

Confidence Intervals and Hypothesis Testing for High-Dimensional Statistical Models

Authored by:

Andrea Montanari
Adel Javanmard

Abstract

Fitting high-dimensional statistical models often requires the use of non-linear parameter estimation procedures. As a consequence, it is generally impossible to obtain an exact characterization of the probability distribution of the parameter estimates. This in turn implies that it is extremely challenging to quantify the ‘uncertainty’ associated with a certain parameter estimate. Concretely, no commonly accepted procedure exists for computing classical measures of uncertainty and statistical significance as confidence intervals or p-values. We consider here a broad class of regression problems, and propose an efficient algorithm for constructing confidence intervals and p-values. The resulting confidence intervals have nearly optimal size. When testing for the null hypothesis that a certain parameter is vanishing, our method has nearly optimal power. Our approach is based on constructing a ‘de-biased’ version of regularized M-estimators. The new construction improves over recent work in the field in that it does not assume a special structure on the design matrix. Furthermore, proofs are remarkably simple. We test our method on a diabetes prediction problem.

1 Paper Body

It is widely recognized that modern statistical problems are increasingly high-dimensional, i.e. require estimation of more parameters than the number of observations/examples. Examples abound from signal processing [16], to genomics [21], collaborative filtering [12] and so on. A number of successful estimation techniques have been developed over the last ten years to tackle these problems. A widely applicable approach consists in optimizing a suitably regularized likelihood function. Such estimators are, by necessity, non-linear and non-explicit (they are solution of certain optimization problems). The use of non-linear parameter estimators comes at a price. In general, it is impossible to characterize the distribution of the estimator. This situation is very different from the one of classical statistics in which either exact characterizations

are available, or asymptotically exact ones can be derived from large sample theory [26]. This has an important and very concrete consequence. In classical statistics, generic and well accepted procedures are available for characterizing the uncertainty associated to a certain parameter estimate in terms of confidence intervals or p-values [28, 14]. However, no analogous procedures exist in high-dimensional statistics. In this paper we develop a computationally efficient procedure for constructing confidence intervals and p-values for a broad class of high-dimensional regression problems. The salient features of our procedure are: (i) Our approach guarantees nearly optimal confidence interval sizes and testing power. (ii) It is the first one that achieves this goal under essentially no assumptions on the population covariance matrix of the parameters, beyond the standard conditions for high-dimensional consistency. (iii) It allows for a streamlined analysis with respect to earlier work in the same area. 1

Table 1: Unbiased estimator for θ_0 in high dimensional linear regression models Input: Measurement vector y , design matrix X , parameter θ . Output: Unbiased estimator $\hat{\theta}_{bu}$. 1: Set $\theta = \theta_0$, and let $\hat{\theta}_{bn}$ be the Lasso estimator as per Eq. (3). 2: Set $\theta = \hat{\theta}_{bn}$. 3: for $i = 1, 2, \dots, p$ do 4: Let m_i be a solution of the convex program: $\min_{\theta} \|y - X\theta\|_2^2$ subject to $\|\theta\|_1 \leq m_i$

minimize
(4)

5: Set $M = (m_1, \dots, m_p)^T$. If any of the above problems is not feasible, then set $M = Ip$. 6: Define the estimator $\hat{\theta}_{bu}$ as follows:

1 $\hat{\theta}_{bu} = \hat{\theta}_{bn} + M^T X^T (Y - X\hat{\theta}_{bn}) / n$
(5)

(iv) Our method has a natural generalization non-linear regression models (e.g. logistic regression, see Section 4). We provide heuristic and numerical evidence supporting this generalization, deferring a rigorous study to future work. For the sake of clarity, we will focus our presentation on the case of linear regression, deferring the generalization to Section 4. In the random design model, we are given n i.i.d. pairs $(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)$, with vectors $X_i \in \mathbb{R}^p$ and response variables Y_i given by $W_i \sim N(0, \Sigma)$.

$Y_i = h^T X_i + W_i$,
(1)

Here $h^T \cdot$ is the standard scalar product in \mathbb{R}^p . In matrix form, letting $Y = (Y_1, \dots, Y_n)^T$ and denoting by X the design matrix with rows X_1^T, \dots, X_n^T , we have $Y = X\theta_0 + W$,

$W \sim N(0, \Sigma \otimes I_n)$.
(2)

The goal is estimate the unknown (but fixed) vector of parameters $\theta_0 \in \mathbb{R}^p$. In the classic setting, $p \ll n$ and the estimation method of choice is ordinary least squares yielding $\hat{\theta}_{OLS} = (X^T X)^{-1} X^T Y$. In particular $\hat{\theta}$ is Gaussian with mean θ_0 and covariance $\Sigma (X^T X)^{-1}$. This directly allows to construct confidence intervals. In the high-dimensional setting where $p \gtrsim n$, the matrix $(X^T X)$ is rank deficient and one has to resort to biased estimators. A particularly successful approach is the Lasso [24, 7] which promotes sparse reconstructions through an ℓ_1 penalty. $\hat{\theta}_{bn}(Y, X; \lambda) = \arg \min_{\theta} \|y - X\theta\|_2^2 + \lambda \|\theta\|_1$. (3)

