

→ Geometry! :-)

→ gl methods à rando + formels

→ gl term "influence"

natural tendency to  
please sparse solutions?

(pouvoir cette nég.)

→ Sparse is an init bien  
cok de due tendency  
→ il n'y a pas trop de  
solution sparse.

Quel enchoir de  
l'algo qu'il ya  
ait)

# 1 Geometry-dependent matching pursuit: a transition phase for convergence of 2 linear regression and LASSO

3 Céline Moucer<sup>†\*</sup>, Francis Bach<sup>†</sup>, and Adrien Taylor<sup>†</sup>

5 Abstract. Greedy first-order methods, such as coordinate descent with Gauss-Southwell rule or matching pur-  
6 suit, have become popular in optimization due to their natural sparse solutions and their refined  
7 convergence guarantees. In this work, we propose a principled approach to generating (regularized)  
8 matching pursuit algorithms adapted to the geometry of the problem at hand, as well as their con-  
9 vergence guarantees. Building on these results, we derive approximate convergence guarantees and  
10 describe a transition phenomenon from underparametrized to overparametrized models.

11 Key words. Optimization, first-order methods, matching pursuit, linear regression, LASSO.

C'est rajouter un peu plus une horloge "en place" ?  
→ "in the conveyor of" ?

12 MSC codes. 90C25, 68Y25, 60B20

natural → par?

13 1. Introduction. Many problems in machine learning and data science take the form of  
14 an  $\ell_1$ -regularized minimization problem:

15 (1.1) 
$$G_* = \min_{\alpha \in \mathbb{R}^d} f(P\alpha) + \lambda \|\alpha\|_1 = F(\alpha) + H(\alpha),$$

16 where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is a smooth strongly convex function,  $P \in \mathbb{R}^{n \times d}$  and  $n, d$  respectively denote  
17 the number of samples and  $d$  the dimension of the problem. Typically, in the vanilla least-  
18 squares regression problem,  $H(\alpha) = 0$  and  $F(\alpha) = f(P\alpha) = \frac{1}{2n} \|P\alpha - y\|_2^2$ , where  $P$  corresponds  
19 to the input data,  $y \in \mathbb{R}^n$  to the labels,  $d$  to the number of features (or parameters) and  $n$  the  
20 number of observations. If in addition  $H(\alpha) = \lambda \|\alpha\|_1$ , Problem (1.1) is exactly the LASSO  
21 problem [57], that belongs to more general variational problems appearing in Fenchel duality  
22 theory [6, Section 15.3]. Problem (1.1) is often compared to its constrained counterpart,

23 (1.2) 
$$\min_{\alpha \in \mathbb{R}^d} F(\alpha), \text{ such that } \|\alpha\|_1 \leq R,$$

24 where  $\lambda$  may be seen as the Lagrange multiplier associated to the constraint  $\|\alpha\|_1 \leq R$  with  
25  $R > 0$ . Problems (1.1) and (1.2) arise when looking for sparsity patterns, such as in signal  
26 processing where we aim for models depending on a small number of variables, or for trace-norm  
27 regularized problems, when looking for low-rank solutions [17]. In particular, Problem (1.1)  
28 is a popular way to induce sparsity on the solution for a well-chosen range of  $\lambda$ . Thus,  $\ell_1$ -  
29 penalization (or constraint) is strongly connected to sparsity and often presented as a convex  
30 substitute for  $\ell_0$ -penalization problems [62, ~~Introduction~~], that correspond to feature selection.

31 First-order methods have become popular to solve optimization Problems (1.1) and (1.2),  
32 due to their low cost per iteration and to the limited accuracy requirements in machine learn-  
33 ing [12, Section 7 and 8]. Within these methods, different algorithms may be used depending

\*Ecole Nationale des Ponts et Chaussées, Marne-la-Vallée, France.

<sup>†</sup>DI ENS, École normale supérieure, Université PSL, CNRS, INRIA, 75005 Paris, France (celine.moucer@inria.fr),  
(adrien.taylor@inria.fr), (francis.bach@inria.fr).

1

[62, Section 1]?

il doit y avoir  
une meilleure w.

→  $\ell_1$  = envol.

convexe de  $f$  sur  $[-1, 1]$

→ Dayan Fogel?

→ T. Hastie ou F. Cucker?

This manuscript is for review purposes only.

34 on the properties of functions  $F$  and  $H$ . For instance, a first-order method may call for the  
 35 gradient, or for the proximal operator [45] as in the proximal gradient method, or the linear  
 36 minimization oracle as in the Frank-Wolfe algorithm [28] for the constrained version (1.2).  
 37 These methods benefit from convergence guarantees, that are mostly helpful to practitioners  
 38 when they can be numerically computed *a priori*.

39 In the context of sparsity, traditional first-order methods, such as the proximal gradient, do  
 40 not always lead to sparse solutions [27]. Boosting strategies (also known as matching pursuit)  
 41 have been developed to ensure sparse representations of approximate solutions [34, 58]. At each  
 42 iteration, a possibly new atom (also referred to as a weak-learner in the boosting literature,  
 43 or a coordinate in the context of coordinate descent) is greedily selected as a best candidate  
 44 among a set of atoms, and combined to past iterates. While boosting benefits from strong  
 45 statistical properties [58], from an optimization perspective, their convergence analyses often  
 46 rely on extra statistical assumptions [63]. More recently, randomized and greedy coordinate  
 47 descent methods have gained interest due to their low-cost per iteration in high dimension [39]  
 48 and to their implicit induced sparsity [7, 18].

49 Correspondences have been highlighted between first-order methods and boosting strate-  
 50 gies for non-regularized minimization problems ( $\lambda = 0$ ), leading to convergence guarantees  
 51 independent of traditional statistical assumptions. For example, coordinate descent has been  
 52 interpreted as matching pursuit [33], as well as Frank-Wolfe algorithms [28, 32] for constrained  
 53 Problems (1.2), by formulating them as minimizers of well-chosen quadratic upper approxi-  
 54 mations. These analyses strongly rely on a well-chosen geometry, characterized by a gauge  
 55 function [22]. To our knowledge, this comparison was only drawn for non-regularized problems  
 56 for which  $\lambda = 0$  [33] and for constrained problems [54].

57 Problem (1.1) in  $\mathbb{R}^d$  can be naturally formulated as an optimization problem in  $\mathbb{R}^n$ , letting  
 58 the gauge geometry appear,

$$\begin{aligned} G_* &= \min_{\alpha \in \mathbb{R}^d, x \in \mathbb{R}^n} G(\alpha) = f(x) + \lambda \|\alpha\|_1, \text{ such that } x = P\alpha, \\ 59 \quad (1.3) \quad &= \min_{x \in \mathbb{R}^n} f(x) + \lambda \inf_{\alpha \in \mathbb{R}^d, x = P\alpha} \|\alpha\|_1, \\ &= \min_{x \in \mathbb{R}^n} f(x) + \lambda \gamma_{\mathcal{P}}(x), \end{aligned}$$

60 where the gauge function is defined by  $\gamma_{\mathcal{P}}(x) = \inf_{\alpha \in \mathbb{R}^d, x = P\alpha} \|\alpha\|_1$  with  $\mathcal{P} = \text{conv}(\{\pm P_{:,i}, i = 1, \dots, d\})$  the centrally symmetric convex hull of the columns of  $P$ . Gauge functions may  
 61 be seen as generalized versions of the  $\ell_1$ -norm, providing a sparse representation  $\alpha \in \mathbb{R}^d$  of  
 62 a vector  $x \in \mathbb{R}^n$  with respect to a set of atoms. Under some assumptions on  $P$ , the gauge  
 63 function may be a norm, as we will see in Section 2. Let us for example take  $\mathcal{P} = \text{conv}(\pm e_i)$ ,  
 64 then  $\gamma_{\mathcal{P}}(x) = \|x\|_1$ .

65 Due to the connection between optimizing in  $\mathbb{R}^d$  and in  $\mathbb{R}^n$ , it is possible to derive al-  
 66 gorithms adapted to one or the other geometry, and to formulate geometry-adapted conver-  
 67 gence guarantees. For the  $\ell_1$ -geometry, Nutini et al. [43, Section 4] analyzed greedy coordi-  
 68 nate descent by considering strong convexity with respect to the  $\ell_1$ -norm, and formulated the  
 69 strong convexity parameter as an optimization problem [43, Appendix 4.1]. More generally,  
 70 D'Aspremont et al. [16, Section 2] extended smoothness and strong convexity with respect  
 71 to the gauge  $\gamma_{\mathcal{P}}$ , which led to formulations of the smoothness parameter as an optimization

73 problem in the work of Sun and Bach [54, Section 2.5]. These optimization problems are often  
 74 hard to solve (yet, they have closed-form reformulation in some cases).

75 The main idea of this work is to propose a principled view on gradient boosting meth-  
 76 ods, that are obtained by minimizing a smoothness upper bound with respect to the  $\ell_1$ -norm.  
 77 This methodology leads to a new boosting strategy for regularized problems, that benefits  
 78 from (sub)linear convergence properties. Unlike former methods, such as orthogonal matching  
 79 pursuit under restricted isometry property (RIP) by [62], convergence analysis is performed  
 80 without statistical assumptions on the data. Convergence guarantees let appear parameters  
 81 characterizing the class of functions and the geometries of optimization problems (1.1) in  
 82  $\mathbb{R}^d$  and (1.3) in  $\mathbb{R}^n$ , but remain mostly intractable. To this end, we compute *a priori* re-  
 83 fined estimates of convergence rates for boosting methods applied to a particular least-squares  
 84 problem. We develop two approaches for computing on the one hand exact estimates using  
 85 SDP relaxations [23], and on the other hand random approximations using random matrix  
 86 theory. As a result, we observe a phase transition in the convergence rate of gradient descent  
 87 (resp. coordinate descent), depending on  $(n, d)$ . Surprisingly, we conclude that for a fixed  
 88 number of samples  $n$ , adding features (dimension  $d$ ) improves their convergence, which may  
 89 be compared to the double descent phenomenon [8] for the generalization error. Building on  
 90 these results, we experimentally highlight a transition phase for the proximal gradient and  
 91 regularized matching pursuit on a LASSO problem, depending on the value for  $\lambda$ . Finally, we  
 92 define an ‘ultimate method’ converging linearly both in the underparametrized ( $n \gg d$ ) and  
 93 in the overparametrized ( $n \ll d$ ) regime, that is nonetheless not a boosting method (it may  
 94 indeed add more than one atom per iteration).

95 **1.1. Prior works.** **Boosting algorithms.** Boosting strategies, also known as matching  
 96 pursuit in signal processing, have been initiated in the context of sparse recovery [34], and  
 97 extended to the fitting of weak-learners with ‘gradient boosting’ techniques such as Adaboost  
 98 by Freund et al. [21]. Matching pursuit (MP) algorithms produce sparse combinations of  
 99 atoms by picking a direction from a set of atoms using information on the gradient. Boosting  
 100 algorithms are suited to both constrained models, with for example orthogonal matching pur-  
 101 suit [51, 58, 62] or greedy algorithms [56], as well as to unconstrained (penalized) optimization  
 102 problems, with for example the vanilla boosting strategy of Zhang et al. [63], that minimizes a  
 103 well-chosen quadratic upper-bound. Recently, Locatello et al. [32] have unified the framework  
 104 for matching pursuit and Frank-Wolfe algorithms [20] leading to non-statistical convergence  
 105 guarantees for matching pursuit.

106 **Coordinate descent.** Coordinate descent has gained interest due to the increasing access to  
 107 large amounts of data, and thereby to the use of large-scale optimization models. Tseng [59]  
 108 opened the path to convergence guarantees for proximal coordinate descent on composite mini-  
 109 mization problems [60]. Nesterov [39] derived the first coordinate gradient descent method with  
 110 global guarantees for convex objectives, paving the way to families of randomized coordinate  
 111 updates [49], and greedy updates [7]. Yet, these analyses often lead to dimension-dependent  
 112 convergence guarantees. Nutini et al. [43] provided the first convergence guarantee of greedy  
 113 coordinate descent (or coordinate descent with Gauss-Southwell rule) without dependence in  
 114 the dimension, formulating the update as the minimization of a smoothness upper bound with  
 115 respect to the  $\ell_1$ -norm. More precisely, they showed a significantly better performance of

↓  
 Sinc de ça ? Tseng court?  
 check ref. dans paper Yui?

*This manuscript is for review purposes only.*

~~Aussi il y en avait certainement~~

ce n'est pas l'analyse qui me pose problème ?  
done ? Nous allons tirer la conclusion !

greedy coordinate descent compared to randomized coordinate descent. However, the analysis did not extend well to proximal coordinate descent, letting a dependence in the dimension appear in the convergence bound. This led to refined techniques such as the greedy update of Karimireddy [30], with dimension-independent convergence guarantees. Finally, these methods often present the benefit of an induced sparsity, that can be linked to the  $\ell_1$ -norm. Locatello et al. [32] interpreted steepest coordinate descent as a matching pursuit algorithm, where the atoms corresponds to the unitary directions. More precisely, steepest coordinate descent may be seen as the minimization of a smoothness upper bound with respect to the  $\ell_1$ -norm. Considering gauge functions, coordinate descent can be extended to producing solutions sparse with respect to atoms, as Sun and Bach [54] did with the generalized conditional gradient method [3].

Approximate convergence guarantees. Sparse optimization often reveals a gap between theoretical convergence guarantees and observed behaviors. The LASSO has been widely studied for statistical recovery. From an optimization point of view, most of the analyses depend on the statistical recovery efficiency. For constrained optimization problems, Zhang [62] proposed a forward-backward greedy algorithm for which he derived convergence guarantees under restricted isometry properties (RIP). Similarly, Agarwal et al. [1] analyzed the proximal gradient and the projected gradient under restricted strong convexity and smoothness, that comes directly from restricted eigenvalue conditions [48], that appear for example for random Gaussian matrices. A recent focus on average-case analysis of optimization methods under random matrices was initiated by Pedregosa et Scieur [46], coming from the convergence analysis of the simplex method [11, 53]. On the contrary, other works improved global convergence guarantees considering well-chosen geometries. For separable quadratics, Nutini et al. [43, Section 4.1] have computed explicitly the strong convexity parameter in the  $\ell_1$ -geometry. Generalizing unitary atoms from the  $\ell_1$ -geometry to atoms, Sun et Bach [54, Section 2.5] formulated smoothness and strong convexity with respect to gauge functions as optimization problems. However in most cases, since these parameters are hard to compute, both strong convexity and smoothness parameters remains formulated in the  $\ell_2$ -norm. This often leads to additional terms in convergence guarantees, coming from the norm equivalence [43, Appendix 4] or from the geometry such as the pyramidal width [31] or the directional width [32] in Frank-Wolfe techniques, or to the Hoffman constant [26] for linear mappings with strongly convex functions [37, 29, 25].

**1.2. Assumptions. Convex optimization framework.** In this work, functions  $f$  into consideration are convex, differentiable and admit at least one global minimizer  $x^* \in \mathbb{R}^n$ . Functions  $F(\cdot) = f(P \cdot) : \mathbb{R}^d \rightarrow \mathbb{R}$  benefit from the same properties. We restrict ourselves to the analysis of first-order methods, made of linear combinations of past iterates and gradients.

In this paper, functions  $f$  may be smooth with respect to a generic norm  $\|\cdot\|_{\mathbb{R}^n}$ , if they verify for all  $x, y \in \mathbb{R}^n$ ,

$$(1.4) \quad f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L^f}{2} \|y - x\|_{\mathbb{R}^n}^2.$$

Functions  $F(\cdot) = f(P \cdot)$  are therefore smooth with respect for any norm  $\|\cdot\|_{\mathbb{R}^d}$  with  $L^F \leq L^f L^P$ , where  $L$  is defined such that for all  $\alpha, \beta \in \mathbb{R}^d$ ,  $\|P(\alpha - \beta)\|_{\mathbb{R}^n}^2 \leq L^P \|\alpha - \beta\|_{\mathbb{R}^d}^2$ , that is

where  $\text{iterative methods}$  ... (see  $\mathcal{CXX}$ , Section  $\mathcal{YY}$ )

This manuscript is for review purposes only.

et pour le  
problème (1.1)  
non?  
Si tu as besoin  
de ce minimum  
pour  $f$ , alors  
 $x^*(f)$  est  
une notation  
+ application.

