



HAL
open science

LASSO and Iterative Feature Selection: Oracle Inequalities and Numerical Performances

Pierre Alquier

► **To cite this version:**

Pierre Alquier. LASSO and Iterative Feature Selection: Oracle Inequalities and Numerical Performances. 2007. hal-00181784v1

HAL Id: hal-00181784

<https://hal.science/hal-00181784v1>

Preprint submitted on 24 Oct 2007 (v1), last revised 25 Nov 2008 (v4)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

LASSO AND ITERATIVE FEATURE SELECTION: ORACLE INEQUALITIES AND NUMERICAL PERFORMANCES

PIERRE ALQUIER

ABSTRACT. We propose a general family of algorithms for regression estimation with quadratic loss. Our algorithms is able to select relevant functions into a large dictionary. We prove that some algorithms that have already been studied (LASSO, by Tibshirani [15], Iterative Feature Selection, in our paper [1], among others) belong to our family. We prove oracle-type inequalities in some particular cases, and compare numerical performances of LASSO and Iterative Feature Selection on a toy example.

CONTENTS

1. Introduction	1
1.1. Setting of the problem	1
1.2. Organization of the paper	3
2. General projection algorithms	4
2.1. Additional notations and hypothesis	4
2.2. General description of the algorithm	5
2.3. First particular case: LASSO	5
2.4. Particular case: Iterative Feature Selection	6
2.5. Particular case: Generalization of the Group LASSO	7
3. Oracle inequalities in some particular cases	7
3.1. Order on the dictionary of functions	7
3.2. Nearly orthogonal dictionary of functions	9
4. Numerical simulations	11
4.1. Motivation	11
4.2. Description of the experiments	11
4.3. Results and comments	12
5. Conclusion	13
References	13

1. INTRODUCTION

1.1. Setting of the problem. Let $n \in \mathbb{N} \setminus \{0\}$. Let P be a probability distribution on

$$\left((\mathcal{X} \times \mathbb{R})^n, (\mathcal{B} \otimes \mathcal{B}_{\mathbb{R}})^{\otimes n} \right)$$

Date: October 24, 2007.

2000 Mathematics Subject Classification. Primary 62G08; Secondary 62J07, 62G15, 68T05.

Key words and phrases. Regression estimation, statistical learning, confidence regions, shrinkage and thresholding methods, LASSO..

I Would like to thank Professors Olivier Catoni and Alexandre Tsybakov for useful remarks.

and

$$\left((X_1, Y_1), \dots, (X_n, Y_n) \right)$$

drawn from P .

For $i \in \{1, \dots, n\}$, let p_i denote the marginal distribution of X_i under P , and let us put:

$$P_X = \frac{1}{n} \sum_{i=1}^n p_i.$$

We assume that P_X is known to the statistician.

Moreover, we put:

$$\bar{P} = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}.$$

The statistician chooses a dictionary of functions: (f_1, \dots, f_m) . For the sake of simplicity we assume that it is such that for any $j \in \{1, \dots, m\}$ we have:

$$P_X [f_j^2] = 1.$$

Definition 1.1. Let us put, for any $\alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m$ and $(x, y) \in \mathcal{X} \times \mathbb{R}$:

$$l_\alpha(x, y) = \left(y - \sum_{j=1}^m \alpha_j f_j(x) \right).$$

We define:

$$r(\alpha) = \bar{P}(l_\alpha) = \frac{1}{n} \sum_{i=1}^n \left[Y_i - \sum_{j=1}^m \alpha_j f_j(X_i) \right]^2$$

and

$$R(\alpha) = P[r(\alpha)].$$

We put:

$$\bar{\alpha} \in \arg \min_{\alpha \in \mathbb{R}^m} R(\alpha).$$

For any $\alpha, \alpha' \in \mathbb{R}^m$ we put:

$$\langle \alpha, \alpha' \rangle_X = P_X \left[\sum_{j=1}^m \sum_{k=1}^m \alpha_j \alpha'_k f_j f_k \right],$$

and

$$\|\alpha\|_X = \sqrt{\langle \alpha, \alpha \rangle}.$$

Finally, we put $e_1 = (1, 0, \dots, 0) \in \mathbb{R}^m$, ..., $e_m = (0, \dots, 0, 1) \in \mathbb{R}^m$ the canonical basis of \mathbb{R}^m .

Let us remark that for any $\alpha \in \mathbb{R}^m$ we have:

$$R(\alpha) - R(\bar{\alpha}) = \|\alpha - \bar{\alpha}\|_X^2.$$

Remark 1.1. We think of two cases of interest. If the pairs (X_i, Y_i) are i. i. d. we have $p_1 = \dots = p_n = P_X$ and so P_X is the marginal distribution of X . It is assumed to be known to the statistician (restrictive hypothesis).

Another case of interest is when the values X_1, \dots, X_n are deterministic. In this case:

$$P_X = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$$

and so we obtain:

$$\langle \alpha, \alpha' \rangle_X = \frac{1}{n} \sum_{i=1}^n \left[\sum_{j,k} \alpha_j \alpha'_k f_j(X_i) f_k(X_i) \right].$$

In this case $\|\cdot\|_X$ is called the empirical norm (usually denoted $\|\cdot\|_n$).

Definition 1.2. Let \mathcal{C} be a closed, convex subset of \mathbb{R}^d . We let $\Pi_{\mathcal{C}}^X(\cdot)$ denote the orthogonal projection on \mathcal{C} with respect to the norm $\|\cdot\|_X$; $\Pi_{\mathcal{C}}(\cdot)$ will denote the orthogonal projection on \mathcal{C} with respect to the euclidian norm $\|\cdot\|$.

1.2. Organization of the paper. The aim of this paper is to propose a method to estimate the real regression function (say f) by selecting a few relevant functions among all the functions in the dictionary.