In case the right hand side has more than one minimizer, one of them can be selected arbitrarily for our purposes. We will often omit the arguments Y, X , as they are clear from the context. We denote by $S = \text{supp}(\beta_0)$ the support of β_0 , and let $s_0 = |S|$. A copious theoretical literature [6, 2, 4] shows that, under suitable assumptions on X , the Lasso is nearly as accurate as if the support S was known a priori. Namely, for $n = (s_0 \log p)$, we have $\|\beta_n - \beta_0\|_2^2 = O(s_0^2 (\log p)/n)$. These remarkable properties come at a price. Deriving an exact characterization for the distribution of β_n is not tractable in general, and hence there is no simple procedure to construct confidence intervals and p-values. In order to overcome this challenge, we construct a de-biased estimator from the Lasso solution. The de-biased estimator is given by the simple formula $\hat{\beta}_{\text{bu}} = \beta_n + (1/n) M^T X^T (Y - X\beta_n)$, as in Eq. (5). The basic intuition is that $X^T (Y - X\beta_n)/(n\sigma)$ is a subgradient of the ‘ ℓ_1 norm at the Lasso solution β_n . By adding a term proportional to this subgradient, our procedure compensates the bias introduced by the ‘ ℓ_1 penalty in the Lasso. For instance, letting $Q = (X^T X/n)^{-1}$, $\hat{\beta}_{\text{biOLS}} \pm 1.96 \sigma \sqrt{\text{diag}(Q)}$ is a 95% confidence interval [28].

$$\hat{\beta}_{\text{biOLS}} \pm 1.96 \sigma \sqrt{\text{diag}(Q)}$$

is a 95% confidence interval

We will prove in Section 2 that $\hat{\beta}_{\text{bu}}$ is approximately Gaussian, with mean β_0 and covariance $\sigma^2 Q$, where $Q = (X^T X/n)^{-1}$ is the empirical covariance of the feature vectors. This result allows to construct confidence intervals and p-values in complete analogy with classical statistics procedures. For instance, letting $Q = (X^T X/n)^{-1}$, $[\hat{\beta}_{\text{biOLS}} \pm 1.96 \sigma \sqrt{\text{diag}(Q)}]$ is a 95% confidence interval. The size of this interval is of order $1/\sqrt{n}$, which is the optimal (minimum) one, i.e. the same that would have been obtained by knowing a priori the support of β_0 . In practice the noise standard deviation is not known, but σ can be replaced by any consistent estimator $\hat{\sigma}$. A key role is played by the matrix $M = R_p^2$ whose function is to ‘decorrelate’ the columns of X . We propose here to construct M by solving a convex program that aims at optimizing two objectives. $\|M\|_1$ (here and below $\|\cdot\|_1$ denotes the entrywise ‘ ℓ_1 norm’) One one hand, we try to control $\|M\|_1$ which ‘as shown in Theorem 2.1’ controls the non-Gaussianity and bias of $\hat{\beta}_{\text{bu}}$. On the other, we bound $\|M\|_2$, for each $i \in [p]$, which controls the variance of $\hat{\beta}_{\text{bu}}$. minimize $\|M\|_1$ The idea of constructing a de-biased estimator of the form $\hat{\beta}_{\text{bu}} = \beta_n + (1/n) M^T X^T (Y - X\beta_n)$ was used by Javanmard and Montanari in [10], that suggested the choice $M = c^{-1} \Sigma^{-1}$, with $\Sigma = E\{X_1 X_1^T\}$ the population covariance matrix and c a positive constant. A simple estimator for Σ^{-1} was proposed for sparse covariances, but asymptotic validity and optimality were proven only for uncorrelated Gaussian designs (i.e. Gaussian X with $\Sigma = I$). Van de Geer, Bühlmann and Ritov [25] used the same construction with M an estimate of Σ^{-1} which is appropriate for sparse inverse covariances. These authors prove semi-parametric optimality in a non-asymptotic setting, provided the sample size is at least $n = (s_0^2 \log p)$. In this paper, we do not assume any sparsity constraint on Σ^{-1} , but still require the sample size scaling $n = (s_0^2 \log p)$.

p). We refer to a forthcoming publication wherein the condition on the sample size scaling is relaxed [11]. From a technical point of view, our proof starts from a simple decomposition of the de-biased estimator $\hat{\beta}_{\text{de-biased}}$ into a Gaussian part and an error term, already used in [25]. However departing radically from earlier work we realize that M need not be a good estimator of Σ in order for the de-biasing procedure to work. We instead set M as to minimize the error term and the variance of the Gaussian term. As a consequence of this choice, our approach applies to general covariance structures Σ . By contrast, earlier approaches applied only to sparse Σ , as in [10], or sparse Σ as in [25]. The only assumptions we make on Σ are the standard compatibility conditions required for high-dimensional consistency [4]. We refer the reader to the long version of the paper [9] for the proofs of our main results and the technical steps. 1.1

