

JOINT SCREENING TESTS FOR LASSO

Cédric Herzet*

Angélique Drémeau*

INRIA Rennes-Bretagne Atlantique
Campus de Beaulieu
35000 Rennes, France

ENSTA Bretagne & Lab-STICC
2 rue Francois Verny
29200 Brest, France

ABSTRACT

This paper focusses on “safe” screening techniques for the LASSO problem. Motivated by the need for low-complexity algorithms, we propose a new approach, dubbed “joint screening test”, allowing to screen a set of atoms by carrying out one single test. The approach is particularized to two different sets of atoms, respectively expressed as sphere and dome regions. After presenting the mathematical derivations of the tests, we elaborate on their relative effectiveness and discuss the practical use of such procedures.

Index Terms— ℓ_1 -norm minimization, LASSO, screening techniques.

1. INTRODUCTION

In the last decade, sparse representations have proven to be powerful tools for solving many problems in signal processing, machine learning, etc. Many central methodologies to find a good sparse representation of an observation vector $\mathbf{y} \in \mathbb{R}^m$ in some dictionary $\mathbf{A} = [\mathbf{a}_1 \dots \mathbf{a}_n] \in \mathbb{R}^{m \times n}$ revolve around the resolution of the so-called (nonnegative) LASSO problem:¹

$$\mathbf{x}_\lambda^* \in \arg \min_{\mathbf{x} \geq \mathbf{0}} P_\lambda(\mathbf{y}, \mathbf{x}), \quad (1)$$

where $P_\lambda(\mathbf{y}, \mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$ and $\mathbf{x} \in \mathbb{R}^n$. Without loss of generality, we will assume hereafter that $\|\mathbf{a}_i\|_2 = 1$ for $i = 1 \dots n$.

Solving (1) may require a heavy computational load when the dimension of \mathbf{x} becomes large. Therefore, the conception of computationally-efficient techniques to solve (1) has become an active field of research [2–4]. One important contribution in this field is the so-called “safe”² screening technique proposed by El Ghaoui *et al.* in [5] and refined in several subsequent works [6–13]. These methodologies aim at

*CH was supported by the French National Research Agency through the GERONIMO and BECOSE projects. AD was supported by the French Defense Procurement Agency (DGA) through the TS-DECO MRIS project. CH is currently working in Lab-STICC (UMR 6285), France.

¹The standard LASSO can be written as a particular case of (1), see [1].

²The term “safe” refers to the fact that the elements identified by the screening method always correspond to zeros in \mathbf{x}_λ^* .

decreasing the dimensionality of the problem to handle by identifying some of the zeros of the target solution \mathbf{x}_λ^* via simple tests. Screening procedures have been shown to allow for a dramatic reduction of the complexity needed to solve (1). Nevertheless, their implementation may still be computationally too-demanding in some applications. More specifically, the complexity of all the screening procedures proposed so far evolves linearly with the number of atoms in the dictionary. Hence, in applications involving very large dictionaries,³ the implementation of these screening tests may lead to unacceptable computational burdens. In this paper, we propose a new screening methodology, dubbed “joint screening”, allowing to circumvent this issue.

2. SOME CONVEX CONSIDERATIONS

We first remind some of the properties of the solutions of (1). Problem (1) is convex and always admits (at least) one solution. The dual problem associated to (1) can be written as (see for example [16])

$$\boldsymbol{\theta}_\lambda^* = \arg \min_{\boldsymbol{\theta} \in \mathcal{D}} D_\lambda(\mathbf{y}, \boldsymbol{\theta}), \quad (2)$$

where

$$D_\lambda(\mathbf{y}, \boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y}\|_2^2 - \frac{1}{2} \|\mathbf{y} - \lambda \boldsymbol{\theta}\|_2^2, \quad (3)$$

$$\mathcal{D} = \{\boldsymbol{\theta} \in \mathbb{R}^m : \langle \mathbf{a}_i, \boldsymbol{\theta} \rangle \leq 1, i = 1 \dots n\}, \quad (4)$$

and $\langle \cdot, \cdot \rangle$ denotes the inner product in \mathbb{R}^m . Since $D_\lambda(\mathbf{y}, \boldsymbol{\theta})$ is a strictly concave and coercive function and \mathcal{D} is a closed set, problem (2) admits a unique solution $\boldsymbol{\theta}_\lambda^*$, see [17, Proposition A.8].

The primal and dual solutions $(\mathbf{x}_\lambda^*, \boldsymbol{\theta}_\lambda^*)$ are related through the Karush-Kuhn-Tucker conditions [17, Proposition 5.1.5]:

$$\mathbf{x}_\lambda^* \geq \mathbf{0}, \langle \mathbf{a}_i, \boldsymbol{\theta}_\lambda^* \rangle \leq 1 \text{ for all } i, \quad (5)$$

$$(\langle \mathbf{a}_i, \boldsymbol{\theta}_\lambda^* \rangle - 1) \mathbf{x}_\lambda^*(i) = 0 \text{ for all } i, \quad (6)$$

$$\mathbf{y} = \lambda \boldsymbol{\theta}_\lambda^* + \mathbf{A} \mathbf{x}_\lambda^*, \quad (7)$$

where $\mathbf{x}_\lambda^*(i)$ denotes the i th component of \mathbf{x}_λ^* .

³As an extreme example, one may think of the “continuous” dictionaries considered in [14, 15], containing an infinite number of atoms.

3. STANDARD SCREENING METHODOLOGIES

The “safe” screening procedures proposed in [5–13] leverage on the following observation: if $\mathcal{R} \subset \mathbb{R}^m$ is a region such that $\boldsymbol{\theta}_\lambda^* \in \mathcal{R}$,⁴ then the following inequality trivially holds

$$\langle \mathbf{a}_i, \boldsymbol{\theta}_\lambda^* \rangle \leq \max_{\boldsymbol{\theta} \in \mathcal{R}} \langle \mathbf{a}_i, \boldsymbol{\theta} \rangle,$$

and from (6), we thus have

$$\max_{\boldsymbol{\theta} \in \mathcal{R}} \langle \mathbf{a}_i, \boldsymbol{\theta} \rangle < 1 \Rightarrow \mathbf{x}_\lambda^*(i) = 0. \quad (8)$$

In other words, if the inequality in the left-hand side of (8) is satisfied, one is ensured that the i th component of the solution vector \mathbf{x}_λ^* is equal to zero.