157  $L^P = \sup_{\|\beta\|_{\mathbb{R}^d} \leq 1} \|P\beta\|_{\mathbb{R}^n}^2$ . For least-squares, functions  $F$  are exactly smooth with  $L^F = L^f L^P$ .  
 158 In addition, functions  $f$  are strongly convex with respect to a norm  $\|\cdot\|_{\mathbb{R}^n}$ , if for all  $x, y \in \mathbb{R}^n$ ,

159 (1.5) 
$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu^f}{2} \|y - x\|_{\mathbb{R}^n}^2.$$

160 Functions  $F(\cdot) = f(P\cdot)$  do not always inherit strong convexity. For example, for least-squares,  
 161 functions  $F$  are not strongly convex as soon as the number of samples  $n$  is lower than the  
 162 dimension  $d$ . The ‘natural’ strong convexity parameter of functions  $F$  is given by  $\mu^F = \mu^f \mu^P$ ,  
 163 with  $\mu^P = \inf_{\|\beta\|_{\mathbb{R}^d} \geq 1} \|P\beta\|_{\mathbb{R}^n}^2$  and may indeed be zero. As we will see in Section 2.4,  $F$  however  
 164 inherits the Łojasiewicz property with parameter  $\mu^{L,F} > 0$ , such that for all  $\beta \in \mathbb{R}^d$ ,

165 (1.6) 
$$\frac{1}{2} \|\nabla F(\beta)\|_{\mathbb{R}^d, *}^2 \geq \mu^{L,F} (F(\beta) - F_*).$$

$\rightarrow$  Say  $\|\cdot\|_{\mathbb{R}^d, *}$  is dual norm center, when def  $\|\cdot\|_{\mathbb{R}^d}$

166 **Random matrices.** A part of this work is devoted to approximating the strong convexity  
 167 and smoothness parameters of  $f$  and  $F$ . We consider on the one hand relaxed formulations  
 168 for strong convexity and smoothness parameters with respect to the data (i.e., geometry)  $P$ .  
 169 On the other hand, we propose random estimates of these parameters, relying on random  
 170 matrix theory. Random matrices often appears in statistical assumptions, such as with re-  
 171 stricted isometry property [14] or the restricted eigenvalue condition [48]. In the machine  
 172 learning literature, random matrices appear in average-case analysis for quadratics [46] with  
 173 the Marchenko-Pastur distribution [35], or when studying the double descent phenomena for  
 174 the generalization error [8, 36, 4] for Gaussian data. Most of the time, these analyses let two  
 175 regimes appear, depending (among others) on the number of samples  $n$  and the dimension  $d$ .  
 176 Throughout this work, we thus consider two regimes depending on the linear mapping struc-  
 177 ture  $P \in \mathbb{R}^{n \times d}$ : the *underparametrized* (respectively *overparametrized*) regime, characterized  
 178 by matrices  $P \in \mathbb{R}^{n \times d}$  for which  $n \geq d$  (resp.  $d \geq n$ ) and  $P^\top P$  (resp.  $PP^\top$ ) is invertible. Note  
 179 that the invertibility of  $PP^\top$  (resp.  $P^\top P$ ) in the overparametrized (resp. underparametrized)  
 180 regime can be obtained by adding sufficiently random noise. More assumptions on  $P$  and  $P^\top P$   
 181 will be made across this study.

182 **2. A phase transition for linear regression.** We begin with the study of a linear regression  
 183 problem, where problem (2.1) is a special case of the optimization Problem (1.1) with  $\lambda = 0$ ,

184 (2.1) 
$$\min_{\alpha \in \mathbb{R}^d} \frac{1}{2n} \|P\alpha - y\|_2^2 = F(\alpha) = f(P\alpha),$$

185 where  $P \in \mathbb{R}^{n \times d}$ , and  $n, d$  respectively denotes the number of samples and the dimension.

186 In this section, we focus on describing the convergence regimes of gradient descent in  
 187 the  $\ell_2$ -geometry and coordinate descent with the Gauss-Southwell (GS) rule [29, 43] in the  
 188  $\ell_1$ -geometry. More precisely, we interpret gradient descent and coordinate descent as ~~the~~  
 189 minimizers of smoothness upper bound with respect to well-chosen norms, which leads to lin-  
 190 ear convergence both in the underparametrized and overparametrized regime. These linear  
 191 guarantees let smoothness and strong convexity parameters appear, that we formulate as opti-  
 192 mization problems in the geometry into account. To ease the computation of these parameters,

↳ phase bigne ...

↳ p.e. dim plot

↳ for characterizing convergence properties ...

This manuscript is for review purposes only.

After we provide  
 estimates of these quantities  
 under some likelihood on the data...?

we derive numerical estimates. A first technique developed in this work is based on an SDP relaxation, and leads to deterministic estimates. A second technique, inspired from statistical assumptions and average-case analysis, leads to random estimates using random matrices. These random estimates let a phase transition appear between the underparametrized and overparametrized regimes, that we illustrate in particular in a random feature experiment. Finally, we interpret coordinate descent as a matching pursuit algorithm depending on the geometry  $P$ .

First, let us formulate smoothness and strong convexity parameters with respect to a generic norm for the least-squares minimization (2.1) as optimization problems. In this context,  $f$  is  $\frac{1}{n}$ -smooth  $\frac{1}{n}$ -strongly convex with respect to the norm  $\|\cdot\|_2$ . Thus,  $F$  is  $L^F$ -smooth with respect to an arbitrary norm  $\|\cdot\|$  in  $\mathbb{R}^d$ , with  $L^F = \frac{1}{n} \sup_{\|\beta\|^2 \leq 1} \|P\beta\|_2^2$ . In addition, the function is (possibly)  $\mu^F$ -strongly convex, with a parameter  $\mu^F$  explicited in Lemma 2.1 and possibly equal to 0 (especially when dimension  $d < n$ ).

**Lemma 2.1.** *Let  $F = \frac{1}{n} \|P\alpha - y\|_2^2$ , where  $P \in \mathbb{R}^{n \times d}$ . Then,  $F$  is  $\mu^F$ -strongly convex with respect to a norm  $\|\cdot\|$  with,*

$$\mu^F = \frac{1}{n} \inf_{\|\beta\|^2 \geq 1} \|P\beta\|_2^2 \quad \text{and} \quad \frac{1}{\mu^F} = n \sup_{\|P\beta\|^2 \leq 1} \|\beta\|_2^2.$$

*Proof.* Let us recall the definition for strong convexity (1.5), for all  $\alpha, \nu \in \mathbb{R}^d$ ,  $F(\alpha) \geq F(\nu) + \langle \nabla F(\nu), \alpha - \nu \rangle + \frac{\mu^F}{2} \|\alpha - \nu\|^2$ . Since  $F$  is a quadratic, the left-hand side of the inequality can be rephrased into, for all  $\beta \in \mathbb{R}^d$ ,  $\|P\beta\|_2^2 \geq \mu^F \|\beta\|^2$ , from which both formulations follow. ■

In Lemma 2.1, we formulate  $\mu^F$ , the strong convexity parameter for  $F$ , as a nonconvex minimization problem, with a convex objective and concave constraints. Such a problem is usually costly to solve. The function  $F$  also verifies the Łojasiewicz inequality (1.6) with  $\mu^{L,F}$ . Again  $\mu^{L,F}$  is formulated as an optimization problem.

**Lemma 2.2.** *Let  $F = \frac{1}{n} \|P\alpha - y\|_2^2$ , where  $P \in \mathbb{R}^{n \times d}$ . Then,  $F$  verifies the Łojasiewicz inequality (1.6) with respect to a (dual) norm  $\|\cdot\|_*$ , with*

$$\mu^{L,F} = \frac{1}{n} \inf_{\|P\beta\|_2^2 \geq 1} \|P^\top P\beta\|_*^2 \quad \text{and} \quad \frac{1}{\mu^{L,F}} = n \sup_{\|P^\top P\beta\|_*^2 \leq 1} \|P\beta\|_2^2.$$

*Proof.* The proof follows from the Łojasiewicz inequality. Given that  $y = P\alpha_*$ , where  $\alpha_*$  is an optimal point for  $F$ , we have for all  $\alpha \in \mathbb{R}^d$ :  $\|\nabla F(\alpha)\|_*^2 = \frac{1}{n} \|P^\top P(\alpha - \alpha_*)\|_*^2$  and  $F(\alpha) - F_* = \frac{1}{2n} \|P(\alpha - \alpha_*)\|_2^2$ . Then, for all  $\beta \in \mathbb{R}^d$ ,  $\|P^\top P(\alpha - \alpha_*)\|_* \geq \mu^{L,F} \|P\beta\|_2^2$ . ■

Again in Lemma 2.2,  $\mu^{L,F}$  is formulated as a (nonconvex) minimization problem. The two quantities  $\mu^F$  and  $\mu^{L,F}$  are compared in Lemma 2.3, with equality in the underparametrized regime in which  $P^\top P$  is invertible.

**Lemma 2.3.** *Let  $F = \frac{1}{2n} \|P\alpha - y\|_2^2$ . Then, we have that  $\mu^{L,F} \geq \mu^F$  for  $\mu^F$  (resp.  $\mu^{L,F}$ ) defined in Lemma 2.1 (resp. Lemma 2.2). If  $P^\top P$  is invertible,  $\mu^{L,F} = \mu^F$ .*

est par vraiment pour "simplifier les calculs" ?  
 Est peut-être pour avoir de bonnes choses générales non?

227 Proof. Let us consider the squared-root formulations of  $\mu^F$  and  $\mu^{L,F}$  given in Lemma 2.1  
 228 and Lemma 2.2.

229  $\frac{1}{\sqrt{n\mu^F}} = \sup_{\|P\beta\|_2 \leq 1} \|\beta\| = \sup_{\|z\|_* \leq 1, \|P\beta\|_2 \leq 1} \langle \beta, z \rangle,$

230  $\frac{1}{\sqrt{n\mu^{L,F}}} = \sup_{\|P^\top P\nu\|_* \leq 1} \|P\nu\|_2 = \sup_{\|P^\top P\nu\|_* \leq 1, \|P\beta\|_2 \leq 1} \langle P\beta, P\nu \rangle = \sup_{\|P^\top P\nu\|_* \leq 1, \|P\beta\|_2 \leq 1} \langle \beta, P^\top P\nu \rangle.$

231 Since  $\text{Im}(P^\top P) \subset \mathbb{R}^d$ , we have  $\frac{1}{\sqrt{n\mu^F}} \geq \frac{1}{\sqrt{n\mu^{L,F}}}$ , and therefore  $\mu^{F,L} \geq \mu^F$ . In the special case  
 232 where  $P^\top P$  is invertible,  $\text{Im}(P^\top P) = \mathbb{R}^d$ , and  $\mu^{F,PL} = \mu^F$ .  $\rightarrow$  relation change. ■

par défaut.

233 In the next sections, we see the role of these parameters in the convergence guarantees  
 234 of gradient descent and steepest coordinate descent, both in the underparametrized and over-  
 235 parametrized regime. To ease the computation of  $\mu^F$  and  $\mu^{L,F}$ , we propose exact estimates of  
 236 these parameters, as well as concentration results based on a simple random model for  $P$ .

237 **2.1. Gradient descent in the  $\ell_2$ -geometry.** We are interested in the (approximate) con-  
 238 vergence of gradient descent in the underparametrized and the overparametrized regimes.  
 239 Assume  $\mathbb{R}^d$  is equipped with the  $\ell_2$ -norm. The function  $F$  is convex,  $L_2^F$ -smooth with respect  
 240 to the norm  $\ell_2$ , with  $L_2^F = \frac{1}{n}\lambda_{\max}(P^\top P)$ . A common interpretation of gradient descent with  
 241 fixed step size  $\gamma = \frac{1}{L_2^F}$  comes from the minimization of the quadratic smoothness upper bound  
 242 on  $F$ :

243 (2.2)  $\alpha_1 = \alpha_0 - \frac{1}{L_2^F} \nabla F(\alpha_0) = \alpha_0 - \frac{1}{L_2^F} P^\top (P\alpha_0 - y).$

244 In the underparametrized regime, the function  $F$  is  $\mu_2^F$ -strongly convex with respect to the  
 245  $\ell_2$ -norm, with  $\mu_2^F = \lambda_{\min}(\frac{P^\top P}{n}) > 0$ . As a result, gradient descent (2.2) converges linearly.  
 246 However, in the overparametrized regime in which  $d \geq n$ ,  $\mu_2^F = 0$ ,  $F$  is not strongly convex.  
 247 Yet, gradient descent still converges linearly [10], since quadratics benefit from the Łojasiewicz  
 248 inequality, with  $\mu_2^{L,F} = \frac{1}{n}\lambda_{\min}(PP^\top) > 0$ .  $\rightarrow \mu_2^{L,F} ?$   $\rightarrow$  value : def de  $\mu_2$ , superscript = objet  
 249 Proposition 2.4. Let  $F$  be convex,  $L_2^F$ -smooth with respect to the norm  $\|\cdot\|_2$ , be  $\mu_2^F$ -strongly  
 250 convex and verify a Łojasiewicz inequality with parameter  $\mu_2^{L,F}$ , with  $0 \leq \mu_2^F \leq \mu_2^{L,F} \leq L_2^F$ .  
 251 Let  $(\alpha_k)_{k \in \mathbb{N}}$  be generated by gradient descent in (2.2) starting from  $\alpha_0 \in \mathbb{R}^d$ . The sequence  
 252 verifies:

253  $F(\alpha_k) - F_* \leq \left(1 - \frac{\max(\mu_2^F, \mu_2^{L,F})}{L_2^F}\right)^{2k} (F(\alpha_0) - F_*),$

254 where  $\mu_2^F = \lambda_{\min}(PP^\top/n)$ ,  $\mu_2^{L,F} = \lambda_{\max}(P^\top P/n)$  and  $L_2^F = \lambda_{\max}(P^\top P/n)$  ■

255 Proof. See appendix A.

256 The convergence speeds obtained in Proposition 2.4 depends on the maximal and minimal  
 257 eigenvalues of  $P^\top P$  and  $PP^\top$ . In the case where  $P$  is generated randomly, we can derive esti-  
 258 mates of these extremal eigenvalues, avoiding a full computation of the extremal eigenvalues,

$0 \leq \mu_2^F \leq L_2^F$  ?  
 $0 \leq \mu_2^{L,F} \leq L_2^F$  ?

question bête ?  
 Est-ce qu'on peut faire identique avec  $L^F$ ?  $\min(\lambda_{\min}(P^\top P), \lambda_{\max}(P^\top P))$

Directement ici ?

Une espèce de  $L$  "loja" ?

295     2.2. Gauss-Southwell coordinate descent in the  $\ell_1$ -geometry. Similar to gradient de-  
 296     scent, we study convergence guarantees of coordinate descent based on the Gauss-Southwell  
 297     (GS) rule. The GS-rule can be obtained from the smoothness upper bound with respect to  
 298     the  $\ell_1$ -norm, as shown by Nutini et al. [43, Section 4]. For all  $\alpha_0, \alpha \in \mathbb{R}^d$ ,

$$299 \quad (2.3) \quad F(\alpha) \leq F(\alpha_0) + \langle \nabla F(\alpha_0), \alpha - \alpha_0 \rangle + \frac{L_1^F}{2} \|\alpha - \alpha_0\|_1^2.$$

300     From this inequality, we compute  $L_1^F = \frac{1}{n} \max_{\alpha \in \mathbb{R}^d, \|\alpha\|_1=1} \|Pz\|_2^2 = \frac{1}{n} \max_{i=1,\dots,d} \|P_{:,i}\|_2^2$  (the  
 301     maximization problem attains its optimum on the extremal point of the simplex). Gauss-  
 302     Southwell coordinate descent follows by minimizing over  $\alpha \in \mathbb{R}^d$ , for a fixed  $\alpha_0 \in \mathbb{R}^d$ ,

$$303 \quad (2.4) \quad i_0 = \arg \max_{k=1,\dots,d} |\nabla_{i_k} F(\alpha_0)|,$$

$$\alpha_1 = \alpha_0 - \frac{1}{L_1^F} \nabla_{i_0} F(\alpha_0) e_{i_0}.$$

304     As for gradient descent, its convergence speed depends on the parametrization regime.  
 305     Depending on the  $(n, d)$ ,  $F$  may be  $\mu_1^F$ -strongly convex, or verify the Łojasiewicz inequality  
 306     with parameter  $\mu_1^{L,F}$ . Both  $\mu_1^F$  and  $\mu_1^{L,F}$  can be formulated as optimization problems. It  
 307     follows from the strong convexity characterization given in Lemma 2.1 with the norm  $\|\cdot\|_1$ ,  
 308     that  $\mu_1^F = \frac{1}{n} \inf_{\|z\|_1^2 \geq 1} \|Pz\|_2^2$ , and from Lemma 2.2 with norm  $\|\cdot\|$  for the Łojasiewicz inequality  
 309      $\mu_1^{L,F} = \inf_{\|P\beta\|_2^2 \geq 1} \|P^\top P\beta\|_\infty^2$ .

310     In the regimes under consideration, Proposition 2.7 states that coordinate descent con-  
 311     verges linearly, as already proven by Karimi et al. [29, Theorem 1].

312     Proposition 2.7. [29, Theorem 1] Let  $F$  be convex,  $L_1^F$ -smooth with respect to the norm  
 313      $\|\cdot\|_1$ , be  $\mu_1^F$ -strongly convex and be  $\mu_1^{L,F}$ -PL with  $0 \leq \mu_1^F \leq \mu_1^{L,F} \leq L_1^F$ . Let  $(\alpha_k)$  be generated  
 314     by coordinate gradient descent (2.4) starting from  $\alpha_0 \in \mathbb{R}^d$ . The sequence verifies:

$$315 \quad F(\alpha_k) - F_* \leq \left(1 - \frac{\max(\mu_1^F, \mu_1^{L,F})}{L_1^F}\right)^k (F(\alpha_0) - F_*).$$

316     Proof. See Appendix A.

317     The convergence guarantee provided in Proposition 2.7 is however not easily computable. Al-  
 318     though  $L_1^F = \frac{1}{n} \max_{i=1,\dots,d} \|P_{:,i}\|_2^2$  has a closed-form solution,  $\mu_1^F = \frac{1}{n} \inf_{\|z\|_1^2 \geq 1} \|Pz\|_2^2$  and  
 319      $\mu_1^{L,F} = \inf_{\|P\beta\|_2^2 \geq 1} \|P^\top P\beta\|_\infty^2$  are formulated as nonconvex minimization problems. We con-  
 320     struct estimates to these quantities, so that they may be computed a priori.