Further related work

The theoretical literature on high-dimensional statistical models is vast and rapidly growing. Restricting ourselves to linear regression, earlier work investigated prediction error [8], model selection properties [17, 31, 27, 5], ℓ_2 consistency [6, 2]. Of necessity, we do not provide a complete set of references, and instead refer the reader to [4] for an in-depth introduction to this area. The problem of quantifying statistical significance in high-dimensional parameter estimation is, by comparison, far less understood. Zhang and Zhang [30], and Bühlmann [3] proposed hypothesis testing procedures under restricted eigenvalue or compatibility conditions [4]. These methods are however effective only for detecting very large coefficients. Namely, they both require $\beta_{0,i} \geq c \max\{\sqrt{s_0 \log p/n}, \sqrt{p/n}\}$, which is s_0 larger than the ideal detection level [10]. In other words, in order for the coefficient $\beta_{0,i}$ to be detectable with appreciable probability, it needs to be larger than the overall ℓ_2 error, rather than the ℓ_2 error per coordinate. Lockart et al. [15] develop a test for the hypothesis that a newly added coefficient along the Lasso regularization path is irrelevant. This however does not allow to test arbitrary coefficients at a given value of λ , which is instead the problem addressed in this paper. It further assumes that the current Lasso support contains the actual support $\text{supp}(\beta_0)$ and that the latter has bounded size. Finally, resampling methods for hypothesis testing were studied in [29, 18, 19]. 1.2

Preliminaries and notations

$\Sigma = \frac{1}{n} X^T X$ be the sample covariance matrix. For $p \leq n$, Σ is always singular. However, We let Σ be nonsingular for a restricted set of directions. 3

Σ and a set S of size s_0 , the compatibility condition is met, if for some Definition 1.1. For a matrix $\Sigma \succeq 0$, and all S satisfying $k_S \leq k_1 \leq 3k_S$, it holds that $k_S \leq k_1 \leq 3k_S$.

$$s_0 \leq \frac{k_1}{k_S} \leq 3s_0$$

Definition 1.2. The sub-gaussian norm of a random variable X , denoted by $\|X\|_{\psi_2}$, is defined as $\|X\|_{\psi_2} = \sup_{p \geq 1} p^{1/2} (\mathbb{E} |X - \mathbb{E} X|^p)^{1/p}$.

The sub-gaussian norm of a random vector $X \in \mathbb{R}^n$ is defined as $\|X\|_{\psi_2} = \sup_{x \in S^{n-1}} \mathbb{E} |x^T X|^2$. Further, for a random variable X , its sub-exponential norm, denoted by $\|X\|_{\psi_1}$, is defined as $\|X\|_{\psi_1} = \sup_{p \geq 1} p^{1/2} (\mathbb{E} |X - \mathbb{E} X|^p)^{1/p}$.

For a matrix A and set of indices I, J , we let $A_{I,J}$ denote the submatrix formed by the rows in I and columns in J . Also, $A_{I,\cdot}$ (resp. $A_{\cdot,I}$) denotes the submatrix containing just the rows (reps. columns) in I . Likewise, for a vector v , v_I is the restriction of v to indices in I . We use the shorthand $\|v\|_1 = \|A_{I,\cdot}\|_1$. In particular, $\|A\|_1 = \sum_j \|A_{\cdot,j}\|_1$. The maximum and the minimum singular values $\sigma_1, \sigma_n = \sqrt{\lambda_1(A^T A)}, \sqrt{\lambda_n(A^T A)}$ of A are respectively denoted by $\sigma_{\max}(A)$ and $\sigma_{\min}(A)$. We write $\|v\|_p$ for the standard ' p ' norm of a vector v and $\|v\|_0$ for the number of nonzero entries of v . For a matrix A , $\|A\|_p$ is the ' p ' operator P norm, and $\|A\|_p$ is the elementwise ' p ' norm, i.e., $\|A\|_p = (\sum_{i,j} |A_{ij}|^p)^{1/p}$. For an integer $p \geq 1$, we let $[p] = \{1, \dots, p\}$. For a vector v , $\text{supp}(v)$ represents the positions of nonzero entries of v . Throughout, (w.h.p) means with probability converging to one as $n \rightarrow \infty$, and $\mathbb{P}(\cdot) \geq 1 - o(1)$ denotes the CDF of the standard normal distribution.

