

Alma Mater Studiorum Università di Bologna
Archivio istituzionale della ricerca

Fast and scalable Lasso via stochastic Frank-Wolfe methods with a convergence guarantee

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version:

Emanuele, F., Ricardo, Ñ., Stefano, L., Claudio, S., Johan, A.K.S. (2016). Fast and scalable Lasso via stochastic Frank-Wolfe methods with a convergence guarantee. MACHINE LEARNING, 104(2-3), 195-221 [10.1007/s10994-016-5578-4].

Availability:

This version is available at: <https://hdl.handle.net/11585/585648> since: 2017-05-10

Published:

DOI: <http://doi.org/10.1007/s10994-016-5578-4>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

This is the final peer-reviewed accepted manuscript of:

Frandi, E., Ñanculef, R., Lodi, S. et al. Fast and scalable Lasso via stochastic Frank–Wolfe methods with a convergence guarantee. Mach Learn 104, 195–221 (2016).

The final published version is available online at: <https://doi.org/10.1007/s10994-016-5578-4>

Rights / License:

The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>)

When citing, please refer to the published version.

Fast and Scalable Lasso via Stochastic Frank-Wolfe Methods with a Convergence Guarantee

Emanuele Frandi¹, Ricardo Nanculef², Stefano Lodi³, Claudio Sartori⁴, and Johan A. K. Suykens⁵

^{1,5}ESAT-STADIUS, KU Leuven, Belgium
`{efrandi,johan.suykens}@esat.kuleuven.be`

²Department of Informatics, Federico Santa María Technical University, Chile
`jnancu@inf.utfsm.cl`

^{3,4}Department of Computer Science and Engineering, University of Bologna, Italy
`{stefano.lodi,claudio.sartori}@unibo.it`

Abstract

Frank-Wolfe (FW) algorithms have been often proposed over the last few years as efficient solvers for a variety of optimization problems arising in the field of Machine Learning. The ability to work with cheap projection-free iterations and the incremental nature of the method make FW a very effective choice for many large-scale problems where computing a sparse model is desirable.

In this paper, we present a high-performance implementation of the FW method tailored to solve large-scale Lasso regression problems, based on a randomized iteration, and prove that the convergence guarantees of the standard FW method are preserved in the stochastic setting. We show experimentally that our algorithm outperforms several existing state of the art methods, including the Coordinate Descent algorithm by Friedman *et al.* (one of the fastest known Lasso solvers), on several benchmark datasets with a very large number of features, without sacrificing the accuracy of the model. Our results illustrate that the algorithm is able to generate the complete regularization path on problems of size up to four million variables in less than one minute.

1 Introduction

Many Machine Learning and Data Mining tasks can be formulated, at some stage, in the form of an optimization problem. As constantly growing amounts of high dimensional data are becoming available in the Big Data era, a fundamental thread in research is the development of high-performance implementations of algorithms tailored to solving these problems in a very large-scale setting. One of the most popular and powerful techniques for high-dimensional data analysis is the *Lasso* [43]. In the last decade there has been intense interest in this method, and several papers describe generalizations and variants of the Lasso

[44]. In the context of supervised learning, it was recently proved that the Lasso problem can be reduced to an equivalent SVM formulation, which potentially allows one to leverage a wide range of efficient algorithms devised for the latter problem [23]. For unsupervised learning, the idea of Lasso regression has been used in [30] for bi-clustering in biological research.

From an optimization point of view, the Lasso can be formulated as an ℓ_1 -regularized least squares problem, and large-scale instances must usually be tackled by means of an efficient first-order algorithm. Several such methods have already been discussed in the literature. Variants of Nesterov’s Accelerated Gradient Descent, for example, guarantee an optimal convergence rate among first-order methods [36]. Stochastic algorithms such as Stochastic Gradient Descent and Stochastic Mirror Descent have also been proposed for the Lasso problem [29, 41]. More recently, Coordinate Descent (CD) algorithms [11, 12], along with their stochastic variants [41, 38], are gaining popularity due to their efficiency on structured large-scale problems. In particular, the CD implementation of Friedman *et al.* mentioned above is specifically tailored for Lasso problems, and is currently recognized as one of the best solvers for this class of problems.

The contribution of the present paper in this context can be summarized as follows:

- We propose a high-performance stochastic implementation of the classical Frank-Wolfe (FW) algorithm to solve the Lasso problem. We show experimentally how the proposed method is able to efficiently scale up to problems with a very large number of features, improving on the performance of other state of the art methods such as the Coordinate Descent algorithm in [12].
- We include an analysis of the complexity of our algorithm, and prove a novel convergence result that yields an $\mathcal{O}(1/k)$ convergence rate analogous (in terms of expected value) to the one holding for the standard FW method.
- We highlight how the properties of the FW method allow to obtain solutions that are significantly more sparse in terms of the number of features compared with those from various competing methods, while retaining the same optimization accuracy.