Recently, a lot of algorithms have been proposed for that purpose, let's cite among others the LASSO by Tibshirani [15] and some variants or generalization like LARS by Efron, Hastie, Johnstone and Tibshirani [11], the Dantzig selector by Candès and Tao [6] and the Group LASSO by Yuan and Lin [17], or Iterative Feature Selection in our paper [1]. This paper proposes a general algorithm that contains LASSO and Iterative Feature Selection as a particular case.

A paper by Butucea, Tsybakov and Wegkamp [5] gives sparsity oracle inequalities for the LASSO, that is inequalities that bounds the risk of the LASSO estimators in terms of the number of selected functions in the dictionary. This paper by Butucea and al. is written in a different context than ours: random design with *unknown* distribution (in the case of a random design, our method require the knowledge of the distribution of the design). Another paper, by Bickel, Ritov and Tsybakov [3] gives sparsity oracle inequalities for the LASSO and the Dantzig selector in the case of the deterministic design. However, in both papers the main results require the assumption $\|f_j\|_{\infty} \leq L$ for some given L that is not necessary in our paper, and that prevents the use of popular basis of functions like wavelets. This is partly due to the use of Hoeffding's inequality.

Our paper uses a geometric point of view that allows to obtain simple sparsity oracle inequalities for the obtained estimator, in both deterministic design case and random design with known distribution. It also uses a deviation inequality proved in a previous work [1] that is sharper than Hoeffding's inequality, and so get rid of the assumption of a (uniform) bound over the functions of the dictionary. Another improvement is that our method is valid for some types of data-dependant of dictionaries of function, for example the case where $m = n$ and:

$$\{f_1(\cdot), \dots, f_m(\cdot)\} = \{K(X_1, \cdot), \dots, K(X_n, \cdot)\}$$

where K is a function $\mathcal{X}^2 \rightarrow \mathbb{R}$.

In section 2, we give the general form for our algorithm and explicit some particular cases (LASSO, Iterative Feature Selection and Group LASSO) under a particular assumption (CRA, Definition 2.2) that says that we are able to build some confidence region for the best value of α in some subspace of \mathbb{R}^m . In section 3, we give some oracle inequalities under some hypothesis on the dictionary of functions.

Finally, section 4 is dedicated to simulations: we compare ordinary least square (OLS), LASSO and Iterative Feature Selection on a toy example. Simulations shows that both particular cases of our family of estimators (LASSO and Iterative Feature Selection) generally outperforms the OLS estimate. Moreover, LASSO performs generally better than Iterative Feature Selection, however, this is not always true: this fact leads to the conclusion that a data-driven choice of a particular algorithm in our general family could lead to optimal results.

2. GENERAL PROJECTION ALGORITHMS

2.1. Additional notations and hypothesis. We choose $M \in \mathbb{N}$ and $S_1 \subset \{1, \dots, m\}, \dots, S_M \subset \{1, \dots, m\}$. We put, for every $S \subset \{1, \dots, m\}$:

$$\mathcal{M}_S = \left\{ \alpha \in \mathbb{R}^m, \quad \ell \notin S \Rightarrow \alpha_\ell = 0 \right\}.$$

So every \mathcal{M}_{S_j} is a submodel of the original model \mathbb{R}^m .

Definition 2.1. We put, for every $S \subset \{1, \dots, m\}$:

$$\bar{\alpha}_S = \arg \min_{\alpha \in \mathcal{M}_S} R(\alpha).$$

Remark that for every $S \subset \{1, \dots, m\}$:

$$\bar{\alpha}_S = \Pi_{\mathcal{M}_S}^X(\bar{\alpha}).$$

Moreover let us put:

$$\hat{\alpha}_S = \arg \min_{\alpha \in \mathcal{M}_S} r(\alpha).$$

Definition 2.2. We say that the confidence region assumption (CRA) is satisfied if $\varepsilon \in [0, 1]$ we have a bound $r(S_j, \varepsilon) \in \mathbb{R}$ such that

$$P \left[\forall j \in \{1, \dots, M\}, \quad \|\bar{\alpha}_{S_j} - \hat{\alpha}_{S_j}\|_X \leq r(S_j, \varepsilon) \right] \geq 1 - \varepsilon.$$

Such confidence regions for $\bar{\alpha}_{S_j}$ can be obtained with standard techniques and various hypothesis on the probability P , we refer the reader to our previous work [1] for example.

Definition 2.3. We define, for any $\varepsilon > 0$ and $j \in \{1, \dots, M\}$, the random set:

$$\mathcal{CR}(j, \varepsilon) = \left\{ \alpha \in \mathbb{R}^m, \quad \left\| \Pi_{\mathcal{M}_{S_j}}^X(\alpha) - \hat{\alpha}_{S_j} \right\|_X^2 \leq r(S_j, \varepsilon) \right\}.$$

We remark that the hypothesis implies that:

$$P \left[\forall j \in \{1, \dots, M\}, \quad \bar{\alpha} \in \mathcal{CR}(j, \varepsilon) \right] \geq 1 - \varepsilon.$$

In our previous work [1] we examined different hypothesis on the probability P such that this hypothesis is satisfied. For example, using inequalities by Catoni [7] and Panchenko [13] we proved the following results (for models of dimension 1, that will be the most used in the sequel of this paper).

