

The LASSO risk: asymptotic results and real world examples

Authored by:

Andrea Montanari
Mohsen Bayati
Jos? Pereira

Abstract

We consider the problem of learning a coefficient vector x_0 from noisy linear observation $y = Ax_0 + w$. In many contexts (ranging from model selection to image processing) it is desirable to construct a sparse estimator. In this case, a popular approach consists in solving an l_1 -penalized least squares problem known as the LASSO or BPDN. For sequences of matrices A of increasing dimensions, with iid gaussian entries, we prove that the normalized risk of the LASSO converges to a limit, and we obtain an explicit expression for this limit. Our result is the first rigorous derivation of an explicit formula for the asymptotic risk of the LASSO for random instances. The proof technique is based on the analysis of AMP, a recently developed efficient algorithm, that is inspired from graphical models ideas. Through simulations on real data matrices (gene expression data and hospital medical records) we observe that these results can be relevant in a broad array of practical applications.

1 Paper Body

Let $x_0 \in \mathbb{R}^N$ be an unknown vector, and assume that a vector $y \in \mathbb{R}^n$ of noisy linear measurements of x_0 is available. The problem of reconstructing x_0 from such measurements arises in a number of disciplines, ranging from statistical learning to signal processing. In many contexts the measurements are modeled by $y = Ax_0 + w$,

(1.1)

where $A \in \mathbb{R}^{n \times N}$ is a known measurement matrix, and w is a noise vector. The LASSO or Basis Pursuit Denoising (BPDN) is a method for reconstructing the unknown vector x_0 given y , A , and is particularly useful when one seeks sparse solutions. For given A , y , one considers the cost functions $C_{A,y} : \mathbb{R}^N \rightarrow \mathbb{R}$ defined by (1.2) $C_{A,y}(x) = \|y - Ax\|_2^2 + \lambda \|x\|_1$, $\lambda \geq 0$. The original signal is estimated by $\hat{x}(\lambda; A, y) = \arg\min_x C_{A,y}(x)$.

(1.3)

In what follows we shall often omit the arguments A, y (and occasionally λ) from the above notations. We will also use $\|b\|_N$ to emphasize the N -dependence. Further $\|v\|_p = (\sum_{i=1}^N |v_i|^p)^{1/p}$ denotes the ℓ_p -norm of a vector $v \in \mathbb{R}^N$ (the subscript p will often be omitted if $p = 2$).

A large and rapidly growing literature is devoted to (i) Developing fast algorithms for solving the optimization problem (1.3); (ii) Characterizing the performances and optimality of the estimator \hat{b} . We refer to Section 1.3 for an unavoidably incomplete overview.

Despite such substantial effort, and many remarkable achievements, our understanding of (1.3) is not even comparable to the one we have of more classical topics in statistics and estimation theory. For instance, the best bound on the mean square error (MSE) of the estimator (1.3), i.e. on the quantity $N^{-1} \text{tr}(\mathbb{E}[\hat{b} \hat{b}^T - b b^T])$, was proved by Candes, Romberg and Tao [CRT06] (who in fact did not consider the LASSO but a related optimization problem). Their result estimates the mean square error only up to an unknown numerical multiplicative factor. Work by Candes and Tao [CT07] on the analogous Dantzig selector, upper bounds the mean square error up to a factor $C \log N$, under somewhat different assumptions. The objective of this paper is to complement this type of ‘rough but robust’ bounds by proving asymptotically exact expressions for the mean square error. Our asymptotic result holds almost surely for sequences of random matrices A with fixed aspect ratio and independent gaussian entries. While this setting is admittedly specific, the careful study of such matrix ensembles has a long tradition both in statistics and communications theory and has spurred many insights [Joh06, Tel99]. Further, our main result provides asymptotically exact expressions for other operating characteristics of the LASSO as well (e.g., False Positive Rate and True positive Rate). We carried out simulations on real data matrices with continuous entries (gene expression data) and binary feature matrices (hospital medical records). The results appear to be quite encouraging. Although our rigorous results are asymptotic in the problem dimensions, numerical simulations have shown that they are accurate already on problems with a few hundreds of variables. Further, they seem to enjoy a remarkable universality property and to hold for a fairly broad family of matrices [DMM10]. Both these phenomena are analogous to ones in random matrix theory, where delicate asymptotic properties of gaussian ensembles were subsequently proved to hold for much broader classes of random matrices. Also, asymptotic statements in random matrix theory have been replaced over time by concrete probability bounds in finite dimensions. Of course the optimization problem (1.2) is not immediately related to spectral properties of the random matrix A . As a consequence, universality and non-asymptotic results in random matrix theory cannot be directly exported to the present problem. Nevertheless, we expect such developments to be foreseeable. Our proof is based on the analysis of an efficient iterative algorithm first proposed by [DMM09], and called AMP, for approximate message passing. The algorithm is inspired by belief-propagation on graphical models, although the resulting iteration is significantly simpler (and scales linearly in the number of nodes). Extensive simulations [DMM10] showed that, in a number of settings, AMP performances