On a broader level, and in continuity with other works from the recent literature [35, 17, 42], the goal of this line of research is to show how FW algorithms provide a general and flexible optimization framework, encompassing several fundamental problems in Machine Learning.

Structure of the Paper

In Section 2 we provide an overview of the Lasso problem and its formulations, and review some of the related literature. Then, in Section 3 we discuss FW optimization and specialize the algorithm for the Lasso problem. The randomized algorithm used in our implementation is discussed in Section 4 and its convergence properties are analyzed. In Section 5 we show several experiments on benchmark datasets and discuss the obtained results. Finally, Section 6 closes the paper with some concluding remarks.

2 The Lasso Problem

Suppose we are given data points (x_ℓ, y_ℓ) , $\ell = 1, 2, \dots, m$, where $x_\ell = (x_{\ell 1}, x_{\ell 2}, \dots, x_{\ell p})^T \in \mathbb{R}^p$ are some predictor variables and y_ℓ the respective responses. A common approach in statistical modeling and Data Mining is the linear regression model, which predicts y_ℓ as a linear combination of the input attributes:

$$\hat{y}_\ell = \sum_{i=1}^p \alpha_i x_{\ell i} + \alpha_0 .$$

In a high-dimensional, low sample size setting ($p \gg m$), estimating the coefficient vector $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_p)^T \in \mathbb{R}^p$ using ordinary least squares leads to an ill-posed problem, i.e. the solution becomes not unique and unstable. In this case, a widely used approach for model estimation is regularization. Among regularization methods, the Lasso is of particular interest due to its ability to perform variable selection and thus obtain more interpretable models [43].

2.1 Formulation

Let X be the $m \times p$ design matrix where the data is arranged row-wise, i.e. $X = [x_1, \dots, x_m]^T$. Similarly, let $y = (y_1, \dots, y_m)^T$ be the m -dimensional response vector. Without loss of generality, we can assume $\alpha_0 = 0$ (e.g. by centering the training data such that each attribute has zero mean) [43]. The Lasso estimates the coefficients as the solution to the following problem

$$\min_{\alpha} f(\alpha) = \frac{1}{2} \|X\alpha - y\|_2^2 \quad \text{s.t.} \quad \|\alpha\|_1 \leq \delta , \quad (1)$$

where the ℓ_1 -norm constraint $\|\alpha\|_1 = \sum_i |\alpha_i| \leq \delta$ has the role of promoting sparsity in the regression coefficients. It is well-known that the constrained problem (1) is equivalent to the unconstrained problem

$$\min_{\alpha} \tilde{f}(\alpha) = \frac{1}{2} \|X\alpha - y\|_2^2 + \lambda \|\alpha\|_1 , \quad (2)$$

in the sense that given a solution α^* of (2) with a certain value for parameter $\bar{\lambda}$, it is possible to find a $\bar{\delta}$ such that α^* is also a solution of (1), and vice versa. Specifically, if one has an exact solution α^* of problem (2), corresponding to a given $\bar{\lambda}$, it is easy to see that α^* is also a solution of (1) for $\bar{\delta} = \|\alpha^*\|_1$. It is immediate to see that $\lambda = 0$ corresponds to the unconstrained solution α^R , i.e. to the plain least-squares regression, which can be obtained by setting $\delta > \|\alpha^R\|_1$ in (1). On the opposite, $\delta = 0$ in (1) corresponds to the null solution, which is obtained for large enough values of λ (specifically, for the case $p > m$, it can be shown that the solution of (2) is $\alpha^* = 0$ whenever $\lambda > \|X^T y\|_\infty$ [47]). Since the optimal tradeoff between the sparsity of the model and its predictive power is not known a-priori, practical applications of the Lasso require to find a solution and the profiles of estimated coefficients for a range of values of the regularization parameter δ (or λ in the penalized formulation). This is known in the literature as the *Lasso regularization path* [12].

