

---

# A unified framework for high-dimensional analysis of $M$ -estimators with decomposable regularizers

---

**Sahand Negahban**  
Department of EECS  
UC Berkeley  
sahand.n@eecs.berkeley.edu

**Pradeep Ravikumar**  
Department of Computer Sciences  
UT Austin  
pradeepr@cs.utexas.edu

**Martin J. Wainwright**  
Department of Statistics  
Department of EECS  
UC Berkeley  
wainwrig@eecs.berkeley.edu

**Bin Yu**  
Department of Statistics  
Department of EECS  
UC Berkeley  
binyu@stat.berkeley.edu

## Abstract

High-dimensional statistical inference deals with models in which the number of parameters  $p$  is comparable to or larger than the sample size  $n$ . Since it is usually impossible to obtain consistent procedures unless  $p/n \rightarrow 0$ , a line of recent work has studied models with various types of structure (e.g., sparse vectors; block-structured matrices; low-rank matrices; Markov assumptions). In such settings, a general approach to estimation is to solve a regularized convex program (known as a regularized  $M$ -estimator) which combines a loss function (measuring how well the model fits the data) with some regularization function that encourages the assumed structure. The goal of this paper is to provide a unified framework for establishing consistency and convergence rates for such regularized  $M$ -estimators under high-dimensional scaling. We state one main theorem and show how it can be used to re-derive several existing results, and also to obtain several new results on consistency and convergence rates. Our analysis also identifies two key properties of loss and regularization functions, referred to as restricted strong convexity and decomposability, that ensure the corresponding regularized  $M$ -estimators have fast convergence rates.

## 1 Introduction

In many fields of science and engineering (among them genomics, financial engineering, natural language processing, remote sensing, and social network analysis), one encounters statistical inference problems in which the number of predictors  $p$  is comparable to or even larger than the number of observations  $n$ . Under this type of high-dimensional scaling, it is usually impossible to obtain statistically consistent estimators unless one restricts to subclasses of models with particular structure. For instance, the data might be sparse in a suitably chosen basis, could lie on some manifold, or the dependencies among the variables might have Markov structure specified by a graphical model.

In such settings, a common approach to estimating model parameters is through the use of a *regularized  $M$ -estimator*, in which some loss function (e.g., the negative log-likelihood of the data) is regularized by a function appropriate to the assumed structure. Such estimators may also be interpreted from a Bayesian perspective as maximum a posteriori estimates, with the regularizer reflecting prior information. In this paper, we study such regularized  $M$ -estimation procedures, and attempt to provide a unifying framework that both recovers some existing results and provides

new results on consistency and convergence rates under high-dimensional scaling. We illustrate some applications of this general framework via three running examples of constrained parametric structures. The first class is that of *sparse vector* models; we consider both the case of “hard-sparse” models which involve an explicit constraint on the number of non-zero model parameters, and also a class of “weak-sparse” models in which the ordered coefficients decay at a certain rate. Second, we consider *block-sparse models*, in which the parameters are matrix-structured, and entire rows are either zero or not. Our third class is that of low-rank matrices, which arise in system identification, collaborative filtering, and other types of matrix completion problems.

To motivate the need for a unified analysis, let us provide a brief (and hence necessarily incomplete) overview of the broad range of past and on-going work on high-dimensional inference. For the case of sparse regression, a popular regularizer is the  $\ell_1$  norm of the parameter vector, which is the sum of the absolute values of the parameters. A number of researchers have studied the Lasso [15, 3] as well as the closely related Dantzig selector [2] and provided conditions on various aspects of its behavior, including  $\ell_2$ -error bounds [7, 1, 21, 2] and model selection consistency [22, 19, 6, 16]. For generalized linear models (GLMs) and exponential family models, estimators based on  $\ell_1$ -regularized maximum likelihood have also been studied, including results on risk consistency [18] and model selection consistency [11]. A body of work has focused on the case of estimating Gaussian graphical models, including convergence rates in Frobenius and operator norm [14], and results on operator norm and model selection consistency [12]. Motivated by inference problems involving block-sparse matrices, other researchers have proposed block-structured regularizers [17, 23], and more recently, high-dimensional consistency results have been obtained for model selection and parameter consistency [4, 8].

In this paper, we derive a single main theorem, and show how we are able to rederive a wide range of known results on high-dimensional consistency, as well as some novel ones, including estimation error rates for low-rank matrices, sparse matrices, and “weakly”-sparse vectors. Due to space constraints, many of the technical details are deferred to the full-length version of this conference paper.

## 2 Problem formulation and some key properties

In this section, we begin with a precise formulation of the problem, and then develop some key properties of the regularizer and loss function. In particular, we define a notion of *decomposability* for regularizing functions  $r$ , and then prove that when it is satisfied, the error  $\hat{\Delta} = \hat{\theta} - \theta^*$  of the regularized  $M$ -estimator must satisfy certain constraints. We use these constraints to define a notion of *restricted strong convexity* that the loss function must satisfy.

### 2.1 Problem set-up

Consider a random variable  $Z$  with distribution  $\mathbb{P}$  taking values in a set  $\mathcal{Z}$ . Let  $Z_1^n := \{Z_1, \dots, Z_n\}$  denote  $n$  observations drawn in an i.i.d. manner from  $\mathbb{P}$ , and suppose  $\theta^* \in \mathbb{R}^p$  is some parameter of this distribution. We consider the problem of estimating  $\theta^*$  from the data  $Z_1^n$ , and in order to do so, we consider the following class of regularized  $M$ -estimators. Let  $\mathcal{L} : \mathbb{R}^p \times \mathcal{Z}^n \mapsto \mathbb{R}$  be some loss function that assigns a cost to any parameter  $\theta \in \mathbb{R}^p$ , for a given set of observations  $Z_1^n$ . Let  $r : \mathbb{R}^p \mapsto \mathbb{R}$  denote a regularization function. We then consider the regularized  $M$ -estimator given by

$$\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \{ \mathcal{L}(\theta; Z_1^n) + \lambda_n r(\theta) \}, \quad (1)$$

where  $\lambda_n > 0$  is a user-defined regularization penalty. For ease of notation, in the sequel, we adopt the shorthand  $\mathcal{L}(\theta)$  for  $\mathcal{L}(\theta; Z_1^n)$ . Throughout the paper, we assume that the loss function  $\mathcal{L}$  is convex and differentiable, and that the regularizer  $r$  is a norm.