are statistically indistinguishable to the ones of LASSO, while its complexity is essentially as low as the one of the simplest greedy algorithms. The proof technique just described is new. Earlier literature analyzes the convex optimization problem (1.3) or similar problems by a clever construction of an approximate optimum, or of a dual witness. Such constructions are largely explicit. Here instead we prove an asymptotically exact characterization of a rather non-trivial iterative algorithm. The algorithm is then proved to converge to the exact optimum. Due to limited space in this paper we only state the main steps of the proof. More details are available in [BM10b].

1.1 Definitions In order to define the AMP algorithm, we denote by $\eta : \mathbb{R} \rightarrow \mathbb{R}$ the soft thresholding function $\eta(x) = \begin{cases} x - \tau & \text{if } x \geq \tau \\ x + \tau & \text{if } x \leq -\tau \\ 0 & \text{otherwise} \end{cases}$ (1.4) $\eta(x) = \begin{cases} x - \tau & \text{if } x \geq \tau \\ x + \tau & \text{if } x \leq -\tau \\ 0 & \text{otherwise} \end{cases}$. The algorithm constructs a sequence of estimates $x_t \in \mathbb{R}^n$, and residuals $z_t \in \mathbb{R}^n$, according to the iteration $x_{t+1} = \eta(A^T z_t + x_t; \tau_t)$, (1.5) $x_{t+1} = \eta(A^T z_t + x_t; \tau_t)$, initialized with $x_0 = 0$. Here A^T denotes the transpose of matrix A , and k_t is number of nonzero entries of x_t . Given a scalar function f and a vector $u \in \mathbb{R}^m$, we let $f(u)$ denote the vector $(f(u_1), \dots, f(u_m)) \in \mathbb{R}^m$ obtained by applying f componentwise. Finally $\bar{u} = \frac{1}{m} \sum_{i=1}^m u_i$ is the average of the vector $u \in \mathbb{R}$. $z_t = y - Ax_t$

2

As already mentioned, we will consider sequences of instances of increasing sizes, along which the LASSO behavior has a non-trivial limit. Definition 1. The sequence of instances $\{x_0(N), w(N), A(N)\}_{N \in \mathbb{N}}$ indexed by N is said to be a converging sequence if $x_0(N) \in \mathbb{R}^n$, $w(N) \in \mathbb{R}^n$, $A(N) \in \mathbb{R}^{n \times n}$ with $n = n(N)$ such that $n(N) \rightarrow \infty$ (0, ∞), and in addition the following conditions hold: (a) The empirical distribution of the entries of $x_0(N)$ converges weakly to a probability measure P on \mathbb{R} with bounded second moment. Further $\frac{1}{n} \sum_{i=1}^n x_{0,i}(N) \rightarrow \int \mathbb{R} P$. (b) The empirical distribution of the entries of $w(N)$ converges weakly to a probability measure P_W on \mathbb{R} with bounded second moment. Further $\frac{1}{n} \sum_{i=1}^n w_i(N) \rightarrow \int \mathbb{R} P_W$.

(c) If $\{e_i\}_{i=1}^n$, $e_i \in \mathbb{R}^n$ denotes the standard basis, then $\max_{i \in [N]} |A(N) e_i|_2 \rightarrow 0$, $\min_{i \in [N]} |A(N) e_i|_2 \rightarrow 1$, as $N \rightarrow \infty$ where $[N] = \{1, 2, \dots, N\}$.