321     SDP relaxations. Building on the formulation of  $\mu_1^F$  and  $\mu_1^{L,F}$  as optimization problems,  
 322     we rephrase them into relaxed SDPs.

323     Proposition 2.8. Gauss-Southwell coordinate descent's linear convergence rate has the fol-  
 324     lowing approximations

- in the underparametrized regime,  $\frac{1}{n\mu_1^F} = \sup_{X \succcurlyeq 0} \text{Tr}((P^\top P)^{-1} X)$ , s.t.  $\text{diag}(X) \leq 1$ ,

$$326 \quad 1 - \frac{\pi}{2} \frac{\tilde{\mu}_1^F}{L_1^F} \leq 1 - \frac{\mu_1^F}{L_1^F} \leq 1 - \frac{\hat{\mu}_1^F}{L_1^F},$$

This manuscript is for review purposes only. *Very close*

→ the explicit having good enough  
 of  $\mu_i$ 's is expected. They stand  
 (in Ch 2.7) relies on such estimate!

évaluer le point  
 à la fois.  
 je pense qu'il faut  
 faire un paragraphe  
 équivalent à la théorie sur  
 comment déterminer les  
 évaluations

- 327 • in the overparametrized regime,  $\frac{1}{n\tilde{\mu}_1^{L,F}} = \sup_{X \succcurlyeq 0} \text{Tr}(P^\top P X)$  s.t.  $\|P^\top P X P^\top P\|_\infty \leq 1$ ,

$$328 \quad 1 - \frac{\mu_1^{L,F}}{L_1^F} \leq 1 - \frac{\tilde{\mu}_1^{L,F}}{L_1^F},$$

329 where  $L_1^F = \frac{1}{n} \max_{i=1,\dots,d} \|P_{:,i}\|_2^2$ . In addition, we still have that  $\tilde{\mu}_1^F \leq \tilde{\mu}_1^{PL}$ .

330 Proof. See Appendix B.1.

331 In Proposition 2.8, we find out SDP relaxations that yield an exact approximation for  $\mu_1^F$ , and  
332 an exact lower bound for  $\mu_1^{L,F}$ . Yet, the larger  $n, d$ , the longer the computation of these SDPs.

333 **Random estimates.** We now assume that  $P$  is randomly generated, as in the  $\ell_2$ -geometry.  
334 Under subgaussian assumptions, we derive in Proposition 2.9 random estimates of  $\mu_1^F, \mu_1^{L,F}$  and  
335  $L_1^F$ . More precisely, we prove that  $L_1^F$  concentrates around the variance  $\sigma^2$ ,  $\mu_1^F$  concentrates  
336 around  $\frac{\sigma^2}{d}$  and  $\mu_1^{L,F}$  concentrates around  $\frac{\sigma^2}{n}$  with subgaussian tails.

337 Proposition 2.9. Let  $P \in \mathbb{R}^{n \times d}$ , with  $P_i \in \mathbb{R}^d$  i.i.d. subgaussian such that  $\mathbb{E}[P_{i,j}] = 0$ ,  
338  $\mathbb{E}[P_{i,j}] = \sigma^2$ . There exists  $C, C_1, C_2, C_3, C_4 > 0$  absolute constants such that,

- 339 • For all  $u \geq 2K^2 \sqrt{\frac{C_1 \log(d)}{n}}$ ,

$$340 \quad \left(1 + C_2 K^2 \frac{1}{\sqrt{n}} - t\right)^2 \leq \frac{L_1^F}{\sigma^2} \leq \left(1 + 2K^2 \sqrt{\frac{C_1 \log(d)}{n}} + t\right)^2,$$

341 holds with probability  $1 - e^{-\frac{C}{\sigma^2 K^4} \min(u_1(t), u_2(t))}$  where  $u_1(t) = \log(d) \sigma^2 (t + \frac{C_2 K^2}{\sqrt{n}})^2$  and

$$342 \quad u_2(t) = d \sigma^2 (t - 2K^2 \sqrt{\frac{C_1 \log(d)}{n}})^2.$$

- 343 • For all  $t \geq 0$ , it holds with probability  $1 - 2 \exp(-t^2)$ ,

$$344 \quad \left(1 - C_3 K^2 \left(\sqrt{\frac{d}{n}} + \frac{t}{\sqrt{n}}\right)\right)^2 \leq \mu_1^F \frac{d}{\sigma^2} \leq \left(1 + C_3 K^2 \left(\sqrt{\frac{1}{n}} + \frac{t}{\sqrt{dn}}\right)\right)^2.$$

- 345 • For all  $t \geq 2\sigma K^2 \sqrt{\frac{C_1 \log(d)}{n}}$ , it holds with probability  $1 - 2 \exp(-\min(t^2, u_2(t)))$ ,

$$346 \quad \left(1 - C_4 K^2 \left(\sqrt{\frac{n}{d}} + \frac{t}{\sqrt{d}}\right)\right)^2 \leq \mu_1^{L,F} \frac{n}{\sigma^2} \leq \left(1 + 2K^2 \sqrt{\frac{C_1 \log(d)}{n}} + t\right)^2.$$

347 The constant  $K > 0$  characterizes subgaussian vectors of  $P$  (and defined in Appendix B.2).

348 Proof. See Appendix B.2. ■

349 Compared with Proposition 2.4, Proposition 2.9 provides concentration inequalities for  
350  $L_1^F, \mu_1^F$  and  $\mu_1^{L,F}$  depending on dimension  $d$ , the variance  $\sigma^2$ , the number of samples  $n$  and  
351 absolute constants. In the overparametrized regime (resp. underparametrized), we conclude  
352 with limiting concentration of the convergence rate for large dimensions (resp. large number  
353 of samples).

295     **2.2. Gauss-Southwell coordinate descent in the  $\ell_1$ -geometry.** Similar to gradient de-  
 296     scent, we study convergence guarantees of coordinate descent based on the Gauss-Southwell  
 297     (GS) rule. The GS-rule can be obtained from the smoothness upper bound with respect to  
 298     the  $\ell_1$ -norm, as shown by Nutini et al. [43, Section 4]. For all  $\alpha_0, \alpha \in \mathbb{R}^d$ ,

299     (2.3)     
$$F(\alpha) \leq F(\alpha_0) + \langle \nabla F(\alpha_0), \alpha - \alpha_0 \rangle + \frac{L_1^F}{2} \|\alpha - \alpha_0\|_1^2.$$

300     From this inequality, we compute  $L_1^F = \frac{1}{n} \max_{\alpha \in \mathbb{R}^d, \|\alpha\|_1=1} \|Pz\|_2^2 = \frac{1}{n} \max_{i=1,\dots,d} \|P_{:,i}\|_2^2$  (the  
 301     maximization problem attains its optimum on the extremal point of the simplex). Gauss-  
 302     Southwell coordinate descent follows by minimizing over  $\alpha \in \mathbb{R}^d$ , for a fixed  $\alpha_0 \in \mathbb{R}^d$ ,

303     (2.4)     
$$\begin{aligned} i_0 &= \arg \max_{k=1,\dots,d} |\nabla_{i_k} F(\alpha_0)|, \\ \alpha_1 &= \alpha_0 - \frac{1}{L_1^F} \nabla_{i_0} F(\alpha_0) e_{i_0}. \end{aligned}$$

304     As for gradient descent, its convergence speed depends on the parametrization regime.  
 305     Depending on the  $(n, d)$ ,  $F$  may be  $\mu_1^F$ -strongly convex, or verify the Łojasiewicz inequality  
 306     with parameter  $\mu_1^{L,F}$ . Both  $\mu_1^F$  and  $\mu_1^{L,F}$  can be formulated as optimization problems. It  
 307     follows from the strong convexity characterization given in Lemma 2.1 with the norm  $\|\cdot\|_1$ ,  
 308     that  $\mu_1^F = \inf_{\|z\|_1^2 \geq 1} \|Pz\|_2^2$ , and from Lemma 2.2 with norm  $\|\cdot\|$  for the Łojasiewicz inequality  
 309      $\mu_1^{L,F} = \inf_{\|P\beta\|_2^2 \geq 1} \|P^\top P\beta\|_\infty^2$ .

310     In the regimes under consideration, Proposition 2.7 states that coordinate descent con-  
 311     verges linearly, as already proven by Karimi et al. [29, Theorem 1].

312     Proposition 2.7. [29, Theorem 1] Let  $F$  be convex,  $L_1^F$ -smooth with respect to the norm  
 313      $\|\cdot\|_1$ , be  $\mu_1^F$ -strongly convex and be  $\mu_1^{L,F}$ -PL with  $0 \leq \mu_1^F \leq \mu_1^{L,F} \leq L_1^F$ . Let  $(\alpha_k)$  be generated  
 314     by coordinate gradient descent (2.4) starting from  $\alpha_0 \in \mathbb{R}^d$ . The sequence verifies:

315     
$$F(\alpha_k) - F_* \leq \left(1 - \frac{\max(\mu_1^F, \mu_1^{L,F})}{L_1^F}\right)^{2k} (F(\alpha_0) - F_*).$$

316     Proof. See Appendix A.

317     The convergence guarantee provided in Proposition 2.7 is however not easily computable. Al-  
 318     though  $L_1^F = \frac{1}{n} \max_{i=1,\dots,d} \|P_{:,i}\|_2^2$  has a closed-form solution,  $\mu_1^F = \frac{1}{n} \inf_{\|z\|_1^2 \geq 1} \|Pz\|_2^2$  and  
 319      $\mu_1^{L,F} = \inf_{\|P\beta\|_2^2 \geq 1} \|P^\top P\beta\|_\infty^2$  are formulated as nonconvex minimization problems. We con-  
 320     struct estimates to these quantities, so that they may be computed a priori.

321     SDP relaxations. Building on the formulation of  $\mu_1^F$  and  $\mu_1^{L,F}$  as optimization problems,  
 322     we rephrase them into relaxed SDPs.

323     Proposition 2.8. Gauss-Southwell coordinate descent's linear convergence rate has the fol-  
 324     lowing approximations

- in the underparametrized regime,  $\frac{1}{n\mu_1^F} = \sup_{X \succcurlyeq 0} \text{Tr}((P^\top P)^{-1} X)$ , s.t.  $\text{diag}(X) \leq 1$ ,

$$1 - \frac{\pi \tilde{\mu}_1^F}{2 L_1^F} \leq 1 - \frac{\mu_1^F}{L_1^F} \leq 1 - \frac{\tilde{\mu}_1^F}{L_1^F},$$

*pas une phrase complète!*

*For computing it, explicit estimate of  $\mu_1^F$  is needed.*

*(X)*

*Si c'est可行的话  
il n'y aurait pas de mon*

*pas ce qui est montré à la page.*

*lettre plus rien pas très formelle (pas une prop.)*

*... the following lemma hold...*

*il y a*

*plusieurs idées*

*différentes*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

*qui sont dans le bon sens*

*mais il y a*

*quelques difficultés*

- 327 • in the overparametrized regime,  $\frac{1}{n\tilde{\mu}_1^{L,F}} = \sup_{X \succcurlyeq 0} \text{Tr}(P^\top P X)$  s.t.  $\|P^\top P X P^\top P\|_\infty \leq 1$ ,

328 
$$1 - \frac{\mu_1^{L,F}}{L_1^F} \leq 1 - \frac{\tilde{\mu}_1^{L,F}}{L_1^F},$$

329 where  $L_1^F = \frac{1}{n} \max_{i=1,\dots,d} \|P_{:,i}\|_2^2$ . In addition, we still have that  $\tilde{\mu}_1^F \leq \tilde{\mu}_1^{PL}$ .

330 Proof. See Appendix B.1.

331 In Proposition 2.8, we find out SDP relaxations that yield an exact approximation for  $\mu_1^F$ , and  
332 an exact lower bound for  $\mu_1^{L,F}$ . Yet, the larger  $n, d$ , the longer the computation of these SDPs.

333 **Random estimates.** We now assume that  $P$  is randomly generated, as in the  $\ell_2$ -geometry.  
334 Under subgaussian assumptions, we derive in Proposition 2.9 random estimates of  $\mu_1^F, \mu_1^{L,F}$  and  
335  $L_1^F$ . More precisely, we prove that  $L_1^F$  concentrates around the variance  $\sigma^2$ ,  $\mu_1^F$  concentrates  
336 around  $\frac{\sigma^2}{d}$  and  $\mu_1^{L,F}$  concentrates around  $\frac{\sigma^2}{n}$  with subgaussian tails.

337 Proposition 2.9. Let  $P \in \mathbb{R}^{n \times d}$ , with  $P_i \in \mathbb{R}^d$  i.i.d. subgaussian such that  $\mathbb{E}[P_{i,j}] = 0$ ,  
338  $\mathbb{E}[P_{i,j}] = \sigma^2$ . There exists  $C, C_1, C_2, C_3, C_4 > 0$  absolute constants such that,

- 339 • For all  $u \geq 2K^2 \sqrt{\frac{C_1 \log(d)}{n}}$ ,

340 
$$\left(1 + C_2 K^2 \frac{1}{\sqrt{n}} - t\right)^2 \leq \frac{L_1^F}{\sigma^2} \leq \left(1 + 2K^2 \sqrt{\frac{C_1 \log(d)}{n}} + t\right)^2,$$

341 holds with probability  $1 - e^{-\frac{C}{\sigma^2 K^4} \min(u_1(t), u_2(t))}$  where  $u_1(t) = \log(d)\sigma^2(t + \frac{C_2 K^2}{\sqrt{n}})^2$  and

342  $u_2(t) = d\sigma^2(t - 2K^2 \sqrt{\frac{C_1 \log(d)}{n}})^2$ .

- 343 • For all  $t \geq 0$ , it holds with probability  $1 - 2 \exp(-t^2)$ ,

344 
$$\left(1 - C_3 K^2 \left(\sqrt{\frac{d}{n}} + \frac{t}{\sqrt{n}}\right)\right)^2 \leq \mu_1^F \frac{d}{\sigma^2} \leq \left(1 + C_3 K^2 \left(\sqrt{\frac{1}{n}} + \frac{t}{\sqrt{dn}}\right)\right)^2.$$

- 345 • For all  $t \geq 2\sigma K^2 \sqrt{\frac{C_1 \log(d)}{n}}$ , it holds with probability  $1 - 2 \exp(-\min(t^2, u_2(t)))$ ,

346 
$$\left(1 - C_4 K^2 \left(\sqrt{\frac{n}{d}} + \frac{t}{\sqrt{d}}\right)\right)^2 \leq \mu_1^{L,F} \frac{n}{\sigma^2} \leq \left(1 + 2K^2 \sqrt{\frac{C_1 \log(d)}{n}} + t\right)^2.$$

347 The constant  $K > 0$  characterizes subgaussian vectors of  $P$  (and defined in Appendix B.2).

348 Proof. See Appendix B.2.

349 Compared with Proposition 2.4, Proposition 2.9 provides concentration inequalities for  
350  $L_1^F, \mu_1^F$  and  $\mu_1^{L,F}$  depending on dimension  $d$ , the variance  $\sigma^2$ , the number of samples  $n$  and  
351 absolute constants. In the overparametrized regime (resp. underparametrized), we conclude  
352 with limiting concentration of the convergence rate for large dimensions (resp. large number  
353 of samples).

354 Corollary 2.10. Under the assumption of Proposition 2.9, if in addition  $n, d \rightarrow \infty$ , coordinate  
 355 descent with GS rule (2.4) converges with a limiting rate,

- 356 • in the underparametrized regime, when  $n \rightarrow \infty$ , the convergence guarantee  $1 - \frac{\mu_1^F}{L_1^F}$   
 357 concentrates in  $1 - \frac{1}{d} + O(\frac{1}{\sqrt{n}})$  with subgaussian tails,
- 358 • in the overparametrized regime, when  $d \rightarrow \infty$  and  $\frac{\log(d)}{n} \rightarrow 0$ , the convergence guarantee  
 359  $1 - \frac{1}{\mu_1^{L,F}}$  concentrates in  $1 - \frac{1}{n} + O(\frac{1}{\sqrt{d}}) + O(\sqrt{\frac{\log(d)}{n}})$  with subgaussian tails.