Our goal is to provide general techniques for deriving bounds on the error  $\hat{\theta} - \theta^*$  in some error metric  $d$ . A common example is the  $\ell_2$ -norm  $d(\hat{\theta} - \theta^*) := \|\hat{\theta} - \theta^*\|_2$ . As discussed earlier, high-dimensional parameter estimation is made possible by structural constraints on  $\theta^*$  such as sparsity, and we will see that the behavior of the error is determined by how well these constraints are captured by the regularization function  $r(\cdot)$ . We now turn to the properties of the regularizer  $r$  and the loss function  $\mathcal{L}$  that underlie our analysis.

## 2.2 Decomposability

Our first condition requires that the regularization function  $r$  be decomposable, in a sense to be defined precisely, with respect to a family of subspaces. This notion is a formalization of the manner in which the regularization function imposes constraints on possible parameter vectors  $\theta^* \in \mathbb{R}^p$ . We begin with some abstract definitions, which we then illustrate with a number of concrete examples. Take some arbitrary inner product space  $\mathcal{H}$ , and let  $\|\cdot\|_2$  denote the norm induced by the inner product. Consider a pair  $(A, B)$  of subspaces of  $\mathcal{H}$  such that  $A \subseteq B^\perp$ . For a given subspace  $A$  and vector  $u \in \mathcal{H}$ , we let  $\pi_A(u) := \operatorname{argmin}_{v \in A} \|u - v\|_2$  denote the orthogonal projection of  $u$  onto  $A$ . We let  $\mathcal{V} = \{(A, B) \mid A \subseteq B^\perp\}$  be a collection of subspace pairs. For a given statistical model, our goal is to construct subspace collections  $\mathcal{V}$  such that for any given  $\theta^*$  from our model class, there exists a pair  $(A, B) \in \mathcal{V}$  with  $\|\pi_A(\theta^*)\|_2 \approx \|\theta^*\|_2$ , and  $\|\pi_B(\theta^*)\|_2 \approx 0$ . Of most interest to us are subspace pairs  $(A, B)$  in which this property holds but the subspace  $A$  is relatively small and  $B$  is relatively large. Note that  $A$  represents the constraints underlying our model class, and imposed by our regularizer. For the bulk of the paper, we assume that  $\mathcal{H} = \mathbb{R}^p$  and use the standard Euclidean inner product (which should be assumed unless otherwise specified).

As a first concrete (but toy) example, consider the model class of all vectors  $\theta^* \in \mathbb{R}^p$ , and the subspace collection  $\mathcal{T}$  that consists of a single subspace pair  $(A, B) = (\mathbb{R}^p, 0)$ . We refer to this choice ( $\mathcal{V} = \mathcal{T}$ ) as the *trivial subspace collection*. In this case, for any  $\theta^* \in \mathbb{R}^p$ , we have  $\pi_A(\theta^*) = \theta^*$  and  $\pi_B(\theta^*) = 0$ . Although this collection satisfies our desired property, it is not so useful since  $A = \mathbb{R}^p$  is a very large subspace. As a second example, consider the class of  $s$ -sparse parameter vectors  $\theta^* \in \mathbb{R}^p$ , meaning that  $\theta_i^* \neq 0$  only if  $i \in S$ , where  $S$  is some  $s$ -sized subset of  $\{1, 2, \dots, p\}$ . For any given subset  $S$  and its complement  $S^c$ , let us define the subspaces

$$A(S) = \{\theta \in \mathbb{R}^p \mid \theta_{S^c} = 0\}, \quad \text{and} \quad B(S) = \{\theta \in \mathbb{R}^p \mid \theta_S = 0\},$$

and the  $s$ -sparse subspace collection  $\mathcal{S} = \{(A(S), B(S)) \mid S \subset \{1, \dots, p\}, |S| = s\}$ . With this set-up, for any  $s$ -sparse parameter vector  $\theta^*$ , we are guaranteed that there exists some  $(A, B) \in \mathcal{S}$  such that  $\pi_A(\theta^*) = \theta^*$  and  $\pi_B(\theta^*) = 0$ . In this case, the property is more interesting, since the subspaces  $A(S)$  are relatively small as long as  $|S| = s \ll p$ .

With this set-up, we say that the regularizer  $r$  is *decomposable* with respect to a given subspace pair  $(A, B)$  if

$$r(u + z) = r(u) + r(z) \quad \text{for all } u \in A \text{ and } z \in B. \quad (2)$$

In our subsequent analysis, we impose the following condition on the regularizer:

**Definition 1.** The regularizer  $r$  is decomposable with respect to a given subspace collection  $\mathcal{V}$ , meaning that it is decomposable for each subspace pair  $(A, B) \in \mathcal{V}$ .

Note that any regularizer is decomposable with respect to the trivial subspace collection  $\mathcal{T} = \{(\mathbb{R}^p, 0)\}$ . It will be of more interest to us when the regularizer decomposes with respect to a larger collection  $\mathcal{V}$  that includes subspace pairs  $(A, B)$  in which  $A$  is relatively small and  $B$  is relatively large. Let us illustrate with some examples.