Since the seminal work by El Ghaoui *et al.*, different screening tests, based on different choices of the region \mathcal{R} , have been proposed in the literature. The most popular ones are probably the “sphere” regions, that is

$$\mathcal{R} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \mathbf{c}\|_2 \leq 1 - \tau\}, \quad (9)$$

for some parameters $\mathbf{c} \in \mathbb{R}^m$ and $\tau \leq 1$. Interestingly, for this particular choice, the general screening test (8) takes the following simple form:

$$\langle \mathbf{a}_i, \mathbf{c} \rangle < \tau \Rightarrow \mathbf{x}_\lambda^*(i) = 0. \quad (10)$$

We find in the literature different definitions for the center \mathbf{c} and the radius $1 - \tau$ (see [5–11]) leading to screening tests of different effectiveness.

It is easy to see that the implementation of the standard screening test (8) necessitates the evaluation of $\max_{\boldsymbol{\theta} \in \mathcal{R}} \langle \mathbf{a}_i, \boldsymbol{\theta} \rangle$ for *each* atom of the dictionary. Hence, the computational load required to implement standard screening tests evolves linearly with the number of atoms in the dictionary. For example, in the case where \mathcal{R} is a “sphere” region, (10) involves the computation of the inner product between the center of the sphere \mathbf{c} and each atom of the dictionary, leading to a complexity scaling as $\mathcal{O}(mn)$. In the next section, we propose a new procedure to reduce this computational load.

4. JOINT SCREENING PROCEDURES

In this section, we introduce a new screening procedure having a complexity not depending on the number of atoms in the dictionary. We dub our methodology “joint screening test” because it allows to screen a *set* of atoms by carrying out one *single* test. In a first subsection, we derive tests allowing to screen any atom belonging to some specific region $\mathcal{G} \subset \mathbb{R}^m$. In a second subsection, we elaborate on the relative effectiveness of the proposed test for different choices of the region \mathcal{G} . Finally, in the last subsection, we leverage on these previous results to propose a novel screening procedure having a low complexity with regard to standard screening tests.

⁴Such a region is commonly referred to as “safe region” in the screening literature.

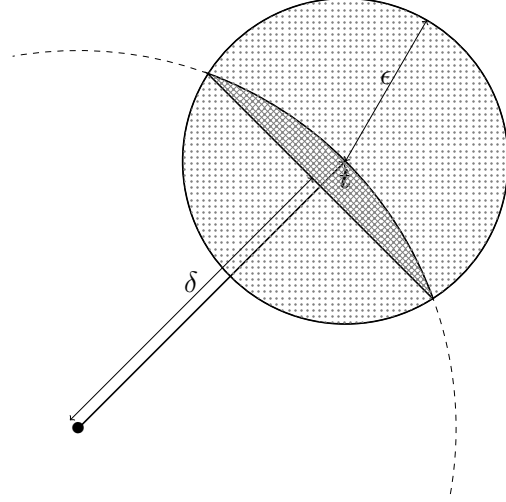


Fig. 1. Illustration of the sphere region $\mathcal{G}^s(\mathbf{t}, \epsilon)$ (dotted area) and the dome region $\mathcal{G}^d(\mathbf{t}, \delta)$ (cross-hatched area). The region inside the dashed curves corresponds to set of atoms with a $\|\cdot\|_2$ -norm smaller than 1.

4.1. Screening the atoms belonging to a region \mathcal{G}

Let $\mathcal{A} = \{\mathbf{a}_i\}_{i=1}^n$ denote the set of atoms of the dictionary and let $\mathcal{R} \subset \mathbb{R}^m$ be a safe region (that is $\boldsymbol{\theta}_\lambda^* \in \mathcal{R}$). The “joint” screening procedure proposed in this paper is a direct consequence of the following observation:⁵

$$\max_{\mathbf{a} \in \mathcal{G}} \max_{\boldsymbol{\theta} \in \mathcal{R}} \langle \mathbf{a}, \boldsymbol{\theta} \rangle < 1 \Rightarrow \mathbf{x}_\lambda^*(i) = 0 \quad \forall i : \mathbf{a}_i \in \mathcal{A} \cap \mathcal{G}. \quad (11)$$

In other words, if the inequality in the left-hand side of (11) is satisfied, all the atoms $\mathbf{a}_i \in \mathcal{A} \cap \mathcal{G}$ can be *safely* and *jointly* screened from problem (1).

In what follows, we will see that the verification of the inequality in the left-hand side of (11) can be done very efficiently for some specific choices of regions \mathcal{R} and \mathcal{G} . First, we will assume that \mathcal{R} is a sphere region (9). The joint screening test (11) then takes the simple form:

$$\max_{\mathbf{a} \in \mathcal{G}} \langle \mathbf{a}, \mathbf{c} \rangle < \tau \Rightarrow \mathbf{x}_\lambda^*(i) = 0 \quad \forall i : \mathbf{a}_i \in \mathcal{A} \cap \mathcal{G}. \quad (12)$$

Moreover, we will consider the two following options for \mathcal{G} :

- “Sphere” : $\mathcal{G}^s(\mathbf{t}, \epsilon) = \{\mathbf{a} : \|\mathbf{a} - \mathbf{t}\|_2 \leq \epsilon\}$,
- “Dome” : $\mathcal{G}^d(\mathbf{t}, \delta) = \{\mathbf{a} : \langle \mathbf{a}, \mathbf{t} \rangle \geq \delta, \|\mathbf{a}\|_2 \leq 1\}$,

where $\mathbf{t} \in \mathbb{R}^m$ and ϵ, δ are some parameters. $\mathcal{G}^s(\mathbf{t}, \epsilon)$ and $\mathcal{G}^d(\mathbf{t}, \delta)$ have some easy geometric interpretations: $\mathcal{G}^s(\mathbf{t}, \epsilon)$ corresponds to the set of vectors located in a ball of radius ϵ centered on \mathbf{t} ; $\mathcal{G}^d(\mathbf{t}, \delta)$ is a dome including all the vectors of norm smaller than one and having an inner product with \mathbf{t} greater than or equal to δ . A graphical representation is given in Fig. 1.