Lemma 2.1. Let us assume that $P = \mathbb{P}_1 \otimes \dots \otimes \mathbb{P}_n$. Let us assume that $Y_i = f(X_i) + \varepsilon_i$ with $\mathbb{P}(\varepsilon_i | X_i) = 0$,

$$\sup_{i \in \{1, \dots, n\}} \mathbb{P}_i(\varepsilon_i^2 | X_i) \leq \sigma^2$$

for some known σ and that $\|f\|_\infty \leq L$ for some known $L > 0$. If we take $S_j = \{j\}$ for any $j \in \{1, \dots, m\}$, assumption CRA is satisfied with:

$$r(\{j\}, \varepsilon) = \frac{4(1 + \log \frac{2m}{\varepsilon})}{n} \left[\frac{1}{n} \sum_{i=1}^n f_j^2(X_i) Y_i^2 + L^2 + \sigma^2 \right].$$

Remark 2.1. It is also shown in [1] that we are allowed to take:

$$\{f_1, \dots, f_m\} = \{K(X_1, \cdot), \dots, K(X_n, \cdot)\}$$

for some function $\mathcal{X}^2 \rightarrow \mathbb{R}$, this being also true in the random design case, but we have to take:

$$r(\{j\}, \varepsilon) = \frac{4(1 + \log \frac{4m}{\varepsilon})}{n} \left[\frac{1}{n} \sum_{i=1}^n f_j^2(X_i) Y_i^2 + L^2 + \sigma^2 \right].$$

Lemma 2.2. *Let us assume that $P = \mathbb{P}_1 \otimes \dots \otimes \mathbb{P}_n$ and that X_1, \dots, X_n are deterministic. Let us assume that there is a $K > 0$ such that $\mathbb{P}_i(|Y_i| \leq K) = 1$ for any i . If we take $S_j = \{j\}$ for any $j \in \{1, \dots, m\}$, assumption CRA is satisfied with:*

$$r(\{j\}, \varepsilon) = \frac{8K^2 \left(1 + \log \frac{2m}{\varepsilon}\right)}{n}.$$

2.2. General description of the algorithm. Now let us choose $N \leq M$ and indices $(j_1, \dots, j_N) \in \{1, \dots, M\}^N$, the region:

$$\bigcap_{\ell=1}^N \mathcal{CR}(j_\ell, \varepsilon)$$

is a closed, convex confidence region for $\bar{\alpha}$.

So we propose the following iterative algorithm.

- Step 0. Choose $\hat{\alpha}^0 = (0, \dots, 0) \in \mathbb{R}^m$. Choose $\varepsilon \in [0, 1]$.
- General Step (k). Choose $N(k) \leq M$ and indices $(j_1^{(k)}, \dots, j_{N(k)}^{(k)}) \in \{1, \dots, M\}^{N(k)}$ and put:

$$\hat{\alpha}^k = \Pi_{\bigcap_{\ell=1}^{N(k)} \mathcal{CR}(j_\ell^{(k)}, \varepsilon)} (\hat{\alpha}^{k-1}).$$

Theorem 2.3. *When the CRA assumption is satisfied we have:*

$$P \left[\forall k \in \mathbb{N}, \quad R(\hat{\alpha}^k) \leq R(\hat{\alpha}^0) - \sum_{j=1}^k \|\hat{\alpha}^j - \hat{\alpha}^{j-1}\|_X^2 \right] \geq 1 - \varepsilon.$$

Proof. Let us choose $k \in \mathbb{N}$.

$$\begin{aligned} R(\hat{\alpha}^k) - R(\bar{\alpha}) &= \|\hat{\alpha}^k - \bar{\alpha}\|_X^2 = \left\| \Pi_{\bigcap_{\ell=1}^{N(k)} \mathcal{CR}(j_\ell^{(k)}, \varepsilon)} (\hat{\alpha}^{k-1}) - \bar{\alpha} \right\|_X^2 \\ &\leq \|\hat{\alpha}^{k-1} - \bar{\alpha}\|_X^2 - \left\| \Pi_{\bigcap_{\ell=1}^{N(k)} \mathcal{CR}(j_\ell^{(k)}, \varepsilon)} (\hat{\alpha}^{k-1}) - \hat{\alpha}^{k-1} \right\|_X^2 \\ &= R(\hat{\alpha}^{k-1}) - R(\bar{\alpha}) - \|\hat{\alpha}^k - \hat{\alpha}^{k-1}\|_X^2. \end{aligned}$$

A recurrence ends the proof. \square

We choose as our estimator $\hat{\alpha} = \hat{\alpha}^k$ for some step $k \in \mathbb{N}$; the choice of the stopping step k will depend of the particular choices of the projections and is detailed in what follows.

2.3. First particular case: LASSO. We first look at the case where $S_j = \{j\}$ for any $j \in \{1, \dots, m\}$ (and so $M = m$). In this case, we only use submodels of dimension 1.

Here, we use only one step where we project 0 onto the intersection of all the confidence regions and so we obtain:

$$\hat{\alpha} = \hat{\alpha}^1 = \Pi_{\bigcap_{\ell=1}^m \mathcal{CR}(\ell, \varepsilon)} (0).$$

Definition 2.4. *Let us put, for any $j \in \{1, \dots, m\}$:*

$$\tilde{\alpha}_j = \frac{1}{n} \sum_{i=1}^n Y_i f_j(X_i).$$

Note that we have:

$$\hat{\alpha}_{S_j} = \hat{\alpha}_{\{j\}} = (0, \dots, 0, \tilde{\alpha}_j, 0, \dots, 0)$$

with the $\tilde{\alpha}_j$ in j -th position, and that:

$$\mathcal{CR}(j, \varepsilon) = \left\{ \alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m, \quad \tilde{\alpha}_j - r(\{j\}, \varepsilon) \leq \langle \alpha, e_j \rangle_X \leq \tilde{\alpha}_j + r(\{j\}, \varepsilon) \right\}.$$

The optimization program to obtain $\hat{\alpha}$ is given by:

$$\begin{cases} \arg \min_{\alpha=(\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m} \|\alpha\|_X^2 \\ \text{s. t. } \alpha \in \bigcap_{\ell=1}^m \mathcal{CR}(\ell, \varepsilon) \end{cases}$$

and so:

$$\begin{cases} \arg \min_{\alpha \in \mathbb{R}^m} \|\alpha\|_X^2 \\ \text{s. t. } \forall j \in \{1, \dots, m\}, \quad |\langle \alpha, e_j \rangle_X - \tilde{\alpha}_j| \leq \sqrt{r(\{j\}, \varepsilon)} \end{cases}$$

