

High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity

Authored by:

Martin J. Wainwright
Po-ling Loh

Abstract

Although the standard formulations of prediction problems involve fully-observed and noiseless data drawn in an i.i.d. manner, many applications involve noisy and/or missing data, possibly involving dependencies. We study these issues in the context of high-dimensional sparse linear regression, and propose novel estimators for the cases of noisy, missing, and/or dependent data. Many standard approaches to noisy or missing data, such as those using the EM algorithm, lead to optimization problems that are inherently non-convex, and it is difficult to establish theoretical guarantees on practical algorithms. While our approach also involves optimizing non-convex programs, we are able to both analyze the statistical error associated with any global optimum, and prove that a simple projected gradient descent algorithm will converge in polynomial time to a small neighborhood of the set of global minimizers. On the statistical side, we provide non-asymptotic bounds that hold with high probability for the cases of noisy, missing, and/or dependent data. On the computational side, we prove that under the same types of conditions required for statistical consistency, the projected gradient descent algorithm will converge at geometric rates to a near-global minimizer. We illustrate these theoretical predictions with simulations, showing agreement with the predicted scalings.

1 Paper Body

In standard formulations of prediction problems, it is assumed that the covariates are fully-observed and sampled independently from some underlying distribution. However, these assumptions are not realistic for many applications, in which covariates may be observed only partially, observed with corruption, or exhibit dependencies. Consider the problem of modeling the voting behavior of politicians: in this setting, votes may be missing due to abstentions, and temporally dependent due to collusion or "tit-for-tat" behavior. Similarly, surveys

often suffer from the missing data problem, since users fail to respond to all questions. Sensor network data also tends to be both noisy due to measurement error, and partially missing due to failures or drop-outs of sensors. There are a variety of methods for dealing with noisy and/or missing data, including various heuristic methods, as well as likelihood-based methods involving the expectation-maximization (EM) algorithm (e.g., see the book [1] and references therein). A challenge in this context is the possible non-convexity of associated optimization problems. For instance, in applications of EM, problems in which the negative likelihood is a convex function often become non-convex with missing or noisy data. Consequently, although the EM algorithm will converge to a local minimum, it is difficult to guarantee that the local optimum is close to a global minimum. In this paper, we study these issues in the context of high-dimensional sparse linear regression, in the case when the predictors or covariates are noisy, missing, and/or dependent. Our main contribution is to develop and study some simple methods for handling these issues, and to prove theoretical results about both the associated statistical error and the optimization error. Like EM-based approaches, our estimators are based on solving optimization problems that may be non-convex; however, despite this non-convexity, we are still able to prove that a simple form of projected gradient descent will produce an output that is sufficiently close to a global minimum.

As a second result, we bound the size of this statistical error, showing that it has the same scaling as the minimax rates for the classical cases of perfectly observed and independently sampled covariates. In this way, we obtain estimators for noisy, missing, and/or dependent data with guarantees similar to the usual fully-observed and independent case. The resulting estimators allow us to solve the problem of high-dimensional Gaussian graphical model selection with missing data. There is a large body of work on the problem of corrupted covariates or errors-in-variables for regression problems (see the papers and books [2, 3, 4, 5] and references therein). Much of the earlier theoretical work is classical in nature, where the sample size n diverges with the dimension p held fixed. Most relevant to this paper is more recent work that has examined issues of corrupted and/or missing data in the context of high-dimensional sparse linear models, allowing for $n \geq p$. Stauder and Bühlmann [6] developed an EM-based method for sparse inverse covariance matrix estimation in the missing data regime, and used this result to derive an algorithm for sparse linear regression with missing data. As mentioned above, however, it is difficult to guarantee that EM will converge to a point close to a global optimum of the likelihood, in contrast to the methods studied here. Rosenbaum and Tsybakov [7] studied the sparse linear model when the covariates are corrupted by noise, and proposed a modified form of the Dantzig selector, involving a convex program. This convexity produces a computationally attractive method, but the statistical error bounds that they establish scale proportionally with the size of the additive perturbation, hence are often weaker than the bounds that can be proved using our methods. The remainder of this paper is organized as follows. We begin in Section 2 with background and a precise description of the problem. We then introduce the class

of estimators we will consider and the form of the projected gradient descent algorithm. Section 3 is devoted to a description of our main results, including a pair of general theorems on the statistical and optimization error, and then a series of corollaries applying our results to the cases of noisy, missing, and dependent data. In Section 4, we demonstrate simulations to confirm that our methods work in practice. For detailed proofs, we refer the reader to the technical report [8].