⁵The validity of (11) straightforwardly follows from (8).

For these two choices of regions, the joint screening test defined in (12) admits the following simple analytical solutions:

- *Joint sphere test:* $\max_{\mathbf{a} \in \mathcal{G}^s(\mathbf{t}, \epsilon)} \langle \mathbf{a}, \mathbf{c} \rangle < \tau$ if and only if
$$\langle \mathbf{t}, \mathbf{c} \rangle < \tau - \epsilon \|\mathbf{c}\|_2. \quad (13)$$

- *Joint dome test:* $\max_{\mathbf{a} \in \mathcal{G}^d(\mathbf{t}, \delta)} \langle \mathbf{a}, \mathbf{c} \rangle < \tau$ if and only if⁶

$$\langle \mathbf{t}, \mathbf{c} \rangle < \tau, \quad (14)$$

and

$$\delta > \frac{\langle \mathbf{t}, \mathbf{c} \rangle \tau + \sqrt{\|\mathbf{c}\|_2^2 - \langle \mathbf{t}, \mathbf{c} \rangle^2} \sqrt{\|\mathbf{c}\|_2^2 - \tau^2}}{\|\mathbf{c}\|_2^2}. \quad (15)$$

A proof of this result can be found in Appendix A.

4.2. Relative effectiveness of the screening tests

It is easy to see from (11) that the choice of region \mathcal{G} is a compromise between the number of atoms that can be jointly screened and the ease of passing the test. Indeed, although large regions allow to screen more atoms, they are also less likely to pass the joint screening test since, for any $\mathcal{G}_1 \subseteq \mathcal{G}_2$, we have

$$\max_{\mathbf{a} \in \mathcal{G}_1} \max_{\boldsymbol{\theta} \in \mathcal{R}} \langle \mathbf{a}, \boldsymbol{\theta} \rangle \leq \max_{\mathbf{a} \in \mathcal{G}_2} \max_{\boldsymbol{\theta} \in \mathcal{R}} \langle \mathbf{a}, \boldsymbol{\theta} \rangle. \quad (16)$$

In particular, letting $\mathcal{G}_1 = \{\mathbf{a}_i\}$ and $\mathcal{G}_2 = \mathcal{G}$ in the above inequality, we see that passing the joint screening test (11) requires in particular that the standard screening test (8) is verified for any atom $\mathbf{a}_i \in \mathcal{A} \cap \mathcal{G}$. Hence, quite logically, joint screening test (11) can only lead to inferior screening performance as compared to standard screening test (8).

Another question of interest is the relative effectiveness of the joint sphere and dome tests proposed in (13) and (14)-(15), respectively. The next lemma provides some insights into this question.

Lemma 1. *The smallest⁷ dome containing a set of unit-norm vectors \mathcal{S} is always contained in the smallest sphere containing \mathcal{S} .*

A proof of this statement can be found in Appendix B. In view of (16), a direct consequence of Lemma 1 is as follows: if one wishes to jointly screen a set of unit-norm atoms, there always exists a joint dome test leading to screening performance at least as good as the “best”⁸ joint sphere test.

⁶We assume that $\|\mathbf{t}\|_2 = 1$.

⁷“Smallest” should be understood as the one with the smallest volume.

⁸“Best” should be understood as the test involving the sphere of smallest volume.

4.3. Low-complexity screening procedures

In this section, we discuss the following screening procedure:

- 1) Select a set of L (dome or sphere) regions $\{\mathcal{G}_l\}_{l=1}^L$,
- 2) Apply screening test (13) or (14)-(15) for $l = 1 \dots L$.

We note that since (13) and (14)-(15) only involve the evaluation of one inner product (namely $\langle \mathbf{t}, \mathbf{c} \rangle$), the screening procedure described above only requires to carry out a total number of L inner products.

On the other hand, the effectiveness of this screening procedure obviously depends on the choice of the regions $\{\mathcal{G}_l\}_{l=1}^L$. In order to find a trade-off between the two conflicting objectives emphasized in the first point of the previous section, we consider hereafter some specific choices of the regions $\{\mathcal{G}_l\}_{l=1}^L$.

Let us focus on one particular region \mathcal{G}_l for a given l . We first assume that \mathcal{G}_l is a sphere, that is $\mathcal{G}_l = \mathcal{G}^s(\mathbf{t}, \epsilon)$. Given the value of \mathbf{t} , we want to tune the value of ϵ so that \mathcal{G} is the largest sphere passing (if possible) the joint screening test (13). We note then that the joint sphere test (13) is satisfied as soon as the radius ϵ verifies

$$\epsilon < \epsilon_{\mathbf{t}, \mathbf{c}} \triangleq \frac{\tau - \langle \mathbf{t}, \mathbf{c} \rangle}{\|\mathbf{c}\|_2}. \quad (17)$$

Now, letting the radius ϵ of \mathcal{G}_l tend to its largest value $\epsilon_{\mathbf{t}, \mathbf{c}}$, we have that any atom $\mathbf{a}_i \in \mathcal{A}$ having a distance to \mathbf{t} strictly smaller than $\epsilon_{\mathbf{t}, \mathbf{c}}$ will be screened by test (13), that is

$$\|\mathbf{a}_i - \mathbf{t}\|_2 < \epsilon_{\mathbf{t}, \mathbf{c}} \Rightarrow \mathbf{x}_\lambda^*(i) = 0. \quad (18)$$