We can write the program in dual form:

$$(2.1) \quad \arg \min_{\alpha \in \mathbb{R}^m} \left\{ \|\alpha\|_X^2 - 2 \sum_{j=1}^m \alpha_j \tilde{\alpha}_j + 2 \sum_{j=1}^m \sqrt{r(\{j\}, \varepsilon)} |\alpha_j| \right\}.$$

If $\|\cdot\|_X$ is the empirical norm we obtain:

$$\begin{aligned} \|\alpha\|_X^2 - 2 \sum_{j=1}^m \alpha_j \tilde{\alpha}_j &= \frac{1}{n} \sum_{i=1}^n \left[\sum_{j=1}^m \alpha_j f_j(X_i) \right]^2 - 2 \frac{1}{n} \sum_{i=1}^n Y_i \left[\sum_{j=1}^m \alpha_j f_j(X_i) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \left[Y_i - \sum_{j=1}^m \alpha_j f_j(X_i) \right]^2 - \frac{1}{n} \sum_{i=1}^n Y_i^2 \\ &= r(\alpha) - \frac{1}{n} \sum_{i=1}^n Y_i^2, \end{aligned}$$

and so the program is equivalent to:

$$\arg \min_{\alpha \in \mathbb{R}^m} \left\{ r(\alpha) + 2 \sum_{j=1}^m \sqrt{r(\{j\}, \varepsilon)} |\alpha_j| \right\}.$$

Note that, if $r(\{j\}, \varepsilon)$ does not depend on j , this is exactly the formulation of the original LASSO algorithm as introduced by Tibshirani [15]. An explicit algorithm to obtain the projection is given by Efron, Hastie, Johnstone and Tibshirani [11].

However, in the cases where $r(\{j\}, \varepsilon)$ is not constant, the difference with the LASSO algorithm is the following: coordinates that are more difficult to estimate (because the confidence interval is larger) are more penalized.

Moreover, note that the program 2.1 gives a form different of the usual LASSO program for the cases where we do not use the empirical norm.

2.4. Particular case: Iterative Feature Selection. Here, we choose general subsets $S_1, \dots, S_m \subset \{1, \dots, N\}$.

Moreover, instead of taking the intersection of every confidence region, we project on each of them iteratively. So the algorithm is the following:

$$\hat{\alpha} = (0, \dots, 0)$$

and at each step k we choose a $j(k) \in \{1, \dots, m\}$ and

$$\hat{\alpha}^k = \Pi_{\mathcal{CR}(j(k), \varepsilon)}^X(\hat{\alpha}^{k-1}).$$

In the case where, as in the LASSO, we actually have $S_j = \{j\}$ for any j , this is exactly the Iterative Feature Selection algorithm that was introduced in Alquier [1], with the choice of $j(k)$:

$$j(k) = \arg \max_j \left\| \hat{\alpha}^{k-1} - \Pi_{\mathcal{CR}(j, \varepsilon)}^X(\hat{\alpha}^{k-1}) \right\|_X,$$

and the suggestion to take as an estimator:

$$\hat{\alpha} = \hat{\alpha}^{\hat{k}}$$

where

$$\hat{k} = \inf \{k \in \mathbb{N}^*, \quad \|\hat{\alpha}^k - \hat{\alpha}^{k-1}\|_X \leq \kappa\}$$

for some small $\kappa > 0$. In [1] is also given the explicit computation of every step of this algorithm.

2.5. Particular case: Generalization of the Group LASSO. Here we choose general subsets $S_1, \dots, S_M \subset \{1, \dots, N\}$.

As in the LASSO algorithm we only use one step where we project 0 onto the intersection of all the confidence regions and so we obtain:

$$\hat{\alpha} = \hat{\alpha}^1 = \Pi_{\bigcap_{\ell=1}^M \mathcal{C}\mathcal{R}(\ell, \varepsilon)}^X(0).$$

The optimization program to obtain $\hat{\alpha}$ is given by:

$$\begin{cases} \arg \min_{\alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m} \|\alpha\|_X^2 \\ \text{s. t. } \forall j \in \{1, \dots, M\}, \quad \left\| \Pi_{\mathcal{M}_{S_j}^X}(\alpha) - \hat{\alpha}_j \right\|_X \leq \sqrt{r(\{j\}, \varepsilon)}. \end{cases}$$

In the case of the empirical norm, this program is equivalent to the following:

$$\arg \min_{\alpha \in \mathbb{R}^m} \left\{ r(\alpha) + \sum_{j=1}^M \sqrt{r(\{j\}, \varepsilon)} \|\Pi_{\mathcal{M}_j} \alpha\|_X \right\},$$

that is a generalization of the Group LASSO algorithm defined by Yuan and Lin [17] in the case of orthogonal basis functions and extended by Chesneau and Hebiri [8] to the general case.

3. ORACLE INEQUALITIES IN SOME PARTICULAR CASES

Some particular assumptions about the dictionary of functions chosen by the statistician allow us to obtain oracle inequalities for some particular order of projection.

3.1. Order on the dictionary of functions. In this subsection, we assume that there is an order on the basis function: in some sense, the statistician knows that a function f_j with a small indice $j \in \{1, \dots, m\}$ is more likely to be useful for his regression problem than another function $f_{j'}$ with a large indice $j' \in \{1, \dots, m\}$.

A usual way to formalize this hypothesis is to make the following regularity assumption.

Definition 3.1. *We say that the ordered regularity assumption with order $\beta > 0$ and constant $C > 0$ is satisfied if, for any $j \in \{1, \dots, m\}$, we have:*

$$\|\bar{\alpha}_{\{1, \dots, j\}} - \bar{\alpha}\|_X \leq Cj^{-\beta}.$$

Remark 3.1. This is a Sobolev-type regularity assumption, see Tsybakov [16] and the references therein for estimation in Sobolev spaces.