For a converging sequence of instances, and an arbitrary sequence of thresholds $\{\tau_t\}_{t=0}^\infty$ (independent of N), the asymptotic behavior of the recursion (1.5) can be characterized as follows. Define the sequence $\{\tau_t\}_{t=0}^\infty$ by setting $\tau_0 = \tau + E\{X_0^2\}^{1/2}$ (for $X_0 \sim P_{X_0}$ and $\tau \geq E\{W^2\}^{1/2}$, $W \sim P_W$) and letting, for all $t \geq 0$: $\tau_{t+1} = F(\tau_t, \tau_t)$ with

$F(\tau, \tau) = E\{[\eta(X_0 + \tau Z)]^2\}$, where $Z \sim N(0, 1)$ is independent of X_0 . Notice that the function F depends on the law P_{X_0} . $F(\tau, \tau) \leq \tau +$

We say a function $\eta : \mathbb{R}^2 \rightarrow \mathbb{R}$ is pseudo-Lipschitz if there exist a constant $L \geq 0$ such that for all $x, y \in \mathbb{R}^2$: $|\eta(x) - \eta(y)| \leq L(1 + |x|^2 + |y|^2)|x - y|$. (This is a special case of the definition used in [BM10a] where such a function is called pseudo-Lipschitz of order 2.) Our next proposition that was conjectured in [DMM09] and proved in [BM10a]. It shows that the behavior of AMP can be tracked by the above one dimensional recursion. We often refer to this prediction by state evolution. Theorem 1 ([BM10a]). Let $\{x_0(N), w(N), A(N)\}_{N \in \mathbb{N}}$ be a converging sequence of instances with the entries of $A(N)$ iid

normal with mean 0 and variance $1/n$ and let $\varphi : \mathbb{R} \rightarrow \mathbb{R} \rightarrow \mathbb{R}$ be a pseudoLipschitz function. Then, almost surely $N \rightarrow \infty$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \varphi(x_{t+1}^i, x_t^i) = E[\varphi(X + \varphi(Z; \varphi), X, 0)] \quad (1.6)$$

where $Z \sim N(0, 1)$ is independent of $X_0 \sim p_{X_0}$.

In order to establish the connection with the LASSO, a specific policy has to be chosen for the thresholds $\{\varphi_t\}_{t \geq 0}$. Throughout this paper we will take $\varphi_t = \varphi_t^*$ with φ is fixed. In other words, the sequence $\{\varphi_t\}_{t \geq 0}$ is given by the recursion $\varphi_{t+1} = F(\varphi_t^2, \varphi_t)$. This choice enjoys several convenient properties [DMM09].

1.2 Main result Before stating our results, we have to describe a calibration mapping between φ and φ^* that was introduced in [DMM10] (Propositions 2, 3 and Corollary 4). Their proofs are presented in [BM10b]. Let us start by stating some convenient properties of the state evolution recursion. Proposition 2 ([DMM09]). Let $\varphi_{\min} = \varphi_{\min}(\varphi)$ be the unique non-negative solution of the equation $\varphi^2 (1 + \varphi^2 R)(\varphi) = \varphi^2(\varphi)$, with $\varphi(z) = e^{z^2/2} / \int_{\mathbb{R}} e^{z^2/2} \varphi(x) dx$.

For any $\varphi \geq 0$, $\varphi \geq \varphi_{\min}(\varphi)$, the fixed point equation $\varphi^2 = F(\varphi^2, \varphi)$ admits a unique solution. Denoting by $\varphi^* = \varphi^*(\varphi)$ this solution, we have $\lim_{t \rightarrow \infty} \varphi_t = \varphi^*(\varphi)$. Further the convergence takes

$$d_F^2$$

place for any initial condition and is monotone. Finally $d_F^2(\varphi, \varphi^*) \leq 1$ at $\varphi = \varphi^* \leq 1$.

The probability distribution that puts a point mass $1/N$ at each of the N entries of the vector.

3

We then define the function $\varphi^* \varphi^*(\varphi)$ on $(\varphi_{\min}(\varphi), \varphi)$, by $\varphi^*(\varphi) = \varphi^* \varphi^* [1 - P\{-X_0 + \varphi^* Z \leq \varphi^* \varphi^*\}]$.

This function defines a correspondence (calibration) between the sequence of thresholds $\{\varphi_t\}_{t \geq 0}$ and the regularization parameter φ . It should be intuitively clear that larger φ corresponds to larger thresholds and hence larger φ^* since both cases yield smaller estimates of x_0 . In the following we will need to invert this function. We thus define $\varphi : (0, \varphi) \rightarrow (\varphi_{\min}, \varphi)$ in such a way that $\varphi^*(\varphi) = \varphi$ $\varphi \in (\varphi_{\min}, \varphi) : \varphi^*(\varphi) = \varphi$.

The next result implies that the set on the right-hand side is non-empty and therefore the function $\varphi^* \varphi^*(\varphi)$ is well defined. Proposition 3 ([DMM10]). The function $\varphi^* \varphi^*(\varphi)$ is continuous on the interval $(\varphi_{\min}, \varphi)$ with $\varphi^*(\varphi_{\min} +) = \varphi_{\min}$ and $\lim_{\varphi \rightarrow \varphi} \varphi^*(\varphi) = \varphi$.