2

An de-biased estimator for β_0

Theorem 2.1. Consider the linear model (1) and let $\hat{\beta}_{\text{bu}}$ be defined as per Eq. (5). Then, $\sqrt{n}(\hat{\beta}_{\text{bu}} - \beta_0) \xrightarrow{d} N(0, \Sigma)$, $\Sigma = n^{-1} \text{Cov}(\sum_{i=1}^n Z_i \epsilon_i) = Z^T Z + \Sigma_\epsilon$, $Z^T Z \xrightarrow{p} N(0, \Sigma_\epsilon)$. Further, suppose that $\sigma_{\min}(\Sigma) = \sigma_{\min}(\Sigma_\epsilon) = O(1)$, and $\sigma_{\max}(\Sigma) = O(1)$. In addition assume the rows of the whitened matrix $X^T X^{-1/2}$ are sub-gaussian, i.e., $\|X^T X^{-1/2} x_k\|_2 = O(1)$. Let E be the event that the p b.i.i = $O(1)$. Then, using $\sigma = O((\log p)/n)$ b and $\max_i \|\epsilon_i\|_p \leq \sigma$ compatibility condition holds for β_0 , β_0 (see inputs in Table 1), the following holds true. On the event E , w.h.p, $\|\hat{\beta}_{\text{bu}} - \beta_0\|_2 = O(\sigma \log p / n)$. Note that compatibility condition (and hence the event E) holds w.h.p. for random design matrices of a general nature. In fact [22] shows that under some general assumptions, the compatibility b w.h.p., when n is sufficiently large. Bounds on condition on β_0 implies a similar condition on β_0 , T b the variances $[M^T M]_{ii}$ will be given in Section 3.2. Finally, the claim of Theorem 2.1 does not rely on the specific choice of the objective function in optimization problem (4) and only uses the optimization constraints. Remark 2.2. Theorem 2.1 does not make any assumption β_0 about the parameter vector β_0 . If we further assume that the support size s_0 satisfies $s_0 = o(n / \log p)$, then we have $\|\hat{\beta}_{\text{bu}} - \beta_0\|_2 = o(1)$, w.h.p. Hence, $\hat{\beta}_{\text{bu}}$ is an asymptotically unbiased estimator for β_0 .

3

Statistical inference

A direct application of Theorem 2.1 is to derive confidence intervals and statistical hypothesis tests β_0 for high dimensional models. Throughout, we make the sparsity assumption $s_0 = o(n / \log p)$. 3.1

Confidence intervals

We first show that the variances of variables $Z_j - X$ are $O(1)$. 4

Lemma 3.1. Let $M = (m_1, \dots, m_p)^T$ be the matrix with rows m_i obtained by solving convex T b b program (4). Then for all $i \in [p]$, $[M^T M]_{ii} = O(1)$. By Remark 2.2 and Lemma 3.1, we have $\|m_i\|_2 = O(1)$.

$0, i \in [p] \Rightarrow \sum_{j=1}^n x_j X_j = \beta(x) + o(1)$, $1/2 \leq T \leq [M^T M]_{ii}$

$\beta(x) \in \mathbb{R}$.

(6)

Since the limiting probability is independent of X , Eq. (6) also holds unconditionally for random design X . For constructing confidence intervals, a consistent estimate of σ is needed. To this end, we use the scaled Lasso [23] given by $\hat{\sigma} = \sqrt{\frac{1}{n} \min_{\beta} \{ \sum_{i=1}^n (y_i - \sum_{j=1}^p \beta_j X_{ij})^2 + \lambda \sum_{j=1}^p |\beta_j| \}}$. This is a joint convex minimization which provides an estimate of the noise level in addition to an ℓ_1 norm, under the assumptions estimate of σ . We use $\hat{\sigma} = c_1 (\log p)/n$ that yields a consistent estimate $\hat{\sigma}$ of Theorem 2.1 (cf. [23]). We hence obtain the following. Corollary 3.2. Let $q = 1 + 1/2 \log T$. $\hat{\sigma}(\beta, n) = \hat{\sigma} (1 + \sqrt{2}) \sqrt{n} [M + M(7)]$. Then $I_i = [\hat{\sigma}(\beta, n), \hat{\sigma}(\beta, n)]$ is an asymptotic two-sided confidence interval for $\beta_{0,i}$ with significance α . Notice that the same corollary applies to any other consistent estimator $\hat{\sigma}$ of the noise standard deviation. 3.2

Hypothesis testing

An important advantage of sparse linear regression models is that they provide parsimonious explanations of the data in terms of a small number of covariates. The easiest way to select the ‘active’ covariates is to choose the indexes i for which $\beta_{0,i} \neq 0$. This approach however does not provide a measure of statistical significance for the finding that the coefficient is non-zero. More precisely, we are interested in testing an individual null hypothesis $H_{0,i} : \beta_{0,i} = 0$ versus the alternative $H_{A,i} : \beta_{0,i} \neq 0$, and assigning p-values for these tests. We construct a p-value P_i for the test $H_{0,i}$ as follows:

$\hat{\sigma}(\beta, n) = \hat{\sigma}(\beta, n)$. (8) $P_i = 1 - \frac{1}{2} \log \frac{\hat{\sigma}(\beta, n)}{\hat{\sigma}(\beta, n)}$ The decision rule is then based on the p-value P_i :

$$\begin{aligned} & 1 \text{ if } P_i \leq \alpha, T_i(X)(y) = 0 \text{ otherwise} \\ & (\text{reject } H_{0,i}), (\text{accept } H_{0,i}). \end{aligned} \quad (9)$$

We measure the quality of the test $T_i(X)(y)$ in terms of its significance level α and statistical power $1 - \beta_i$. Here α is the probability of type I error (i.e. of a false positive at i) and β_i is the probability of type II error (i.e. of a false negative at i). Note that it is important to consider the tradeoff between statistical significance and power. Indeed any significance level α can be achieved by randomly rejecting $H_{0,i}$ with probability α . This test achieves power $1 - \beta_i = \alpha$. Further note that, without further assumption, no nontrivial power can be achieved. In fact, choosing $\beta_{0,i} \neq 0$ arbitrarily close to zero, $H_{0,i}$ becomes indistinguishable from its alternative. We will therefore assume that, whenever $\beta_{0,i} \neq 0$, we have $|\beta_{0,i}| \geq \delta$ as well. We take a minimax perspective and require the test to behave uniformly well over s_0 -sparse vectors. Formally, for $\delta > 0$ and $i \in [p]$, define $\alpha_i(n) = \sup_{\beta} P(T_i(X)(y) = 1) : \beta_{0,i} \neq 0, \|\beta\|_0 \leq s_0(n), |\beta_{0,i}| \geq \delta$. $\beta_i(n; \delta) = \sup_{\beta} P(T_i(X)(y) = 0) : \beta_{0,i} = 0, \|\beta\|_0 \leq s_0(n), |\beta_{0,i}| \geq \delta$.

5

Here, we made dependence on n explicit. Also, $P_i(\beta)$ is the induced probability for random design X and noise realization w , given the fixed parameter vector β . Our next theorem establishes bounds on $\alpha_i(n)$ and $\beta_i(n; \delta)$. Theorem 3.3. Consider a random design model that satisfies the conditions of Theorem 2.1. Under the sparsity assumption $s_0 = o(n/\log p)$, the following holds true

for any fixed sequence of integers $i = i(n)$: $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \frac{1}{i} = \frac{1}{2}$.

$$\begin{array}{l} (10) \\ n?? \\ \lim \\ n?? \\ 1\ ?\ ?i\ (?\ ;\ n)\ ?\ 1,\ 1\ ?\ ?i? \ (?\ ;\ n) \\ 1\ ?\ ?i? \ (?\ ;\ n)\ ?\ G\ ?, \\ ? \\ n? \ 1/2\ ?[??1\ i,i] \end{array}$$

(11)

where, for $\alpha \in [0, 1]$ and $u \in \mathbb{R}^+$, the function $G(\alpha, u)$ is defined as follows:

$$G(\alpha, u) = \frac{1}{2} \left(\frac{1}{1 + u} + \frac{1}{1 + u} \right) = \frac{1}{2} \left(\frac{1}{1 + u} + \frac{1}{1 + u} \right).$$

It is easy to see that, for any $\alpha \in [0, 1]$, $G(\alpha, u)$ is continuous and monotone increasing. Moreover, $G(\alpha, 0) = \alpha$ which is the trivial power obtained by randomly rejecting H_0 with probability α . As α deviates from zero, we obtain nontrivial tests to achieve a specific power β . Notice that in order to achieve a power β , our scheme requires $n = O(1/\alpha)$, since $\frac{1}{1 + u} = O(1)$. i.e. $\frac{1}{1 + u} = O(1)$.

Minimax optimality

The authors of [10] prove an upper bound for the minimax power of tests with a given significance level α , under the Gaussian random design models (see Theorem 2.6 therein). This bound is obtained by considering an oracle test that knows all the active parameters except i , i.e., $S \setminus \{i\}$. To state the bound formally, for a set $S \subseteq [p]$ and $i \in S^c$, define $\beta_{i|S} = \beta_{i, S \setminus \{i\}}$ ($\beta_{i,S} = \beta_{i, S \cup \{i\}}$), and let $\alpha_{i|S} = \min_{\beta \in S : \beta \in [p] \setminus \{i\}} \|\beta - \beta_{i|S}\|_2$.