Note that this is equivalent to:

$$R(\bar{\alpha}_{\{1, \dots, j\}}) - R(\bar{\alpha}) \leq Cj^{-2\beta}.$$

Let us put $S_1 = \{1\}$, $S_2 = \{1, 2\}$, ..., $S_m = \{1, \dots, m\}$ and so $M = m$ and follow the following iterative projection scheme:

$$\hat{\alpha} = \Pi_{\mathcal{C}\mathcal{R}(m, \varepsilon)}^X \dots \Pi_{\mathcal{C}\mathcal{R}(2, \varepsilon)}^X \Pi_{\mathcal{C}\mathcal{R}(1, \varepsilon)}^X(0).$$

Then we have the following result.

Theorem 3.1. *Let us assume that the CRA assumption is satisfied. Then we have:*

$$(3.1) \quad P \left\{ R(\hat{\alpha}) \leq R(\bar{\alpha}) + \inf_{1 \leq j \leq m} [R(\bar{\alpha}_{\{1, \dots, j\}}) - R(\bar{\alpha}) + 4r(\{1, \dots, j\}, \varepsilon)] \right\} \geq 1 - \varepsilon.$$

If we assume moreover that there is a $k > 0$ such that:

$$r(\{1, \dots, j\}, \varepsilon) \leq \frac{jk \log \frac{m}{\varepsilon}}{n}$$

and that the ordered regularity assumption is satisfied with regularity β and constant C then we have:

$$(3.2) \quad P \left\{ R(\hat{\alpha}) \leq R(\bar{\alpha}) + (2\beta + 1) C^{\frac{1}{2\beta+1}} \left(\frac{2k \log \frac{m}{\varepsilon}}{\beta n} \right)^{\frac{2\beta}{2\beta+1}} + \left(\frac{4k \log \frac{m}{\varepsilon}}{n} \right) \right\} \geq 1 - \varepsilon.$$

Proof. Let us assume that the event:

$$\left\{ \forall j \in \{1, \dots, M\}, \quad \|\bar{\alpha}_{S_j} - \hat{\alpha}_{S_j}\|_X \leq r(S_j, \varepsilon) \right\}$$

is satisfied (this is true with probability at least $1 - \varepsilon$ thanks to assumption CRA). For any $j \in \{1, \dots, m\}$ we have:

$$\hat{\alpha}^j = \Pi_{\mathcal{C}\mathcal{R}(j, \varepsilon)}^X \dots \Pi_{\mathcal{C}\mathcal{R}(2, \varepsilon)}^X \Pi_{\mathcal{C}\mathcal{R}(1, \varepsilon)}^X(0),$$

and $\hat{\alpha} = \hat{\alpha}^m$. It is evident that $\hat{\alpha}^j \in \mathcal{C}\mathcal{R}(j, \varepsilon)$. Moreover, we can prove that $\hat{\alpha}^j \in \mathcal{M}_{\{1, \dots, j\}}$ by recurrence. So we have:

$$\|\hat{\alpha}^j - \bar{\alpha}_{\{1, \dots, j\}}\|_X^2 \leq 4r(S_j, \varepsilon),$$

which means that:

$$R(\hat{\alpha}^j) \leq R(\bar{\alpha}_{\{1, \dots, j\}}) + 4r(S_j, \varepsilon).$$

Now, Theorem 2.3 ensures that:

$$R(\hat{\alpha}) = R(\hat{\alpha}^m) \leq R(\hat{\alpha}^j),$$

this proves Equation 3.1:

$$R(\hat{\alpha}) \leq R(\bar{\alpha}) + \inf_{1 \leq j \leq m} [R(\bar{\alpha}_{\{1, \dots, j\}}) - R(\bar{\alpha}) + 4r(S_j, \varepsilon)].$$

For Equation 3.2 note that:

$$R(\hat{\alpha}) \leq R(\bar{\alpha}) + \inf_{1 \leq j \leq m} \left[Cj^{-2\beta} + \frac{4jk \log \frac{m}{\varepsilon}}{n} \right],$$

and take:

$$j = \left\lceil \left(\frac{\beta C n}{2k \log \frac{m}{\varepsilon}} \right)^{\frac{1}{2\beta+1}} \right\rceil + 1$$

to conclude. \square

Remark 3.2. Note that, as soon as $\beta > 1/2$, we reach the minimax rate of convergence up to a $\log m$ term. In general, as β is unknown to the statistician, we propose to take $m = n$ so we miss the optimal rate of convergence by a $\log n$ term. In this case, we propose the following modification (inspired by the block method, see Tsybakov [16] and the references therein for example). We take $S_h = \{1, \dots, 2^h\}$ with $h \in \{1, \dots, m\}$ and

$$m = \left\lfloor \frac{\log n}{\log 2} \right\rfloor$$

as this choice ensures that $2^{m-1} < n \leq 2^m$ so we obtain:

$$\begin{aligned} R(\hat{\alpha}) &\leq R(\bar{\alpha}) + \inf_{1 \leq h \leq m} [R(\bar{\alpha}_{\{1, \dots, 2^h\}}) - R(\bar{\alpha}) + r(S_h, \varepsilon)] \\ &\leq R(\bar{\alpha}) + \inf_{1 \leq h \leq m} \left[2^{-2\beta h} + \frac{2^h 4k \log \frac{\log \lfloor \frac{\log n}{\log 2} \rfloor}{\varepsilon}}{n} \right] \end{aligned}$$

and so we can conclude:

$$R(\hat{\alpha}) \leq R(\bar{\alpha}) + \mathcal{O} \left[\left(\frac{\log \frac{\log n}{\varepsilon}}{n} \right)^{\frac{2\beta}{2\beta+1}} \right],$$

so in this case we reach the minimax rate of convergence up to a $\log \log n$ term.