Therefore the function $\varphi^* \varphi^*(\varphi)$ satisfying $\varphi^*(\varphi) = \varphi$ $\varphi \in (\varphi_{\min}, \varphi) : \varphi^*(\varphi) = \varphi$ exists. We will denote by $A = \varphi^*((0, \varphi))$ the image of the function φ . Notice that the definition of φ is a priori not unique. We will see that uniqueness follows from our main theorem. Examples of the mappings $\varphi \mapsto \varphi^* F(\varphi^2, \varphi)$, $\varphi \mapsto \varphi^* \varphi^*(\varphi)$ and $\varphi \mapsto \varphi^* \varphi^*(\varphi)$ are presented in [BM10b]. We can now state our main result. Theorem 2. Let $\{x_0(N), w(N), A(N)\}_{N \in \mathbb{N}}$ be a converging sequence of instances with the entries of $A(N)$ iid normal with mean 0 and variance $1/n$. Denote by $x_b(\varphi; N)$ the LASSO estimator for instance $(x_0(N), w(N), A(N))$.

)), with $\lambda \geq 0$, $P\{X_0 \neq 0\}$ and let $\gamma : \mathbb{R} \rightarrow \mathbb{R}$ be a pseudo-Lipschitz function. Then, almost surely $N \rightarrow \infty$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \gamma(x_i) = E[\gamma(X_0 + \sqrt{\lambda} Z)], \quad (1.7)$$

where $Z \sim N(0, 1)$ is independent of $X_0 \sim p_{X_0}$, $\gamma = \gamma(\lambda)$ and $\gamma = \gamma(\lambda)$. As a corollary, the function $\gamma(\lambda)$ is indeed uniquely defined. Corollary 4. For any $\lambda \geq 0$ there exists a unique λ_{\min} such that $\gamma(\lambda) = \gamma$ (with the function $\gamma(\lambda)$ defined by $\gamma(\lambda) = E[\gamma(X_0 + \sqrt{\lambda} Z)]$). Hence the function $\gamma(\lambda)$ is continuous non-decreasing with $\gamma(0) = \gamma$. The assumption of a converging problem-sequence is important for the result to hold, while the hypothesis of gaussian measurement matrices $A(N)$ is necessary for the proof technique to be correct. On the other hand, the restrictions $\lambda \geq 0$, and $P\{X_0 \neq 0\} > 0$ (whence $\gamma \neq 0$ using $\gamma(\lambda) = E[\gamma(X_0 + \sqrt{\lambda} Z)]$) are made in order to avoid technical complications due to degenerate cases. Such cases can be resolved by continuity arguments. 1.3 Related work The LASSO was introduced in [Tib96, CD95]. Several papers provide performance guarantees for the LASSO or similar convex optimization methods [CRT06, CT07], by proving upper bounds on the resulting mean square error. These works assume an appropriate ‘isometry’ condition to hold for A . While such condition hold with high probability for some random matrices, it is often difficult to verify them explicitly. Further, it is only applicable to very sparse vectors x_0 . These restrictions are intrinsic to the worst-case point of view developed in [CRT06, CT07]. Guarantees have been proved for correct support recovery in [ZY06], under an appropriate ‘irrepresentability’ assumption on A . While support recovery is an interesting conceptualization for some applications (e.g. model selection), the metric considered in the present paper (mean square error) provides complementary information and is quite standard in many different fields. Closer to the spirit of this paper [RFG09] derived expressions for the mean square error under the same model considered here. Similar results were presented recently in [KWT09, GBS09]. These papers argue that a sharp asymptotic characterization of the LASSO risk can provide valuable

guidance in practical applications. For instance, it can be used to evaluate competing optimization methods on large scale applications, or to tune the regularization parameter λ . Unfortunately, these results were non-rigorous and were obtained through the famously powerful ‘replica method’ from statistical physics [MM09]. Let us emphasize that the present paper offers two advantages over these recent developments: (i) It is completely rigorous, thus putting on a firmer basis this line of research; (ii) It is algorithmic in that the LASSO mean square error is shown to be equivalent to the one achieved by a low-complexity message passing algorithm.