? In asymptotic regime and under our sparsity assumption $s_0 = o(n/\log p)$, the bound of [10] simplifies to \lim

$$\begin{aligned} & n^{??} \\ & 1 \text{ ? } ?_{\text{iopt}}(?, ?) \text{ ? } 1, G(?, ?/\text{eff}) \\ & ?_{\text{eff}} = ? \\ & ? , n^{??}, s_0 \\ & (12) \end{aligned}$$

Using the bound of (12) and specializing the result of Theorem 3.3 to Gaussian design \mathbf{X} , we obtain that our scheme achieves a near optimal minimax power for a broad class of covariance matrices. We can compare our test to the optimal test by computing how much β must be increased in order to achieve the minimax optimal power. It follows from the above that β must be increased to β^* , with the two differing by a factor: $q/p = \lambda_{\max}(\Sigma)/\lambda_{\min}(\Sigma)$, $\lambda_{\max}(\Sigma)$ and $\lambda_{\min}(\Sigma)$ being the largest and the smallest eigenvalues of Σ , respectively.

4

General regularized maximum likelihood

In this section, we generalize our results beyond the linear regression model to general regularized maximum likelihood. Here, we only describe the debiasing method. Formal guarantees can be obtained under suitable restricted

strong convexity assumptions [20] and will be the object of a forthcoming publication. For univariate Y , and vector $X \in \mathbb{R}^p$, we let $\{f_\theta(Y - X)\}_{\theta \in \mathbb{R}^p}$ be a family of conditional probability densities parameterized by θ , that are absolutely continuous with respect to a common measure $\nu(dy)$, and suppose that the gradient $\nabla_\theta f_\theta(Y - X)$ exists and is square integrable. As in for linear regression, we assume that the data is given by n i.i.d. pairs $(X_1, Y_1), \dots, (X_n, Y_n)$, where conditional on X_i , the response variable Y_i is distributed as $Y_i \sim f_{\theta_0}(Y - X_i)$. 6

for some parameter vector $\theta_0 \in \mathbb{R}^p$. Let $L_i(\theta) = -\log f_\theta(Y_i - X_i)$ be the normalized negative P_n log-likelihood corresponding to the observed pair (Y_i, X_i) , and define $L(\theta) = \frac{1}{n} \sum_{i=1}^n L_i(\theta)$. We consider the following regularized estimator:

$$(13) \quad \hat{\theta}_\lambda = \arg \min_{\theta \in \mathbb{R}^p} L(\theta) + \lambda R(\theta), \quad \lambda \in \mathbb{R}_+$$

where λ is a regularization parameter and $R : \mathbb{R}^p \rightarrow \mathbb{R}_+$ is a norm. Let $I_i(\theta)$ be the Fisher information of $f_\theta(Y - X_i)$, defined as $I_i(\theta) = \mathbb{E}[\nabla_\theta \log f_\theta(Y - X_i) \nabla_\theta \log f_\theta(Y - X_i)^T]$. We next generalize the definition of I_i .

Let

$I(\theta) = \mathbb{E}[\nabla_\theta \log f_\theta(Y - X) \nabla_\theta \log f_\theta(Y - X)^T]$ where the second identity holds under suitable regularity conditions [13], and \mathbb{H}_θ denotes the Hessian of $L(\theta)$ as follows: $\mathbb{H}_\theta = -\mathbb{E}[\nabla_\theta^2 \log f_\theta(Y - X)]$. We assume \mathbb{H}_θ is positive definite for all θ .

$$(14) \quad \mathbb{H}_\theta = \mathbb{H}_\theta^T, \quad \mathbb{H}_\theta \succ 0 \quad \text{for all } \theta \in \mathbb{R}^p$$

Finally, the de-biased estimator $\hat{\theta}_{\text{bu}}$ is defined by $\hat{\theta}_{\text{bu}} = \hat{\theta}_\lambda + \lambda^{-1} \mathbb{H}_\lambda^{-1} \nabla L(\hat{\theta}_\lambda)$. Note that (in general) $\hat{\theta}_{\text{bu}} \neq \hat{\theta}_\lambda$, with M given again by the solution of the convex program (4), and the definition of b provided here. Notice that this construction is analogous to the one in [25] (although the present setting is somewhat more general) with the crucial difference of the construction of M . A simple heuristic derivation of this method is the following. By Taylor expansion of $L(\theta)$ around θ_0 we get $L(\theta) \approx L(\theta_0) + \nabla L(\theta_0)^T (\theta - \theta_0) + \frac{1}{2} (\theta - \theta_0)^T \mathbb{H}_{\theta_0} (\theta - \theta_0)$. Approximating \mathbb{H}_{θ_0} by $\mathbb{H}_{\hat{\theta}_\lambda}$ (which amounts to taking expectation with respect to the response variables y_i), we get $\mathbb{H}_{\hat{\theta}_\lambda} \approx \mathbb{H}_{\theta_0}$. Conditionally on $\{X_i\}_{i=1}^n$, the first term has zero expectation $\mathbb{E}[\nabla L(\theta_0)^T (\hat{\theta}_\lambda - \theta_0)] = 0$. Further, by central limit theorem, its low-dimensional marginals are asymptotically Gaussian. The bias term $\mathbb{E}[\nabla L(\theta_0)^T (\hat{\theta}_\lambda - \theta_0)]$ can be bounded as in the linear regression proximately Gaussian. The bias term $\mathbb{E}[\nabla L(\theta_0)^T (\hat{\theta}_\lambda - \theta_0)]$ is given by $I_i = [\text{var}(Y_i | X_i), \text{cov}(Y_i, X_i)]$, where $\text{var}(Y_i | X_i) = \mathbb{E}[Y_i^2 | X_i] - (\mathbb{E}[Y_i | X_i])^2$. Moreover, an asymptotically valid p -value P_i for testing null hypothesis $H_{0,i}$ is constructed as:

$P_i = 2 \min\{1, \sqrt{1 - \frac{1}{n} \nabla L(\hat{\theta}_\lambda)^T \mathbb{H}_{\hat{\theta}_\lambda}^{-1} \nabla L(\hat{\theta}_\lambda)}\}$. In the next section, we shall apply the general approach presented here to L_1 -regularized logistic regression. In this case, the binary response $Y_i \in \{0, 1\}$ is distributed as $Y_i \sim f_{\theta_0}(Y - X_i)$ where $f_{\theta_0}(1 - x) = (1 + e^{hx, \theta_0})^{-1}$ and $f_{\theta_0}(0 - x) = (1 + e^{-hx, \theta_0})^{-1}$. It is easy to see that in this case $b_i = \text{qbi}(1 - \text{qbi}) X_i^T$, with $\text{qbi} = (1$

+ $e^{h^*} X_i$ (?) , and thus in
 $X = \frac{1}{\sqrt{p}} \sum_{i=1}^n q_i X_i$
 5

Diabetes data example

We consider the problem of estimating relevant attributes in predicting type-2 diabetes. We evaluate the performance of our hypothesis testing procedure on the Practice Fusion Diabetes dataset [1]. This dataset contains de-identified medical records of 10000 patients, including information on diagnoses, medications, lab results, allergies, immunizations, and vital signs. From this dataset, we extract p numerical attributes resulting in a sparse design matrix X_{tot} of size $n \times p$, with $n = 10000$, $p = 7$

0.3
 ZSc
 ZS $N(0, 1)$
 0.1
 0.2
 Density
 2 0 -2
 0.0
 -4
 Sample Quantiles of Z
 4
 0.4
 Histograms of Z
 -3
 -2
 -1
 0
 1
 2
 3
 -10
 -5
 0
 5
 10
 Standard normal quantiles
 (a)
 (b)
 Q-Q plot of Z

Figure 1: Q-Q plot of Z and normalized histograms of Z (in red) and Z (in blue) for one realization. No fitting of the Gaussian mean and variance was done in panel (b).

and $p = 805$ (only 5.9% entries of X_{tot} are non-zero). Next, we standardize the columns of X to have mean 0 and variance 1. The attributes consist of:

(i) Transcript records: year of birth, gender and BMI; (ii) Diagnoses informations: 80 binary attributes corresponding to different ICD-9 codes. (iii) Medications: 80 binary attributes indicating the use of different medications. (iv) Lab results: For 70 lab test observations, we include attributes indicating patients tested, abnormality flags, and the observed values. We also bin the observed values into 10 quantiles and make 10 binary attributes indicating the bin of the corresponding observed value. We consider logistic model as described in the previous section with a binary response identifying the patients diagnosed with type-2 diabetes. For the sake of performance evaluation, we need to know the true significant attributes. Letting $L(\lambda)$ be the logistic loss corresponding to the design X_{tot} and response vector $Y \in \mathbb{R}^{n_{\text{tot}}}$, we take λ_0 as the minimizer of $L(\lambda)$. Notice that here, we are in the low dimensional regime ($n_{\text{tot}} \ll p$) and no regularization is needed. Next, we take random subsamples of size $n = 500$ from the patients, and examine the performance of our testing procedure. The experiment is done using glmnet-package in R that fits the entire path of the regularized logistic estimator. We then choose the value of λ that yields maximum AUC (area under ROC curve), approximated by a 5-fold cross validation. Results: Type I errors and powers of our decision rule (9) are computed by comparing to λ_0 . The average error and power (over 20 random subsamples) and significance level $\alpha = 0.05$ are respectively, 0.0319 and 0.818. Let $Z = (z_i)_{i=1}^n$ denote the vector with $z_i = n(\lambda_i - \lambda_0)/[M \sum_{i=1}^n \lambda_i]$. In Fig. 1(a), sample quantiles of Z are depicted versus the quantiles of a standard normal distribution. The plot clearly corroborates our theoretical result regarding the limiting distribution of Z . In order to build further intuition about the proposed p-values, let $Z^* = (z_i^*)_{i=1}^n$ be the vector with $z_i^* = b_i - \lambda_0$. In Fig. 1(b), we plot the normalized histograms of Z^* (in red) and Z (in blue). As the plot showcases, Z^* has roughly standard normal distribution, and the entries of Z appear as distinguishable spikes. The entries of Z with larger magnitudes are easier to be marked off from the normal distribution tail.