360 Proof. See the proof in Appendix B.4. ■

361 For large overparametrized models (resp. underparametrized), the convergence guarantee of  
 362 coordinate descent with GS rule concentrates to  $(1 - \frac{1}{n})$  (resp.  $(1 - \frac{1}{d})$ ), that is independent  
 363 of the dimension  $d$  (resp. of the number of samples). Note that the condition  $\log(d) \ll n$   
 364 is indeed reasonable, since  $e^n$  grows quickly (when  $n = 50$ ,  $e^n \approx 5 \times 10^{21}$ ). A numerical  
 365 comparison for the expected and exact lower bounds for  $\mu_1^F$  and  $\mu_1^{L,F}$  is provided in SM2.  
 366 Unlike the approximate convergence guarantees for gradient descent in the underparametrized  
 367 regime (resp. overparametrized) detailed in Corollary 2.6, coordinate descent with GS-rule  
 368 does not improve when adding samples (resp. features).

369 As for gradient descent in the  $\ell_2$ -geometry, we have formulated coordinate descent with  
 370 GS rule as the minimization of the smoothness upper bound with respect to the  $\ell_1$ -norm,  
 371 leading its linear convergence in both the underparametrized and overparametrized regime.  
 372 For a linear regression problem, nor the strong convexity parameter neither the Łojasiewicz in  
 373 the  $\ell_1$ -geometry benefit from a closed-form formulation (but it did in the  $\ell_2$ -geometry). In a  
 374 first approach, we approximate these quantities by SDPs, that may take longer computation in  
 375 large models (either in the number of samples or the dimension). Instead of that, we consider  
 376 randomly generated matrices  $P$  to approximate these parameters. Under subgaussian data, it  
 377 appears  $\mu_1^F, \mu_1^{L,F}$  and  $L_1^F$  concentrate to there expectation with subgaussian tails. In the next  
 378 section, we perform numerical experiments that let a transition phenomenon appear.

379 **2.3. A phase transition phenomenon: experimental results.** We compare the approx-  
 380 imated convergence guarantees to numerical experimental convergence for gradient descent  
 381 from Corollary 2.6 and of coordinate descent from Propositions B.1–B.9. More precisely, we  
 382 verify the expected phase transition in  $(n, d)$ : in the overparametrized (respectively under-  
 383 parametrized) regime, the larger the dimension (resp. the number of samples), the better the  
 384 convergence. To this end, we perform experiments on several datasets: several least-squares  
 385 problems obtained with synthetic quadratics, quadratics from the Leukemia dataset and a  
 386 random features experiment, that we described below.

387 **Synthetic quadratics.** We consider several least-squares problems (2.1), where the num-  
 388 ber of samples  $n = 50$  is fixed, and the number of dimensions varies so that both the over-  
 389 parametrized and the underparametrized regimes are explored. In this model, the feature  
 390 matrix  $P$  into account is generated such that  $P_{:,i} \sim \mathcal{N}(0, I_d)$  are i.i.d.,  $\alpha_* \in \{-1, 1\}^d$  has a  
 391 sparsity equal to  $s = 8 < d$ , and  $y = P\alpha_* + \epsilon$ , where  $\epsilon_i \sim \mathcal{N}(0, \sigma)$ .

392 **Leukemia dataset.** We consider the standard Leukemia dataset [24], where  $n = 72$  and  
 393  $d = 7129$ . Again, we consider submatrices, so that the dimensions vary from both the un-  
 394 derparametrized regime to the overparametrized. For each model,  $P$  has zero mean and the

under consideration

one

\* missing parenthesis

showing the transition phase between the two regimes

a few experiments with

395 features have a variance equal to one,

396 **Random features.** Let us consider the example of random features for a fixed prediction  
 397 model. We consider the regression model  $\hat{a} = \arg \min_{a \in \mathbb{R}^m} \frac{1}{2n} \|y - f(P, a, \theta)\|_2^2$ , where the family  
 398 of models is given by  $\mathcal{F}(\theta) = \{f(P, a, \theta) = \sum_{i=1}^m a_i \sigma(\langle \theta_{:, i}, P_{:, j} \rangle) = \phi_P(\theta)^\top a, a \in \mathbb{R}^N\}$ , where  
 399  $\theta \in \mathbb{R}^{m \times d} \sim \mathcal{N}(0, \nu^2)$ , and  $\sigma(\cdot) = \max(0, \cdot)$ . In this experiment, we increase the number of  
 400 features  $m$  (from 10 to 1000) while the initial data taken from the leukemia dataset is such that  
 401  $n = 72$ ,  $d = 200$ , and  $\theta_i \sim \mathcal{N}(0, I_d)$ . Compared to the experiments on synthetic quadratics  
 402 and on the leukemia dataset, the model does not vary in random features: all models converge  
 403 to the same optimal solution  $y$ .

404 In Figure 1, we plot the iteration number  $k(\epsilon)$  at which a certain accuracy  $\epsilon$  is reached  
 405 for the three models described above, both for gradient descent and coordinate descent with  
 406 GS-rule. We consider precisions  $\epsilon = \{10^0, \dots, 10^{-10}\}$  and we refer to these curves by  $\epsilon$ -curves.  
 In Figure 1, both steepest coordinate descent and gradient descent converge faster for the

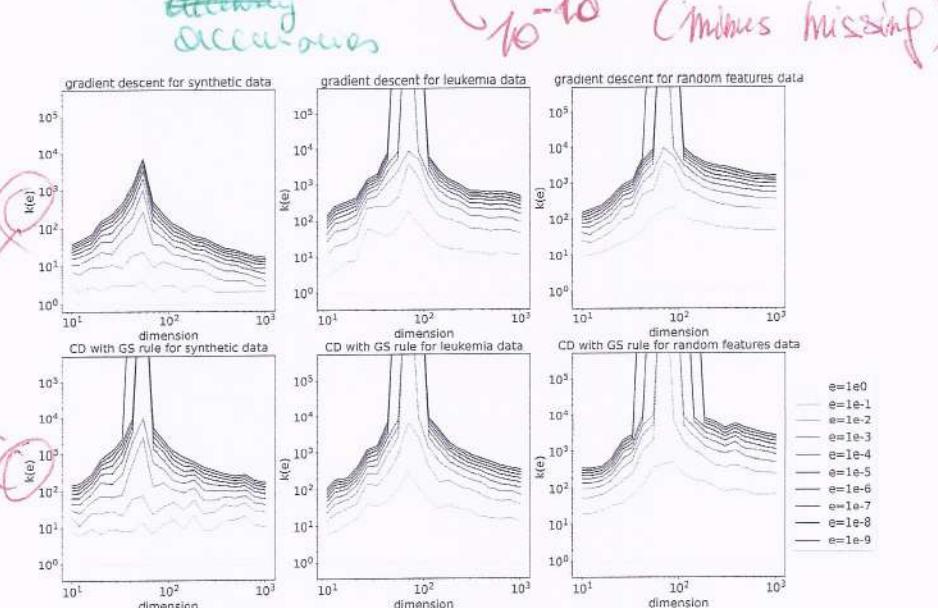


Figure 1:  $\epsilon$ -curve for gradient descent (top) and coordinate descent with the GS-rule (bottom), for the three models: synthetic quadratics (on the left) with  $n = 50$ , the leukemia dataset (in the middle) with  $n = 72$ , a random feature model (on the right) with  $n = 72$ .

407  
 408 three models when  $n \gg d$  (resp.  $d \gg n$ ) in the underparametrized (resp. overparametrized)  
 409 regime. For  $n \approx d$ , convergence slows down and tends to be sublinear, as expected from the  
 410 theory for smooth convex functions. In other words, we observe a transition phenomenon for  
 411 dimensions  $d \approx n$ . For the random feature models, a double descent phenomena was empirically  
 412 highlighted by Belkin et al. [8], and formalized by Mei and Montanari [36]. For a fixed  
 413 prediction model, as the number of features increase, the excess risk follows is U-shaped for un-

*ausi: longer font*

*This manuscript is for review purposes only.*

*(ideally: labels et ticks don't z in table que  
 stan teste)*

derparametrized optimization models and goes down for overparametrized models. As for the excess risk, we observe a transition at  $m \approx n$  as well as a better precision for overparametrized models. Contrary to the generalization error, underparametrized models ( $d \ll n$ ) performs well even when  $\frac{d}{n} \rightarrow 0$  and are not U-shaped. We thus rather speak of phase transition for gradient and coordinate descent.

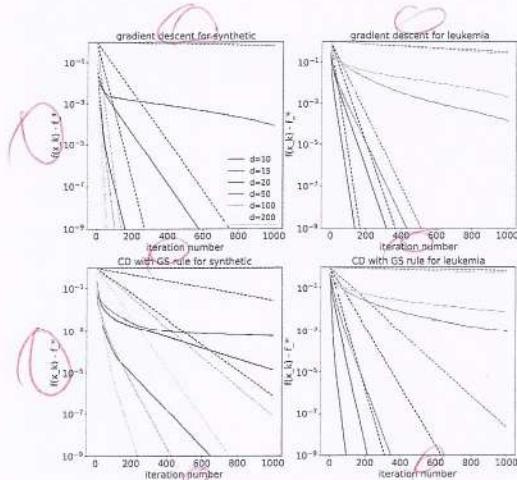


Figure 2: Convergence in function value for gradient descent and coordinate descent with GS rule, on synthetic quadratics ( $n = 20$ ) and on the leukemia dataset ( $n = 72$ ), for several values of dimensions. Dashed line: comparison to the exact upper bound for synthetic quadratics, and to the approximated upper bound for the leukemia dataset.

In Figure 2, we compare the exact and approximated upper bound to the convergence guarantee in function value for gradient descent and coordinate descent with the GS rule. For gradient descent, the theoretical approximation guarantee from Corollary 2.6 matches the real convergence behavior of gradient descent. For steepest coordinate descent, we compare its convergence in function values to the exact upper bound obtained from the SDP relaxation in Proposition 2.8 for ‘small’ values of  $d$  and  $n$ , and to its random estimate otherwise. In both cases, they do not match exactly. However, we numerically recover that convergence is improved as the dimension increases.

**2.4. Coordinate descent is an instance of matching pursuit.** Coordinate descent, as well as gradient descent, converges linearly both in the underparametrized and overparametrized regime, as proven in Proposition 2.7, due respectively to strong convexity and the Łojasiewicz property. Given a certain structure on  $f$ , we prove that  $F(\cdot) = f(P\cdot)$  inherits some convex properties from  $f$  and that coordinate descent can be interpreted as a matching pursuit algorithm in either  $\mathbb{R}^d$  or  $\mathbb{R}^n$ . We now go back to the more general regime,

$$\min_{\alpha \in \mathbb{R}^d} F(\alpha) = f(P\alpha),$$

provided by

Bef. some  
regularity properties?

Consider Now, let us consider  
the more general formulation?

434 where  $f$  is  $L^f$ -smooth,  $\mu^f$ -strongly convex and  $F$  is  $L_1^F$ -smooth with respect to the  $\ell_1$ -norm.

435 In the underparametrized regime,  $F$  is  $\mu_1^F$ -strongly convex (since  $P^\top P$  is invertible).  
 436 The connection between coordinate descent with GS-rule (2.4) and matching pursuit was  
 437 highlighted by Locatello et al [33]. Considering the set of unitary direction  $\mathcal{A} = \text{conv}(\{\pm e_i, i = 1, \dots, d\})$ , that is the unit ball  $L_1$ , coordinate descent may be rewritten as a matching pursuit:

439  $e_{i_0} \in -\text{LMO}_{\mathcal{A}}(\nabla F(\alpha_0)), \rightarrow e_{i_0} = \inf_{z \in \mathcal{A}} \dots$   
 440  $\alpha_1 = \alpha_0 - \frac{1}{L_1^F} \nabla_{i_0} F(\alpha_0) e_{i_0},$

441 where  $\text{LMO}_{\mathcal{A}}(\nabla F(\alpha_0)) = \inf_{z \in \mathbb{R}^d} \nabla F(\alpha_0)^\top z$ . Steepest coordinate descent converges linearly  
 442 from Proposition 2.7, with the same convergence guarantee as in the context of matching  
 443 pursuit [33, Theorem 5].

444 In the overparametrized regime,  $F$  does not inherit strong convexity. Yet, for least-  
 445 squares,  $F$  but does inherit some structure from  $f$  (see Lemma 2.2). We prove that coordinate  
 446 descent can be interpreted as a matching pursuit algorithm in  $\mathbb{R}^n$ .

447 Recall the gauge function, for  $x \in \mathbb{R}^n$ ,  $\gamma_P(x) = \inf_{\alpha \in \mathbb{R}^d, x = P\alpha} \|\alpha\|_1$ . Lemma 2.11 ensures  
 448  $\gamma_P(\cdot)$  is a norm in the overparametrized regime.

449 Lemma 2.11. Let  $\alpha \in \mathbb{R}^d \rightarrow P\alpha \in \mathbb{R}^n$  be a surjection in  $\mathbb{R}^d$ , and  $\mathcal{P} = \text{conv}(P)$  be centrally  
 450 symmetric. The function  $\gamma_P(\cdot)$  is a norm, and its dual norm is  $\gamma_P^*(\cdot) = \sup_{s \in \mathcal{P}} \langle s, \cdot \rangle = \|P^\top \cdot\|_\infty$ .

451 Proof. See appendix C.

452 Let  $f$  be convex,  $L_P^f$ -smooth and  $\mu_P^f$ -strongly convex with respect to the norm  $\gamma_P(\cdot)$ . We  
 453 compute in Lemma 2.12 its smoothness and strong convexity parameters with respect to a  
 454 gauge  $\gamma_P(\cdot)$  (that is a thus a norm).

455 Lemma 2.12. Let  $L_P^f$  (resp.  $L_2^f$ ) be the smoothness parameter of  $f$  with respect to  $\gamma_P(\cdot)$   
 456 (resp.  $\ell_2$ ), and  $\mu_P^f$  (resp.  $\mu_2^f$ ) be the strong convexity of  $f$  with respect to  $\gamma_P$  (resp.  $\ell_2$ ). Then,  
 457  $L_P^f = L_2^f \sup_{j=1,\dots,d} \|P_j\|_2^2$  and  $\mu_P^f = \mu_2^f \inf_{z \in \mathbb{R}^n} \|P^\top z\|_\infty^2$ , such that  $\|z\|_2^2 = 1$ .

458 Proof.  $L_P^f$  is the tightest constant such that for all  $x \in \mathbb{R}^n$ ,  $L_P^f \|x\|_2^2 \leq L_P^f \gamma_P(x)^2$  and  
 459 we have that  $L_P^f = L_2^f \sup_{u,\beta,u=P\beta,\|\beta\|_1 \leq 1} \|u\|_2^2 = L_2^f \sup_{\beta,\|\beta\|_1 \leq 1} \|P\beta\|_2^2 = L_2^f \sup_{j=1,\dots,d} \|P_j\|_2^2$ .  
 460 The same reasoning leads to the result for  $\mu_P^f$ .

461 In Lemma 2.12, the smoothness (resp. strong convexity) parameter is formulated as optimization  
 462 problems, multiplied by the smoothness (resp. strong convexity) in the  $\ell_2$ -norm, which are  
 463 closely related to the parameters for least-squares from Lemma 2.1 and 2.2. In the context of  
 464 least-squares where  $f(x) = \frac{1}{2n} \|x - y\|_2^2$  for  $x \in \mathbb{R}^n$ , we indeed have that  $L_2^f = \mu_2^f = \frac{1}{n}$ . For the  
 465  $\ell_1$ -norm, that  $L_P^f = L_1^F$  as defined in (2.3), and  $\mu_P^f = \mu_1^F$  as soon as  $P P^\top$  is invertible (which  
 466 is true here). Multiplying (2.4) by  $P$ , noticing that  $\min_{e, \|e\|_1=1} \langle \nabla F(\alpha), e \rangle = \min_{p \in \mathcal{P}} \langle \nabla f(x), p \rangle$   
 467 and using Lemma 2.12, coordinate descent with the GS-rule on  $F$  can be formulated as match-

beau de dire ça ? Au moins vraiment pas tout le  
 temps? Juste une forme sup claire suffisante pour  
 This manuscript is for review purposes only.

en arguments? Bon sûr, mais si tu as  
 une meilleure offre. Le problème si tu incorpores "the tightest"  
 dans la définition c'est que les règles de comparaison ne tiennent pas

468 ing pursuit on  $f$ ,

$$469 \quad (2.5) \quad z_0 \in \text{LMO}_{\mathcal{P}}(\nabla f(x_0)), \\ x_1 = x_0 - \frac{1}{L_{\mathcal{P}}^f} \langle \nabla f(x_0), z_0 \rangle z_0.$$

470 Let  $x_k$  be generated by matching pursuit (2.5), starting from  $x_0 \in \mathbb{R}^n$  for  $L_{\mathcal{P}}^f$ -smooth and  
 471  $\mu_{\mathcal{P}}^f$ -strongly convex functions, then, Locatello et al [33, Theorem 5] proved linear convergence  
 472 of the sequence with

$$473 \quad (2.6) \quad f(x_k) - f_* \leq \left(1 - \frac{\mu_{\mathcal{P}}^f}{L_{\mathcal{P}}^f}\right) (f(x_0) - f_*).$$

474 By construction, since  $x_k = P\alpha_k$ , we have that  $F(\alpha_k) - F_* \leq (1 - \frac{\mu_{\mathcal{P}}^f}{L_{\mathcal{P}}^f})(F(\alpha_0) - F_*) =$   
 475  $(1 - \frac{\mu_{\mathcal{P}}^f}{L_{\mathcal{P}}^f})(F(\alpha_0) - F_*)$ . The same result could have been derived from Proposition 2.7  
 476 and the observation that strongly convex functions composed with a linear mapping verify  
 477 a Łojasiewicz-inequality.