We note that (18) defines a screening test on *all* the elements of the dictionary (it may be applied to any atom $\mathbf{a}_i \in \mathcal{A}$) although it only requires the evaluation of one inner product $\langle \mathbf{t}, \mathbf{c} \rangle$ (to evaluate $\epsilon_{\mathbf{t}, \mathbf{c}}$). Let us moreover mention that the quantities $\{\|\mathbf{a}_i - \mathbf{t}\|_2\}_{i=1}^n$ can be precomputed and sorted once for all in advance, so that the identification of the atoms verifying (18) can be done very efficiently (one may for example achieve a complexity scaling as $\mathcal{O}(\log_2 n)$). As a consequence, the “on-line” complexity associated to the implementation of (18) is of the order of $\mathcal{O}(m + \log_2 n)$.

We can apply the same kind of reasoning when \mathcal{G}_l is a dome region, that is $\mathcal{G}_l = \mathcal{G}^d(\mathbf{t}, \delta)$. If we assume that $\mathbf{t} \in \mathbb{R}^m$ is given, a tight lower bound on the value of δ verifying the joint dome test is trivially given by the right-hand side of (15), that is

$$\delta > \delta_{\mathbf{t}, \mathbf{c}}, \quad (19)$$

where

$$\delta_{\mathbf{t}, \mathbf{c}} \triangleq \frac{\langle \mathbf{t}, \mathbf{c} \rangle \tau + \sqrt{\|\mathbf{c}\|_2^2 - \langle \mathbf{t}, \mathbf{c} \rangle^2} \sqrt{\|\mathbf{c}\|_2^2 - \tau^2}}{\|\mathbf{c}\|_2^2}. \quad (20)$$

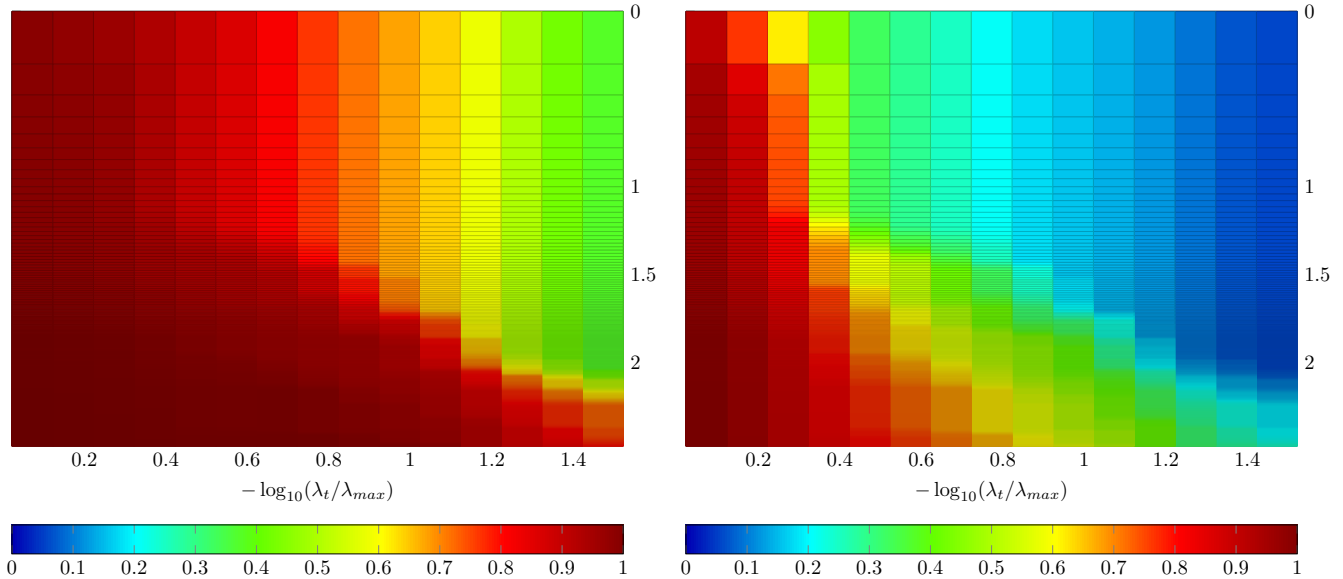


Fig. 2. Proportion of zeros identified by the screening procedures as a function of $-\log_{10}(\lambda_t/\lambda_{max})$ (horizontal axis) and the (\log_{10} of the) number of iterations (vertical axis): GAP sphere test (left) and proposed procedure with a dome region (right)

Hence, provided that $\langle \mathbf{t}, \mathbf{c} \rangle < \tau$, letting parameter δ tend to its smallest value $\delta_{\mathbf{t}, \mathbf{c}}$ will lead to the screening of any atom $\mathbf{a}_i \in \mathcal{A}$ having an inner product with \mathbf{t} strictly greater than $\delta_{\mathbf{t}, \mathbf{c}}$, that is:

$$\begin{cases} \langle \mathbf{t}, \mathbf{c} \rangle < \tau \\ \langle \mathbf{t}, \mathbf{a}_i \rangle > \delta_{\mathbf{t}, \mathbf{c}} \end{cases} \Rightarrow \mathbf{x}_\lambda^*(i) = 0. \quad (21)$$

Again, the quantities $\{\langle \mathbf{t}, \mathbf{a}_i \rangle\}_{i=1}^n$ can be precomputed and sorted once for all in advance, so that a complexity scaling as $\mathcal{O}(m + \log_2 n)$ can also be achieved here.

