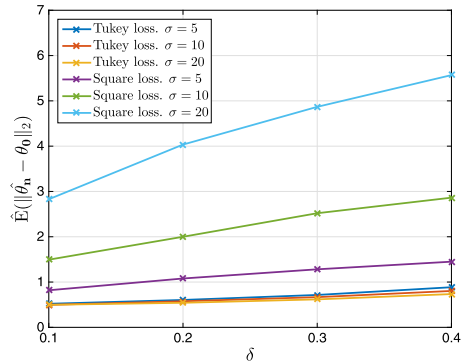
FIG. 7. Robust regression. The standard deviation of each iteration point with respect to random initialization.

sponses  $Y_i = \langle \theta_0, \mathbf{X}_i \rangle + \varepsilon_i$ , where  $\|\theta_0\|_2 = 1$ . Again, we used projected gradient descent to solve the optimization problem (4.4) with  $r = 10$ . For the loss function, we used Tukey’s loss (4.5) with  $t_0 = 4.685$ .

In Figure 7, we plot the standard deviation of the iterates  $\text{std}(\hat{\theta}_n(i)) = \text{Tr}(\widehat{\text{Var}}(\hat{\theta}_n(i)))^{1/2}$  over random initializations  $\theta_s \sim \mathcal{N}(\mathbf{0}, 25\mathbf{I}_{d \times d}/d)$ . In this case  $\varepsilon_i \sim \mathcal{N}(0, 1)$ . Again, this standard deviation converges exponentially fast to 0 sup-



(a)



(b)

FIG. 8. Robust regression: (a) The standard deviation of each iteration point with respect to random initialization, for different proportion of contamination. (b) The robustness of the global minimum between linear regression and Tukey regression.



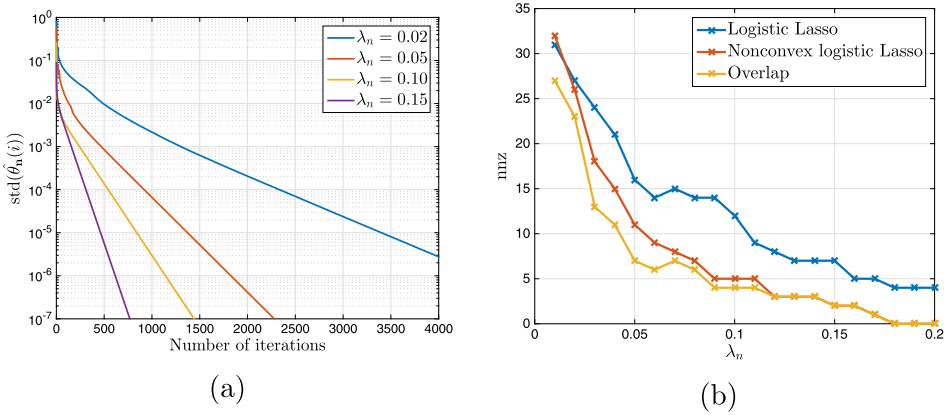


FIG. 9. Colon cancer data: (a) The standard deviation of each iteration point with respect to random initialization, for different regularization parameter. (b) Number of nonzero elements of logistic Lasso and nonconvex logistic Lasso.

porting the claim that proximal gradient descent converges to a unique global minimum irrespective of the initialization.

In Figure 8(a), (b), we study the a contaminated model for the noise, namely  $\varepsilon_i \sim (1 - \delta)N(0, 1) + \delta N(0, \sigma^2)$ . In Figure 8(a), we plot the standard deviation of the estimates obtained with random initializations  $\theta_s \sim N(\mathbf{0}, 25\mathbf{I}_{d \times d}/d)$ , for  $n = 480, d = 80$ . Convergence rate remains exponential even for large contamination fraction. In Figure 8(b), we investigated the dependence of the estimation error on the contamination fraction, and the scale of outliers. Tukey’s regression is fairly insensitive to outliers, while the least squares regression deteriorates as expected.