- *Sparse vectors and  $\ell_1$  norm regularization.* Consider a model involving  $s$ -sparse regression vectors  $\theta^* \in \mathbb{R}^p$ , and recall the definition of the  $s$ -sparse subspace collection  $\mathcal{S}$  discussed above. We claim that the  $\ell_1$ -norm regularizer  $r(u) = \|u\|_1$  is decomposable with respect to  $\mathcal{S}$ . Indeed, for any  $s$ -sized subset  $S$  and vectors  $u \in A(S)$  and  $v \in B(S)$ , we have  $\|u + v\|_1 = \|u\|_1 + \|v\|_1$ , as required.
- *Group-structured sparse matrices and  $\ell_{1,q}$  matrix norms.* Various statistical problems involve matrix-valued parameters  $\Theta \in \mathbb{R}^{k \times m}$ ; examples include multivariate regression problems or (inverse) covariance matrix estimation. We can define an inner product on such matrices via  $\langle\langle \Theta, \Sigma \rangle\rangle = \operatorname{trace}(\Theta^T \Sigma)$  and the induced (Frobenius) norm  $\sum_{i=1}^k \sum_{j=1}^m \Theta_{i,j}^2$ . Let us suppose that  $\Theta$  satisfies a group sparsity condition, meaning that the  $i^{\text{th}}$  row, denoted  $\Theta_i$ , is non-zero only if  $i \in S \subseteq \{1, \dots, k\}$  and the cardinality of  $S$  is controlled. For a given subset  $S$ , we can define the subspace pair

$$B(S) = \{\Theta \in \mathbb{R}^{k \times m} \mid \Theta_i = 0 \text{ for all } i \in S^c\}, \quad \text{and} \quad A(S) = (B(S))^\perp,$$

For some fixed  $s \leq k$ , we then consider the collection

$$\mathcal{V} = \{(A(S), B(S)) \mid S \subset \{1, \dots, k\}, |S| = s\},$$

which is a group-structured analog of the  $s$ -sparse set  $\mathcal{S}$  for vectors. For any  $q \in [1, \infty]$ , now suppose that the regularizer is the  $\ell_1/\ell_q$  matrix norm, given by  $r(\Theta) = \sum_{i=1}^k [\sum_{j=1}^m |\Theta_{ij}|^q]^{1/q}$ , corresponding to applying the  $\ell_q$  norm to each row and then taking the  $\ell_1$ -norm of the result. It can be seen that the regularizer  $r(\Theta) = \|\Theta\|_{1,q}$  is decomposable with respect to the collection  $\mathcal{V}$ .

- *Low-rank matrices and nuclear norm.* The estimation of low-rank matrices arises in various contexts, including principal component analysis, spectral clustering, collaborative filtering, and matrix completion. In particular, consider the class of matrices  $\Theta \in \mathbb{R}^{k \times m}$  that have rank  $r \leq \min\{k, m\}$ . For any given matrix  $\Theta$ , we let  $\text{row}(\Theta) \subseteq \mathbb{R}^m$  and  $\text{col}(\Theta) \subseteq \mathbb{R}^k$  denote its row space and column space respectively. For a given pair of  $r$ -dimensional subspaces  $U \subseteq \mathbb{R}^k$  and  $V \subseteq \mathbb{R}^m$ , we define a pair of subspaces  $A(U, V)$  and  $B(U, V)$  of  $\mathbb{R}^{k \times m}$  as follows:

$$A(U, V) := \{\Theta \in \mathbb{R}^{k \times m} \mid \text{row}(\Theta) \subseteq V, \text{col}(\Theta) \subseteq U\}, \quad \text{and} \quad (3a)$$

$$B(U, V) := \{\Theta \in \mathbb{R}^{k \times m} \mid \text{row}(\Theta) \subseteq V^\perp, \text{col}(\Theta) \subseteq U^\perp\}. \quad (3b)$$

Note that  $A(U, V) \subseteq B^\perp(U, V)$ , as is required by our construction. We then consider the collection  $\mathcal{V} = \{(A(U, V), B(U, V)) \mid U \subseteq \mathbb{R}^k, V \subseteq \mathbb{R}^m\}$ , where  $(U, V)$  range over all pairs of  $r$ -dimensional subspaces. Now suppose that we regularize with the nuclear norm  $r(\Theta) = \|\Theta\|_1$ , corresponding to the sum of the singular values of the matrix  $\Theta$ . It can be shown that the nuclear norm is decomposable with respect to  $\mathcal{V}$ . Indeed, since any pair of matrices  $M \in A(U, V)$  and  $M' \in B(U, V)$  have orthogonal row and column spaces, we have  $\|M + M'\|_1 = \|M\|_1 + \|M'\|_1$  (e.g., see the paper [13]).

Thus, we have demonstrated various models and regularizers in which decomposability is satisfied with interesting subspace collections  $\mathcal{V}$ . We now show that decomposability has important consequences for the error  $\widehat{\Delta} = \widehat{\theta} - \theta^*$ , where  $\widehat{\theta} \in \mathbb{R}^p$  is any optimal solution of the regularized  $M$ -estimation procedure (1). In order to state a lemma that captures this fact, we need to define the dual norm of the regularizer, given by  $r^*(v) := \sup_{u \in \mathbb{R}^p} \frac{\langle u, v \rangle}{r(u)}$ . For the regularizers of interest, the dual norm can be obtained via some easy calculations. For instance, given a vector  $\theta \in \mathbb{R}^p$  and  $r(\theta) = \|\theta\|_1$ , we have  $r^*(\theta) = \|\theta\|_\infty$ . Similarly, given a matrix  $\Theta \in \mathbb{R}^{k \times m}$  and the nuclear norm regularizer  $r(\Theta) = \|\Theta\|_1$ , we have  $r^*(\Theta) = \|\Theta\|_2$ , corresponding to the operator norm (or maximal singular value).

**Lemma 1.** *Suppose  $\widehat{\theta}$  is an optimal solution of the regularized  $M$ -estimation procedure (1), with associated error  $\Delta = \widehat{\theta} - \theta^*$ . Furthermore, suppose that the regularization penalty is strictly positive with  $\lambda_n \geq 2r^*(\nabla \mathcal{L}(\theta^*))$ . Then for any  $(A, B) \in \mathcal{V}$*

$$r(\pi_B(\widehat{\Delta})) \leq 3r(\pi_{B^\perp}(\widehat{\Delta})) + 4r(\pi_{A^\perp}(\theta^*)).$$

This property plays an essential role in our definition of restricted strong convexity and subsequent analysis.