Notation. For a matrix M , we write $\|M\|_{\max} := \max_{i,j} |m_{ij}|$ to be the elementwise ‘ ∞ ’-norm of M . Furthermore, $\|M\|_1$ denotes the induced ‘1’-operator norm (maximum absolute column sum) of M , and $\|M\|_{\text{op}}$ is the (M) , the condition number of M . induced ‘2’-operator norm (spectral norm) of M . We write $\kappa(M) := \frac{\|M\|_{\text{op}}}{\lambda_{\min}(M)}$

2

Background and problem set-up

In this section, we provide a formal description of the problem and motivate the class of estimators studied in the paper. We then describe a class of projected gradient descent algorithms to be used in the sequel.

Observation model and high-dimensional framework

Suppose we observe a response variable $y_i \in \mathbb{R}$ that is linked to a covariate vector $x_i \in \mathbb{R}^p$ via the linear model $y_i = \beta^T x_i + \epsilon_i$, for $i = 1, 2, \dots, n$.

(1)

β

Here, the regression vector $\beta \in \mathbb{R}$ is unknown, and $\epsilon_i \in \mathbb{R}$ is observation noise, independent of x_i . Rather than directly observing each $x_i \in \mathbb{R}^p$, we observe a vector $z_i \in \mathbb{R}^p$ linked to x_i via some conditional distribution: $z_i \sim Q(\cdot | x_i)$, for $i = 1, 2, \dots, n$.

(2)

This setup allows us to model various types of disturbances to the covariates, including (a) Additive noise: We observe $z_i = x_i + w_i$, where $w_i \in \mathbb{R}^p$ is a random vector independent of x_i , say zero-mean with known covariance matrix Σ_w . (b) Missing data: For a fraction $\rho \in [0, 1)$, we observe a random vector $z_i \in \mathbb{R}^p$ such that independently for each component j , we observe $z_{ij} = x_{ij}$ with probability $1 - \rho$, and $z_{ij} = ?$ with probability ρ . This model can also be generalized to allow for different missing probabilities for each covariate. Our first set of results is deterministic, depending on specific instantiations of the observed variables $\{(y_i, z_i)\}_{i=1}^n$. However, we are also interested in proving results that hold with high probability when the x_i ’s and z_i ’s are drawn at random from some distribution. We develop results for both the i.i.d. setting and the case of dependent covariates, where the x_i ’s are generated according to a stationary vector autoregressive (VAR) process. Furthermore, we work within a high-dimensional framework where $n \gg p$, and assume β has at most k non-zero parameters, where the sparsity k is also allowed to increase to infinity with the sample size n . We assume the scaling $k/n^2 = O(1)$, which is reasonable in order to have a non-diverging signal-to-noise ratio.

2.2

M-estimators for noisy and missing covariates

We begin by examining a simple deterministic problem. Let $\text{Cov}(X) = \Sigma$, and consider the program

(3)
$$\hat{b} = \arg \min_{\|b\|_1 \leq R} \|y - Xb\|_2^2$$
 As long as the constraint radius R is at least $\sqrt{\lambda_{\min}(\Sigma)}$, the unique solution to this convex program is $\hat{b} = \Sigma^{-1}X^T y$. This idealization suggests various estimators based on the plug-in principle. We form unbiased estimates of Σ and b and $\hat{\Sigma}$, denoted by $\hat{\Sigma}$, \hat{b} , respectively, and consider the modified program and its regularized version:

(4)
$$\hat{b} = \arg \min_{\|b\|_1 \leq R} \|y - Xb\|_2^2 + \lambda \|b\|_1$$

(5)
$$\hat{b} = \arg \min_{\|b\|_1 \leq R} \|y - Xb\|_2^2 + \lambda \|b\|_1 + \frac{\lambda}{2} \|b\|_2^2$$
 where $\lambda \geq 0$ is the regularization parameter. The Lasso [9, 10] is a special case of these programs, where $\hat{b}_{\text{Lasso}} := (X^T X + \lambda I)^{-1} X^T y$, (6) $\hat{\Sigma} = X^T X$ and we have introduced the shorthand $y = (y_1, \dots, y_n)^T \in \mathbb{R}^n$, and $X \in \mathbb{R}^{n \times p}$, with x_i as its i th row. In this paper, we focus on more general instantiations of the programs (4) and (5), involving different choices of $\hat{\Sigma}$ and \hat{b} . \hat{b}_{Lasso} is positive the pair $(\hat{\Sigma}, \hat{b})$ that are adapted to the cases of noisy and/or missing data. Note that the matrix $\hat{\Sigma}$ is semidefinite, so the Lasso program is convex. In sharp contrast, for the cases of noisy or missing data, the most \hat{b} is not positive semidefinite, hence the loss functions appearing in the problems (4) natural choice of the matrix $\hat{\Sigma}$ and (5) are non-convex. It is generally impossible to provide a polynomial-time algorithm that converges to a (near) global optimum of a non-convex problem. Remarkably, we prove that a simple projected gradient descent algorithm still converges with high probability to a vector close to any global optimum in our setting. Let us illustrate these ideas with some examples: Example 1 (Additive noise). Suppose we observe the $n \times p$ matrix $Z = X + W$, where W is a random matrix independent of X , with rows w_i drawn i.i.d. from a zero-mean distribution with known covariance Σ_w . Consider the pair $\hat{b}_{\text{add}} := (Z^T Z)^{-1} Z^T y$ and $\hat{\Sigma}_{\text{add}} := Z^T Z$, (7) $\hat{\Sigma} = Z^T Z$ which correspond to unbiased estimators of Σ and b , respectively. Note that when $\Sigma_w = 0$ (corresponding \hat{b}_{add} is to the noiseless case), the estimators reduce to the standard Lasso. However, when $\Sigma_w \neq 0$, the matrix $\hat{\Sigma}$ is not positive semidefinite in the high-dimensional regime ($n > p$) of interest. Indeed, since the matrix $Z^T Z$ has rank at most n , the subtracted matrix Σ_w may cause $\hat{\Sigma}$ to have a large number of negative eigenvalues. Example 2 (Missing data). Suppose each entry of X is missing independently with probability $\gamma \in [0, 1)$, and we observe the matrix $Z \in \mathbb{R}^{n \times p}$ with entries

$Z_{ij} = X_{ij}$ with probability $1 - \gamma$, $Z_{ij} = 0$ otherwise. Given the observed matrix $Z \in \mathbb{R}^{n \times p}$, consider an estimator of the general form (4), based on the choices $\hat{\Sigma} = Z^T Z$ and $\hat{b}_{\text{mis}} := (Z^T Z)^{-1} Z^T y$, (8) $\hat{\Sigma} = Z^T Z$ and $\hat{b}_{\text{mis}} := (Z^T Z)^{-1} Z^T y$ where $Z_{\text{mis}} = Z$ reduces to the pair $(\hat{\Sigma}, \hat{b}_{\text{Lasso}})$ for the standard Lasso when $\gamma = 0$, corresponding to no missing data. In the more interesting case when $\gamma \in (0, 1)$, the matrix $\hat{\Sigma}$ in equation (8) has rank at most n , so the subtracted diagonal matrix may cause the matrix $\hat{\Sigma}_{\text{mis}}$ to have a large number of negative eigenvalues when $n > p$, and the associated quadratic function is not convex.

2.3 Restricted eigenvalue conditions

there are various ways to assess its closeness to β . We focus on the ℓ_2 -norm $\|\hat{\beta} - \beta\|_2$, as well as the closely related ℓ_1 -norm $\|\hat{\beta} - \beta\|_1$. When the covariate matrix X is fully observed (so that the Lasso $\hat{\beta}_{\text{Las}} = (X^T X)^{-1} X^T y$ can be applied), it is well understood that a sufficient condition for ℓ_2 -recovery is that the matrix Σ_n satisfy a restricted eigenvalue (RE) condition (e.g., [11, 12, 13]). In this paper, we use the following condition:

Σ_n satisfies a lower restricted eigenvalue condition with curvature 1 (Lower-RE condition). The matrix Σ_n is ℓ_2 -bounded and tolerance $\ell_2(n, p) \leq 0$ if $\|\beta\|_2 \leq k^{-1/2} \ell_2(n, p) k^{1/2}$ for all $\beta \in \mathbb{R}^p$.