478 Lemma 2.13. *Let  $f$  be  $\mu_{\mathcal{P}}^f$ -strongly convex with respect to the norm  $\gamma_{\mathcal{P}}(\cdot)$ . Then,  $F$  verifies  
 479 a Łojasiewicz inequality with parameters  $\mu_{\mathcal{P}}^f$ , that is for all  $\alpha \in \mathbb{R}^d$ ,*

$$480 \quad \text{pos de } \| \cdot \| \quad \frac{1}{2} \|\nabla F(\alpha)\|_{\infty} \geq \mu_{\mathcal{P}}^f (F(\alpha) - F_*).$$

481 Proof. Let  $x \in \mathbb{R}^n$ ,  $f_* \geq f(x) - \sup_z \langle -\nabla f(x), y - x \rangle - \frac{\mu_{\mathcal{P}}^f}{2} \gamma_{\mathcal{P}}^2(y - x) \geq f(x) -$   
 482  $(\frac{\mu_{\mathcal{P}}^f}{2} \gamma_{\mathcal{P}}^2(\cdot))^*(-\nabla f(x)) \geq f(x) - \frac{1}{2\mu_{\mathcal{P}}^f} \|P^\top \nabla f(x)\|_{\infty}^2$ . Since  $F(\cdot) = f(P\cdot)$ , the inequality is ob-  
 483 tained by taking  $x = P\alpha$  and since  $\nabla F(\alpha) = P^\top \nabla f(P\alpha)$ . ■

484 Lemma 2.13 corresponds exactly to the result of Karimi et al. [29, Appendix B], that let a  
 485 Hoffman constant appear (that is in their context equal to the smallest non-zero eigenvalue of  
 486  $P$ ), as defined in [37, Section 3 and 4.1] by  $\theta(P) = \max_{z, \|P^\top z\|_{\infty}=1} \|z\|_2^2$ . They indeed proved  
 487 that  $F$  verifies the Łojasiewicz inequality, for all  $\alpha \in \mathbb{R}^d$ ,  $\frac{1}{2} \|\nabla F(\alpha)\|_2^2 \geq \theta(P) \mu_{\mathcal{P}}^f (F(\alpha) - F_*)$ .

488 Depending on the parametrization regime, we have proven that coordinate descent may  
 489 be formulated as a (possibly rebased) matching pursuit method. In the underparametrized  
 490 regime on the one hand, since  $F$  inherits all convex properties from  $f$ , the atoms are defined  
 491 by the Euclidean basis and the matching pursuit is formulated in  $\mathbb{R}^d$ . On the other hand  
 492 in the overparametrized regime, the introduction of a well-chosen gauge function  $\gamma_{\mathcal{P}}$  allows  
 493 to formulate coordinate descent as a matching pursuit algorithm in  $\mathbb{R}^n$ , and to perform a  
 494 convergence analysis using the strong convexity assumption on  $f$ . Again, smoothness and  
 495 strong convexity can be formulated as optimization problems depending on the gauge. The  
 496 gauge let also appear how  $F$  inherits some structure from  $f$ . In the next sections, we will see  
 497 how to generalize this framework for analyzing penalized linear models.

global values of the smoothness  
and strong convexity parameters  
can be ...

498     3. Phase transition for penalized linear models. We now consider the penalized linear  
 499 model,

500 (3.1)                    $\min_{\alpha \in \mathbb{R}^d} f(P\alpha) + \lambda \|\alpha\|_1 = F(\alpha) + H(\alpha) = G(\alpha),$

501 where  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $L_2^f$ -smooth,  $\mu_2^f$ -strongly convex, (and thus,  $F : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L_2^F$  smooth),  
 502  $P \in \mathbb{R}^{n \times d}$ ,  $\lambda > 0$  and where  $H(\alpha) = \lambda \|\alpha\|_1$  is closed convex and proper. In this section,  
 503 we derive a new matching pursuit algorithm for a  $\ell_1$ -regularized model, that we compare to  
 504 proximal coordinate descent with GS rule and to the proximal gradient descent. Building on  
 505 the result of Section 2, we derive convergence guarantees depending on the properties of  $f$   
 506 and  $P$ , and notice a strong connection to the proximal coordinate descent with GS rule. Yet,  
 507 in the overparametrized regime, neither the proximal gradient nor the regularized matching  
 508 pursuit benefits from linear convergence. Instead of that, we describe experimentally the role  
 509 of  $\lambda$  in the LASSO, as a continuous mapping between low-rank solutions and full-rank solution  
 510 to the least-squares.

511 **Proximal gradient descent.** Proximal gradient descent, a.k.a. forward-backward (see e.g.  
 512 [15]) was developed for such ‘composite’ convex optimization problem. Given a starting point  
 513  $\alpha_0 \in \mathbb{R}^d$ , each iterate is obtained by minimizing a smooth quadratic upper bound on  $F$ ,

514 (3.2)                    $G(\alpha) \leq F(\alpha_k) + \langle \nabla F(\alpha_k), \alpha - \alpha_k \rangle + \frac{L_2^F}{2} \|\alpha_k - \alpha\|_2^2 + \lambda \|\alpha\|_1.$

515 Minimizing the right side of the inequality leads to the proximal gradient method as follows:

516                        $\alpha_{k+1} = \arg \min_{\alpha \in \mathbb{R}^d} \langle \nabla F(\alpha_k), \alpha - \alpha_k \rangle + \frac{L_2^F}{2} \|\alpha - \alpha_k\|_2^2 + \lambda \|\alpha\|_1.$

517 The proximal gradient method converges sublinearly if  $F$  is not strongly convex, and linearly  
 518 if  $F$  is in addition  $\mu_2^F$ -strongly convex, such as in the underparametrized regime. Then, the  
 519 sequence  $\alpha_k$  starting from  $\alpha_0 \in \mathbb{R}^d$  verifies ([55, Theorem 2.1]),

520                        $G(\alpha_k) - G_* \leq \left(1 - \frac{\mu_2^F}{L_2^F}\right)^k (G(\alpha_0) - G_*).$

521 **Coordinate descent.** In practice, (randomized) coordinate gradient descent is widely used  
 522 to avoid computing the full gradient (that costs  $O(d)$ ), and is particularly suited to sparse  
 523 regression problems. Nutini et al. [43] analyzed coordinate descent with the Gauss-Southwell  
 524 selection rule, that tends to perform better than randomized coordinate descent.

525 (3.3)                    $\alpha_{k+1} = \arg \min_{\alpha \in \mathbb{R}^d} \nabla_{i_k} F(\alpha_k)(\alpha^{(i_k)} - \alpha_k^{(i_k)}) + \frac{L_2^F}{2} (\alpha^{(i_k)} - \alpha_k^{(i_k)})^2 + \lambda |\alpha^{(i_k)}|,$

526 where  $i_k = \arg \min_l \min_{t \in \mathbb{R}} \nabla_l F(P\alpha_k)(t - \alpha^{(l)}) + \frac{L_2^F}{2} (t - \alpha^{(l)})^2 + \lambda |t|$  corresponds to the GS  
 527 rule. Nutini et al. [43, Appendix 8] proved that coordinate descent with the Gauss-Southwell  
 528 rule makes at least as much progress as randomized coordinate descent,

529                        $G(\alpha_{k+1}) - G_* \leq \left(1 - \frac{\mu_2^F}{dL_2^F}\right) (G(\alpha_k) - G_*).$

530 A refinement, that let a sublinear dependence in the parameter  $\mu_1^F$  appear, is mentioned in [43,  
 531 Appendix 8]. Coordinate descent with GS-rule is closely related to matching pursuit as for  
 532 nonpenalized models, as detailed in Supplementary Materials SM3, where we formulate this  
 533 method as a ‘nearly’ matching pursuit algorithm.

534 In the following, we derive a matching pursuit procedure for  $\ell_1$ -regularized problems (3.1),  
 535 that we compare to classical boosting algorithm and coordinate descent with GS-rule. Af-  
 536 ter that, we compute convergence guarantees for smooth (possibly strong) convex functions.  
 537 Finally, we propose an interpretation of the convergence regimes depending on the penalty  $\lambda$ .  
on a function of

538 **3.1. Regularized matching pursuit.** We propose a new regularized matching pursuit algo-  
 539 rithm based on the  $\ell_1$ -geometry. The main idea is to replace the  $\ell_2$ -norm in the minimization  
 540 Problem (3.2) leading to the proximal gradient by  ~~$\ell_1$ -norm~~. Let  $F$  be convex,  $L_1^F$ -smooth, as  
 541 for coordinate descent with GS rule in the linear regression problem from Section 2. We define  
 542 the penalized matching pursuit method starting from  $\alpha_0 \in \mathbb{R}^d$  as the sequence minimizing  
 543 smoothness with respect to the  $\ell_1$ -norm at each iteration:

$$544 \quad (3.4) \quad \alpha_{k+1} = \arg \min_{\alpha \in \mathbb{R}^d} \langle P^\top \nabla f(P\alpha), \alpha - \alpha_k \rangle + \frac{L_1^F}{2} \|\alpha - \alpha_k\|_1^2 + \lambda \|\alpha\|_1. \quad \text{coord-wise}$$

545 Whereas the optimization steps in proximal gradient descent (3.2) and proximal coordinate  
 546 descent with the GS rule (3.3) can be decomposed per coordinate, the function  $\alpha \rightarrow \|\alpha\|_1^2$ ,  
 547 that is not decomposable. Based on the same upper bound (3.4), Song et al. [52, Algorithm  
 548 1] generalized greedy coordinate descent with the “SOft ThresOlding PrOjection” (SOTOP)  
 549 algorithm using a reweighted least-squares formulation (SM1.1). However, their methods is  
 550 neither a coordinate-based method, nor a boosting method. We propose instead a regularized  
 551 matching pursuit algorithm that draws a clean connection to boosting and proximal coordinate  
 552 descent.

553 In the following, we formulate this optimization step (3.4) as a matching pursuit algorithm,  
 554 that only calls for a linear minimization oracle. Using a variational trick detailed in SM1.1 to  
 555 approach  $\|\beta\|_1^2$ , let us begin with formulating the problem (3.4) starting from  $\alpha_k \in \mathbb{R}^d$  as a  
 556 decomposable optimization problem,

$$\begin{aligned} 557 \quad V_* &= \min_{\beta \in \mathbb{R}^d} \langle \nabla F(\alpha_k), \beta \rangle + \frac{L_1^F}{2} \|\beta\|_1^2 + \lambda \|\beta + \alpha_k\|_1, \\ 558 &= \min_{\beta \in \mathbb{R}^d} \max_{z \geq 0} \langle \nabla F(\alpha_k), \beta \rangle - \frac{z^2}{2L_1^F} + z\|\beta\|_1 + \lambda \|\beta + \alpha_k\|_1, \\ 559 &= \max_{z \geq 0} -\frac{z^2}{2L_1^F} + \min_{\beta \in \mathbb{R}^d} \sum_{i=1}^d \{\nabla_i F(\alpha_k) \beta^{(i)} + z|\beta^{(i)}| + \lambda|\beta^{(i)} + \alpha_k^{(i)}|\}. \end{aligned}$$

560 At the optimum,  $z = L_1^F \|\beta\|_1$ . The problem is now decomposable for each coordinate  $\beta^{(i)}$ ,  
 561 and can be reduced to an optimization problem in  $z \geq 0$  in Lemma 3.1.

562 **Lemma 3.1.** *The optimization step (3.4) can be reformulated as*

$$563 \quad (3.5) \quad V_* = \max_{z_{\min} \leq z} h(z) \triangleq -\frac{z^2}{2L_1^F} + \sum_{i \in I} \min \left( \lambda |\alpha_k^{(i)}|, -\nabla_i F(\alpha_k) \alpha_k^{(i)} + z |\alpha_k^{(i)}| \right),$$

564 where  $z_{\min} = (\max_i |\nabla_i F(\alpha_k)| - \lambda)_+$  and where  $I = \{i, \alpha_k^{(i)} \neq 0\}$  is the set of active atoms.

565 *Proof.* The function  $\phi_i(z, \beta^{(i)}) = \nabla_i F(\alpha_k) \beta^{(i)} + z|\beta^{(i)}| + \lambda|\beta^{(i)} + \alpha_k^{(i)}|$  is lower bounded if  
 566  $z \geq |\nabla_i F(\alpha_k)| - \lambda$  for all  $i \in I$ . Then,  $\phi_i(z, \cdot)$  attains its minima in  $\beta^{(i)} = 0$  with  $\phi_i(z, 0) =$   
 567  $\lambda|\alpha_k^{(i)}|$  or in  $\beta^{(i)} = -\alpha_k^{(i)}$  with  $\phi_i(z, -\alpha_k^{(i)}) = -\nabla_i F(\alpha_k) \alpha_k^{(i)} + z|\alpha_k^{(i)}|$  or in every possible value  
 568 for  $\beta^{(i)}$  if  $z = \pm \nabla_i F(\alpha_k) - \lambda$  with  $\phi_i(z, 0) = \lambda|\alpha_k^{(i)}|$ . ■

569 Lemma 3.1 leads to a convex constrained optimization problem in  $\mathbb{R}^+$ , whose objective is  
 570 piecewise quadratic with slope coefficients changing at  $z_i = \lambda + \nabla_i F(\alpha_k) \alpha_k^{(i)} / |\alpha_k^{(i)}|$ . Its con-  
 571 straints in  $z_{\min}$  includes the LMO. By construction, the LMO given by  $z_{\min}$  may correspond  
 572 to several atoms  $\beta_j$  such that  $j \in \arg \min_i |\nabla_i F(\alpha_k)|$ . Since we aim at solving Problem (3.4)  
 573 by constructing a solution as sparse as possible, we introduce Assumption 3.2.

574 *Assumption 3.2.* We consider only one atom corresponding to the LMO, that is  $i_{\min} \in$   
 575  $\arg \max_i |\nabla_i F(\alpha_k)|$ , such that  $z_{\min} = (\max_i |\nabla_i F(\alpha_k)| - \lambda)_+$ .

576 In the following ~~lemmas~~, we compute explicitly the minimum of  $h$ . Under Assumption 3.2,  
 577 we first deal in ~~Lemma 3.3~~ with the situation in which the objective is a pure quadratic, that  
 578 is for all  $i \in I$ ,  $z_i \geq z_{\min}$ .

579 *Lemma 3.3.* Let  $(z_*, \beta_*)$  be a solution of (3.5), assume  $\{i, z_i \geq z_{\min}\} = \emptyset$  and verify Assumption 3.2.  
 580 Then  $z_* = z_{\min}$ ,  $\beta_*^{(i_{\min})} = -\text{sign}(\nabla_{i_{\min}} F(\alpha_k)) \frac{z_{\min}}{L_1^F}$  and  $\beta_*^{(i)} = 0$  for  $i \neq i_{\min}$ .

581 *Proof.* The objective is quadratic and attains its minimum at  $z_{\min} = L_1^F |\beta^{(i_{\min})}|$ . ■

582 In the context of Lemma 3.3 and Assumption 3.2, only the atom given by the LMO in  $z_{\min} =$   
 583  $(\max_i |\nabla_i F(\alpha_k)| - \lambda)_+$  can be added to the set of active atoms. Now, we assume the objective  
 584 is piecewise quadratic, that is  $\mathcal{S} = \{i, z_i \geq z_{\min}\} \neq \emptyset$ .

585 *Lemma 3.4.* Let  $z_*, \beta_*$  be a solution of (3.5) and assume  $\mathcal{S} = \{i, z_i \geq z_{\min}\} \neq \emptyset$  and verify Assumption 3.2. There are four possible solutions to the Problem (3.5),

- If  $h'(z_{\min}) \leq 0$ , then  $z_* = z_{\min}$ .