2 References

- [1] Practice Fusion Diabetes Classification. <http://www.kaggle.com/c/pf2012-diabetes>, 2012. Kaggle competition dataset.
- [2] P. J. Bickel, Y. Ritov, and A. B. Tsybakov. Simultaneous analysis of Lasso and Dantzig selector. *Amer. J. of Mathematics*, 37:1705–1732, 2009.
- [3] P. Bühlmann. Statistical significance in high-dimensional linear models. *arXiv:1202.1377*, 2012.
- [4] P. Bühlmann and S. van de Geer. *Statistics for high-dimensional data*. Springer-Verlag, 2011.
- [5] E. Candès and Y. Plan. Near-ideal model selection by ℓ_1 minimization. *The Annals of Statistics*, 37(5A):2145–2177, 2009.
- [6] E. J. Candès and T. Tao. Decoding by linear programming. *IEEE Trans. on Inform. Theory*, 51:4203–4215, 2005.
- [7] S. Chen and D. Donoho. Examples of basis pursuit. In *Proceedings of Wavelet Applications in Signal and Image Processing III*, San Diego, CA, 1995.
- [8] E. Greenshtein and Y.

Ritov. Persistence in high-dimensional predictor selection and the virtue of overparametrization. *Bernoulli*, 10:971?988, 2004. [9] A. Javanmard and A. Montanari. Confidence Intervals and Hypothesis Testing for High-Dimensional Regression. *arXiv:1306.3171*, 2013. [10] A. Javanmard and A. Montanari. Hypothesis testing in high-dimensional regression under the gaussian random design model: Asymptotic theory. *arXiv:1301.4240*, 2013. [11] A. Javanmard and A. Montanari. Nearly Optimal Sample Size in Hypothesis Testing for High-Dimensional Regression. *arXiv:1311.0274*, 2013. [12] Y. Koren, R. Bell, and C. Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30?37, August 2009. [13] E. Lehmann and G. Casella. Theory of point estimation. Springer, 2 edition, 1998. [14] E. Lehmann and J. Romano. Testing statistical hypotheses. Springer, 2005. [15] R. Lockhart, J. Taylor, R. Tibshirani, and R. Tibshirani. A significance test for the lasso. *arXiv preprint arXiv:1301.7161*, 2013. [16] M. Lustig, D. Donoho, J. Santos, and J. Pauly. Compressed sensing mri. *IEEE Signal Processing Magazine*, 25:72?82, 2008. [17] N. Meinshausen and P. B?uhlmann. High-dimensional graphs and variable selection with the lasso. *Ann. Statist.*, 34:1436?1462, 2006. [18] N. Meinshausen and P. B?uhlmann. Stability selection. *J. R. Statist. Soc. B*, 72:417?473, 2010. [19] J. Minnier, L. Tian, and T. Cai. A perturbation method for inference on regularized regression estimates. *Journal of the American Statistical Association*, 106(496), 2011. [20] S. N. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of m-estimators with decomposable regularizers. *Statistical Science*, 27(4):538?557, 2012. [21] J. Peng, J. Zhu, A. Bergamaschi, W. Han, D.-Y. Noh, J. R. Pollack, and P. Wang. Regularized multivariate regression for identifying master predictors with application to integrative genomics study of breast cancer. *The Annals of Applied Statistics*, 4(1):53?77, 2010. [22] M. Rudelson and S. Zhou. Reconstruction from anisotropic random measurements. *IEEE Transactions on Information Theory*, 59(6):3434?3447, 2013. [23] T. Sun and C.-H. Zhang. Scaled sparse linear regression. *Biometrika*, 99(4):879?898, 2012. [24] R. Tibshirani. Regression shrinkage and selection with the Lasso. *J. Royal. Statist. Soc B*, 58:267?288, 1996. [25] S. van de Geer, P. B?uhlmann, and Y. Ritov. On asymptotically optimal confidence regions and tests for high-dimensional models. *arXiv:1303.0518*, 2013. [26] A. W. Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000. [27] M. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ‘1 -constrained quadratic programming. *IEEE Trans. on Inform. Theory*, 55:2183?2202, 2009. [28] L. Wasserman. All of statistics: a concise course in statistical inference. Springer Verlag, 2004. [29] L. Wasserman and K. Roeder. High dimensional variable selection. *Annals of statistics*, 37(5A):2178, 2009. [30] C.-H. Zhang and S. Zhang. Confidence Intervals for Low-Dimensional Parameters in High-Dimensional Linear Models. *arXiv:1110.2563*, 2011. [31] P. Zhao and B. Yu. On model selection consistency of Lasso. *The Journal of Machine Learning Research*, 7:2541?2563, 2006.