### 2.3 Restricted Strong Convexity

Next we state our assumption on the loss function  $\mathcal{L}$ . In general, guaranteeing that  $\mathcal{L}(\widehat{\theta}) - \mathcal{L}(\theta^*)$  is small is *not sufficient* to show that  $\widehat{\theta}$  and  $\theta^*$  are close. (As a trivial example, consider a loss function that is identically zero.) The standard way to ensure that a function is “not too flat” is via the notion of strong convexity—in particular, by requiring that there exist some constant  $\gamma > 0$  such that  $\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) - \langle \nabla \mathcal{L}(\theta^*), \Delta \rangle \geq \gamma d^2(\Delta)$  for all  $\Delta \in \mathbb{R}^p$ . In the high-dimensional setting, where the number of parameters  $p$  may be much larger than the sample size, the strong convexity assumption need not be satisfied. As a simple example, consider the usual linear regression model  $y = X\theta^* + w$ , where  $y \in \mathbb{R}^n$  is the response vector,  $\theta^* \in \mathbb{R}^p$  is the unknown parameter vector,  $X \in \mathbb{R}^{n \times p}$  is the design matrix, and  $w \in \mathbb{R}^n$  is a noise vector, with i.i.d. zero mean elements. The least-squares loss is given by  $\mathcal{L}(\theta) = \frac{1}{2n} \|y - X\theta\|_2^2$ , and has the Hessian  $H(\theta) = \frac{1}{n} X^T X$ . It is easy to check that the  $p \times p$  matrix  $H(\theta)$  will be rank-deficient whenever  $p > n$ , showing that the least-squares loss cannot be strongly convex (with respect to  $d(\cdot) = \|\cdot\|_2$ ) when  $p > n$ .

Herein lies the utility of Lemma 1: it guarantees that the error  $\widehat{\Delta}$  must lie within a restricted set, so that we only need the loss function to be strongly convex for a limited set of directions. More precisely, we have:

**Definition 2.** Given some subset  $\mathcal{C} \subseteq \mathbb{R}^p$  and error norm  $d(\cdot)$ , we say that the loss function  $\mathcal{L}$  satisfies *restricted strong convexity* (RSC) (with respect to  $d(\cdot)$ ) with parameter  $\gamma(\mathcal{L}) > 0$  over  $\mathcal{C}$  if

$$\mathcal{L}(\theta^* + \Delta) - \mathcal{L}(\theta^*) - \langle \nabla \mathcal{L}(\theta^*), \Delta \rangle \geq \gamma(\mathcal{L}) d^2(\Delta) \quad \text{for all } \Delta \in \mathcal{C}. \quad (4)$$

In the statement of our results, we will be interested in loss functions that satisfy RSC over sets  $\mathcal{C}(A, B, \epsilon)$  that are indexed by a subspace pair  $(A, B)$  and a tolerance  $\epsilon \geq 0$  as follows:

$$\mathcal{C}(A, B, \epsilon) := \{ \Delta \in \mathbb{R}^p \mid r(\pi_B(\Delta)) \leq 3r(\pi_{B^\perp}(\Delta)) + 4r(\pi_{A^\perp}(\theta^*)), \quad d(\Delta) \geq \epsilon \}. \quad (5)$$

In the special case of least-squares regression with hard sparsity constraints, the RSC condition corresponds to a lower bound on the sparse eigenvalues of the Hessian matrix  $X^T X$ , and is essentially equivalent to a restricted eigenvalue condition introduced by Bickel et al. [1].

### 3 Convergence rates

We are now ready to state a general result that provides bounds and hence convergence rates for the error  $d(\hat{\theta} - \theta^*)$ . Although it may appear somewhat abstract at first sight, we illustrate that this result has a number of concrete consequences for specific models. In particular, we recover the best known results about estimation in  $s$ -sparse models with general designs [1, 7], as well as a number of new results, including convergence rates for estimation under  $\ell_q$ -sparsity constraints, estimation in sparse generalized linear models, estimation of block-structured sparse matrices and estimation of low-rank matrices.

In addition to the regularization parameter  $\lambda_n$  and RSC constant  $\gamma(\mathcal{L})$  of the loss function, our general result involves a quantity that relates the error metric  $d$  to the regularizer  $r$ ; in particular, for any set  $A \subseteq \mathbb{R}^p$ , we define

$$\Psi(A) := \sup_{\{u \in \mathbb{R}^p \mid d(u)=1\}} r(u), \quad (6)$$

so that  $r(u) \leq \Psi(A)d(u)$  for  $u \in A$ .

**Theorem 1 (Bounds for general models).** *For a given subspace collection  $\mathcal{V}$ , suppose that the regularizer  $r$  is decomposable, and consider the regularized  $M$ -estimator (1) with  $\lambda_n \geq 2r^*(\nabla \mathcal{L}(\theta^*))$ . Then, for any pair of subspaces  $(A, B) \in \mathcal{V}$  and tolerance  $\epsilon \geq 0$  such that the loss function  $\mathcal{L}$  satisfies restricted strong convexity over  $\mathcal{C}(A, B, \epsilon)$ , we have*

$$d(\hat{\theta} - \theta^*) \leq \max \left\{ \epsilon, \frac{1}{\gamma(\mathcal{L})} \left[ 2\Psi(B^\perp) \lambda_n + \sqrt{2\lambda_n \gamma(\mathcal{L}) r(\pi_{A^\perp}(\theta^*))} \right] \right\}. \quad (7)$$

The proof is motivated by arguments used in past work on high-dimensional estimation (e.g., [9, 14]); we provide the details in the full-length version. The remainder of this paper is devoted to illustrations of the consequences of Theorem 1 for specific models. In all of these uses of Theorem 1, we choose the regularization parameter as small as possible—namely,  $\lambda_n = 2r^*(\nabla \mathcal{L}(\theta^*))$ . Although Theorem 1 allows for more general choices, in this conference version, we focus exclusively on the case when  $d(\cdot)$  to be the  $\ell_2$ -norm. In addition, we choose a tolerance parameter  $\epsilon = 0$  for all of the results except for the weak-sparse models treated in Section 3.1.2.