(9)

$\hat{\beta}_{\text{Las}} = (X^T X)^{-1} X^T y$ satisfies this RE condition (9), the Lasso estimate $\hat{\beta}_{\text{Las}}$ can be shown that when the Lasso matrix Σ_n has low ℓ_2 -error for any vector β supported on any subset of size at most k . $\ell_2(n, p)$. Moreover, it is known $\hat{\beta}_{\text{Las}}$ will satisfy such an RE condition that for various random choices of the design matrix X , the Lasso matrix Σ_n with high probability (e.g., [14]). We also make use of the analogous upper restricted eigenvalue condition: Σ_n satisfies an upper restricted eigenvalue condition with Definition 2 (Upper-RE condition). The matrix Σ_n smoothness $\ell_2(n, p) \leq 0$ and tolerance $\ell_2(n, p) \leq 0$ if $\|\beta\|_2 \leq k^{-1/2} \ell_2(n, p) k^{1/2}$ for all $\beta \in \mathbb{R}^p$.

(10)

In recent work on high-dimensional projected gradient descent, Agarwal et al. [15] use a more general form of bounds (9) and (10), called the restricted strong convexity (RSC) and restricted smoothness (RSM) conditions.

2.4 Projected gradient descent

In addition to proving results about the global minima of programs (4) and (5), we are also interested in polynomial-time procedures for approximating such optima. We show that the simple projected gradient descent algorithm can be used to solve the program (4). The algorithm generates a sequence of iterates β_t according to

$\beta_{t+1} = \Pi(\beta_t - \eta \nabla f(\beta_t))$, (11) where $\eta \geq 0$ is a stepsize parameter, and Π denotes the ℓ_2 -projection onto the ℓ_1 -ball of radius R . This projection can be computed rapidly in $O(p)$ time, for instance using a procedure due to Duchi et al. [16]. Our analysis shows that under a reasonable set of conditions, the iterates for the family of programs (4) converges to a point extremely close to any global optimum in both ℓ_1 -norm and ℓ_2 -norm, even for the non-convex program.

3

Main results and consequences

We provide theoretical guarantees for both the constrained estimator (4) and the regularized variant (5)

$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \|\beta\|_1 + \lambda \sum_{i=1}^n \frac{1}{2} \|\beta_i\|_2^2$ (12)

for a constant $b_0 \leq k^2$, which is a hybrid between the constrained (4) and regularized (5) programs. 3.1

Statistical error

\mathbf{b} satisfies a lower-RE condition with curvature In controlling the statistical error, we assume that the matrix \mathbf{Q} and vector \mathbf{y} and tolerance δ (n, p), as previously defined (9). In addition, recall that the matrix \mathbf{Q} serve as surrogates to the deterministic quantities $\mathbf{x} \in \mathbb{R}^p$ and $\mathbf{x} \in \mathbb{R}^p$, respectively. We assume there is a function $\psi(Q, \delta)$, depending on the standard deviation σ of the observation noise vector from equation (1) and the conditional distribution Q from equation (2), such that the following deviation conditions are satisfied: $r \leq \log p \leq \psi(\mathbf{x}) \leq k \psi(Q, \delta) \leq \log p \leq kb \leq \psi(\mathbf{x}) \leq k \psi(Q, \delta)$ and $k \leq (13) \leq n \leq q$ The following result applies to any global optimum \mathbf{b} of the program (12) with $n \leq 4 \psi(Q, \delta) \leq \log p \leq b$ Theorem 1 (Statistical error). Suppose the surrogates (\mathbf{Q}, \mathbf{b}) satisfy the deviation bounds (13), and the matrix \mathbf{b} satisfies the lower-RE condition (9) with parameters (δ, δ) such that $r \leq \delta \psi(Q, \delta) \leq \log p \leq k \delta (n, p) \leq \min \{ \dots \}$ (14) $2 b_0 n \leq 128 k^4$

Then for any vector \mathbf{b} with sparsity at most k , there is a universal positive constant c_0 such that any global optimum \mathbf{b} satisfies the bounds $r \leq$

$$k \log p \leq c_0 \max \{ \psi(Q, \delta), n \}, \text{ and (15a) } k \mathbf{b} \leq k^2 \delta \leq n r$$

$\log p \leq c_0 k \leq b \max \{ \psi(Q, \delta), n \}$ (15b) $k \leq k_1 \leq \delta \leq n$ The same bounds (without n) also apply to the constrained program (4) with radius choice $R = k \leq k_1 \leq b$ Las, Remarks: Note that for the standard Lasso pair (\mathbf{Q}, \mathbf{b}) Las, bounds of the form (15) for sub-Gaussian noise are well-known from past work (e.g., [12, 17, 18, 19]). The novelty of Theorem 1 is in allowing for general pairs of such surrogates, which can lead to non-convexity in the underlying M -estimator. 3.2