588 In addition,  $\beta_*^{(i)} = \begin{cases} -\alpha_k^{(i)} & \text{if } z_i \geq z_*, \\ 0 & \text{if } z_i \leq z_*, \\ -\text{sign}(\nabla_{i_{\min}} F(\alpha_k)) \frac{z_{\min} - \sum_{i \in \mathcal{S}} |\alpha_k^{(i)}|}{L_1^F} & \text{if } i = i_{\min}. \end{cases}$

- If there exists  $k \in \mathcal{S}$  such that  $h'(z_k^+) \geq 0$  and  $h'(z_{k+1}^-) \leq 0$ , then  $z_* \in ]z_k, z_{k+1}[$ . In

590 addition,  $\beta_*^{(i)} = \begin{cases} -\alpha_k^{(i)} & \text{if } z_i \geq z_*, \\ 0 & \text{if } z_i \leq z_*. \end{cases}$

- If there exists  $k \in \mathcal{S}$  such that  $h'(z_k^-) \geq 0$  and  $h'(z_k^+) \leq 0$  then  $z_* = z_k$ . In addition,

592  $\beta_*^{(i)} = \begin{cases} -\alpha_k^{(i)} & \text{if } z_i > z_*, \\ 0 & \text{if } z_i < z_*, \\ -\text{sign}(\alpha_k^{(i)}) \left( \frac{z_k}{L_1^F} - \sum_{i, z_i > z_k} |\alpha_k^{(i)}| \right) & \text{if } i = k. \end{cases}$

- If  $h'(z_{|I|}) > 0$ , then for all  $i \in I$ ,  $z_* > z_i$  and  $\beta_*^{(i)} = \begin{cases} -\text{sign}(\alpha_k^{(i)}) \frac{z_i}{L_1^F} & \text{if } i = |I|, \\ 0 & \text{otherwise.} \end{cases}$

594 Proof. The function  $h$  is strictly concave and quadratic by part on  $[z_i, z_{i+1}]$ . The solution  
 595 to the optimization Problem (3.5) is thus obtained by studying the sign of  $h'(\cdot)$  at  $z_i^-$  and  $z_i^+$ .  
 596 By construction of the solution given in the proof of Lemma 3.1, for all  $i$  such that  $z_* > z_i$   
 597 (resp.  $z_* < z_i$ ), then  $\beta^{(i)} = 0$  (resp.  $\beta^{(i)} = -\alpha_k^{(i)}$ ). Finally, we have  $z_* = L_1^F \|\beta_*\|_1$  which  
 598 gives the solution for  $z_* = z_{\min}$  or  $z_* = z_k$ . *ok partout? ou bien mille part!!*

599 Lemmas 3.3 and 3.4 provides a closed form solution by calling only for the linear mini-  
 600 mization oracle  $\min_i |\nabla_i F(\alpha_k)|$ , and performing  $O(|I|)$  operations on the active atoms. From  
 601 that, we deduce Algorithm 3.1.

**Algorithm 3.1** Regularized matching pursuit (RMP)?

```

 $\alpha \in \mathbb{R}^d, N \in \mathbb{N}$ 
for  $k \in [0, \dots, N]$  do
   $z_{\min} = (\max_i |\nabla_i F(\alpha_k)| - \lambda)_+$  and  $i_{\min} = \arg \max_i |\nabla_i F(\alpha_k)|$ 
  For  $\alpha_k^{(i)} \neq 0$ , compute  $z_i = \lambda + \frac{\alpha_k^{(i)}}{|\alpha_k^{(i)}|} \nabla_i F(\alpha_k)$  such that  $z_{i+1} \geq z_i$ 
  if  $\{i, z_i \geq z_{\min}\} = \emptyset$  then
     $\beta_{i_{\min}} = -\text{sign}(\nabla_{i_{\min}} F(\alpha_k)) \frac{z_{\min}}{L_1^F}$ 
  else
    Compute  $u = \arg \min_i \{z_i \geq z_{\min}\}$  and for  $i \in [u, v]$ , compute  $h'(z_i)$ 
    if  $h'(z_{\min}) \leq 0$  or  $h'(z_u) \leq 0$  then
      For  $i \in [u, v]$ ,  $\beta^{(i)} = -\alpha_k^{(i)}$ 
      If  $h'(z_{\min}) \leq 0$ , then  $\beta_{i_{\min}} = -\text{sign}(\nabla_{i_{\min}} F(\alpha_k)) (\frac{z_{\min}}{L_1^F} - \sum_{i=u}^v |\alpha_k^{(i)}|)$ 
    else
       $n = \arg \max\{i, i \in [u, v-1], h'(z_i^+), h'(z_{i+1}^-) \geq 0\}$ 
      if  $n = v-1$  then
         $\beta_v = -\text{sign}(\alpha_v) \frac{1}{L_1^F} (\lambda + \frac{\alpha_v}{|\alpha_v|} \nabla_v F(\alpha_k))$ 
      else
        For  $i \in [n+1, v]$ ,  $\beta^{(i)} = -\alpha_k^{(i)}$ 
        If  $h'(z_n^+) \leq 0$ , then  $\beta^{(n)} = -\text{sign}(\alpha_k^{(n)}) (\frac{1}{L_1^F} \left( \lambda + \frac{\alpha_k^{(n)}}{|\alpha_k^{(n)}|} \nabla_n F(\alpha_k) \right) - \sum_{i=n+1}^v |\alpha_k^{(i)}|)$ 
      end if
    end if
  end if
   $\alpha_{k+1} = \alpha_k + \beta$ 
end for

```

*Notation* ↓  
*Set index, not superindex,*  
*Mais pas de changement de variable par la suite !*

*at each iteration, Algo 3.1 performs one of these..?*

602 In short, Algorithm 3.1 performs three possible actions: either one new atom is added (at  
 603 most) by calling the LMO  $\arg \max_i |\nabla_i F(\alpha_k)| = \arg \max_{p \in \mathcal{P}} p^\top \nabla f(P\alpha_k)$  while some active  
 604 atoms may be set to zero, or one active atom may be optimized while some active atoms may  
 605 be set to zero, or some active atoms are set to zero (but none is added nor optimized). To  
 606 sum it up, at each iteration, it constructs the next iterate using only past active atoms plus  
 607 possibly a new one generated by the LMO. Therefore, Algorithm 3.1 belongs to the family of  
 608 boosting algorithms. We refer to it as the regularized matching pursuit.

y  
 (RMP)?

609 Compared to the boosting approach of Zhang et al. [64] for metric-norm regularization or  
 610 to the generalized conditional gradient [3, 54], active atoms are not modified uniformly since  
 611 only some of them may be reduced to zero. The SOTOPPO method of Song [52] minimizes  
 612 the same upper bound with respect to the  $\ell_1$ -norm (3.4). Yet, it is resolved with a different  
 613 variational formulation, that does not let a linear minimization oracle appear. Compared to  
 614 proximal coordinate descent with GS-rule (3.3) applied with  $L_1^F$  (instead of  $L_2^F$ ), the regularized  
 615 matching pursuit happens to often follow exactly the same path when starting from zero  
 616 (but not when starting from a nonzero point), as observed in Figure 3. This suggests a  
 617 connection between regularized matching pursuit and proximal coordinate descent, as proven  
 618 by Locatello et al. [33] for gradient descent and steepest coordinate descent.

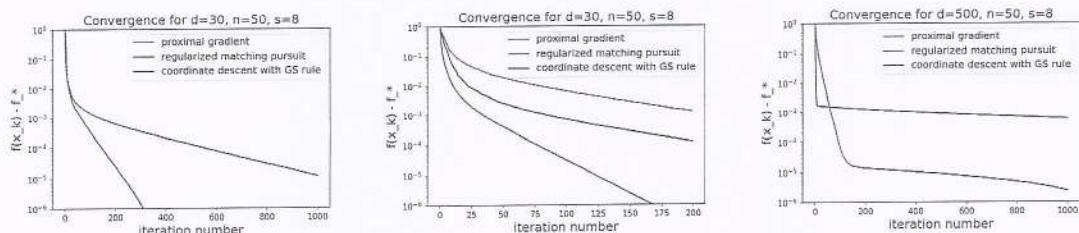


Figure 3: Convergence in function value of the proximal gradient descent, coordinate descent with Gauss-Southwell rule and with  $L = L_1^F$  (instead of  $L = L_2^F$ ) and of the regularized matching pursuit, for  $n = 50$ ,  $s = 8$ ,  $\lambda = 0.001$ ,  $\sigma = 0.5$  and for  $d = 30$  starting from zero on the left (underparametrized regime), from a non zero point in the middle (underparametrized regime) and for  $d = 500$  on the right (overparametrized regime). RMP and coordinate descent with GS-rule match exactly in these examples.

619 In Figure 3, the regularized matching pursuit seems to converge linearly in the under-  
 620 parametrized regime, and sublinearly in the overparametrized regime. We will compute its  
 621 theoretical convergence guarantees in the next section.

622 *3.2. Convergence guarantee.* We now establish convergence guarantees of the regularized  
 623 matching pursuit, both for strongly convex and non-strongly convex functions  $F$ . We  
 624 consider a more general composite minimization problem,

625 (3.6) 
$$\min_{\alpha \in \mathbb{R}^d} G(\alpha) = F(\alpha) + H(\alpha)$$

626 where  $H$  is closed, convex, proper, and where  $F$  is  $L_1^F$ -smooth and (possibly)  $\mu_1^F$ -strongly  
 627 convex with respect to the  $\ell_1$ -norm. If in addition,  $F$  is a linear mapping, and  $H(\cdot) = \|\cdot\|_1$ , this  
 628 is exactly the original optimization Problem (1.1). We evaluate the convergence guarantee a  
 629 generalized version of the regularized matching pursuit (that is not always a boosting method).

630 (3.7) 
$$\alpha_{k+1} = \arg \min_{\alpha \in \mathbb{R}^d} \langle \nabla F(\alpha_k), \alpha - \alpha_k \rangle + \frac{L_1^F}{2} \|\alpha - \alpha_k\|_1^2 + H(\alpha).$$

631 As we will see, our proofs are closely related to those for randomized coordinate descent [49,  
 632 Theorem 5, 7].

633     **3.2.1. Strongly convex functions**  $\text{R}$ . Let us assume that  $F$  is  $L_1^F$ -smooth and  $\mu_1^F$ -strongly  
 634 convex, typically in the underparametrized regime. Similarly to coordinate gradient descent  
 635 with GS rule which converges linearly in this context [43], regularized matching pursuit is  
 636 formulated as the minimization of the smoothness upper bound with respect to the  $\ell_1$ -norm.  
 637 Therefore, it benefits from linear convergence guarantees, detailed below.

638     Proposition 3.5. [41, Appendix 8] *If  $F$  be convex,  $L_1^F$ -smooth with respect to the  $\ell_1$ -norm,  
 639 and  $\mu_1^F$ -strongly convex with respect to the  $\ell_1$ -norm. Then, the sequence  $(\alpha_k)$  generated by (3.7)  
 640 verifies,*

$$641 \quad G(\alpha_{k+1}) - G_* \leq \left(1 - \frac{\mu_1^F}{L_1^F}\right) (G(\alpha_k) - G_*).$$

642     *Proof.* The proof consists in an optimization step over all trajectories, and not on a ran-  
 643 domized step (see [49]), as it was done by [41]. See SM4.2. ■

644     The regularized matching pursuit is a special case of method (3.7), and verifies the conver-  
 645 gence guarantee of Proposition 3.5. As a conclusion, it beats traditional boosting techniques  
 646 converging sublinearly, such as coordinate descent with GS rule (with  $1 - \frac{\mu_2^F}{dL_2^F} \leq 1 - \frac{\mu_1^F}{L_1^F}$ ), or  
 647 the generalized conditional gradient that is also appropriate to a gauge geometry. In addition,  
 648 its linear guarantee only depends on the strong convexity and smoothness parameters of  $F$   
 649 with respect to the  $\ell_1$ -norm. In the special case of the LASSO, the estimates established in  
 650 Proposition 2.8 and 2.9 still apply. In the overparametrized regime however, Figure 3 suggests  
 651 that the method does not converge linearly (since it is stuck at a precision around  $10^{-5}$ ). adopted?

652     **3.2.2. Smooth convex functions**  $\text{R}$ . Let now  $F$  be  $L_1^F$ -smooth, convex, but not strongly  
 653 convex (which is verified in the overparametrized regime). Usually, guarantees for splitting  
 654 methods, such as proximal gradient, states a sublinear convergence guarantee. Similarly in  
 655 Proposition 3.6, we prove sublinear convergence of the (generalized) regularized matching  
 656 pursuit. To our knowledge, there is no such theorem for sublinear convergence of SOTOP [52]  
 657 or of coordinate descent with the Gauss-Southwell rule, which is very close to the generalized  
 658 regularized matching pursuit. GRMP?

659     Proposition 3.6. *Let  $(\alpha_k)$  be generated by the generalized regularized matching pursuit (3.7),  
 660 starting from  $\alpha_0 \in \mathbb{R}^d$*

$$661 \quad G(\alpha_k) - G_* \leq \frac{2L_1^F \mathcal{R}_{\alpha_0}^2}{k+1},$$

662 where  $\mathcal{R}_{\alpha_0}^2 = \max_{\alpha \in \mathbb{R}^d} \max_{\alpha_* \in \mathbb{R}^d} \{\|\alpha - \alpha_*\|_1^2, \text{ s.t. } G(\alpha) \leq G(\alpha_0)\}$ .

663     *Proof.* This technique is inspired from the proof for sublinear convergence for randomized  
 664 proximal coordinate descent established by Richtarik and Takac [49, Theorem 5]. Let  $\alpha_{k+1} \in$

665  $\mathbb{R}^d$  be a minimizer of the smooth upper bound:

$$\begin{aligned}
 666 \quad G(\alpha_{k+1}) &\leq \inf_{\alpha \in \mathbb{R}^d} F(\alpha_k) + \langle \nabla F(\alpha_k), \alpha - \alpha_k \rangle + \frac{L_1^F}{2} \|\alpha - \alpha_k\|_1^2 + H(\alpha), \\
 667 \quad &\leq \inf_{\alpha \in \mathbb{R}^d} F(\alpha) + H(\alpha) + \frac{L_1^F}{2} \|\alpha - \alpha_k\|_1^2 (= G(\alpha) + \frac{L_1^F}{2} \|\alpha - \alpha_k\|_1^2) \text{ (} F \text{ convex),} \\
 668 \quad &\leq \inf_{t \in [0,1]} G(t\alpha_* + (1-t)\alpha_k) + \frac{L_1^F t^2}{2} \|\alpha_k - \alpha_*\|_1^2, \\
 669 \quad &\leq \inf_{t \in [0,1]} G(\alpha_k) - t(G(\alpha_k) - G_*) + \frac{L_1^F t^2}{2} \|\alpha_k - \alpha_*\|_1^2 \text{ (convexity of } H, F\text{),} \\
 670 \quad G(\alpha_{k+1}) - G_* &\leq \inf_{t \in [0,1]} (1-t)(G(\alpha_k) - G_*) + \frac{L_1^F t^2}{2} \|\alpha_k - \alpha_*\|_1^2.
 \end{aligned}$$

671 The solution of this minimization problem is given by  $t_* = \min(1, \frac{G(\alpha_k) - G_*}{L_1^F \|\alpha_k - \alpha_*\|_1^2})$ . We conclude  
672 the minimization bound, depending on the sign of  $G(\alpha_k) - G_* - L_1^F \|\alpha_k - \alpha_*\|_1^2$ :

$$673 \quad \text{Explain} \quad G(\alpha_{k+1}) - G_* \leq \max \left( 1 - \frac{G(\alpha_k) - G_*}{2L_1^F \|\alpha_k - \alpha_*\|_1^2}, \frac{1}{2} \right) (G(\alpha_k) - G_*).$$

674 As a first conclusion, notice that  $G(\alpha_k) - G_*$  is nonincreasing. Recall now that  $\mathcal{R}_{\alpha_0}^2 =$   
675  $\max_{\alpha \in \mathbb{R}^d} \max_{\alpha_* \in \mathbb{R}^d} \{\|\alpha - \alpha_*\|_1^2, \text{ s.t. } G(\alpha) \leq G(\alpha_0)\}$ . Then, using the notation  $\delta_k = G(\alpha_k) - G_*$ ,  
676 an upper bound for  $\delta_{k+1}$  is given by  $\delta_{k+1} \leq \max \left( 1 - \frac{\delta_k}{2L_1^F \mathcal{R}_{\alpha_0}^2}, \frac{1}{2} \right) \delta_k$ . Assume now that  
677  $\delta_0 \leq L_1^F \mathcal{R}_{\alpha_0}^2$  and notice that  $\delta_k \leq L_1^F \mathcal{R}_{\alpha_0}^2$  since  $\delta_k$  is nonincreasing. If not, notice that  
678 the inequality satisfied at the next iteration  $\delta_1 \leq \frac{1}{2} \mathcal{R}_{\alpha_0}^2$ . Then, we have for  $\omega = \frac{1}{2L_1^F \mathcal{R}_{\alpha_0}^2}$ ,  
679  $\delta_{k+1} \leq (1 - \delta_k \omega) \delta_k$ . Following the same argument than in the proof for sublinear convergence  
680 of steepest coordinate descent, detailed in SM4.1, it leads to the convergence guarantee  
681  $G(\alpha_k) - G_* \leq \frac{2L_1^F \mathcal{R}_{\alpha_0}^2}{k+1}$ .

682 In Proposition 3.6, we have obtained a sublinear convergence guarantee for the regularized  
683 matching pursuit for non-strongly convex functions. To our knowledge, this is the first sublin-  
684 ear guarantee for a boosting algorithm under classical assumptions from convex optimization.  
685 This method does not benefit from linear convergence guarantee. Yet, as we will see in nu-  
686 matical experiments in the next section, in the case of  $\ell_1$ -regularized model, the regularized  
687 matching pursuit converges linearly in certain regimes.