5.4. Colon cancer data. In Figure 9(a), (b), we consider a gene-expression dataset from [2]. The data set contains expression levels of of 2000 genes in 22 normal and 40 tumor colon tissues, hence  $n = 62$  data points. Expression levels are normalized as in [2] to have zero mean and unit standard deviation. We use the expression levels to form feature vectors  $\mathbf{x}_i \in \mathbb{R}^d, d = 2000$  and encode the type of tissue using a binary label  $y_i = 1$  (tumor) or  $y_i = 0$  (no tissue).

We fit a model of the form  $\mathbb{P}(Y_i = 1 | \mathbf{X}_i = \mathbf{x}) = \sigma(\langle \theta_0, \mathbf{x} \rangle)$  with  $\sigma(u) = \sigma_L(u)$  the logistic function, by using the nonconvex approach (4.3) and proximal gradient. We also used  $\ell_1$ -regularized logistic regression, for comparison. Let us emphasize here that our focus here is not on the accuracy of the predictive model, but rather on showing that the nonconvex approach is a viable alternative to the more standard regularized logistic regression.

In Figure 9(a), we plot the standard deviation of the estimate  $\hat{\theta}_n(i)$ , over random initializations  $\theta_s \sim N(\mathbf{0}, \mathbf{I}_{d \times d}/d)$ . The standard deviation decreases exponentially fast, suggesting that indeed the optimization problem has a unique local minimum. In Figure 9(b), we compare the model selected by the nonconvex approach (4.3)

to the one from  $\ell_1$ -regularized logistic regression, and also plot the number of overlaps of their selected variables. Note that most of the covariates selected by the nonconvex regression method also appear in logistic regression. This suggests that the model produced by the nonconvex approach is comparable to that produced by  $\ell_1$ -regularized logistic regression.

## SUPPLEMENTARY MATERIAL

**Supplement: Proofs and simulations** (DOI: [10.1214/17-AOS1637SUPP](https://doi.org/10.1214/17-AOS1637SUPP); .pdf). The supplement provides some technical background lemmas and gives all the proofs of the theorems, and additional numerical simulations.

## REFERENCES

- [1] AI, A., LAPANOWSKI, A., PLAN, Y. and VERSHYNIN, R. (2014). One-bit compressed sensing with non-Gaussian measurements. *Linear Algebra Appl.* **441** 222–239. [MR3134344](#)
- [2] ALON, U., BARKAI, N., NOTTERMAN, D. A., GISH, K., YBARRA, S., MACK, D. and LEVINE, A. J. (1999). Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proc. Natl. Acad. Sci. USA* **96** 6745–6750.
- [3] ANANDKUMAR, A., GE, R. and JANZAMIN, M. (2015). Learning overcomplete latent variable models through tensor methods. In *Proceedings of the Conference on Learning Theory (COLT), Paris, France*.
- [4] ATIA, G. K. and SALIGRAMA, V. (2012). Boolean compressed sensing and noisy group testing. *IEEE Trans. Inform. Theory* **58** 1880–1901. [MR2932872](#)
- [5] BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of lasso and Dantzig selector. *Ann. Statist.* **37** 1705–1732. [MR2533469](#)
- [6] BOUCHERON, S., LUGOSI, G. and MASSART, P. (2013). *Concentration Inequalities: A Nonasymptotic Theory of Independence*. OUP Oxford.
- [7] CANDÈS, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *Ann. Statist.* **35** 2313–2351. [MR2382644](#)
- [8] CANDÈS, E. J. and TAO, T. (2005). Decoding by linear programming. *IEEE Trans. Inform. Theory* **51** 4203–4215.
- [9] CHAPPELLE, O., DO, C. B., TEO, C. H., LE, Q. V. and SMOLA, A. J. (2009). Tighter bounds for structured estimation. In *Advances in Neural Information Processing Systems* 281–288.
- [10] CHEN, Y. and CANDÈS, E. (2015). Solving random quadratic systems of equations is nearly as easy as solving linear systems. In *Advances in Neural Information Processing Systems* 739–747.
- [11] DASKALAKIS, C., TZAMOS, C. and ZAMPETAKIS, M. (2016). Ten steps of EM suffice for mixtures of two Gaussians. Available at [arXiv:1609.00368](https://arxiv.org/abs/1609.00368).
- [12] DONOHO, D. L. (2006). Compressed sensing. *IEEE Trans. Inform. Theory* **52** 1289–1306. [MR2241189](#)
- [13] DUBROVIN, B. A., FOMENKO, A. T. and NOVIKOV, S. P. (2012). *Modern Geometry—Methods and Applications: Part II: The Geometry and Topology of Manifolds* **104**. Springer, Berlin.
- [14] FISHER, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* **222** 309–368.