Optimization error

Although Theorem 1 provides guarantees that hold uniformly for any choice of global minimizer, it does not provide any guidance on how to approximate such a global minimizer using a polynomial-time algorithm. Nonetheless, we are able to show that for the family of programs (4), under reasonable conditions on δ isified in various settings, a simple projected gradient method will converge geometrically fast to a very good approximation of any global optimum. Theorem 2 (Optimization error). Consider the program (4) with any choice of radius R for which the constraint \mathbf{b} satisfies the lower-RE (9) and upper-RE (10) conditions with δ is active. Suppose that the surrogate matrix \mathbf{Q} $\log p \leq u \leq n$, and that we apply projected gradient descent (11) with constant stepsize $\eta = 2u$. Then as long as $n \geq k \log p$, there is a contraction coefficient $\rho \in (0, 1)$ independent of (n, p, k) and universal positive b the gradient descent iterates $\{\mathbf{b}^t\}$ satisfy the bound constants (c_1, c_2) such that for any global optimum \mathbf{b}^* , $t \geq 0$ $\log p \leq b^2 \leq t k^2 \leq 0 \leq \mathbf{b}^2 + c_1 k^2 \leq t \leq k \mathbf{b}^2 \leq k^2 + c_2 k \mathbf{b}^2 \leq k^2$ for all $t = 0, 1, 2, \dots$ (16) $2 \leq n$ In addition, we have the ‘1-bound’ $\mathbf{b}^1 \leq 2 k k^2 \leq t \leq k \mathbf{b}^2 + 2 k k \mathbf{b}^2 \leq k^2 + 2 k \mathbf{b}^2 \leq k_1$ for all $t = 0, 1, 2, \dots$ $k^2 \leq t \leq k$ (17) Note that the bound (16) controls the ‘2-distance between the iterate \mathbf{b}^t at time t , which is easily computed in polynomial-time, and any global optimum \mathbf{b}^* of the program (4), which may be

difficult to compute. Since $\epsilon \in (0, 1)$, the first term in the bound vanishes as t increases. Together with Theorem 4.1, equations (16) and (17) imply that the ℓ_2 - and ℓ_1 -optimization error are bounded as $O(\frac{1}{n})$ and $O(\frac{1}{k \log p})$, respectively. 3.3

Some consequences

Both Theorems 1 and 2 are deterministic results; applying them to specific models requires additional work to establish the stated conditions. We turn to the statements of some consequences of these theorems for different cases of noisy, missing, and dependent data. A zero-mean random variable Z is sub-Gaussian with parameter σ if $E(e^{tZ}) \leq \exp(\frac{1}{2} \sigma^2 t^2)$ for all $t \in \mathbb{R}$. We say that a random matrix $X \in \mathbb{R}^{n \times p}$ is sub-Gaussian with parameters (σ, κ) if each row $x_i^T \in \mathbb{R}^p$ is sampled independently from a zero-mean distribution with covariance Σ , and for any unit vector $u \in \mathbb{R}^p$, the random variable $u^T x_i$ is sub-Gaussian with parameter at most σ . We begin with the case of i.i.d. samples with additive noise, as described in Example 1. Corollary 1. Suppose we observe $Z = X + W$, where the random matrices $X, W \in \mathbb{R}^{n \times p}$ are sub-Gaussian with parameters (σ_x, κ_x) and (σ_w, κ_w) , respectively, and the sample size is lower-bounded as $n \geq \frac{1}{\epsilon} \frac{1}{\max(\sigma_x, \sigma_w)} \log \frac{1}{\epsilon}$. Then for the M -estimator based on the surrogates (ϕ_{badd}) , the results of Theorems 1 and 2 hold with parameters

$\sigma = \max(\sigma_x, \sigma_w)$ and $\kappa = \max(\kappa_x, \kappa_w) = c_0 \sigma_x^2 + \sigma_w^2 + \epsilon \sigma_x^2 + \sigma_w^2$ with probability at least $1 - c_1 \exp(-c_2 \log p)$. 5

For i.i.d. samples with missing data, we have the following: Corollary 2. Suppose $X \in \mathbb{R}^{n \times p}$ is a sub-Gaussian matrix with parameters (σ_x, κ_x) , and Z is the missing