#### 3.1 Bounds for linear regression

Consider the standard linear regression model  $y = X\theta^* + w$ , where  $\theta^* \in \mathbb{R}^p$  is the regression vector,  $X \in \mathbb{R}^{n \times p}$  is the design matrix, and  $w \in \mathbb{R}^n$  is a noise vector. Given the observations  $(y, X)$ , our goal is to estimate the regression vector  $\theta^*$ . Without any structural constraints on  $\theta^*$ , we can apply Theorem 1 with the trivial subspace collection  $\mathcal{T} = \{(\mathbb{R}^p, 0)\}$  to establish a rate  $\|\hat{\theta} - \theta^*\|_2 = \mathcal{O}(\sigma\sqrt{p/n})$  for ridge regression, which holds as long as  $X$  is full-rank (and hence requires  $n > p$ ). Here we consider the sharper bounds that can be obtained when it is assumed that  $\theta^*$  is an  $s$ -sparse vector.

### 3.1.1 Lasso estimates of hard sparse models

More precisely, let us consider estimating an  $s$ -sparse regression vector  $\theta^*$  by solving the Lasso program  $\hat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|y - X\theta\|_2^2 + \lambda_n \|\theta\|_1 \right\}$ . The Lasso is a special case of our  $M$ -estimator (1) with  $r(\theta) = \|\theta\|_1$ , and  $\mathcal{L}(\theta) = \frac{1}{2n} \|y - X\theta\|_2^2$ . Recall the definition of the  $s$ -sparse subspace collection  $\mathcal{S}$  from Section 2.2. For this problem, let us set  $\epsilon = 0$  so that the restricted strong convexity set (5) reduces to  $\mathcal{C}(A, B, 0) = \{\Delta \in \mathbb{R}^p \mid \|\Delta_{S^c}\|_1 \leq 3\|\Delta_S\|_1\}$ . Establishing restricted strong convexity for the least-squares loss is equivalent to ensuring the following bound on the design matrix:

$$\|X\theta\|_2^2/n \geq \gamma(\mathcal{L}) \|\theta\|_2^2 \quad \text{for all } \theta \in \mathbb{R}^p \text{ such that } \|\theta_S\|_1 \leq 3\|\theta_{S^c}\|_1. \quad (8)$$

As mentioned previously, this condition is essentially the same as the restricted eigenvalue condition developed by Bickel et al. [1]. In very recent work, Raskutti et al. [10] show that condition (8) holds with high probability for various random ensembles of Gaussian matrices with non-i.i.d. elements.

In addition to the RSC condition, we assume that  $X$  has bounded column norms (specifically,  $\|X_i\|_2 \leq 2\sqrt{n}$  for all  $i = 1, \dots, p$ ), and that the noise vector  $w \in \mathbb{R}^n$  has i.i.d. elements with zero-mean and sub-Gaussian tails (i.e., there exists some constant  $\sigma > 0$  such that  $\mathbb{P}[|w_i| > t] \leq \exp(-t^2/2\sigma^2)$  for all  $t > 0$ ). Under these conditions, we recover as a corollary of Theorem 1 the following known result [1, 7].

**Corollary 1.** *Suppose that the true vector  $\theta^* \in \mathbb{R}^p$  is exactly  $s$ -sparse with support  $S$ , and that the design matrix  $X$  satisfies condition (8). If we solve the the Lasso with  $\lambda_n^2 = \frac{16\sigma^2 \log p}{n}$ , then with probability at least  $1 - c_1 \exp(-c_2 n \lambda_n^2)$ , the solution satisfies*

$$\|\hat{\theta} - \theta^*\|_2 \leq \frac{8\sigma}{\gamma(\mathcal{L})} \sqrt{\frac{s \log p}{n}}. \quad (9)$$

*Proof.* As noted previously, the  $\ell_1$ -regularizer is decomposable for the sparse subspace collection  $\mathcal{S}$ , while condition (8) ensures that RSC holds for all sets  $\mathcal{C}(A, B, 0)$  with  $(A, B) \in \mathcal{S}$ . We must verify that the given choice of regularization satisfies  $\lambda_n \geq 2r^*(\nabla \mathcal{L}(\theta^*))$ . Note that  $r^*(\cdot) = \|\cdot\|_\infty$ , and moreover that  $\nabla \mathcal{L}(\theta^*) = X^T w/n$ . Under the column normalization condition on the design matrix  $X$  and the sub-Gaussian nature of the noise, it follows that  $\|X^T w/n\|_\infty \leq \sqrt{4\sigma^2 \frac{\log p}{n}}$  with high probability. The bound in Theorem 1 is thus applicable, and it remains to compute the form that its different terms take in this special case. For the  $\ell_1$ -regularizer and the  $\ell_2$  error metric, we have  $\Psi(A_S) = \sqrt{|S|}$ . Given the hard sparsity assumption,  $r(\theta_{S^c}^*) = 0$ , so that Theorem 1 implies that  $\|\hat{\theta} - \theta^*\|_2 \leq \frac{2}{\gamma(\mathcal{L})} \sqrt{s} \lambda_n = \frac{8\sigma}{\gamma(\mathcal{L})} \sqrt{\frac{s \log p}{n}}$ , as claimed.  $\square$