Going back to the screening procedure advocated at the beginning of this section, we see that adapting the parameter ϵ (resp. δ) for each region \mathcal{G}_l as discussed above is equivalent to applying test (18) (resp. (21)) for each of the L different test vectors \mathbf{t} specifying the regions $\{\mathcal{G}_l\}_{l=1}^L$. The overall complexity of this procedure thus scales as $\mathcal{O}(Lm + L \log_2 n)$. This is to compare to the complexity in $\mathcal{O}(mn)$ of the standard screening tests. The joint screening procedures proposed in this paper will then be of particular interest when dealing with high-dimensional dictionaries.

5. NUMERICAL EXPERIMENTS

In this section we perform some numerical experiments to evaluate the behavior of the proposed method. We confront our screening procedure to the standard one (8) within the following setup. We consider a dictionary $\mathbf{A} \in \mathbb{R}^{100 \times 2000}$ made up of $L = 100$ clusters of 20 atoms. For each cluster, a “seed” atom is created as the m -dimensional realization of a zero-mean circular Gaussian with covariance matrix $m^{-1} \mathbf{I}_m$, where \mathbf{I}_m is the $m \times m$ identity matrix. The other atoms of the

cluster are generated so that their inner products with the seed atom are not smaller than 0.9. We assume that \mathbf{y} is a linear combination of 10 columns (chosen randomly) of the dictionary. The nonzero coefficients are generated as independent realizations of a zero-mean Gaussian with variance equal to 1.

Considering this data set, we target the solution of problem (1) for a decreasing sequence of λ going from $\lambda_{max} \triangleq \|\mathbf{A}^T \mathbf{y}\|_\infty$ to $10^{-1.5} \lambda_{max}$. The solution is searched via the FISTA algorithm [4]. At each iteration of FISTA, the screening procedure is implemented as follows. A “safe” sphere \mathcal{R} is computed according to the GAP procedure proposed in [10]. \mathcal{R} is then used to implement both a standard screening test (8) and the reduced-complexity screening test presented in Section 4.3. Due to space limitation, we only consider the results obtained with the test based on the dome region,⁹ see (21). The L test vectors appearing in (21) are set to be equal to the “seed” vectors used to generate each cluster of the dictionary.

We evaluate the performance of the methods as the good detection rate of zeros in the solution vector \mathbf{x}_λ^* . Fig. 2 presents this figure of merit as a function of λ (horizontal axis) and the iteration number (vertical axis). As expected (see Section 4.2), the standard screening test (left figure) presents better performance than the proposed methodology (right figure). However, this performance must be weighed against the complexity required to perform the tests: the standard procedure requires 2000 scalar products for each test, whereas the proposed method involves only 100 scalar products. We thus see that the proposed procedure allows for a good compromise between computational complexity and

⁹The results corresponding to the test based on the sphere region (18) are however sensibly similar.

the ability to identify zeros of \mathbf{x}_λ^* .

6. CONCLUSIONS

In this paper, we proposed a new screening methodology for the (nonnegative) LASSO problem. Our procedure aims to jointly screen a set of similar atoms by carrying out one single test. We considered two instances of such test (focussing on sets of atoms belonging to either a sphere or a dome region of \mathbb{R}^m) and showed that the latter take a very simple form. In particular, both tests only require the evaluation of *one* inner product in \mathbb{R}^m . Leveraging on this result, we showed that screening procedures for the entire dictionary can be devised by considering an arbitrary number, say L , of regions. The resulting screening procedure has a complexity scaling as $\mathcal{O}(Lm + L \log_2 n)$, where n is the number of atoms in the dictionary. This has to be compared to the complexity of standard screening procedures scaling as $\mathcal{O}(mn)$.

Our future avenues of research include designing new joint tests (based on more refined regions \mathcal{G}) and devising low-complexity methodologies to identify safe regions \mathcal{R} .

A. DERIVATION OF THE JOINT SCREENING TESTS

In this appendix, we provide a the detailed derivations leading to the joint screening tests stated in (13) and (14)-(15). Some steps of our reasoning are based on the technical lemmas stated in Appendix C.

A.1. Sphere joint test (13)

Let us first notice that the sphere region $\mathcal{G}^s(\mathbf{t}, \epsilon)$ can be written as

$$\mathcal{G}^s(\mathbf{t}, \epsilon) = \{\mathbf{a} = \mathbf{t} + \mathbf{z} : \|\mathbf{z}\|_2 \leq \epsilon\}. \quad (22)$$

Therefore, we have

$$\begin{aligned} \max_{\mathbf{a} \in \mathcal{G}^s(\mathbf{t}, \epsilon)} \langle \mathbf{a}, \mathbf{c} \rangle &= \langle \mathbf{t}, \mathbf{c} \rangle + \max_{\|\mathbf{z}\|_2 \leq \epsilon} \langle \mathbf{z}, \mathbf{c} \rangle \\ &= \langle \mathbf{t}, \mathbf{c} \rangle + \epsilon \|\mathbf{c}\|_2, \end{aligned} \quad (23)$$

where the last equality is a consequence of the tightness of the Cauchy-Schwarz inequality. The joint screening test (13) then straightforwardly follows from (23).

A.2. Dome joint test (14)-(15)

We assume $\|\mathbf{t}\|_2 = 1$. We first note that the dome region $\mathcal{G}^d(\mathbf{t}, \delta)$ can be written as

$$\mathcal{G}^d(\mathbf{t}, \delta) = \{\alpha \mathbf{t} + \mathbf{z} : \mathbf{z} \in \mathcal{Z}_\alpha, \delta \leq \alpha \leq 1\}, \quad (24)$$

where

$$\mathcal{Z}_\alpha = \left\{ \mathbf{z} : \mathbf{z} \in (\text{span}[\mathbf{t}])^\perp, \|\mathbf{z}\|_2 \leq \sqrt{1 - \alpha^2} \right\}. \quad (25)$$