σ_x , κ_x , $1/k \log p$, then Theorems 1 and 2 hold with data matrix with parameter σ . If $n \geq \frac{1}{\epsilon} \frac{1}{\max(\sigma_x, \sigma_w)} \log \frac{1}{\epsilon}$

probability at least $1 - c_1 \exp(-c_2 \log p)$ for $\sigma = \max(\sigma_x, \sigma_w)$ and $\kappa = \max(\kappa_x, \kappa_w) = c_0$

$\sigma_x \sigma_w + \epsilon \sigma_x^2$

Now consider the case where the rows of X are drawn from a vector autoregressive (VAR) process according to $x_{i+1} = Ax_i + v_i$,

for $i = 1, 2, \dots, n-1$,

(18)

where $v_i \in \mathbb{R}^p$ is a zero-mean noise vector with covariance matrix Σ_v , and $A \in \mathbb{R}^{p \times p}$ is a driving matrix with spectral norm $\|A\|_2 \leq 1$. We assume the rows of X are drawn from a Gaussian distribution with covariance Σ_x , such that $\Sigma_x = A \Sigma_x A^T + \Sigma_v$, so the process is stationary. Corollary 3 corresponds to the case of additive noise for a Gaussian VAR process. A similar result can be derived in the missing data setting. Corollary 3. Suppose the rows of X are drawn according to a Gaussian VAR process with driving matrix A .

4 Suppose the additive noise matrix W is i.i.d. with Gaussian rows. If $n \geq \frac{1}{\epsilon} \frac{1}{\max(\sigma_x, \sigma_w)} \log \frac{1}{\epsilon}$, with \min

$\sigma^2 = \sigma_w^2 + \sigma_x^2$

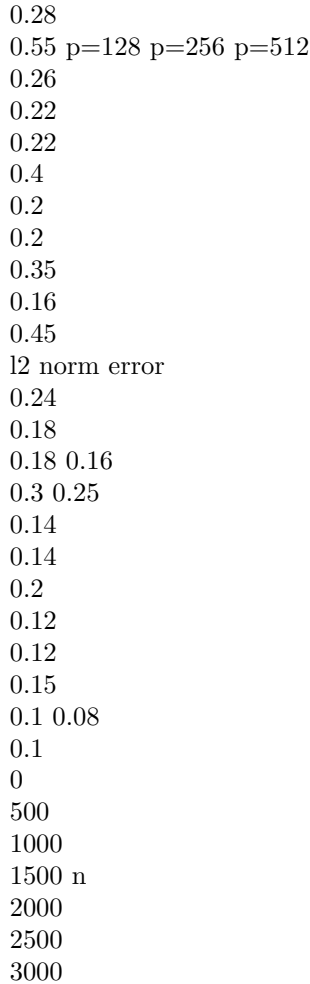
$\frac{1}{2} \sigma_x^2 + \frac{1}{2} \|A\|_2^2 \sigma_x^2$

there exists a universal constant c_0 such that $\| \hat{Q} - Q \|_F \leq c_0 \sqrt{\frac{p \log p}{n}}$

4 Simulations

In this section, we provide simulation results to confirm that the scalings predicted by our theory are sharp. In Figure 1, we plot the results of simulations under the additive noise model described in Example 1, using $\Sigma_x = I$ and $\Sigma_w = \Sigma_x$ with $\sigma_w = 0.2$. Panel (a) provides plots of ℓ_2 -error versus the sample size n , for $p \in \{128, 256, 512\}$. For all three choices of dimensions, the error decreases to zero as the sample size n increases, showing consistency of the method. If we plot the ℓ_2 -error versus the rescaled sample size $n/(k \log p)$, as depicted in panel (b), the curves roughly align for different values of p , agreeing with Theorem 1. Panel (c) shows analogous curves for VAR data with additive noise, using a driving matrix A with $\|A\|_F = 0.2$. Additive noise

Additive noise



0.08
 p=128 p=256 p=512
 0.5
 0.24
 l2 norm error
 l2 norm error
 Additive noise with autoregressive data
 0.28 p=128 p=256 p=512
 0.26
 0.1
 2
 4
 6
 8
 10 12 $n/(k \log p)$
 (a)
 (b)
 14
 16
 18
 20
 0.05
 0
 2
 4
 6
 8
 10 12 $n/(k \log p)$
 14
 16
 18
 20
 (c)

Figure 1. Plots of the error $\| \hat{\beta} - \beta^* \|_2^2$ after running projected gradient descent on the non-convex objective, with sparsity $k \leq p$. Plot (a) is an error plot for i.i.d. data with additive noise, and plot (b) shows ℓ_2 -error versus the rescaled sample size $n/(k \log p)$. Plot (c) depicts a similar (rescaled) plot for VAR data with additive noise. As predicted by Theorem 1, the curves align for different values of p in the rescaled plot.