3.2. Nearly orthogonal dictionary of functions.

Definition 3.2. For a pair $(\delta, D) \in (\mathbb{R}_+)^2$, we say that the condition $\text{NO}(\delta, D)$ is satisfied if:

$$\forall \alpha \in \mathbb{R}^m, \quad \delta \|\alpha\| \leq \|\alpha\|_X \leq D \|\alpha\|.$$

Remark 3.3. This condition was given by Kerkycharian and Picard [12] to study the statistical properties of algorithms generalizing the idea of thresholding. The meaning of condition $\text{NO}(\delta, D)$ when $\delta \simeq 1$ and $D \simeq 1$ is clear: the norms $\|\cdot\|$ and $\|\cdot\|_X$ have the same behaviour, in other words the dictionary of functions $(f_j)_{j \in \{1, \dots, m\}}$ is "nearly orthogonal". Of course, as we are in a space of finite dimension m , we can always find a small δ (of the order $1/m$) and a large D (of the order m) such that condition $\text{NO}(\delta, D)$ is satisfied. However, Theorem 3.2 clearly shows that the results are interesting for values of δ and D that does not depend on m .

We will also use the following regularity hypothesis.

Definition 3.3. We say that the general regularity assumption with order $\beta > 0$ and constant $C > 0$ if, for any $j \in \{1, \dots, m\}$, we have:

$$\inf_{\substack{S \subset \{1, \dots, m\} \\ |S| \leq j}} \|\bar{\alpha}_S - \bar{\alpha}\|_X \leq C j^{-\beta}.$$

Remark 3.4. This is the type of regularity assumption used to define weak Besov spaces, see Cohen [9] and the references therein.

We still take $S_1 = \{1\}, \dots, S_m = \{m\}$, and

$$\hat{\alpha} = \Pi_{\bigcap_{j=1}^m \mathcal{C}\mathcal{R}(j, \varepsilon)}^X(0)$$

the LASSO estimator.

Theorem 3.2. Let us assume that the CRA assumption is satisfied. For any $(\delta, D) \in (\mathbb{R}_+)^2$ such that condition $\text{NO}(\delta, D)$ is satisfied we have:

$$P \left\{ R(\hat{\alpha}) \leq R(\bar{\alpha}) + \frac{1}{\delta^2} \inf_{S \subset \{1, \dots, m\}} \left[D^2 [R(\bar{\alpha}_S) - R(\bar{\alpha})] + 4 \sum_{j \in S} r(\{j\}, \varepsilon) \right] \right\} \geq 1 - \varepsilon,$$

and so, if moreover the general regularity assumption is satisfied with regularity $\beta > 0$ and constant $C > 0$ and if there is a $k > 0$ such that:

$$r(\{j\}, \varepsilon) \leq \frac{k \log \frac{m}{\varepsilon}}{n}$$

then we have:

$$P \left\{ R(\hat{\alpha}) \leq R(\bar{\alpha}) + \frac{(2\beta+1)C^{\frac{1}{2\beta+1}}}{\delta^2} \left(\frac{2Dk \log \frac{m}{\varepsilon}}{\beta n} \right)^{\frac{2\beta}{2\beta+1}} + \frac{1}{\delta^2} \left(\frac{4k \log \frac{m}{\varepsilon}}{n} \right) \right\} \geq 1 - \varepsilon.$$

Proof. In this proof, we adopt the notation α' for the transposed vector of α , and let M denote the matrix:

$$M = (\langle e_j, e_k \rangle_X)_{(j,k) \in \{1, \dots, m\}^2} = (P_X(f_j f_k))_{(j,k) \in \{1, \dots, m\}^2}.$$

So note that for any $(\alpha, \beta) \in (\mathbb{R}^m)^2$ we have:

$$\langle \alpha, \beta \rangle = \alpha' \beta$$

and

$$\langle \alpha, \beta \rangle_X = \alpha' M \beta.$$

We have:

$$\begin{aligned} R(\hat{\alpha}) - R(\bar{\alpha}) &= \|\hat{\alpha} - \bar{\alpha}\|_X^2 = \left\| M^{\frac{1}{2}} (\hat{\alpha} - \bar{\alpha}) \right\|^2 \\ &\leq \frac{1}{\delta^2} \left\| M^{\frac{1}{2}} (\hat{\alpha} - \bar{\alpha}) \right\|_X^2 = \frac{1}{\delta^2} \|M (\hat{\alpha} - \bar{\alpha})\|^2 = \frac{1}{\delta^2} (\hat{\alpha} - \bar{\alpha})' M M (\hat{\alpha} - \bar{\alpha}) \\ &= \frac{1}{\delta^2} \text{Tr} [(\hat{\alpha} - \bar{\alpha})' M M (\hat{\alpha} - \bar{\alpha})] = \frac{1}{\delta^2} \text{Tr} [M (\hat{\alpha} - \bar{\alpha}) (\hat{\alpha} - \bar{\alpha})' M] \\ &= \frac{1}{\delta^2} \sum_{j=1}^m e_j' M (\hat{\alpha} - \bar{\alpha}) (\hat{\alpha} - \bar{\alpha})' M e_j = \frac{1}{\delta^2} \sum_{j=1}^m \langle \hat{\alpha} - \bar{\alpha}, e_j \rangle_X^2. \end{aligned}$$

Now, note that the fact that for any $j \in \{1, \dots, m\}$, $\hat{\alpha} \in \mathcal{CR}(j, \varepsilon)$ implies that:

$$\forall j \in \{1, \dots, m\}, \quad \langle \hat{\alpha} - \bar{\alpha}, e_j \rangle_X^2 \leq 4r(\{j\}, \varepsilon).$$

Moreover, the by the definition of $\hat{\alpha}$ we have:

$$\forall j \in \{1, \dots, m\}, \quad \langle \hat{\alpha} - \bar{\alpha}, e_j \rangle_X^2 \leq \langle \bar{\alpha}, e_j \rangle_X^2.$$

So, for any $S \subset \{1, \dots, m\}$ we have:

$$R(\hat{\alpha}) - R(\bar{\alpha}) \leq \frac{1}{\delta^2} \left[4 \sum_{j \in S} r(\{j\}, \varepsilon) + \sum_{j \notin S} \langle \bar{\alpha}, e_j \rangle_X^2 \right].$$