Therefore, we have

$$\begin{aligned} \max_{\mathbf{a} \in \mathcal{G}^d(\mathbf{t}, \delta)} (\langle \mathbf{a}, \mathbf{c} \rangle) &= \max_{\delta \leq \alpha \leq 1} \left(\alpha \langle \mathbf{t}, \mathbf{c} \rangle + \max_{\mathbf{z} \in \mathcal{Z}_\alpha} (\langle \mathbf{z}, \mathbf{c} \rangle) \right) \\ &= \max_{\delta \leq \alpha \leq 1} \left(\alpha \langle \mathbf{t}, \mathbf{c} \rangle + \max_{\mathbf{z} \in \mathcal{Z}_\alpha} (\langle \mathbf{z}, P_{\mathbf{t}}^\perp(\mathbf{c}) \rangle) \right) \\ &= \max_{\delta \leq \alpha \leq 1} \underbrace{\left(\alpha \langle \mathbf{t}, \mathbf{c} \rangle + \sqrt{1 - \alpha^2} \|P_{\mathbf{t}}^\perp(\mathbf{c})\|_2 \right)}_{\triangleq g(\delta)}, \end{aligned}$$

where $P_{\mathbf{t}}^\perp(\cdot)$ denotes the projector onto $(\text{span}[\mathbf{t}])^\perp$. The last equality is a consequence of the tightness of the Cauchy-Schwarz inequality.

The joint dome test can thus simply be rewritten as

$$g(\delta) < \tau. \quad (26)$$

Using Lemma 2 with $A = \frac{\langle \mathbf{t}, \mathbf{c} \rangle}{\|\mathbf{c}\|_2}$,¹⁰ we obtain that

$$g(\delta) = \begin{cases} \|\mathbf{c}\|_2 & \text{if } A < \frac{\langle \mathbf{t}, \mathbf{c} \rangle}{\|\mathbf{c}\|_2} \\ \delta \langle \mathbf{t}, \mathbf{c} \rangle + \sqrt{1 - \delta^2} \sqrt{\|\mathbf{c}\|_2^2 - \langle \mathbf{t}, \mathbf{c} \rangle^2} & \text{otherwise} \end{cases} \quad (27)$$

The joint group test (14)-(15) can then be shown as follows. First, satisfying (26) necessarily requires that

$$\tau > \min_{\tilde{\delta} \in [-1, 1]} g(\tilde{\delta}) = g(1) = \langle \mathbf{t}, \mathbf{c} \rangle. \quad (28)$$

This corresponds to the condition enforced by (14).

Moreover, we have from Lemma 3 that

$$\tau \leq \|\mathbf{c}\|_2$$

provided that τ is associated to the radius of a safe sphere. Therefore, if (14) holds, owing to the continuity of g and the fact that it is strictly decreasing over $[\frac{\langle \mathbf{t}, \mathbf{c} \rangle}{\|\mathbf{c}\|_2}, 1]$ (see Lemma 2), there exists $\delta_{\mathbf{t}, \mathbf{c}} \in [\frac{\langle \mathbf{t}, \mathbf{c} \rangle}{\|\mathbf{c}\|_2}, 1]$ such that $g(\delta_{\mathbf{t}, \mathbf{c}}) = \tau$. Using the expression of g over $[\frac{\langle \mathbf{t}, \mathbf{c} \rangle}{\|\mathbf{c}\|_2}, 1]$ in (27), we find

$$\delta_{\mathbf{t}, \mathbf{c}} = \frac{\langle \mathbf{t}, \mathbf{c} \rangle \tau + \sqrt{\|\mathbf{c}\|_2^2 - \langle \mathbf{t}, \mathbf{c} \rangle^2} \sqrt{\|\mathbf{c}\|_2^2 - \tau^2}}{\|\mathbf{c}\|_2^2}. \quad (29)$$

Invoking the strict decrease of g over $[\frac{\langle \mathbf{t}, \mathbf{c} \rangle}{\|\mathbf{c}\|_2}, 1]$ (see Lemma 2), we have that $\tau = g(\delta_{\mathbf{t}, \mathbf{c}}) > g(\delta)$ if and only if

$$\delta > \delta_{\mathbf{t}, \mathbf{c}}.$$

Combining this condition with the expression of $\delta_{\mathbf{t}, \mathbf{c}}$ in (29), we obtain (15).

¹⁰We use the fact that $\|P_{\mathbf{t}}^\perp(\mathbf{c})\|_2 = \sqrt{\|\mathbf{c}\|_2^2 - \langle \mathbf{t}, \mathbf{c} \rangle^2}$.

B. PROOF OF LEMMA 1

On the one hand, the dome of smallest volume including all the elements of \mathcal{S} is given by $\mathcal{G}^d(\mathbf{t}^*, \delta^*)$ with

$$\mathbf{t}^* = \arg \max_{\mathbf{t}: \|\mathbf{t}\|_2=1} \min_{\mathbf{a} \in \mathcal{S}} \langle \mathbf{t}, \mathbf{a} \rangle, \quad (30)$$

$$\delta^* = \min_{\mathbf{a} \in \mathcal{S}} \langle \mathbf{t}^*, \mathbf{a} \rangle. \quad (31)$$

On the other hand, the minimum-volume sphere covering all the elements of \mathcal{S} is given by $\mathcal{G}^s(\tilde{\mathbf{t}}^*, \epsilon^*)$ with

$$\tilde{\mathbf{t}}^* = \arg \min_{\tilde{\mathbf{t}}} \max_{\mathbf{a} \in \mathcal{S}} \|\mathbf{a} - \tilde{\mathbf{t}}\|_2, \quad (32)$$

$$\epsilon^* = \max_{\mathbf{a} \in \mathcal{S}} \|\mathbf{a} - \tilde{\mathbf{t}}^*\|_2. \quad (33)$$