We also verified the results of Theorem 2 empirically. Figure 2 shows the results of applying projected gradient descent to solve the optimization problem (4) in the cases of additive noise and missing data. We first applied b then reapplied projected gradient descent 10 times, tracking projected gradient to obtain an initial estimate $\hat{\beta}_0$, then the optimization error $\| \hat{\beta}_t - \beta^* \|_2^2$ (in blue) and statistical error $\| \beta^* - \beta^* \|_2^2$ (in red). As predicted by Theorem 2, the iterates exhibit geometric convergence to roughly the same fixed point, regardless of

starting point. Finally, we simulated the inverse covariance matrix estimation algorithm on three types of graphical models: (a) Chain-structured graphs. In this case, all nodes are arranged in a line. The diagonal entries of Σ are equal to 1, and entries corresponding to links in the chain are equal to 0.1. Then Σ is rescaled so $\text{tr}(\Sigma) = 1$. (b) Star-structured graphs. In this case, all nodes are connected to a central node, which has degree $k = 0.1p$. All other nodes have degree 1. The diagonal entries of Σ are set equal to 1, and all entries corresponding to edges in the graph are set equal to 0.1. Then Σ is rescaled so $\text{tr}(\Sigma) = 1$. (c) Erdős-Rényi graphs. As in Rothman et al. [22], we first generate a matrix B with diagonal entries 0, and all other entries independently equal to 0.5 with probability k/p , and 0 otherwise. Then Σ is chosen so $\Sigma = B + I$ has condition number p , and Σ is rescaled so $\text{tr}(\Sigma) = 1$. 7

Log error plot: additive noise case

Log error plot: missing data case

0.5

0.5 Stat error Opt error

0

0.5

1

$\log(\text{tr}(\Sigma_t) / \text{tr}(\Sigma))$

$\log(\text{tr}(\Sigma_t) / \text{tr}(\Sigma))$

0.5

1.5 2 2.5

1 1.5 2 2.5

3 3.5

Stat error Opt error

0

3

0

20

40 60 Iteration count

80

3.5

100

0

20

40 60 Iteration count

(a)

80

100

(b)

Figure 2. Plots of the optimization error $\log(\text{tr}(\Sigma_t) / \text{tr}(\Sigma))$ and statistical error $\log(k / \text{tr}(\Sigma))$ versus iteration number t , generated by running projected gradient descent on the non-convex objective. As predicted by Theorem 2, the optimization error decreases geometrically.

After generating the matrix X of n i.i.d. samples from the appropriate graphical model, with covariance matrix $\Sigma_X = I$, we generated the corrupted matrix $Z = X + W$ with $W = (0.2)I$. Figure 3 shows the rescaled ℓ_2 error plotted against the sample size n for a chain-structured graph, with panel (a) showing the original plot and panel (b) plotting against the rescaled sample size. We obtained qualitatively similar results for the star and Erdős-Rényi graphs, in the presence of missing and/or dependent data.

Chain graph

Chain graph

0.7

0.7 $p=64$ $p=128$ $p=256$

0.5

0.4

0.3

0.2

0.1

0

$p=64$ $p=128$ $p=256$

0.6 $1/\sqrt{k}$ * ℓ_2 operator norm error

$1/\sqrt{k}$ * ℓ_2 operator norm error

0.6

0.5

0.4

0.3

0.2

0.1

0

100

200

300

400

500

600

0

700

n

(a) ℓ_2 error plot for chain graph, additive noise

10

20

30 40 $n/(k \log p)$

50

60

(b) rescaled plot

Figure 3. (a) Plots of the rescaled error ℓ_2 error after running projected gradient descent for a chain-structured Gaussian graphical model with additive noise. As predicted by Theorems 1 and 2, all curves align

when the rescaled n , as shown in (b). Each point represents the average over 50 trials. error is plotted against the ratio $k \log p$