2 Numerical illustrations Theorem 2 assumes that the entries of matrix A are iid gaussians. We expect however that our predictions to be robust and hold for much larger family of matrices. Rigorous evidence in this direction is presented in [KM10] where the normalized cost $C(b, x)/N$ is shown to have a limit as $N \rightarrow \infty$

λ which is universal with respect to random matrices A with iid entries. (More precisely, it is universal if $E\{A_{ij}\} = 0$, $E\{A_{2ij}\} = 1/n$ and $E\{A_{6ij}\} \leq C/n^3$ for a uniform constant C .) Further, our result is asymptotic, while one might wonder how accurate it is for instances of moderate dimensions. Numerical simulations were carried out in [DMM10] and suggest that the result is robust and relevant already for N of the order of a few hundreds. As an illustration, we present in Figures 1-3 the outcome of such simulations for four types of real data and random matrices. We generated the signal vector randomly with entries in $\{+1, 0, -1\}$ and $P(x_{0,i} = +1) = P(x_{0,i} = -1) = 0.05$. The noise vector w was generated by using i.i.d. $N(0, 0.2)$ entries. We obtained the optimum estimator x by using OWLQN and l1 ls, packages for solving large-scale l1-regularized regressions [KKL+ 07], [AJ07]. We used 40 values of λ between .05 and 2 and N equal to 500, 1000, and 2000. For each case, the point (λ, MSE) was plotted and the results are shown in the figures. Continuous lines corresponds to the asymptotic prediction by Theorem 2 for

$\lambda^2(a, b) = (a - b)^2$, namely $\text{MSE} = \lim_{N \rightarrow \infty} \frac{1}{N} \|x - x_0\|_2^2 = E \|X_0 + \lambda Z; \lambda\|^2$ where $X_0 = (x_{0,1}, \dots, x_{0,N})^T$.

The agreement is remarkably good already for N, n of the order of a few hundreds, and deviations are consistent with statistical fluctuations. The four figures correspond to measurement matrices A :

Figure 1(a): Data consist of 2253 measurements of expression level of 7077 genes (this data is provided to us by Broad Institute). From this matrix we took sub-matrices A of aspect ratio λ for each N . The entries were continuous variables. We standardized all columns of A to have mean 0 and variance 1. Figure 1(b): From a data set of 1932 patient records we extracted 4833 binary features describing demographic information, medical history, lab results, medications etc. The 0-1 matrix was sparse (with only 3.1% non-zero entries). Similar to genes data, for each N , the sub-matrices A with aspect ratio λ were selected and standardized. Figure 2(a): Random ± 1 matrices with aspect ratio λ . Each entry is independently equal to $+1/\sqrt{n}$ or $-1/\sqrt{n}$ with equal probability.

Figure 2(b): Random gaussian matrices with aspect ratio λ and iid $N(0, 1/n)$ entries (as in Theorem 2). Notice the behavior appears to be essentially indistinguishable. Also the asymptotic prediction has a minimum as a function of λ . The location of this minimum can be used to select the regularization parameter. Further empirical analysis is presented in [BBM10]. For the second data set (patient records) we repeated the simulation 20 times (each time with fresh x_0 and w) and obtained the average and standard error for MSE, False Positive Rate (FPR) and True Positive Rate (TPR). The results with error bars are shown in Figure 3. The length of each error bar is 5

(a) Gene expression data
(b) Hospital records
0.4
0.4 N=500 N=1000 N=2000 Prediction
0.3
0.3

0.25
 0.25
 0.2
 0.2
 0.15
 0.15
 0.1
 0.1
 0.05 0
 0.5
 1
 ?
 1.5
 2
 N=500 N=1000 N=2000 Prediction
 0.35
 MSE
 MSE
 0.35
 0.05 0
 2.5
 0.5
 1
 ?
 1.5
 2
 2.5

Figure 1: Mean square error (MSE) as a function of the regularization parameter λ compared to the asymptotic prediction for $\lambda = .5$ and $\lambda^2 = .2$. In plot (a) the measurement matrix A is a real valued (standardized) matrix of gene expression data and in plot (b) A is a (standardized) 0-1 feature matrix of hospital records. Each point in these plots is generated by finding the LASSO predictor \hat{x} using a measurement vector $y = Ax_0 + w$ for an independent signal vector x_0 and an independent noise vector w . (a) λ_1 matrices

(b) Gaussian matrices
 0.4
 0.4 N=500 N=1000 N=2000 Prediction
 0.3
 0.3
 0.25
 0.25
 0.2
 0.2
 0.15
 0.15
 0.1

0.1
0.05 0
0.5
1
?
1.5
2
N=500 N=1000 N=2000 Prediction
0.35
MSE
MSE
0.35
0.05 0
2.5
0.5
1
?
1.5
2
2.5