688 **3.3. A phase transition depending on  $\lambda$ : experimental results.** The regularized match-  
689 ing pursuit algorithm benefits from convergence guarantees similar to those for the proximal  
690 gradient: these methods converge linearly under strong convexity assumptions (under-  
691 parametrized regime) but sublinearly for non-strongly convex problems (overparametrized  
692 regime). In the context of sparsity though, the proximal gradient descent benefits from  
693 linear convergence under additional assumptions on the problem classes such as restricted  
694 eigenvalue properties [48], and for a well-chosen parameter  $\lambda$  [1, Theorem 2]. In this sec-  
695 tion, our experiments reveal a transition phenomenon driven by  $\lambda$  on a LASSO problem

*Re: non → regularized.*

*This manuscript is for review purposes only.*

*Only sublinear converg-  
guarantees when no  
strong-convexity - like  
property ...?*

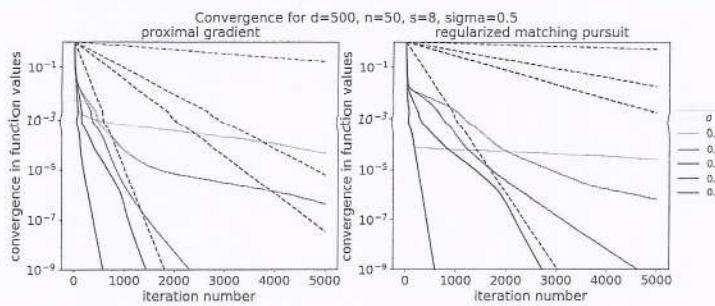


Figure 4: Convergence in function values for the proximal gradient on the left and the regularized matching pursuit on the right for  $n = 50$ ,  $d = 500$  and a sparsity  $s = 8$  and for several penalty  $\lambda$ . Convergence is compared in dashed lines to local convergence guarantee, taken on the support  $S$  on the last iterates and the SDP relaxation from Proposition 2.8.

696  $F(\alpha) = \frac{1}{2n} \|P\alpha - y\|_2^2 + \lambda \|\alpha\|_1$ , where  $P$  are synthetic Gaussian data in the overparametrized  
697 setting as in Section 2.3. pas

698 The convergence behavior of proximal gradient descent follows a transition phase, that can  
699 be divided into three phases: first, the method converges linearly according to the nonregu-  
700 larized trajectory, then it converges sublinearly, and it converges linearly once the support is  
701 identified. Iutzeler and Malick [27, Theorem 1] always identifies the structure of the solution  
702 (described by manifolds, or sparsity patterns) after a certain number of steps. For strongly  
703 convex functions, Nutini et al. [42] bounded the ‘active-set’-complexity of the proximal gra-  
704 dient method. The regularized matching pursuit follows the same behavior. It appears that  
705 the sparsity of the solution, and the sparsity identification highly depends on the value of  
706  $\lambda$ : the larger  $\lambda$ , the sparser the solution and the quicker the identification. Thanks to this  
707 observation, we derive a posteriori guarantee in Figure 4, based on the sparsity of the solution  
708 to the optimization problem. In Figure 4, local strong convexity parameters are given by  
709 the estimated of Corollary 2.6 and Proposition 2.8. We recover that large parameter  $\lambda$  both  
710 induces a stricter sparsity on the solution and a better convergence.

711 We describe this transition phase numerically in Figure 5 by plotting the  $\epsilon$ -curve (see  
712 Section 2.3) as a function of  $\lambda$ . For  $\lambda = 0$ , both methods converge linearly (as expected in the  
713 overparametrized regime for gradient descent and coordinate descent with the GS-rule). For  
714 ‘large’ values of  $\lambda$  for which the support is quickly identified, the convergence is linear too. In  
715 the middle phase however, convergence depends on the effective dimension of the trajectory,  
716  $d_{eff} \approx n$  by construction (and thus, on the effective strong convexity of  $f$  long the trajectory).  
717 The  $\epsilon$ -curves can be seen as equivalent in the optimization perspective with the regularization  
718 path usually drawn in the context of statistical recovery.

719 The regularized matching pursuit Algorithm 3.1 formulation allows some intuition re-  
720 garding the interplay between  $\lambda$  and the sparsity of the solution. Let  $\mathcal{A} = \{i, \alpha_k^{(i)} > 0\}$   
721 be the set of active atoms. Algorithm 3.1 may reduce an active atom  $i \in \mathcal{A}$  to zero if  
722  $\|\nabla F(\alpha)\|_\infty - \frac{\alpha_k^{(i)}}{|\alpha_k^{(i)}|} \nabla_i F(\alpha) \leq 2\lambda$ . The larger  $\lambda$ , the more active directions may be canceled out.

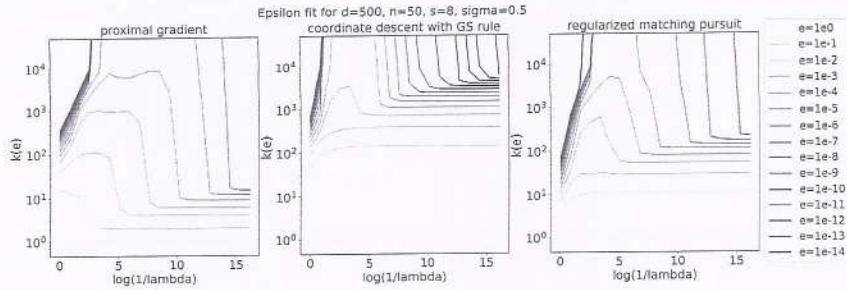


Figure 5:  $\epsilon$ -curve of the proximal gradient, coordinate descent with the GS rule and regularized matching pursuit for a LASSO problem with  $d = 500$ ,  $n = 50$ , a sparsity level  $s = 8$ ,  $\sigma = 0.5$ , after  $k = 10000$  iterations for several values of  $\lambda$ .

723 The smaller  $\lambda$  ( $\lambda \ll \|\nabla F(\alpha)\|_\infty$ ), the closer is regularized matching pursuit to coordinate  
 724 descent with GS rule (on the right in Figure 5): indeed, only the linear minimization oracle may  
 725 be added to the set of atoms without modifying other active atoms ( $z_i \lesssim z_{\min}$ ). For  $\lambda \approx 0$ , the  
 726 regularized matching pursuit thus converges linearly up to a certain iteration number, which  
 727 appears with the parallel level lines in Figure 5.

728 Based on the minimization of a smoothness upper bound with respect to the  $\ell_1$ -norm, we  
 729 have developed a regularized matching pursuit algorithm, that benefits from linear convergence  
 730 in the underparametrized regime (where  $F$  is strongly convex), and sublinear convergence in  
 731 the overparametrized regime (where  $F$  is not strongly convex). Thanks to the  $\epsilon$ -curve, we  
 732 numerically described the role of  $\lambda$  on the convergence of the method (and on the sparsity).  
 733 In the following section, we propose to develop a method suited to the gauge geometry in the  
 734 overparametrized regime.

735 **3.4. An ultimate method adapted to the geometry of regularized models.** The regular-  
 736 ized matching pursuit 3.1 was derived from the  $\ell_1$ -geometry. In Section 2.4 for non-regularized  
 737 models, coordinate descent with GS-rule was interpreted as a matching pursuit algorithm in the  
 738  $\ell_1$ -geometry in the underparametrized regime, but in the  $\gamma_P$ -geometry in the overparametrized  
 739 regime. We will see that the regularized matching pursuit as developed above does not benefit  
 740 from this formulation in the overparametrized regime. Instead, we will propose an ‘ultimate  
 741 method’ for the gauge geometry, that benefits from linear convergence in the overparametrized  
 742 regime but lacks a simple formulation.

743 Recall the equivalent regularized minimization problems (1.1) and (1.3),

744 *cont.* 
$$\min_{\alpha \in \mathbb{R}^d} f(P\alpha) + \lambda \|\alpha\|_1 = \min_{x \in \mathbb{R}^n} f(x) + \lambda \gamma_P(x),$$

745 where  $\gamma_P$  is a gauge function as defined in Section 2.4, and  $f$  is  $L_{\gamma_P}^f$ -smooth and  $\mu_{\gamma_P}^f$ -  
 746 strongly convex with respect to the gauge. The problem in  $\mathbb{R}^d$  is reformulated as  $\mathbb{R}^n$ , of lower  
 747 dimension.

748 *Remark 3.7.* This reformulation only requires  $\gamma_P$  to be a gauge function, but not specifically to be a norm. Reversely, minimizing a function penalized by a gauge function (or a

*fine 2 phon?*

*linear function**Mapping is**e.g.  $\mathbb{R}^d \rightarrow \mathbb{R}^d$ .*

750 semi-norm) can be reformulated as minimizing a linear mapping penalized by and  $\ell_1$ -norm.

751 As for the regularized matching pursuit, we formulate an optimization method as the  
752 minimization of the smoothness upper bound with respect to the gauge function, starting  
753 from  $x_0 \in \mathbb{R}^n$ .

$$754 \quad (3.8) \quad x_{k+1} = \arg \min_{x \in \mathbb{R}^n} \langle \nabla f(x_k), x - x_k \rangle + \frac{L_f^f}{2} \gamma_P(x - x_k)^2 + \lambda \gamma_P(x).$$

755 We refer to this method as the ultimate method for the gauge  $\gamma_P$ , that is adapted to the  
756 geometry of the regularized problem 1.1. Let us reformulate the minimization Problem (3.8)  
757 on  $\mathbb{R}^n$  into a minimization problem in  $\mathbb{R}^d$ . Let  $x_k = P\alpha_k$  with  $\alpha_k \in \mathbb{R}^d$ , then

$$\begin{aligned} 758 \quad & \min_{x \in \mathbb{R}^n} \langle \nabla f(x_k), x - x_k \rangle + \frac{L_f^f}{2} \gamma_P(x - x_k)^2 + \lambda \gamma_P(x), \\ 759 \quad &= \min_{\alpha, \nu \in \mathbb{R}^d} \langle \nabla f(P\alpha_k), P(\alpha - \alpha_k) \rangle + \frac{L_f^f}{2} \|\alpha - \alpha_k\|_1^2 + \lambda \|\nu\|_1, \text{ s.t. } x = P\alpha = P\nu, \\ 760 \quad &= \min_{\alpha, \nu \in \mathbb{R}^d} \langle \nabla F(\alpha_k), \alpha - \alpha_k \rangle + \frac{L_f^f}{2} \|\alpha - \alpha_k\|_1^2 + \lambda \|\nu\|_1, \text{ s.t. } P\alpha = P\nu. \end{aligned}$$

761 When  $P\alpha = P\nu$  implies  $\alpha = \nu$ , such as in the underparametrized regime where  $P^\top P$   
762 is invertible, the ultimate method for the gauge is equivalent with the regularized matching  
763 pursuit (3.1). However, in the overparametrized regime,  $P\alpha = P\nu$  does not imply  $\alpha = \nu$  in  
764 general. We believe this method does not belong to boosting algorithms due to the evaluation  
765 of the gauge function in  $x$  and in  $x - x_k$  in (3.8). In addition, this minimization problem admits  
766 neither a simple closed-form solution in general nor a solution based on the KKT conditions  
767 (as we did for regularized matching pursuit). While not directly computable in general, the  
768 minimization step (3.8) converges linearly to the optimum, as proven below in Proposition 3.8.

769 *Proposition 3.8.* *Let  $f$  be  $L_f^f$ -smooth and  $\mu_f^f$ -strongly convex with respect to the norm  $\gamma_P(\cdot)$ .  
770 The ultimate matching pursuit (3.8) ( $x_k$ ) converges linearly with*

$$771 \quad f(x_k) - f_* \leq \left(1 - \frac{\mu_f^f}{L_f^f}\right)^k (f(x_0) - f_*).$$

772 *Proof.* The proof follows exactly the proof for Theorem 3.5, replacing the function  $F$  by  $f$   
773 and the norm  $\|\cdot\|_1$  by  $\gamma_P(\cdot)$ . *qed*

774 In Proposition 3.8, we prove linear convergence for the ultimate matching pursuit algorithm.  
775 We recover the convergence guarantee of coordinate descent with GS rule in the non regularized  
776 model (2.6).

777 In Figure 6, we solve the optimization step for the ultimate method for the gauge with  
778 the solver MOSEK [2] on a LASSO problem. It converges linearly in the overparametrized  
779 regime, while the other method are first stuck in a sublinear phase. Compared to the proximal  
780 gradient, proximal coordinate descent with GS rule, the regularized matching pursuit and  
781 the ultimate method starts with sparse solution, and differs after a small number of iteration

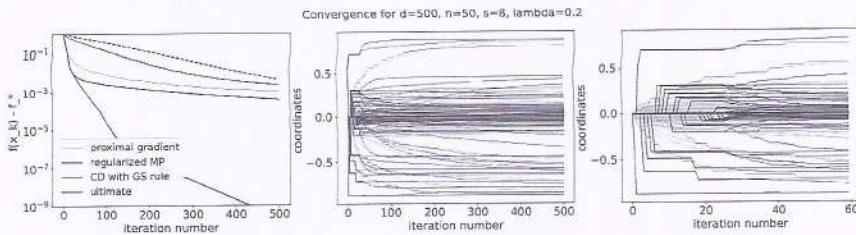


Figure 6: Convergence in function value and coordinates as a function of the iteration number for the proximal gradient descent, the proximal coordinate descent with GS rule, the regularized matching pursuit and the ultimate method on a LASSO problem, with  $d = 500$ ,  $n = 50$ ,  $s = 8$ ,  $\lambda = 0.2$ . The approximate guarantee is provided in dashed lines.

(about 30 here). In the special case of the LASSO, it is possible to approximate its convergence guarantee as for the linear regression problem. Noticing that  $L_{\gamma_P}^F = L_1^F$  and  $\mu_{\gamma_P}^F = \mu_1^F$ , the estimate of the convergence guarantee of coordinate descent with GS rule from Proposition 2.8 apply here. In the Supplementary Material ??, we propose an inner loop strategy to avoid the use of an optimization solver, together with the convergence analysis of the outer loop given the precision of the inner loop.



*Unclear*

**Conclusion and future works.** In this paper, we developed a principled view for generating optimization algorithms from the minimization of a smoothness upper bound with respect to a well-chosen norm. For non-regularized models, this procedure leads to coordinate descent with GS-rule, that can be interpreted as a matching pursuit algorithm both in the  $\ell_1$ -geometry for underparametrized models, and in the  $\gamma_P$ -geometry for overparametrized models. Building on these results, we derive a new regularized matching pursuit algorithm based on the minimization of smoothness with respect to the  $\ell_1$ -norm (whose counterpart is proximal gradient descent in the  $\ell_2$ -geometry). While being strongly connected to proximal coordinate descent with GS-rule, the regularized matching pursuit cannot be interpreted as a matching pursuit algorithm in the gauge geometry for overparametrized models, and thus, does not converge linearly in this regime. We finally formulate the ultimate method adapted to overparametrized geometries. Yet, this method lacks a closed-form formulation, that we approximate thanks to an inner-loop strategy.

*an*

From this approach, we obtain refined convergence guarantees for (resp. regularized) matching pursuit (resp. coordinate descent with GS rule), that are adapted to the geometry of the problem under consideration. For linear regression and the LASSO, we derive approximate convergence guarantees using SDP relaxations and random matrix theory. As a byproduct, convergence guarantees of both gradient descent and steepest coordinate descent applied to least-squares follow a transition phase from the underparametrized to the overparametrized regime. For  $\ell_1$ -regularized models, a similar transition phase for  $\lambda$  appears, and allows to interpret it as a measure of the sparsity of the solution.

*In manifold  $X_P$ , we approx it using ..*

Building on these results, it could be of interest to extend this principled approach to accelerated matching pursuit algorithm (and thus, to accelerated coordinate descent algo-

*We believe in  
would be ..*

*work "unfortunate"  
working.*

811 rithms). Some accelerated techniques have already been developed relying on randomly se-  
 812 lected coordinates, such as those of Nesterov et al. [40] for nonregularized minimization and  
 813 Fercoq and Richtarik [19] or Locatello et al. [33, Section 3] for composite minimization prob-  
 814 lems, but. Another interesting line of research could be to understand the connection between  
 815 the observed phase transition for optimization methods and the double descent phenomenon  
 816 observed for the generalization error in machine learning.

817 **Acknowledgments.** This work was funded by MTE and the Agence Nationale de la  
 818 Recherche as part of the “Investissements d’avenir” program, reference ANR-19-P3IA-0001  
 819 (PRAIRIE 3IA Institute). We also acknowledge support from the European Research Council  
 820 (grant SEQUOIA 724063).