5

Discussion

In this paper, we formulated an ℓ_1 -constrained minimization problem for sparse linear regression on corrupted data. The source of corruption may be additive noise or missing data, and although the resulting objective is not generally convex, we showed that projected gradient descent is guaranteed to converge to a point within statistical precision of the optimum. In addition, we established ℓ_1 - and ℓ_2 -error bounds that hold with high probability when the data are drawn i.i.d. from a sub-Gaussian distribution, or drawn from a Gaussian VAR process. Finally, we used our results on linear regression to perform sparse inverse covariance estimation for a Gaussian graphical model, based on corrupted data. The bounds we obtain for the spectral norm of the error are of the same order as existing bounds for inverse covariance matrix estimation with uncorrupted, i.i.d. data. Acknowledgments PL acknowledges support from a Hertz Foundation Fellowship and an NDSEG Fellowship; MJW and PL were also partially supported by grants NSF-DMS-0907632 and AFOSR-09NL184. The authors thank Alekh Agarwal, Sahand Negahban, and John Duchi for discussions and guidance.

8

2 References

- [1] R. Little and D. B. Rubin. Statistical analysis with missing data. Wiley, New York, 1987.
- [2] J. T. Hwang. Multiplicative errors-in-variables models with applications to recent data released by the U.S. Department of Energy. *Journal of the American Statistical Association*, 81(395):pp. 680?688, 1986.
- [3] R. J. Carroll, D. Ruppert, and L. A. Stefanski. *Measurement Error in Nonlinear Models*. Chapman and Hall, 1995.
- [4] S. J. Iturria, R. J. Carroll, and D. Firth. Polynomial regression and estimating functions in the presence of multiplicative measurement error. *Journal of the Royal Statistical Society Series B - Statistical Methodology*, 61:547?561, 1999.
- [5] Q. Xu and J. You. Covariate selection for linear errors-in-variables regression models. *Communications in Statistics - Theory and Methods*, 36(2):375?386, 2007.
- [6] N. St?adler and P. B?uhlmann. Missing values: Sparse inverse covariance estimation and an extension to sparse regression. *Statistics and Computing*, pages 1?17, 2010.
- [7] M. Rosenbaum and A. B. Tsybakov. Sparse recovery under matrix uncertainty. *Annals of Statistics*, 38:2620?2651, 2010.
- [8] P. Loh and M.J. Wainwright. High-dimensional regression with noisy and missing data: Provable guarantees with non-convexity. Technical report, UC Berkeley, September 2011. Available at <http://arxiv.org/abs/1109.3714>.
- [9] R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58(1):267?288, 1996.
- [10] S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Journal on Sci-*

entific Computing, 20(1):33?61, 1998. [11] S. van de Geer. The deterministic Lasso. In Proceedings of Joint Statistical Meeting, 2007. [12] P. J. Bickel, Y. Ritov, and A. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics*, 37(4):1705?1732, 2009. [13] S. van de Geer and P. Bühlmann. On the conditions used to prove oracle results for the Lasso. *Electronic Journal of Statistics*, 3:1360?1392, 2009. [14] G. Raskutti, M. J. Wainwright, and B. Yu. Restricted eigenvalue properties for correlated Gaussian designs. *Journal of Machine Learning Research*, 11:2241?2259, 2010. [15] A. Agarwal, S. Negahban, and M.J. Wainwright. Fast global convergence of gradient methods for highdimensional statistical recovery. Technical report, UC Berkeley, April 2011. Available at <http://arxiv.org/abs/1104.4824>. [16] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra. Efficient projections onto the ℓ_1 -ball for learning in high dimensions. In International Conference on Machine Learning, pages 272?279, 2008. [17] C. H. Zhang and J. Huang. The sparsity and bias of the Lasso selection in high-dimensional linear regression. *Annals of Statistics*, 36(4):1567?1594, 2008. [18] N. Meinshausen and B. Yu. Lasso-type recovery of sparse representations for high-dimensional data. *Annals of Statistics*, 37(1):246?270, 2009. [19] S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for the analysis of regularized M-estimators. In Advances in Neural Information Processing Systems, 2009. [20] N. Meinshausen and P. Bühlmann. High-dimensional graphs and variable selection with the Lasso. *Annals of Statistics*, 34:1436?1462, 2006. [21] M. Yuan. High-dimensional inverse covariance matrix estimation via linear programming. *Journal of Machine Learning Research*, 99:2261?2286, August 2010. [22] A. J. Rothman, P. J. Bickel, E. Levina, and J. Zhu. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494?515, 2008.