2.2 Relevance and Applications

The Lasso is part of a powerful family of regularized linear regression methods for high-dimensional data analysis, which also includes ridge regression (RR) [19]

[18], the ElasticNet [53], and several recent extensions thereof [52, 54, 45]. From a statistical point of view, they can be viewed as methods for trading off the bias and variance of the coefficient estimates in order to find a model with better predictive performance. From a Machine Learning perspective, they allow to adaptively control the capacity of the model space in order to prevent overfitting. In contrast to RR, which is obtained by substituting the ℓ_1 norm in (2) by the squared ℓ_2 norm $\sum_i |\alpha_i|^2$, it is well-known that the Lasso does not only reduce the variance of coefficient estimates but is also able to perform variable selection by shrinking many of these coefficients to zero. Elastic-net regularization trades off ℓ_1 and ℓ_2 norms using a “mixed” penalty $\Omega(\alpha) = \gamma \|\alpha\|_1 + (1 - \gamma) \|\alpha\|_2^2$ which requires tuning the additional parameter γ [53]. ℓ_p norms with $p \in [0, 1)$ can enforce a more aggressive variable selection, but lead to computationally challenging non-convex optimization problems. For instance, $p = 0$, which corresponds to “direct” variable selection, leads to an NP-hard problem [49].

Thanks to its ability to perform variable selection and model estimation simultaneously, the Lasso is used in many fields involving high-dimensional data. These scenarios have become of increasing importance in the last decades since, as technologies for collecting and processing data evolve, classification and regression problems with a large number of candidate predictors have become ubiquitous. Advances in molecular technologies, for example, enable scientists to measure the status of thousands to millions of biomolecules simultaneously [16]. In text analysis, vector space models for representing documents easily leads to several thousands or even millions of document-term counts that correspond to potentially informative variables [25, 21]. Similarly, in the analysis of functional magnetic resonance imaging (fMRI) data, one can easily obtain datasets with millions of *voxels* representing the activity of particular portions of the brain [31]. In all these cases, the number dimensions or attributes p can far exceed the number of data instances m . It is also worth mentioning that popular data analysis tools such as ℓ_1 -regularized logistic regression and penalized Cox regression models can be implemented using iterative algorithms which solve Lasso problems at each iteration [12].

2.3 Related Work

Problem (1) is a quadratic programming problem with a convex constraint, which in principle may be solved using standard techniques such as interior-point methods, guaranteeing convergence in few iterations. However, the computational work required *per iteration* as well as the memory demanded by these approaches make them practical only for small and medium-sized problems. A faster specialized interior point method for the Lasso was proposed in [24], which however compares unfavorably with the baseline used in this paper [12].

One of the first efficient algorithms proposed in the literature for finding a solution of (2) is the Least Angle Regression (LARS) by Efron *et al.* [4]. As its main advantage, LARS allows to generate the entire Lasso regularization path with the same computational cost as standard least-squares via QR decomposition, i.e. $\mathcal{O}(mp^2)$, assuming $m < p$ [18]. At each step, it identifies the variable most correlated with the residuals of the model obtained so far and includes it in the set of “active” variables, moving the current iterate in a direction equiangular with the rest of the active predictors. It turns out that the algorithm we propose makes the same choice but updates the solution differently using

cheaper computations. A similar homotopy algorithm to calculate the regularization path has been proposed in [47], which differs slightly from LARS in the choice of the search direction.

More recently, it has been shown by Friedman *et al.* that a careful implementation of the Coordinate Descent method (CD) provides an extremely efficient solver [11], [12], [10], which also applies to more general models such as the ElasticNet proposed by Zou and Hastie [53]. In contrast to LARS, this method cyclically chooses one variable at a time and performs a simple analytical update. The full regularization path is built by defining a sensible range of values for the regularization parameter and taking the solution for a given value as warm-start for the next. This algorithm has been implemented into the Glmnet package and can be considered the current standard for solving this class of problems. Recent works have also advocated the use of Stochastic Coordinate Descent (SCD) [41], where the order of variable updates is chosen randomly instead of cyclically. This strategy can prevent the adverse effects caused by possibly unfavorable orderings of the coordinates, and allows to prove stronger theoretical guarantees compared to the plain CD [38].

Other methods for ℓ_1 -regularized regression may be considered. For instance, Zhou *et al.* recently proposed a geometrical approach where the Lasso is reformulated as a nearest point problem and solved using an algorithm inspired by the classical Wolfe method [51]. However, the popularity and proved efficiency of Glmnet on high-dimensional problems make it the chosen baseline in this work.