We first show that the optimal parameters (\mathbf{t}^*, δ^*) and $(\tilde{\mathbf{t}}^*, \epsilon^*)$ are related as follows

$$\tilde{\mathbf{t}}^* = \bar{\delta} \mathbf{t}^*, \quad (34)$$

$$\epsilon^* = \sqrt{1 - \bar{\delta}^2}, \quad (35)$$

where

$$\bar{\delta} = \max(0, \delta^*). \quad (36)$$

Indeed, setting $\tilde{\mathbf{t}} = \beta \mathbf{u}$ with $\beta \geq 0$ and $\|\mathbf{u}\|_2 = 1$, (32) can also be rewritten as

$$\begin{aligned} (\beta^*, \mathbf{u}^*) &= \arg \min_{\beta \geq 0, \|\mathbf{u}\|_2=1} \max_{\mathbf{a} \in \mathcal{S}} \|\mathbf{a} - \beta \mathbf{u}\|_2^2, \\ &= \arg \min_{\beta \geq 0, \|\mathbf{u}\|_2=1} \left(\beta^2 - 2\beta \min_{\mathbf{a} \in \mathcal{S}} \langle \mathbf{u}, \mathbf{a} \rangle \right). \end{aligned} \quad (37)$$

From (37), we clearly have that

$$\mathbf{u}^* = \arg \max_{\|\mathbf{u}\|_2=1} \min_{\mathbf{a} \in \mathcal{S}} \langle \mathbf{u}, \mathbf{a} \rangle.$$

In view of (30), we thus have $\mathbf{u}^* = \mathbf{t}^*$. Taking this fact into account, we deduce

$$\beta^* = \arg \min_{\beta \geq 0} (\beta^2 - 2\beta \delta^*) = \bar{\delta},$$

where $\bar{\delta}$ is defined in (36), and thus $\tilde{\mathbf{t}}^* = \bar{\delta} \mathbf{t}^*$. Plugging $\tilde{\mathbf{t}}^* = \bar{\delta} \mathbf{t}^*$ in (33) and using the definition of δ^* in (31), we find

$$\epsilon^* = \sqrt{1 - \bar{\delta}^2}.$$

This shows (34)-(35).

We now prove the result of the lemma, that is

$$\mathcal{G}^d(\mathbf{t}^*, \delta^*) \subseteq \mathcal{G}^s(\tilde{\mathbf{t}}^*, \epsilon^*).$$

This statement can equivalently be rewritten as

$$\forall \mathbf{a} \in \mathcal{G}^d(\mathbf{t}^*, \delta^*) : \|\mathbf{a} - \bar{\delta} \mathbf{t}^*\|_2^2 \leq 1 - \bar{\delta}^2.$$

If $\bar{\delta} = 0$, the inequality is true since $\|\mathbf{a}\|_2 \leq \forall \mathbf{a} \in \mathcal{G}^d(\mathbf{t}^*, \delta^*)$. If $\bar{\delta} = \delta^*$, the inequality is also verified because

$$\begin{aligned} \|\mathbf{a} - \delta^* \mathbf{t}^*\|_2^2 &= 1 + (\delta^*)^2 - 2\delta^* \langle \mathbf{t}^*, \mathbf{a} \rangle \\ &\leq 1 - (\delta^*)^2 \end{aligned}$$

where the last inequality follows from the fact that

$$\forall \mathbf{a} \in \mathcal{G}^d(\mathbf{t}^*, \delta^*) : \langle \mathbf{t}^*, \mathbf{a} \rangle \geq \delta^*.$$

C. MISCELLANEOUS TECHNICAL LEMMAS

In this section, we state and prove two technical lemmas which are useful for the derivation of the dome sphere test in Appendix A. The first lemma (Lemma 2) characterizes the properties of a particular function. The second lemma (Lemma 3) establishes a relationship between the radius parameter τ and the center \mathbf{c} of a safe sphere.

Lemma 2. *If $A \in [-1, 1]$ (resp. $A \in (-1, 1)$), the function*

$$f(\xi) = A\xi + \sqrt{1 - A^2} \sqrt{1 - \xi^2} \quad (38)$$

is concave (resp. strictly concave) over the interval $[-1, 1]$. Moreover, the function $g(\xi) \triangleq \max_{\xi \leq \xi' \leq 1} f(\xi')$ can be written as follows in the interval $\xi \in [-1, 1]$:

$$g(\xi) = \begin{cases} 1 & \text{if } \xi < A \\ A\xi + \sqrt{1 - A^2} \sqrt{1 - \xi^2} & \text{otherwise} \end{cases}. \quad (39)$$

$g(\xi)$ is concave and non-increasing over the interval $[-1, 1]$. If $A \in (-1, 1)$, it is strictly concave and strictly decreasing over the interval $[A, 1]$.

Proof: The strict concavity of $f(\xi)$ over $[-1, 1]$ for $A \in (-1, 1)$ can be proved by showing that its second derivative is strictly negative over this interval. Now, we have

$$f''(\xi) = -\sqrt{\frac{1 - A^2}{(1 - \xi^2)^3}} < 0 \quad \forall \xi \in [-1, 1].$$

The concavity of $f(\xi)$ over $[-1, 1]$ for $A \in [-1, 1]$ then follows by noticing that $f(\xi)$ is a linear function when $A = \pm 1$.

The expression of $g(\xi)$ in (39) can be found as follows. If $A = 1$, we have

$$\max_{\xi \leq \xi' \leq 1} f(\xi') = \max_{\xi \leq \xi' \leq 1} \xi' = 1 \quad \forall \xi \in [-1, 1]. \quad (40)$$

This corresponds to the definition of $g(\xi)$ in (39). If $A = -1$, we have

$$\max_{\xi \leq \xi' \leq 1} f(\xi') = \max_{\xi \leq \xi' \leq 1} -\xi' = -\xi \quad \forall \xi \in [-1, 1]. \quad (41)$$

This again corresponds to the definition of $g(\xi)$ in (39). Let us finally consider the case where $A \in (-1, 1)$. Using the first part of the lemma, we have that $f(\xi')$ is strictly concave

over $\xi' \in [-1, 1]$. Simple calculus shows that $\xi' = A$ cancels out the first derivative of $f(\xi')$. We thus have

$$1 = f(A) = \max_{-1 \leq \xi' \leq 1} f(\xi'). \quad (42)$$

Using standard optimality conditions for concave problems, we then obtain (39).