So, in order to prove Theorem 3.2 we just have to upper bound the second sum, using the same techniques:

$$\begin{aligned} \sum_{j \notin S} \langle \bar{\alpha}, e_j \rangle_X^2 &= \sum_{j=1}^m \langle \bar{\alpha} - \bar{\alpha}_S, e_j \rangle_X^2 = \sum_{j=1}^m e_j' M (\bar{\alpha} - \bar{\alpha}_S) (\bar{\alpha} - \bar{\alpha}_S)' M e_j \\ &= \text{Tr} [M (\bar{\alpha} - \bar{\alpha}_S) (\bar{\alpha} - \bar{\alpha}_S)' M] = \text{Tr} [(\bar{\alpha} - \bar{\alpha}_S)' M M (\bar{\alpha} - \bar{\alpha}_S)] \\ &= \|M (\bar{\alpha} - \bar{\alpha}_S)\|^2 = \left\| M^{\frac{1}{2}} (\bar{\alpha} - \bar{\alpha}_S) \right\|_X^2 \leq D^2 \left\| M^{\frac{1}{2}} (\bar{\alpha} - \bar{\alpha}_S) \right\|^2 \\ &= D^2 \|\bar{\alpha} - \bar{\alpha}_S\|_X^2 = D^2 [R(\bar{\alpha}_S) - R(\bar{\alpha})]. \end{aligned}$$

□

Remark 3.5. A particularly interesting case is the orthogonal case, where condition NO(1, 1) is satisfied. In this case:

$$\hat{\alpha} = \Pi_{\mathcal{CR}(m, \varepsilon)}^X \Pi_{\mathcal{CR}(m-1, \varepsilon)}^X \dots \Pi_{\mathcal{CR}(1, \varepsilon)}^X (0) = \Pi_{\bigcap_{j=1}^m \mathcal{CR}(j, \varepsilon)}^X (0)$$

so in this case, LASSO and Iterative Feature Selection (with $\kappa = 0$) are equivalent and actually we have:

$$\hat{\alpha} = \sum_{j=1}^m \operatorname{sgn}(\tilde{\alpha}_j) \left(|\tilde{\alpha}_j| - \sqrt{r(j, \varepsilon)} \right)_+ e_j$$

that is a soft-thresholded estimator. See our previous work on Iterative Feature Selection [1] for a proof. Soft-thresholding is now a standard way to deal with selection of functions in an orthogonal family, see for example the seminal paper by Donoho and Johnstone [10] in the case of a wavelet basis. The bound just becomes:

$$P \left\{ R(\hat{\alpha}) \leq R(\bar{\alpha}) + \inf_{S \subset \{1, \dots, m\}} \left[R(\bar{\alpha}_S) - R(\bar{\alpha}) + 4 \sum_{j \in S} r(\{j\}, \varepsilon) \right] \right\} \geq 1 - \varepsilon.$$

4. NUMERICAL SIMULATIONS

4.1. Motivation. We compare here LASSO and Iterative Feature Selection on a toy example, introduced by Tibshirani [15]. We also compare their performances to the ordinary least square (OLS) estimate as a benchmark. Note that we will not propose a very fine choice for the $r(\{j\}, \varepsilon)$. The idea of these simulations is not to identify a good choice for the penalization in practice. The idea is to observe the similarity and differences between different order in projections in our general algorithm, using the same confidence regions.

4.2. Description of the experiments. The model defined by Tibshirani [15] is the following. We have:

$$\forall i \in \{1, \dots, 20\}, \quad Y_i = \langle \beta, X_i \rangle + \varepsilon_i$$

with $X_i \in \mathcal{X} = \mathbb{R}^8$, $\beta \in \mathbb{R}^8$ and the ε_i are i. i. d. from a gaussian distribution with mean 0 and standard device σ .

The X_i 's are i. i. d. too, and each X_i comes from a gaussian distribution with mean $(0, \dots, 0)$ and with variance-covariance matrix:

$$\Sigma(\rho) = \left(\rho^{|i-j|} \right)_{\substack{i \in \{1, \dots, 8\} \\ j \in \{1, \dots, 8\}}}$$

for $\rho \in [0, 1[$.

We will use the three particular values for β taken by Tibshirani [15]:

$$\beta^1 = (3, 1.5, 0, 0, 2, 0, 0, 0),$$

$$\beta^2 = (1.5, 1.5, 1.5, 1.5, 1.5, 1.5, 1.5, 1.5),$$

$$\beta^3 = (5, 0, 0, 0, 0, 0, 0, 0),$$

corresponding to a "sparse" situation (β^1), a "non-sparse" situation (β^2) and a "very sparse" situation (β^3).

We use two values for σ : 1 (the "low noise case") and 3 (the "noisy case").

Finally, we use two values for ρ : 0.1 ("weakly correlated variables") and 0.5 ("highly correlated variables").

We run each example (corresponding to a given value of β , σ and ρ) 250 times. We use the software R [14] for simulations. We implement Iterative Feature Selection as described in subsection 2.4 page 6, while using the standard OLS estimate and the LASSO estimator given by the LARS package described in [11]. The choice:

$$r(\{j\}, \varepsilon) = \frac{\sigma}{3} \sqrt{\frac{\log m}{n}} = \frac{\sigma}{3} \sqrt{\frac{\log 8}{20}}$$

was not motivated by theoretical considerations but seems to perform well in practice.

4.3. **Results and comments.** The results are reported in Table 1.

TABLE 1. Results of the Simulations. For each possible combination of β , σ and ρ , we report in a column the mean empirical loss over the 250 simulations, the standard deviation of this quantity over the simulations and finally the mean number of non-zero coefficients in the estimate, this for each estimate, ordinary least square (OLS), LASSO and Iterative Feature Selection (IFS).