3 Frank-Wolfe Optimization

One of the earliest constrained optimization approaches [9, 50], the Frank-Wolfe algorithm has recently seen a sudden resurgence in interest from the Optimization and Machine Learning communities, due to its powerful theoretical properties and proved efficiency in the context of large-scale problems [3, 22, 17]. On the theoretical side, FW methods come with iteration complexity bounds that are independent of the number of variables in the problem, and sparsity guarantees that hold during the whole execution of the algorithm [3, 22]. In addition, several variants of the basic procedure have been analyzed, which can improve the convergence rate and practical performance of the basic FW iteration [15, 35, 26, 6]. From a practical point of view, they have emerged as efficient alternatives to traditional methods in several contexts, such as large-scale SVM classification [7, 8, 35, 6] and nuclear norm-regularized matrix recovery [22, 42]. In view of these developments, FW algorithms have come to be regarded as a suitable approach to large-scale optimization in various areas of Machine Learning, statistics, bioinformatics and related fields [1, 27].

Overall, though, the number of works showing experimental results for FW on practical applications is limited compared to that of the theoretical studies which have appeared in the literature. In the context of problems with ℓ_1 -regularization or sparsity constraints, the use of FW has been discussed in [40], but no experiments are provided. A closely related algorithm has been proposed in [51], however its implementation has a high computational cost in terms of time and memory requirements, and is not suitable for solving large-scale problems on a standard desktop or laptop machine. As such, the current

literature does not provide many examples of efficient FW-based software for large-scale Lasso or l_1 -regularized optimization. In this work, we aim to fill this gap by showing how a properly implemented stochastic FW method can best the current state of the art solvers on Lasso problems with a very large number of features.

3.1 The Standard Frank-Wolfe Algorithm

The FW algorithm is a general method to solve problems of the form

$$\min_{\alpha \in \Sigma} f(\alpha), \quad (3)$$

where $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is a convex differentiable function, and $\Sigma \subset \mathbb{R}^p$ is a compact convex set. Given an initial guess $\alpha^{(0)} \in \Sigma$, the standard FW method consists of the steps outlined in Algorithm 1. From an implementation point of view, a

Algorithm 1 The standard Frank-Wolfe algorithm

- 1: **Input:** an initial guess α^0 .
- 2: **for** $k = 0, 1, \dots$ **do**
- 3: Define a search direction $d^{(k)}$ by optimizing a linear model:

$$u^{(k)} \in \operatorname{argmin}_{u \in \Sigma} (u - \alpha^{(k)})^T \nabla f(\alpha^{(k)}), \quad d^{(k)} = u^{(k)} - \alpha^{(k)}. \quad (4)$$

- 4: Choose a stepsize $\lambda^{(k)}$, e.g. via line-search:

$$\lambda^{(k)} \in \operatorname{argmin}_{\lambda \in [0,1]} f(\alpha^{(k)} + \lambda d^{(k)}). \quad (5)$$

- 5: Update: $\alpha^{(k+1)} = \alpha^{(k)} + \lambda^{(k)} d^{(k)}$.
 - 6: **end for**
-

fundamental advantage of FW is that the most demanding part of the iteration, i.e. the solution of the linear subproblem (4), has a computationally convenient solution for several problems of practical interest, mainly due to the particular form of the feasible set. The key observation is that, when Σ is a polytope (e.g. the unit simplex for L_2 -SVMs [35], the ℓ_1 -ball of radius δ for the Lasso problem [1], a spectrahedron in nuclear norm for matrix recovery [14]), the search in step 3 can be reduced to a search among the vertices of Σ . This allows to devise cheap analytical formulas to find $u^{(k)}$, ensuring that each iteration has an overall cost of $\mathcal{O}(p)$. The fact that FW methods work with *projection-free* iterations is also a huge advantage on many practical problems, since a projection step to maintain the feasibility of the iterates (as needed by classical approaches such as proximal methods for Matrix Recovery problems) generally has a super-linear complexity, making the solution of large-scale problems difficult in practice [1].

Another distinctive feature of the algorithm is the fact that the solution at a given iteration K can be expressed as a convex combination of the vertices $u^{(k)}$, $k = 0, \dots, K-1$. Due to the incremental nature of the FW iteration, at most one new extreme point of Σ is discovered at each iteration, implying that at most k of such points are active at iteration k . Furthermore, this sparsity bound holds

for the entire run of the algorithm, effectively allowing to control the sparsity of the solution as it is being computed. This fact carries a particular relevance in the context of sparse approximation, and generally in all applications where it is crucial to find models with a small number of features. It also represents, as we will show in our experiments in Section 5, one of the major differences between incremental, forward approximation schemes and more general solvers for ℓ_1 -regularized optimization, which in general cannot guarantee to find sparse solutions along the regularization path.

3.2 Theoretical Properties

We summarize here some well-known theoretical results for the FW algorithm which are instrumental in understanding the behaviour of the method. We refer the reader to [22] for the proof of the following proposition. To prove the result, it is sufficient to assume that f has bounded curvature, which, as explained in [22], is roughly equivalent to the Lipschitz continuity of ∇f .