Finally, the concavity and the non-increasing (resp. the strict concavity and the strictly decreasing) nature of $g(\xi)$ over $[-1, 1]$ (resp. $[A, 1]$) for $A \in [-1, 1]$ (resp. $A \in (-1, 1)$) follows from its definition of $g(\xi)$ and the concavity (resp. strict concavity) of $f(\xi)$. \square

Lemma 3. Let $\mathcal{R} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \mathbf{c}\|_2 \leq 1 - \tau\}$ be a safe sphere for problem (2) with $\lambda < \lambda_{\max}$ (that is $\boldsymbol{\theta}_\lambda^* \in \mathcal{R}$). Then, we have

$$\tau \leq \|\mathbf{c}\|_2. \quad (43)$$

Proof: If $\mathcal{R} = \{\boldsymbol{\theta} : \|\boldsymbol{\theta} - \mathbf{c}\|_2 \leq 1 - \tau\}$ is a safe region, then

$$\|\boldsymbol{\theta}_\lambda^* - \mathbf{c}\|_2 \leq 1 - \tau, \quad (44)$$

which leads, by using a triangle inequality, to

$$\|\boldsymbol{\theta}_\lambda^*\|_2 - 1 + \tau \leq \|\mathbf{c}\|_2. \quad (45)$$

Now, if $\lambda < \lambda_{\max}$, we have

$$\|\boldsymbol{\theta}_\lambda^*\|_2 \geq 1. \quad (46)$$

The latter claim follows from the following arguments: if $\lambda < \lambda_{\max}$, we necessarily have $\mathbf{x}_\lambda^*(i) > 0$ for some $i \in [1 \dots n]$. From the optimality condition (7), we have for such i : $\langle \mathbf{a}_i, \boldsymbol{\theta}_\lambda^* \rangle = 1$. Hence, we obtain (46) by using the Cauchy-Schwarz inequality:

$$\|\boldsymbol{\theta}_\lambda^*\|_2 \geq \langle \mathbf{a}_i, \boldsymbol{\theta}_\lambda^* \rangle = 1.$$

Finally, we obtain the main result (43) by combining (45) and (46). \square

D. REFERENCES

- [1] Mario A. T. Figueiredo, “Teaching a new trick to an old dog: Revisiting the quadratic programming formulation of sparse recovery using ADMM,” in *Acoustics, Speech and Signal Processing (ICASSP), 2014 IEEE International Conference on*. May 2014, pp. 1512–1516, IEEE.
- [2] M. R. Osborne, B. Presnell, and B. A. Turlach, “A new approach to variable selection in least squares problems,” *IMA Journal of Numerical Analysis*, vol. 20, no. 3, pp. 389–403, July 2000.
- [3] Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani, “Least angle regression,” *The Annals of Statistics*, vol. 32, no. 2, pp. 407–451, 2004.
- [4] Amir Beck and Marc Teboulle, “A fast iterative shrinkage-thresholding algorithm for linear inverse problems,” *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, Jan. 2009.
- [5] Laurent El Ghaoui, Vivian Viallon, and Tarek Rabbani, “Safe feature elimination in sparse supervised learning,” Tech. Rep. UC/EECS-2010-126, EECS Dept., University of California at Berkeley, Sept. 2010.
- [6] Zhen J. Xiang, Hao Xu, and Peter J. Ramadge, “Learning sparse representations of high dimensional data on large scale dictionaries,” in *Advances in Neural Information Processing Systems 24*, J. Shawe-taylor, Zemel, P. Bartlett, Pereira, and Weinberger, Eds., pp. 900–908, 2011.
- [7] Liang Dai and Kristiaan Pelckmans, “An ellipsoid based, two-stage screening test for BPDN,” in *Signal Processing Conference (EUSIPCO), 2012 Proceedings of the 20th European*. Aug. 2012, pp. 654–658, IEEE.
- [8] Jie Wang, Peter Wonka, and Jieping Ye, “Lasso screening rules via dual polytope projection,” *Journal of Machine Learning Research*, 2015.
- [9] Antoine Bonnefoy, Valentin Emiya, Liva Ralaivola, and Remi Gribonval, “Dynamic screening: Accelerating first-order algorithms for the Lasso and group-Lasso,” *Signal Processing, IEEE Transactions on*, vol. 63, no. 19, pp. 5121–5132, Oct. 2015.
- [10] O. Fercoq, A. Gramfort, and J. Salmon, “Mind the duality gap: safer rules for the Lasso,” in *ICML*, 2015.
- [11] C. Herzet and A. Malti, “Safe screening tests for LASSO based on firmly non-expansiveness,” in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2016, pp. 4732–4736.
- [12] Zhen J. Xiang and Peter J. Ramadge, “Fast Lasso screening tests based on correlations,” in *Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on*. Mar. 2012, pp. 2137–2140, IEEE.
- [13] Zhen J. Xiang, Yun Wang, and Peter J. Ramadge, “Screening tests for Lasso problems,” arXiv:1405.4897v1, May 2014.
- [14] Emmanuel J. Candès and Carlos Fernandez-Granda, “Towards a mathematical theory of super-resolution,” *Comm. Pure Appl. Math*, vol. 67, no. 6, pp. 906–956, June 2014.

- [15] Vincent Duval and Gabriel Peyré, “Exact support recovery for sparse spikes deconvolution,” *Foundations of Computational Mathematics*, vol. 15, no. 5, pp. 1315–1355, 2015.
- [16] S. Foucart and H. Rauhut, *A mathematical introduction to compressive sensing.*, Applied and Numerical Harmonic Analysis. Birkhäuser, 2013.
- [17] D. P. Bertsekas, *Nonlinear Programming*, Athena Scientific, USA, 2003.
- [18] Cédric Herzet and Angélique Drémeau, “Joint screening tests for Lasso,” *arXiv:1710.09809*, 2017.