β	σ	ρ	OLS	LASSO	IFS
β^1 (sparse)	3	0.5	3.67 1.84 8	1.64 1.25 4.64	1.56 1.20 4.62
	1	0.5	0.40 0.22 8	0.29 0.19 5.42	0.36 0.23 5.70
	3	0.1	3.75 1.86 8	2.72 1.50 5.70	2.85 1.58 5.66
	1	0.1	0.40 0.19 8	0.30 0.19 5.92	0.31 0.19 5.96
β^2 (non sparse)	3	0.5	3.54 8 1.82	3.36 7.08 1.64	4.90 6.57 1.58
	1	0.5	0.41 8 0.21	0.54 7.94 0.93	0.84 7.89 0.36
	3	0.1	3.78 8 1.78	3.82 7.06 1.51	4.50 7.03 1.59
	1	0.1	0.40 8 0.20	0.42 7.98 0.29	0.71 7.98 0.32
β^3 (very sparse)	3	0.5	3.55 8 1.79	1.65 4.48 1.28	1.59 4.49 1.27
	1	0.5	0.40 8 0.21	0.18 4.46 0.14	0.17 4.48 0.14
	3	0.1	3.46 8 1.74	1.69 4.92 1.29	1.62 4.92 1.18
	1	0.1	0.40 8 0.20	0.20 4.98 0.14	0.19 4.91 0.14

The following remarks can easily be made in view of the results:

- both methods based on projection on random confidence regions clearly outperforms the OLS in the sparse cases, moreover they present the advantage of giving sparse estimates;

- in the non-sparse case, the OLS performs generally better than the other methods, but LASSO is very close, it is known that a better choice for the value $r(\{j\}, \varepsilon)$ would lead to a better result (see Tibshirani [15]);
- LASSO seems to be the best method on the whole set of experiments. In every case, it is never the worst method, and always performs almost as well as the best method;
- in the "sparse case" (β^1), note that IFS and LASSO are very close for the small value of ρ . This is coherent with the previous theory, see remark 3.5 page 10;
- IFS gives very bad results in the non-sparse case (β^2), but is the best method in the sparse case (β^3). This last point tends to indicate that different situations should lead to a different choice for the confidence regions we are to project on. However, theoretical results leading on that choice are missing.

5. CONCLUSION

This paper provides a simple interpretation of well-known algorithms of statistical learning theory in terms of orthogonal projections on confidence regions. This very intuitive approach provides a very simple way to prove oracle inequalities.

Also note that this approach can be easily extended into general statistical problems with quadratic loss: in our paper [2], the Iterative Feature Selection method is generalized to the density estimation with quadratic loss problem, leading to a proposition of a LASSO-like program for density estimation, that have also been proposed and studied by Bunea, Tsybakov and Wegkamp [4] under the name SPADES.

Simulations shows that methods based on confidence regions clearly outperforms the OLS estimate in most examples. However, theoretical results leading the statistician to a particular choice for the order of the successive projections are still missing.

REFERENCES

- [1] ALQUIER, P. Iterative feature selection in regression estimation. Preprint Laboratoire de Probabilités et Modèles Aléatoires (submitted to the Annales de l'IHP, accepted in 2007), 2005.
- [2] ALQUIER, P. Density estimation with quadratic loss: A confidence intervals method. Preprint Laboratoire de Probabilités et Modèles Aléatoires (submitted to ESAIM PS, accepted in 2007), 2006.
- [3] BICKEL, P. J., RITOV, Y., AND TSYBAKOV, A. Simultaneous analysis of lasso and dantzig selector. Preprint Submitted to the Annals of Statistics, 2007.
- [4] BUNEA, F., TSYBAKOV, A., AND WEGKAMP, M. Sparse density estimation with ℓ_1 penalties. In *Proceedings of 20th Annual Conference on Learning Theory (COLT 2007)* (2007), Springer-Verlag, pp. 530–543.
- [5] BUNEA, F., TSYBAKOV, A., AND WEGKAMP, M. Sparsity oracle inequalities for the lasso. *Electronic Journal of Statistics* 1 (2007), 169–194.
- [6] CANDÈS, E., AND TAO, T. The dantzig selector: statistical estimation when p is much larger than n . *The Annals of Statistics* 35 (2007).
- [7] CATONI, O. A pac-bayesian approach to adaptive classification. Preprint Laboratoire de Probabilités et Modèles Aléatoires, 2003.
- [8] CHESNEAU, C., AND HEBIRI, M. Some theoretical results on the grouped variable lasso. Preprint Laboratoire de Probabilités et Modèles Aléatoires (submitted), 2007.
- [9] COHEN, A. *Handbook of Numerical Analysis*, vol. 7. North-Holland, Amsterdam, 2000.
- [10] DONOHO, D., AND JOHNSTONE, I. Ideal spatial adaptation by wavelet shrinkage. *Biometrika* 81 (1994), 425–455.
- [11] EFRON, B., HASTIE, T., JOHNSTONE, I., AND TIBSHIRANI, R. Least angle regression. *The Annals of Statistics* 32, 2 (2004), 407–499.
- [12] KERKYCHARIAN, G., AND PICARD, D. Thresholding in learning theory. Preprint Laboratoire de Probabilités et Modèles Aléatoires, 2005.

- [13] PANCHENKO, D. Symmetrization approach to concentration inequalities for empirical processes. *The Annals of Probability* 31, 4 (2003), 2068–2081.
- [14] TEAM, R. D. C. R: A language and environment for statistical computing. Vienna, Austria. URL: <http://www.R-project.org/>, 2004.
- [15] TIBSHIRANI, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society B* 58, 1 (1996), 267–288.
- [16] TSYBAKOV, A. *Introduction à l'Estimation Non-Paramétrique*. Springer, 2004.
- [17] YUAN, M., AND LIN, Y. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society B* 68, 1 (2006), 49–67.

CREST, AND, LABORATOIRE DE PROBABILITÉS ET MODÈLES ALÉATOIRES (UNIVERSITÉ PARIS 7),,
175, RUE DU CHEVALERET, 75252 PARIS CEDEX 05, FRANCE.

URL: <http://www.crest.fr/pageperso/alquier/alquier.htm>

E-mail address: alquier@ensae.fr