? Figure 2: As in Figure 1, but the measurement matrix A has iid entries that are equal to $1/n$ with equal probabilities in plot (a), and has iid $N(0, 1/n)$ entries in plot (b). Additionally, each point in these plots uses an independent matrix A . is equal to twice the standard error (in each direction). FPR and TPR are calculated using $\frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\hat{x}_i \neq 0\}$ $\frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\hat{x}_{i,0} \neq 0\}$ $\frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\hat{x}_i = 0\}$ $\frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\hat{x}_{i,0} = 0\}$

(2.1)
where $\mathbb{I}\{S\} = 1$ if statement S holds and $\mathbb{I}\{S\} = 0$ otherwise. The predictions for FPR and TPR are obtained by applying Theorem 2 to $\text{fpr}(a, b) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{a_i \neq 0\} \mathbb{I}\{b_i = 0\}$ and $\text{tpr}(a, b) = \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{a_i = 0\} \mathbb{I}\{b_i \neq 0\}$, which yields $\lim_{N \rightarrow \infty} \text{FPR} = 2\alpha(1-\alpha)$, $\lim_{N \rightarrow \infty} \text{TPR} = \alpha(1-\alpha) + \alpha(1-\alpha)$ (2.2) $\lim_{N \rightarrow \infty} \text{FPR} = 2\alpha(1-\alpha)$, $\lim_{N \rightarrow \infty} \text{TPR} = \alpha(1-\alpha) + \alpha(1-\alpha)$ where α is defined in Proposition 2. Note that functions $\text{fpr}(a, b)$ and $\text{tpr}(a, b)$ are not pseudoLipschitz but the limits (2.2) follow from Theorem 2 via standard weak-convergence arguments.

3 A structural property and proof of the main theorem We will prove the following theorem which implies our main result, Theorem 2. Theorem 3. Assume the hypotheses of Theorem 2. Denote by $\{\hat{x}_t(N)\}_{t \geq 0}$ the sequence of estimates produced by AMP. Then $\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \mathbb{I}\{\hat{x}_t(N) \neq x_i\} = 0$, almost surely. 6

0.4
0.5 N500 N1000 N2000 Prediction
0.35
0.4 False Positive Rate
0.3
MSE
N=500 N=1000 N=2000 Prediction

0.45
 0.25 0.2
 0.35 0.3 0.25 0.2 0.15
 0.15
 0.1 0.1 0.05 0.05 0
 0.5
 1
 ?
 1.5
 2
 0 0
 2.5
 0.5
 1
 ?
 1.5
 2
 2.5
 0.8 N=500 N=1000 N=2000 Prediction
 0.7
 True Positive Rate
 0.6 0.5 0.4 0.3 0.2 0.1 0 0
 0.5
 1
 ?
 1.5
 2
 2.5

Figure 3: Average of MSE, FPR and TPR versus γ for medical data, using 20 samples per γ and N . All parameters are similar to Figure 1(b). Error bars are twice the standard errors (in each direction). The rest of the paper is devoted to the proof of this theorem. Section 3.1 proves a structural property that is the key tool in this proof. Section 3.2 uses this property together with a few lemmas to prove Theorem 3. Proofs of lemmas and more details can be found in [BM10b]. The proof of Theorem 2 follows immediately. Since when γ is Lipschitz there is a constant B where $\| \sum_{i=1}^N (\mathbf{x}_{t+1} - \mathbf{x}_{0,i}) - \sum_{i=1}^N \mathbf{x}_{i,1} \| \leq B \sqrt{\sum_{i=1}^N \|\mathbf{x}_{t+1} - \mathbf{x}_{0,i}\|^2}$. We then obtain $\sum_{i=1}^N \|\mathbf{x}_{i,1} - \mathbf{x}_{0,i}\| \leq B \sqrt{N}$.

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_{i,1} - \mathbf{x}_{0,i}\| = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_{t+1} - \mathbf{x}_{0,i}\| = \mathbb{E} \{ \|\mathbf{x}_0 + \gamma \mathbf{Z}; \gamma\|, \mathbf{x}_0 \}$$

where we used Theorem 1 and Proposition 2. The case of pseudo-Lipschitz γ is a straightforward generalization. Some notations. For any non-empty subset S of $[m]$ and any $k \times m$ matrix M we refer by M_S to the k by $|S|$ sub-matrix of M that contains only the columns of M corresponding to S . Also define $\langle u, v \rangle = \sum_{i=1}^m u_i v_i$ for $u, v \in \mathbb{R}^m$. Finally, the subgradient of a convex function $f : \mathbb{R}^m \rightarrow \mathbb{R}$ at point $x \in \mathbb{R}^m$ is denoted by $\partial f(x)$. In

particular, remember that the subgradient of the ℓ_1 norm, $\|x\|_1$ is given by

$\|x\|_1 = v \cdot R_m$ such that $-v_i \leq 1$ if $x_i \neq 0$ and $v_i = \text{sign}(x_i)$.