### 3.1.2 Lasso estimates of weak sparse models

We now consider models that satisfy a weak sparsity assumption. More concretely, suppose that  $\theta^*$  lies in the  $\ell_q$ -“ball” of radius  $R_q$ —namely, the set  $\mathbb{B}_q(R_q) := \{\theta \in \mathbb{R}^p \mid \sum_{i=1}^p |\theta_i|^q \leq R_q\}$  for some  $q \in (0, 1]$ . Our analysis exploits the fact that any  $\theta^* \in \mathbb{B}_q(R_q)$  can be well approximated by an  $s$ -sparse vector (for an appropriately chosen sparsity index  $s$ ). It is natural to approximate  $\theta^*$  by a vector supported on the set  $S = \{i \mid |\theta_i^*| \geq \tau\}$ . For any choice of threshold  $\tau > 0$ , it can be shown that  $|S| \leq R_q \tau^{-q}$ , and it is optimal to choose  $\tau$  equal to the same regularization parameter  $\lambda_n$  from Corollary 1 (see the full-length version for details). Accordingly, we consider the  $s$ -sparse subspace collection  $\mathcal{S}$  with subsets of size  $s = R_q \lambda_n^{-q}$ . We assume that the noise vector  $w \in \mathbb{R}^n$  is as defined above and that the columns are normalized as in the previous section. We also assume that the matrix  $X$  satisfies the condition

$$\|Xv\|_2 \geq \kappa_1 \|v\|_2 - \kappa_2 \left(\frac{\log p}{n}\right)^{\frac{1}{2}} \|v\|_1 \quad \text{for constants } \kappa_1, \kappa_2 > 0. \quad (10)$$

Raskutti et al. [10] show that this property holds with high probability for suitable Gaussian random matrices. Under this condition, it can be verified that RSC holds with  $\gamma(\mathcal{L}) = \kappa_1/2$  over the set  $\mathcal{C}(A(S), B(S), \epsilon_n)$ , where  $\epsilon_n = (4/\kappa_1 + \sqrt{4/\kappa_1}) R_q^{\frac{1}{2}} \left(\sqrt{\frac{16\sigma^2 \log p}{n}}\right)^{1-q/2}$ . The following result, which we obtain by applying Theorem 1 in this setting, is new to the best of our knowledge:

**Corollary 2.** *Suppose that the true vector  $\theta^* \in \mathbb{B}_q(R_q)$ , and the design matrix  $X$  satisfies condition (10). If we solve the Lasso with  $\lambda_n^2 = \frac{16\sigma^2 \log p}{n}$ , then with probability  $1 - c_1 \exp(-c_2 n \lambda_n^2)$ , the solution satisfies*

$$\|\widehat{\theta} - \theta^*\|_2 \leq R_q^{\frac{1}{2}} \left( \sqrt{\frac{16\sigma^2 \log p}{n}} \right)^{1-q/2} \left[ \frac{2}{\gamma(\mathcal{L})} + \frac{\sqrt{2}}{\sqrt{\gamma(\mathcal{L})}} \right]. \quad (11)$$

We note that both of the rates—for hard-sparsity in Corollary 1 and weak-sparsity in Corollary 2—are known to be optimal<sup>1</sup> in a minimax sense [10].

### 3.2 Bounds for generalized linear models

Our next example is a generalized linear model with canonical link function, where the distribution of response  $y \in \mathcal{Y}$  based on a predictor  $x \in \mathbb{R}^p$  is given by  $p(y \mid x; \theta^*) = \exp(y(\theta^*, x) - a(\langle \theta^*, X \rangle) + d(y))$ , for some fixed functions  $a : \mathbb{R} \mapsto \mathbb{R}$  and  $d : \mathcal{Y} \mapsto \mathbb{R}$ , where  $\|x\|_\infty \leq A$ , and  $|y| \leq B$ . We consider estimating  $\theta^*$  from observations  $\{(x_i, y_i)\}_{i=1}^n$  by  $\ell_1$ -regularized maximum likelihood  $\widehat{\theta} \in \arg \min_{\theta \in \mathbb{R}^p} \left\{ -\frac{1}{n} \langle \theta, \sum_{i=1}^n y_i x_i \rangle + \frac{1}{n} \sum_{i=1}^n a(\langle \theta, x_i \rangle) + \|\theta\|_1 \right\}$ . This is a special case of our  $M$ -estimator (1) with  $\mathcal{L}(\theta) = -\langle \theta, \left(\frac{1}{n} \sum_{i=1}^n y_i x_i\right) \rangle + \frac{1}{n} \sum_{i=1}^n a(\langle \theta, x_i \rangle)$ , and  $r(\theta) = \|\theta\|_1$ . Let  $X \in \mathbb{R}^{n \times p}$  denote the matrix with  $i^{\text{th}}$  row  $x_i$ . For analysis, we again use the  $s$ -sparse subspace collection  $\mathcal{S}$  and  $\epsilon = 0$ . With these choices, it can be verified that an appropriate version of the RSC will hold if the second derivative  $a''$  is strongly convex, and the design matrix  $X$  satisfies a version of the condition (8).

**Corollary 3.** *Suppose that the true vector  $\theta^* \in \mathbb{R}^p$  is exactly  $s$ -sparse with support  $S$ , and the model  $(a, X)$  satisfies an RSC condition. Suppose that we compute the  $\ell_1$ -regularized MLE with  $\lambda_n^2 = \frac{32A^2 B^2 \log p}{n}$ . Then with probability  $1 - c_1 \exp(-c_2 n \lambda_n^2)$ , the solution satisfies*

$$\|\widehat{\theta} - \theta^*\|_2 \leq \frac{16AB}{\gamma(\mathcal{L})} \sqrt{\frac{s \log p}{n}}. \quad (12)$$

We defer the proof to the full-length version due to space constraints.