- [15] FISHER, R. A. (1925). Theory of statistical estimation. In *Mathematical Proceedings of the Cambridge Philosophical Society* **22** 700–725. Cambridge Univ Press, Cambridge.
- [16] GE, R., HUANG, F., JIN, C. and YUAN, Y. (2015). Escaping from saddle points—online stochastic gradient for tensor decomposition. In *Conference on Learning Theory* 797–842.
- [17] GUILLEMIN, V. and POLLACK, A. (2010). *Differential Topology* **370**. AMS, Providence.
- [18] HUBER, P. J. (1973). Robust regression: Asymptotics, conjectures and Monte Carlo. *Ann. Statist.* **1** 799–821. [MR0356373](#)
- [19] KESHAVAN, R. H., OH, S. and MONTANARI, A. (2009). Matrix completion from a few entries. In *IEEE International Symposium on Information Theory, 2009. ISIT 2009* 324–328. IEEE, New York.
- [20] LASKA, J. N. and BARANIUK, R. G. (2012). Regime change: Bit-depth versus measurement-rate in compressive sensing. *IEEE Trans. Signal Process.* **60** 3496–3505.
- [21] LASKA, J. N., WEN, Z., YIN, W. and BARANIUK, R. G. (2011). Trust, but verify: Fast and accurate signal recovery from 1-bit compressive measurements. *IEEE Trans. Signal Process.* **59** 5289–5301.
- [22] LECUN, Y., BENGIO, Y. and HINTON, G. (2015). Deep learning. *Nature* **521** 436–444.
- [23] LOH, P.-L. (2017). Statistical consistency and asymptotic normality for high-dimensional robust  $M$ -estimators. *Ann. Statist.* **45** 866–896. [MR3650403](#)
- [24] LOH, P.-L. and WAINWRIGHT, M. J. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *Ann. Statist.* **40** 1637–1664. [MR3015038](#)
- [25] LOH, P.-L. and WAINWRIGHT, M. J. (2013). Regularized  $m$ -estimators with nonconvexity: Statistical and algorithmic theory for local optima. In *Advances in Neural Information Processing Systems* 476–484.
- [26] LOZANO, A. C. and MEINSHAUSEN, N. (2013). Minimum distance estimation for robust high-dimensional regression. Available at [arXiv:1307.3227](#).
- [27] MEI, S., BAI, Y. and MONTANARI, A. (2018). Supplement to “The landscape of empirical risk for nonconvex losses.” DOI:[10.1214/17-AOS1637SUPP](#).
- [28] MILNOR, J. (1963). *Morse Theory* **51**. Princeton Univ. Press, Princeton.
- [29] MILNOR, J. W. (1997). *Topology from the Differentiable Viewpoint*. Princeton Univ. Press, Princeton, NH.
- [30] MONTANARI, A. and RICHARD, E. (2014). A statistical model for tensor PCA. In *Advances in Neural Information Processing Systems* 2897–2905.
- [31] 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. *Statist. Sci.* **27** 538–557. [MR3025133](#)
- [32] NESTEROV, Y. (2013). *Introductory Lectures on Convex Optimization: A Basic Course* **87**. Springer Science & Business Media, New York.
- [33] NESTEROV, Y. (2013). Gradient methods for minimizing composite functions. *Math. Program.* **140** 125–161.
- [34] NGUYEN, T. and SANNER, S. (2013). Algorithms for direct 0–1 loss optimization in binary classification. In *Proceedings of the 30th International Conference on Machine Learning* 1085–1093.
- [35] PENG, J., ZHU, J., BERGAMASCHI, A., HAN, W., NOH, D.-Y., POLLACK, J. R. and WANG, P. (2010). Regularized multivariate regression for identifying master predictors with application to integrative genomics study of breast cancer. *Ann. Appl. Stat.* **4** 53–77. [MR2758084](#)
- [36] PLAN, Y. and VERSHYNIN, R. (2013). One-bit compressed sensing by linear programming. *Comm. Pure Appl. Math.* **66** 1275–1297. [MR3069959](#)