(3.1) 3.1 A structural property of the LASSO cost function

One main challenge in the proof of Theorem 2 lies in the fact that the function $\|x\|_1 + C_{A,y}(x)$ is not in general strictly convex. Hence there can be, in principle, vectors x of cost very close to the optimum and nevertheless far from the optimum. The following Lemma provides conditions under which this does not happen. Lemma 1. There exists a function $\phi(c_1, \dots, c_5)$ such that the following happens. If $x, r \in \mathbb{R}^N$ satisfy the following conditions:

- (1) $\|r\|_2 \leq c_1 N$;
- (2) $C(x+r) \leq C(x)$;
- (3) There exists a subgradient $\text{sg}(C, x) = \phi(x)$ with $\|\text{sg}(C, x)\|_2 \leq N$;
- (4) Let $v = (1/\phi)[A^T(y - Ax) + \text{sg}(C, x)]$ with $\|v\|_2 \leq \phi$, and $S = \{i \in [N] : -v_i \leq 1 - c_2\}$. Then, for any $S' \subseteq [N]$, $|S'| \leq c_3 N$, we have $\min_{i \in S'} (AS(c_2)) \geq c_4$;
- (5) The maximum and minimum non-zero singular value of A satisfy $c_1 \leq \min(A) \leq \max(A) \leq c_5$. Then $\|r\|_2 \leq N \phi(c_1, \dots, c_5)$. Further for any $c_1, \dots, c_5 \geq 0$, $\phi(c_1, \dots, c_5) = 0$ as $\phi \rightarrow 0$.

Further, if $\ker(A) = \{0\}$, the same conclusion holds under conditions 1, 2, 3, and 5. 3.2 Proof of Theorem 3

The proof is based on a series of Lemmas that are used to check the assumptions of Lemma 1. The next lemma implies that submatrices of A constructed using the first t iterations of the AMP algorithm are non-singular (more precisely, have singular values bounded away from 0). Lemma 2. Let $S \subseteq [N]$ be measurable on the σ -algebra \mathcal{S}_t generated by $\{z_0, \dots, z_{t-1}\}$ and $\{x_0 + A^T z_0, \dots, x_{t-1} + A^T z_{t-1}\}$ and assume $|S| \leq N \phi(c)$ for some $c \geq 0$. Then there exists $a_1 = a_1(c) \geq 0$ (independent of t) and $a_2 = a_2(c, t) \geq 0$ (depending on t and c) such that $\min_{i \in S} \{ \min_{j \in S} (AS(c)) : S' \subseteq [N], |S'| \leq a_1 N \} \geq a_2$, with probability converging to 1 as $N \rightarrow \infty$. We will apply this lemma to a specific choice of the set S . Namely, defining $v_t =$

$$\frac{1}{t} \sum_{i=0}^{t-1} (x_{t-1} + A^T z_{t-1} - x_t), \quad (3.2)$$

our last lemma shows convergence of a particular sequence of sets provided by v_t .

Lemma 3. Fix $\phi \in (0, 1)$ and let the sequence $\{S_t(\phi)\}_{t \geq 0}$ be defined by $S_t(\phi) = \{i \in [N] :$

$-v_{it} \leq 1 - \phi\}$. For any $\epsilon \geq 0$ there exists $t_\epsilon = t_\epsilon(\phi, \epsilon) \geq 0$ such that, for all $t_2 \geq t_1 \geq t_\epsilon$: $\lim_{N \rightarrow \infty} \mathbb{P}(|S_{t_2}(\phi) \cap S_{t_1}(\phi)| \leq \epsilon N) = 0$. The last two lemmas imply the following.

Proposition 5. There exist constants $\phi_1 \in (0, 1)$, $\phi_2, \phi_3 \geq 0$ and $t_{\min} \geq 0$ such that, for any $t \geq t_{\min}$, $\min_{i \in S} \{ \min_{j \in S} (AS(\phi_1)) : S' \subseteq [N], |S'| \leq \phi_2 N \} \geq \phi_3$ with probability converging to 1 as $N \rightarrow \infty$. Proof of Theorem 3. We apply Lemma 1 to $x = x_t$, the AMP estimate and $r = x - x_t$ the

distance from the LASSO optimum. The thesis follows by checking conditions 1-5. Namely we need to show that there exists constants $c_1, \dots, c_5 \geq 0$ and, for each $\epsilon \geq 0$ some $t = t(\epsilon)$ such that 1-5 hold with probability going to 1 as $N \rightarrow \infty$. Condition 1 holds since $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \|x_i - b\|_1$ and $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \|x_i - b\|_2$ for all t are finite. Condition 2 is immediate since $x + r = x - b$ minimizes $C(\cdot)$.