*ANR-19-*

821

## REFERENCES

- 822 [1] A. AGARWAL, S. NEGAHBAN, AND M. J. WAINWRIGHT, *Fast global convergence rates of gradient meth-  
 823 ods for high-dimensional statistical recovery*, in Advances in Neural Information Processing Systems,  
 824 vol. 23, 2010.
- 825 [2] M. AP'S, *The MOSEK optimization toolbox for MATLAB manual. Version 10.0.*, 2022, <http://docs.mosek.com/9.0/toolbox/index.html>.
- 826 [3] F. BACH, *Duality between subgradient and conditional gradient methods*, SIAM Journal on Optimization, 25 (2015), pp. 115–129.
- 827 [4] F. BACH, *High-dimensional analysis of double descent for linear regression with random projections*. Technical report, arXiv:2303.01372, 2023.
- 828 [5] Z. BAI AND J. W. SILVERSTEIN, *Spectral analysis of large dimensional random matrices*, Springer Series  
 829 in Statistics, (2010).
- 830 [6] H. H. BAUSCHKE AND P. L. COMBETTES, *Convex Analysis and Monotone Operator Theory in Hilbert  
 831 Spaces*, 2017.
- 832 [7] A. BECK AND L. TETRUASHVILI, *On the convergence of block coordinate descent type methods*, SIAM  
 833 Journal on Optimization, 23 (2013), pp. 2037–2060.
- 834 [8] M. BELKIN, D. HSU, S. MA, AND S. MANDAL, *Reconciling modern machine learning and the bias-  
 835 variance trade-off*, Proceedings of the National Academy of Sciences, 32 (2019), pp. 15849–15854.
- 836 [9] R. BERTHIER, F. BACH, AND P. GAILLARD, *Accelerated gossip in networks of given dimension using  
 837 Jacobi polynomial iterations*, SIAM Journal on Mathematics of Data Science, 2 (2020), pp. 24–47.
- 838 [10] J. BOLTE, A. DANILINIS, O. LEY, AND L. MAZET, *Characterization of Lojasiewicz inequalities: sub-  
 839 gradient flows, talweg, convexity*, Transactions of the American Mathematical Society, 362 (2010),  
 840 pp. 3319–3363.
- 841 [11] K. BORGWARDT, *The average number of pivot steps required by the simplex-method is polynomial*, Zeitschrift  
 842 für Operations Research, 26 (1987), pp. 157–177.
- 843 [12] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine learning*,  
 844 SIAM Review, 60 (2018), pp. 223–311.
- 845 [13] S. BOUCHERON, G. LUGOSI, AND P. MASSART, *Concentration inequalities. A nonasymptotic theory of  
 846 independence*, Oxford University Press, 2013.
- 847 [14] E. CANDES AND T. TAO, *Decoding by linear programming*, IEEE Transactions on Information Theory, 51 (2005), pp. 4203–4215.
- 848 [15] P. L. COMBETTES AND V. R. WAJS, *Signal recovery by proximal forward-backward splitting*, Multiscale  
 849 Modeling & Simulation, 4 (2005), pp. 1168–1200.
- 850 [16] A. D'ASPREMONT, C. GUZMÁN, AND M. JAGGI, *Optimal affine-invariant smooth minimization algo-  
 851 rithms*, SIAM Journal on Optimization, 28 (2018), pp. 2384–2405.
- 852 [17] M. DUDIK, Z. HARCHAOUI, AND J. MALICK, *Lifted coordinate descent for learning with trace-norm  
 853 regularization*, in Proceedings of the Fifteenth International Conference on Artificial Intelligence and  
 854 Statistics, 2012, pp. 327–336.

- 859 [18] H. FANG, Z. FAN, Y. SUN, AND M. P. FRIEDLANDER, *Greed meets sparsity: Understanding and improving greedy coordinate descent for sparse optimization*, in Proceedings of the International Conference  
860 on Artificial Intelligence and Statistics, 2020.
- 861 [19] O. FERCOQ AND P. RICHTÁRIK, *Accelerated, parallel, and proximal coordinate descent*, SIAM Journal on  
863 Optimization, 25 (2015), pp. 1997–2023.
- 864 [20] M. FRANK AND P. WOLFE, *An algorithm for quadratic programming*, Naval Research Logistics Quarterly,  
865 3 (1956), pp. 95–110.
- 866 [21] Y. FREUND AND R. E. SCHAPIRE, *A short introduction to boosting*, Japanese Society For Artificial  
867 Intelligence, 14 (1999), pp. 771–780.
- 868 [22] M. P. FRIEDLANDER, I. MACÉDO, AND T. K. PONG, *Gauge optimization and duality*, SIAM Journal  
869 on Optimization, 24 (2014), pp. 1999–2022.
- 870 [23] M. X. GOEMANS AND D. P. WILLIAMSON, *Improved approximation algorithms for maximum cut and  
871 satisfiability problems using semidefinite programming*, Journal of the ACM, 42 (1995), pp. 1115–1145.
- 872 [24] T. GOLUB, D. K. SLONIM, P. TAMAYO, C. HUARD, M. GAASENBEEK, J. P. MESIROV, H. COLLER,  
873 M. LOH, J. R. DOWNING, M. A. CALIGIURI, C. D. BLOOMFIELD, AND E. S. LANDER, *Molecular  
874 classification of cancer: Class discovery and class prediction by gene expression monitoring*, Science,  
875 286 (1999), pp. 531–537.
- 876 [25] C. GUILLE-ESCIRET, B. GOUJAUD, M. GIROTTI, AND I. MITLIAGKAS, *A study of condition numbers  
877 for first-order optimization*, in Proceedings of the International Conference on Artificial Intelligence  
878 and Statistics, 2021, pp. 1261–1269.
- 879 [26] A. J. HOFFMAN, *On approximate solutions of systems of linear inequalities*, Journal Research of the  
880 National Bureau of Standards, 49 (1957), pp. 263–264.
- 881 [27] F. IUTZELER AND J. MALICK, *Nonsmoothness in machine learning: specific structure, proximal identifi-  
882 cation, and applications*, Set-Valued and Variational Analysis, 28 (2020), pp. 661–678.
- 883 [28] M. JAGGI, *Revisiting Frank-Wolfe: Projection-free sparse convex optimization*, in Proceedings of the  
884 International Conference on Machine Learning, no. 1, 2013, pp. 427–435.
- 885 [29] H. KARIMI, J. NUTINI, AND M. SCHMIDT, *Linear convergence of gradient and proximal-gradient methods  
886 under the Polyak-Lojasiewicz condition*, in Machine Learning and Knowledge Discovery in Databases,  
887 2016, pp. 795–811.
- 888 [30] S. P. KARIMIREDDY, A. KOLOSKOVA, S. U. STICH, AND M. JAGGI, *Efficient greedy coordinate descent  
889 for composite problems*, in Proceedings of the International Conference on Artificial Intelligence and  
890 Statistics, 2019.
- 891 [31] S. LACOSTE-JULIEN AND M. JAGGI, *On the global linear convergence of frank-wolfe optimization variants*,  
892 in Advances in Neural Information Processing Systems, 2015.
- 893 [32] F. LOCATELLO, R. KHANNA, M. TSCHANNEN, AND M. JAGGI, *A unified optimization view on general-  
894 ized matching pursuit and Frank-Wolfe*, in Proceedings of the International Conference on Artificial  
895 Intelligence and Statistics, 2017, pp. 860–868.
- 896 [33] F. LOCATELLO, A. RAJ, S. P. KARIMIREDDY, G. RAETSCH, B. SCHÖLKOPF, S. STICH, AND M. JAGGI,  
897 *On matching pursuit and coordinate descent*, in Proceedings of the International Conference on Ma-  
898 chine Learning, 2018, pp. 3198–3207.
- 899 [34] S. G. MALLAT AND Z. ZHANG, *Matching pursuits with time-frequency dictionaries*, IEEE Transactions  
900 on Signal Processing, 41 (1993), p. 3397–3415.
- 901 [35] V. A. MARCHENKO AND L. A. PASTUR, *Distribution of eigenvalues for some sets of random matrices*,  
902 Matematicheskii Sbornik, 1 (1967), pp. 457 – 483.
- 903 [36] S. MEI AND A. MONTANARI, *The generalization error of random features regression: Precise asymptotics  
904 and the double descent curve*, Communications on Pure and Applied Mathematics, 75 (2022), pp. 667–  
905 766.
- 906 [37] I. NECOARA, Y. NESTEROV, AND F. GLINEUR, *Linear convergence of first order methods for non-strongly  
907 convex optimization*, Mathematical Programming, 175 (2019), p. 69–107.
- 908 [38] Y. NESTEROV, *A method for solving the convex programming problem with convergence rate  $o(1/k^2)$* ,  
909 Proceedings of the USSR Academy of Sciences, 269 (1983), pp. 543–547.
- 910 [39] Y. NESTEROV, *Efficiency of coordinate descent methods on huge-scale optimization problems*, SIAM Jour-  
911 nal on Optimization, 22 (2012), pp. 341–362.
- 912 [40] Y. NESTEROV AND S. U. STICH, *Efficiency of the accelerated coordinate descent method on structured*

(in Johnson)  
parfois pages,  
parfois pas  
(en cours):  
taut le  
c'est à dire  
je pourrai

et tout de même idem

- 913 optimization problems, SIAM Journal on Optimization, 27 (2017), pp. 110–123.
- 914 [41] J. NUTINI, *Greed is good: greedy optimization methods for large-scale structured problems*, PhD thesis,  
915 University of British Columbia, 2018.
- 916 [42] J. NUTINI, M. SCHMIDT, AND W. HARE, "active-set complexity" of proximal gradient: How long does it  
917 take to find the sparsity pattern?, Optimization Letters, 13 (2018), pp. 645–655.
- 918 [43] J. NUTINI, M. SCHMIDT, I. LARADJI, M. FRIEDLANDER, AND H. KOEPKE, Coordinate descent con-  
919 verges faster with the Gauss-Southwell rule than random selection, in Proceedings of the International  
920 Conference on Machine Learning, 2018, pp. 1632–1641.
- 921 [44] C. PAQUETTE, B. VAN MERRIËNBOER, E. COURTNEY, AND F. PEGREGOSA, Halting time is predictable  
922 for large models: A universality property and average-case analysis, Foundations of Computational  
923 Mathematics, 23 (2023), p. 597–673.
- 924 [45] N. PARikh AND S. P. BOYD, Proximal algorithms, Foundations and Trends Optimization, 1 (2013),  
925 pp. 127–239.
- 926 [46] F. PEDREGOSA AND D. SCIEUR, Acceleration through spectral density estimation, in Proceedings of the  
927 37th International Conference on Machine Learning, 2020, pp. 7553–7562.
- 928 [47] F. PEDREGOSA AND D. SCIEUR, Acceleration through spectral density estimation, in Proceedings of the  
929 37th International Conference on Machine Learning, 2020, pp. 7553–7562.
- 930 [48] G. RASKUTTI, M. J. WAINWRIGHT, AND B. YU, Restricted eigenvalue properties for correlated gaussian  
931 designs, Journal of Machine Learning Research, 11 (2010), pp. 2241–2259.
- 932 [49] P. RICHTÁŘIK AND M. TAKÁČ, Iteration complexity of randomized block-coordinate descent methods for  
933 minimizing a composite function, Mathematical Programming, 144 (2014).
- 934 [50] D. SCIEUR AND F. PEDREGOSA, Universal asymptotic optimality of Polyak momentum, in Proceedings  
935 of the 37th International Conference on Machine Learning, 2020, pp. 8565–8572.
- 936 [51] S. A. B. SHENG CHEN AND W. LUO, Orthogonal least squares methods and their application to non-linear  
937 system identification, International Journal of Control, 50 (1989), pp. 1873–1896.
- 938 [52] C. SONG, S. CUI, Y. JIANG, AND S.-T. XIA, Accelerated stochastic greedy coordinate descent by soft  
939 thresholding projection onto simplex, in Advances in Neural Information Processing Systems, 2017.
- 940 [53] D. SPIELMAN AND S.-H. TENG, Smoothed analysis of algorithms: Why the simplex algorithm usually  
941 takes polynomial time, in Proceedings of the Thirty-Third Annual ACM Symposium on Theory of  
942 Computing, 2001, p. 296–305.
- 943 [54] Y. SUN AND F. BACH, Safe screening for the generalized conditional gradient method, Open Journal of  
944 Mathematical Optimization, 3 (2022), pp. 1–35.
- 945 [55] A. B. TAYLOR, J. HENDRICKX, AND F. A. GLINEUR, Exact worst-case convergence rates of the proximal  
946 gradient method for composite convex minimization, Journal of Optimization Theory and Applications, 178 (2018), pp. 455–476.
- 947 [56] A. TEWARI, P. RAVIKUMAR, AND I. DHILLON, Greedy algorithms for structurally constrained high di-  
948 mensional problems, in Advances in Neural Information Processing Systems, 2011.
- 949 [57] R. TIBSHIRANI, Regression shrinkage and selection via the Lasso, Journal of the Royal Statistical Society.  
950 Series B, 58 (1996), pp. 267–288.
- 951 [58] J. A. TROPP, Greed is good: algorithmic results for sparse approximation, IEEE Transactions on Infor-  
952 mation Theory, 50 (2004), pp. 2231–2242.
- 953 [59] P. TSENG, Convergence of a block coordinate descent method for nondifferentiable minimization, Journal  
954 of Optimization Theory and Applications, 109 (2001), p. 475–494.
- 955 [60] P. TSENG AND S. YUN, A coordinate gradient descent method for nonsmooth separable minimization,  
956 Mathematical Programming, B (2009), pp. 387–423.
- 957 [61] R. VERSHYNIN, *High-Dimensional Probability: An Introduction with Applications in Data Science*, Cam-  
958 bridge Series in Statistical and Probabilistic Mathematics, 2018.
- 959 [62] T. ZHANG, Adaptive forward-backward greedy algorithm for learning sparse representations, IEEE Trans-  
960 actions on Information Theory, 57 (2011), pp. 4689–4708.
- 961 [63] T. ZHANG, Sparse recovery with orthogonal matching pursuit under RIP, IEEE Transactions on Informa-  
962 tion Theory, 57 (2011), pp. 6215–6221.
- 963 [64] X. ZHANG, D. SCHUERMANS, AND Y.-L. YU, Accelerated training for matrix-norm regularization: A  
964 boosting approach, in Advances in Neural Information Processing Systems, vol. 25, 2012.

2 X

Some: am for

966 Appendix.

967 **Appendix A. Linear convergence Propositions 2.4, 2.7.**

968 Let us prove linear convergence of gradient descent with fixed step size and coordinate  
969 descent with GS rule in a more general framework. This proof leads to the results of Proposi-  
970 tions 2.4, 2.7 for the  $\ell_2$ ,  $\ell_1$  norm respectively.

971 Let  $F$  be  $L^F$ -smooth with respect to a norm  $\|\cdot\|$ , (possibly)  $\mu$ -strongly convex with respect  
972 to  $\|\cdot\|$  and verify the Łojasiewicz inequality with parameter  $\mu^{L,F}$ . We consider a method  $(\alpha_k)$   
973 starting from  $\alpha_0 \in \mathbb{R}^d$ , and obtained by minimizing the smoothness quadratic upper bound:

974 
$$F(\alpha_{k+1}) \leq F(\alpha_k) + \min_{\alpha \in \mathbb{R}^d} \left( \langle \nabla F(\alpha_k), \alpha - \alpha_k \rangle + \frac{L^F}{2} \|\alpha - \alpha_k\|^2 \right) \leq F(\alpha_k) - \frac{1}{2L^F} \|\nabla F(\alpha_k)\|_*^2.$$

975 If  $F$  is  $\mu^F$ -strongly convex, then  $F$  verifies the Łojasiewicz inequality with parameter  $\mu^F$ :  
976 for all  $\alpha \in \mathbb{R}^d$ ,  $\mu^F(F(\alpha) - F_*) \leq \frac{1}{2} \|\nabla F(\alpha)\|_*^2$ . Thus, subtracting  $F_*$  on each side of the  
977 smoothness inequality, we have

978 
$$F(\alpha_{k+1}) - F_* \leq \left(1 - \frac{\mu^F}{L^F}\right) (F(\alpha_k) - F_*).$$

979 Similarly, if  $F$  satisfies the Łojasiewicz inequality with parameter  $\mu^{L,F}$ , but is not strong convex  
980 ( $\mu^F = 0$ ). *→ dire que la which juste (4) qui est Loya.*

981 **Appendix B. Approximate linear convergence of coordinate descent for least-squares.**

982 In Proposition 2.7, a sequence  $(\alpha_k)$  generated by steepest coordinate descent for a linear  
983 regression problem has a linear convergence rate in function values,

984 
$$F(\alpha_k) - F_* \leq \left(1 - \frac{\max(\mu_1^F, \mu_1^{L,F})}{L_1^F}\right)^{2k} (F(\alpha_0) - F_*).$$

985 We first derive inequalities connecting  $\mu_1^{L,F}$  and  $\mu_1^F$  to  $L_1^F$ . Then, we prove concentration  
986 inequalities for  $\mu_1^{L,F}$ ,  $\mu_1^F$  and  $L_1^F$ , and derive approximate convergence guarantees of steepest  
987 coordinate descent.

988 **B.1. SDP relaxations for  $\mu_1^F$  and  $\mu_1^{PL}$ :** *proof for Proposition 2.8.* We look for exact  
989 lower bounds for  $\mu_1^F$  and  $\mu_1^{PL}$ . Both in the overparametrized and underparametrized regime,  
990 we are going to reformulate the optimization problems defining  $\mu_1, \mu_1^{L,F}$  into SDPs, and relax  
991 some rank constraints.

992 **SDP relaxation for  $\mu_1^F$  in the underparametrized regime.** Recall from Lemma 2.3 that  
993  $\mu_1^F$  is non zero in this regime and given by  $\frac{1}{\sqrt{n\mu_1^F}} = \max_{\alpha, \|\alpha\|_\infty \leq 1} \max_{\nu, \|P\nu\|_2 \leq 1} \alpha^\top \nu$ . Since