### 3.3 Bounds for sparse matrices

In this section, we consider some extensions of our results to estimation of regression matrices. Various authors have proposed extensions of the Lasso based on regularizers that have more structure than the  $\ell_1$  norm (e.g., [17, 20, 23, 5]). Such regularizers allow one to impose various types of block-sparsity constraints, in which groups of parameters are assumed to be active (or inactive) simultaneously. We assume that the observation model takes on the form  $Y = X\Theta^* + W$ , where  $\Theta^* \in \mathbb{R}^{k \times m}$  is the unknown fixed set of parameters,  $X \in \mathbb{R}^{n \times k}$  is the design matrix, and  $W \in \mathbb{R}^{n \times m}$  is the noise matrix. As a loss function, we use the Frobenius norm  $\frac{1}{n} \mathcal{L}(\Theta) = \|Y - X\Theta\|_F^2$ , and as a regularizer, we use the  $\ell_{1,q}$ -matrix norm for some  $q \geq 1$ , which takes the form  $\|\Theta\|_{1,q} = \sum_{i=1}^k \|(\Theta_{i1}, \dots, \Theta_{im})\|_q$ . We refer to the resulting estimator as the  $q$ -group Lasso. We define the quantity  $\eta(m; q) = 1$  if  $q \in (1, 2]$  and  $\eta(m; q) = m^{1/2-1/q}$  if  $q > 2$ . We then set the regularization parameter as follows:

$$\lambda_n = \begin{cases} \frac{4\sigma}{\sqrt{n}} [\eta(m; q) \sqrt{\log k} + C_q m^{1-1/q}] & \text{if } q > 1 \\ 4\sigma \sqrt{\frac{\log(km)}{n}} & \text{for } q = 1. \end{cases}$$

**Corollary 4.** *Suppose that the true parameter matrix  $\Theta^*$  has non-zero rows only for indices  $i \in S \subseteq \{1, \dots, k\}$  where  $|S| = s$ , and that the design matrix  $X \in \mathbb{R}^{n \times k}$  satisfies condition (8). Then with probability at least  $1 - c_1 \exp(-c_2 n \lambda_n^2)$ , the  $q$ -block Lasso solution satisfies*

$$\|\widehat{\Theta} - \Theta^*\|_F \leq \frac{2}{\gamma(\mathcal{L})} \Psi(S) \lambda_n. \quad (13)$$

<sup>1</sup>Raskutti et al. [10] show that the rate (11) is achievable by solving the computationally intractable problem of minimizing  $\mathcal{L}(\theta)$  over the  $\ell_q$ -ball.

The proof is provided in the full-length version; here we consider three special cases of the above result. A simple argument shows that  $\Psi(S) = \sqrt{s}$  if  $q \geq 2$ , and  $\Psi(S) = m^{1/q-1/2} \sqrt{s}$  if  $q \in [1, 2]$ . For  $q = 1$ , solving the group Lasso is identical solving a Lasso problem with sparsity  $sm$  and ambient dimension  $km$ , and the resulting upper bound  $\frac{8\sigma}{\gamma(\mathcal{L})} \sqrt{\frac{sm \log(km)}{n}}$  reflects this fact (compare to Corollary 1). For the case  $q = 2$ , Corollary 4 yields the upper bound  $\frac{8\sigma}{\gamma(\mathcal{L})} [\sqrt{\frac{s \log k}{n}} + \sqrt{\frac{sm}{n}}]$ , which also has a natural interpretation: the term  $\frac{s \log k}{n}$  captures the difficulty of finding the  $s$  non-zero rows out of the total  $k$ , whereas the term  $\frac{sm}{n}$  captures the difficulty of estimating the  $sm$  free parameters in the matrix (once the non-zero rows have been determined). We note that recent work by Lounici et al. [4] established the bound  $\mathcal{O}(\frac{\sigma}{\gamma(\mathcal{L})} \sqrt{\frac{c\sqrt{m} s \log k}{n} + \frac{sm}{n}})$ , which is equivalent apart from a term  $\sqrt{m}$ . Finally, for  $q = \infty$ , we obtain the upper bound  $\frac{8\sigma}{\gamma(\mathcal{L})} [\sqrt{\frac{s \log k}{n}} + m \sqrt{\frac{s}{n}}]$ .

### 3.4 Bounds for estimating low rank matrices

Finally, we consider the implications of our main result for the problem of estimating low-rank matrices. This structural assumption is a natural variant of sparsity, and has been studied by various authors (see the paper [13] and references therein). To illustrate our main theorem in this context, let us consider the following instance of low-rank matrix learning. Given a low-rank matrix  $\Theta^* \in \mathbb{R}^{k \times m}$ , suppose that we are given  $n$  noisy observations of the form  $Y_i = \langle X_i, \Theta^* \rangle + W_i$ , where  $W_i \sim N(0, 1)$  and  $\langle A, B \rangle := \text{trace}(A^T B)$ . Such an observation model arises in system identification settings in control theory [13]. The following regularized  $M$ -estimator can be considered in order to estimate the desired low-rank matrix  $\Theta^*$ :

$$\min_{\Theta \in \mathbb{R}^{m \times p}} \frac{1}{2n} \sum_{i=1}^n |Y_i - \langle X_i, \Theta \rangle|^2 + \|\Theta\|_1, \quad (14)$$

where the regularizer,  $\|\Theta\|_1$ , is the nuclear norm, or the sum of the singular values of  $\Theta$ . Recall the rank- $r$  collection  $\mathcal{V}$  defined for low-rank matrices in Section 2.2. Let  $\Theta^* = U\Sigma W^T$  be the singular value decomposition (SVD) of  $\Theta^*$ , so that  $U \in \mathbb{R}^{k \times r}$  and  $W \in \mathbb{R}^{m \times r}$  are orthogonal, and  $\Sigma \in \mathbb{R}^{r \times r}$  is a diagonal matrix. If we let  $A = A(U, W)$  and  $B = B(U, W)$ , then,  $\pi_B(\Theta^*) = 0$ , so that by Lemma 1 we have that  $\|\pi_B(\Delta)\|_1 \leq 3 \|\pi_{B^\perp}(\Delta)\|_1$ . Thus, for restricted strong convexity to hold it can be shown that the design matrices  $X_i$  must satisfy

$$\frac{1}{n} \sum_{i=1}^n |\langle X_i, \Delta \rangle|^2 \geq \gamma(\mathcal{L}) \|\Delta\|_F^2 \quad \text{for all } \Delta \text{ such that } \|\pi_B(\Delta)\|_1 \leq 3 \|\pi_{B^\perp}(\Delta)\|_1, \quad (15)$$

and satisfy the appropriate analog of the column-normalization condition. As with analogous conditions for sparse linear regression, these conditions hold w.h.p. for various non-i.i.d. Gaussian random matrices.<sup>2</sup>