Conditions 3-4. Take $v = v_t$ as defined in Eq. (3.2). Using the definition (1.5), it is easy to check that $\|v_t\|_2 \leq \epsilon$. Further it can be shown that $v_t = (1/\epsilon)[A^T(y - Ax_t) + \text{sg}(C, x_t)]$, with $\text{sg}(C, x_t)$ a subgradient satisfying $\lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \|\text{sg}(C, x_t)\|_2^2 = 0$. This proves condition 3 and condition 4 holds by Proposition 5. Condition 5 follows from standard limit theorems on the singular values of Wishart matrices.

Acknowledgement This work was partially supported by a Terman fellowship, the NSF CAREER award CCF-0743978, the NSF grant DMS-0806211 and a Portuguese Doctoral FCT fellowship.

8

2 References

[AJ07]

G. Andrew and G. Jianfeng, Scalable training of l_1 -regularized log-linear models, Proceedings of the 24th international conference on Machine learning, 2007, pp. 33-40.

[BBM10] M. Bayati, J. A. Bento, and A. Montanari, The LASSO risk: asymptotic results and real world examples, Long version (in preparation), 2010.

[BM10a]

[BM10b]

M. Bayati and A. Montanari, The dynamics of message passing on dense graphs, with applications to compressed sensing, Proceedings of IEEE International Symposium on Inform. Theory (ISIT), 2010, Longer version in <http://arxiv.org/abs/1001.3448>, The LASSO risk for gaussian matrices, 2010, preprint available in <http://arxiv.org/abs/1008.2581>.

[CD95]

S.S. Chen and D.L. Donoho, Examples of basis pursuit, Proceedings of Wavelet Applications in Signal and Image Processing III (San Diego, CA), 1995.

[CRT06]

E. Candes, J. K. Romberg, and T. Tao, Stable signal recovery from incomplete and inaccurate measurements, Communications on Pure and Applied Mathematics 59 (2006), 1207-1223.

[CT07]

E. Candes and T. Tao, The Dantzig selector: statistical estimation when p is much larger than n , Annals of Statistics 35 (2007), 2313-2351.

[DMM09] D. L. Donoho, A. Maleki, and A. Montanari, Message Passing Algorithms for Compressed Sensing, Proceedings of the National Academy of Sciences 106 (2009), 18914-18919. [DMM10] D.L. Donoho, A. Maleki, and A. Montanari, The Noise Sensitivity Phase Transition in Compressed Sensing, Preprint, 2010. [GBS09]

- D. Guo, D. Baron, and S. Shamai, A single-letter characterization of optimal noisy compressed sensing, 47th Annual Allerton Conference (Monticello, IL), September 2009.
- [Joh06] I. Johnstone, High Dimensional Statistical Inference and Random Matrices, Proc. International Congress of Mathematicians (Madrid), 2006.
- [KKL+ 07] S. J. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky, A method for large-scale ℓ_1 -regularized least squares., IEEE Journal on Selected Topics in Signal Processing 1 (2007), 606–617. [KM10]
- S. Korada and A. Montanari, Applications of Lindeberg Principle in Communications and Statistical Learning, preprint available in <http://arxiv.org/abs/1004.0557>, 2010.
- [KWT09] Y. Kabashima, T. Wadayama, and T. Tanaka, A typical reconstruction limit for compressed sensing based on ℓ_p -norm minimization, J.Stat. Mech. (2009), L09003. [MM09] M. Mezard and A. Montanari, Information, Physics and Computation, Oxford University Press, Oxford, 2009. [RFG09] S. Rangan, A. K. Fletcher, and V. K. Goyal, Asymptotic analysis of map estimation via the replica method and applications to compressed sensing, PUT NIPS REF, 2009. [Tel99] E. Telatar, Capacity of Multi-antenna Gaussian Channels, European Transactions on Telecommunications 10 (1999), 585–595. [Tib96] R. Tibshirani, Regression shrinkage and selection with the lasso, J. Royal. Statist. Soc B 58 (1996), 267–288. [ZY06]
- P. Zhao and B. Yu, On model selection consistency of Lasso, The Journal of Machine Learning Research 7 (2006), 2541–2563.