- [37] PLAN, Y. and VERSHYNIN, R. (2013). Robust 1-bit compressed sensing and sparse logistic regression: A convex programming approach. *IEEE Trans. Inform. Theory* **59** 482–494.
- [38] PLAN, Y., VERSHYNIN, R. and YUDOVINA, E. (2014). High-dimensional estimation with geometric constraints. Preprint. Available at [arXiv:1404.3749](https://arxiv.org/abs/1404.3749).
- [39] ROBBINS, H. and MONRO, S. (1951). A stochastic approximation method. *Ann. Math. Stat.* **22** 400–407. [MR0042668](https://doi.org/10.2307/2283373)
- [40] SERDOBOLSKII, V. I. (2013). *Multivariate Statistical Analysis: A High-Dimensional Approach* **41**. Springer, New York.
- [41] SUN, J., QU, Q. and WRIGHT, J. (2016). A geometric analysis of phase retrieval. Available at [arXiv:1602.06664](https://arxiv.org/abs/1602.06664).
- [42] TSITSIKLIS, J. N., BERTSEKAS, D. P. and ATHANS, M. (1984). Distributed asynchronous deterministic and stochastic gradient optimization algorithms. In *1984 American Control Conference* 484–489.
- [43] VAN DE GEER, S. A. (2000). *Applications of Empirical Process Theory. Cambridge Series in Statistical and Probabilistic Mathematics* **6**. Cambridge Univ. Press, Cambridge. [MR1739079](https://doi.org/10.1017/CBO9780511526458)
- [44] VAPNIK, V. N. (1998). *Statistical Learning Theory*. Wiley, New York. [MR1641250](https://doi.org/10.1007/978-1-4939-9829-9)
- [45] WU, Y. and LIU, Y. (2012). Robust truncated hinge loss support vector machines. *J. Amer. Statist. Assoc.* **102** 974–983. [MR2411659](https://doi.org/10.1198/01621451101659)
- [46] XU, J., HSU, D. and MALEKI, A. (2016). Global analysis of expectation maximization for mixtures of two Gaussians. Preprint. Available at [arXiv:1608.07630](https://arxiv.org/abs/1608.07630).
- [47] YANG, Z., WANG, Z., LIU, H., ELDAR, Y. C. and ZHANG, T. (2015). Sparse nonlinear regression: Parameter estimation and asymptotic inference. Available at [arXiv:1511.04514](https://arxiv.org/abs/1511.04514).

S. MEI  
STANFORD UNIVERSITY  
HUANG BUILDING 475 VIA ORTEGA  
STANFORD, CALIFORNIA 94305  
USA  
E-MAIL: [songmei@stanford.edu](mailto:songmei@stanford.edu)

Y. BAI  
A. MONTANARI  
STANFORD UNIVERSITY  
390 SERRA MALL  
STANFORD, CALIFORNIA 94305  
USA  
E-MAIL: [yub@stanford.edu](mailto:yub@stanford.edu)  
[montanar@stanford.edu](mailto:montanar@stanford.edu)