**Corollary 5.** *Suppose that the true matrix  $\Theta^*$  has rank  $r \ll \min(k, m)$ , and that the design matrices  $\{X_i\}$  satisfy condition (15). If we solve the regularized  $M$ -estimator (14) with  $\lambda_n = 16 \frac{\sqrt{k+\sqrt{m}}}{\sqrt{n}}$ , then with probability at least  $1 - c_1 \exp(-c_2(k+m))$ , we have*

$$\|\hat{\Theta} - \Theta^*\|_F \leq \frac{16}{\gamma(\mathcal{L})} \left[ \sqrt{\frac{rk}{n}} + \sqrt{\frac{rm}{n}} \right]. \quad (16)$$

*Proof.* Note that if  $\text{rank}(\Theta^*) = r$ , then  $\|\Theta^*\|_1 \leq \sqrt{r} \|\Theta^*\|_F$  so that  $\Psi(B^\perp) = \sqrt{2r}$ , since the subspace  $B(U, V)^\perp$  consists of matrices with rank at most  $2r$ . All that remains is to show that  $\lambda_n \geq 2r^*(\nabla \mathcal{L}(\Theta^*))$ . Standard analysis gives that the dual norm to  $\|\cdot\|_1$  is the operator norm,  $\|\cdot\|_2$ . Applying this observation we may construct a bound on the operator norm of  $\nabla \mathcal{L}(\Theta^*) = \frac{1}{n} \sum_{i=1}^n X_i W_i$ . Given unit vectors  $u \in \mathbb{R}^k$  and  $v \in \mathbb{R}^m$ ,  $\frac{1}{n} \sum_{i=1}^n |\langle X_i, v u^T \rangle|^2 \leq \|v u^T\|_F^2 = 1$ . Therefore,  $\frac{1}{n} \sum_{i=1}^n (u^T X_i v) W_i \sim N(0, \frac{1}{n})$ . A standard argument shows that the supremum over all unit vectors  $u$  and  $v$  is bounded above by  $8 \frac{\sqrt{k+\sqrt{m}}}{\sqrt{n}}$  with probability at least  $1 - c_1 \exp(-c_2(k+m))$ , verifying that  $\lambda_n \geq 2r^*(\nabla \mathcal{L}(\Theta^*))$  with high probability.  $\square$

<sup>2</sup>This claim involves some use of concentration of measure and Gaussian comparison inequalities analogous to arguments in Raskutti et al. [10]; see the full-length length version for details.



## References

- [1] P. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. Submitted to *Annals of Statistics*, 2008.
- [2] E. Candes and T. Tao. The Dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *Annals of Statistics*, 35(6):2313–2351, 2007.
- [3] S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM J. Sci. Computing*, 20(1):33–61, 1998.
- [4] K. Lounici, M. Pontil, A. B. Tsybakov, and S. van de Geer. Taking advantage of sparsity in multi-task learning. *Arxiv*, 2009.
- [5] L. Meier, S. Van de Geer, and P. Bühlmann. The group lasso for logistic regression. *Journal of the Royal Statistical Society, Series B*, 70:53–71, 2008.
- [6] N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34:1436–1462, 2006.
- [7] N. Meinshausen and B. Yu. Lasso-type recovery of sparse representations for high-dimensional data. *Annals of Statistics*, 37(1):246–270, 2009.
- [8] G. Obozinski, M. J. Wainwright, and M. I. Jordan. Union support recovery in high-dimensional multivariate regression. Technical report, Department of Statistics, UC Berkeley, August 2008.
- [9] S. Portnoy. Asymptotic behavior of M-estimators of  $p$  regression parameters when  $p^2/n$  is large: I. consistency. *Annals of Statistics*, 12(4):1296–1309, 1984.
- [10] G. Raskutti, M. J. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over  $\ell_q$ -balls. Technical Report arXiv:0910.2042, UC Berkeley, Department of Statistics, 2009.
- [11] P. Ravikumar, M. J. Wainwright, and J. Lafferty. High-dimensional Ising model selection using  $\ell_1$ -regularized logistic regression. *Annals of Statistics*, 2008. To appear.
- [12] P. Ravikumar, M. J. Wainwright, G. Raskutti, and B. Yu. High-dimensional covariance estimation by minimizing  $\ell_1$ -penalized log-determinant divergence. Technical Report 767, Department of Statistics, UC Berkeley, September 2008.
- [13] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *Allerton Conference*, 2007.
- [14] A.J. Rothman, P.J. Bickel, E. Levina, and J. Zhu. Sparse permutation invariant covariance estimation. *Electron. J. Statist.*, 2:494–515, 2008.
- [15] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58(1):267–288, 1996.
- [16] J. Tropp. Just relax: Convex programming methods for identifying sparse signals in noise. *IEEE Trans. Info Theory*, 52(3):1030–1051, March 2006.
- [17] B. Turlach, W.N. Venables, and S.J. Wright. Simultaneous variable selection. *Technometrics*, 27:349–363, 2005.
- [18] S. Van de Geer. High-dimensional generalized linear models and the lasso. *Annals of Statistics*, 36(2):614–645, 2008.
- [19] M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (Lasso). *IEEE Trans. Information Theory*, 55:2183–2202, May 2009.
- [20] M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society B*, 1(68):49, 2006.
- [21] C. Zhang and J. Huang. Model selection consistency of the lasso selection in high-dimensional linear regression. *Annals of Statistics*, 36:1567–1594, 2008.
- [22] P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2567, 2006.
- [23] P. Zhao, G. Rocha, and B. Yu. Grouped and hierarchical model selection through composite absolute penalties. *Annals of Statistics*, 37(6A):3468–3497, 2009.