Proposition 1 (Sublinear convergence, [22]). *Let α^* be an optimal solution of problem (3). Then, for any $k \geq 1$, the iterates of Algorithm 1 satisfy*

$$f(\alpha^{(k)}) - f(\alpha^*) \leq \frac{4C_f}{k+2},$$

where C_f is the curvature constant of the objective function.

An immediate consequence of Proposition 1 is an upper bound on the iteration complexity: given a tolerance $\varepsilon > 0$, the FW algorithm finds an ε -approximate solution, i.e. an iterate $\alpha^{(k)}$ such that $f(\alpha^{(k)}) - f(\alpha^*) \leq \varepsilon$, after $\mathcal{O}(1/\varepsilon)$ iterations. Besides giving an estimate on the total number of iterations which has been shown experimentally to be quite tight in practice [5, 6], this fact tells us that the tradeoff between sparsity and accuracy can be partly controlled by appropriately setting the tolerance parameter. Recently, Garber and Hazan showed that under certain conditions the FW algorithm can obtain a convergence rate of $\mathcal{O}(1/k^2)$, comparable to that of first-order algorithms such as Nesterov’s method [13]. However, their results require strong convexity of the objective function and of the feasible set, a set of hypotheses which is not satisfied for several important ML problems such as the Lasso or the Matrix Recovery problem with trace norm regularization.

Another possibility is to employ a *Fully-Corrective* variant of the algorithm, where at each step the solution is updated by solving the problem restricted to the currently active vertices. The algorithm described in [51], where the authors solve the Lasso problem via a nearest point solver based on Wolfe’s method, operates with a very similar philosophy. A similar case can be made for the LARS algorithm of [4], which however updates the solution in a different way. The Fully-Corrective FW also bears a resemblance to the Orthogonal Matching Pursuit algorithms used in the Signal Processing literature [46], a similarity which has already been discussed in [3] and [22]. However, as mentioned in [3], the increase in computational cost is not paid off by a corresponding improvement in terms of complexity bounds. In fact, the work in [28] shows that the result in Proposition 1 cannot be improved for any first-order method based on solving linear subproblems without strengthening the assumptions. Greedy approximation techniques based on both the vanilla and the Fully-Corrective FW have

also been proposed in the context of approximate risk minimization with an ℓ_0 constraint by Shalev-Shwartz *et al.*, who proved several strong runtime bounds for the sparse approximations of arbitrary target solutions [40].

Finally, it is worth mentioning that the result of Proposition 1 can indeed be improved by using variants of FW that employ additional search directions, and allow under suitable hypotheses to obtain a linear convergence rate [35, 26]. It should be mentioned, however, that such rates only hold in the vicinity of the solution and that, as shown in [6], a large number of iterations might be required to gain substantial advantages. For this reason, we choose not to pursue this strategy in the present paper.

4 Randomized Frank-Wolfe for Lasso Problems

A specialized FW algorithm for problem (1) can be obtained straightforwardly by setting Σ equal to the ℓ_1 -ball of radius δ , hereafter denoted as \odot_δ . In this case, the vertices of the feasible set (i.e., the candidate points among which $u^{(k)}$ is selected in the FW iterations) are $\mathcal{V}(\odot_\delta) = \{\pm \delta \mathbf{e}_i : i = 1, 2, \dots, p\}$, where \mathbf{e}_i is the i -th element of the canonical basis. It is easy to see that the linear subproblem (4) in Algorithm 1 has a closed-form solution, given by:

$$\begin{aligned} u^{(k)} &= -\delta \operatorname{sign} \left(\nabla f(\alpha^{(k)})_{i_*^{(k)}} \right) \mathbf{e}_{i_*^{(k)}} \equiv \tilde{\delta}^{(k)} \mathbf{e}_{i_*^{(k)}}, \\ i_*^{(k)} &= \operatorname{argmax}_{i=1, \dots, p} \left| \nabla f(\alpha^{(k)})_i \right|. \end{aligned} \quad (6)$$

In order to efficiently execute the iteration, we can exploit the form of the objective function to obtain convenient expressions to update the function value and the gradient after each FW iteration. The gradient of $f(\cdot)$ in (1) is

$$\nabla f(\alpha) = -X^T (y - X\alpha) = -X^T y + X^T X \alpha.$$

There are two possible ways to compute $\nabla f(\alpha^{(k)})_i$ efficiently. One is to keep track of the vector of residuals $R^{(k)} = (y - X\alpha^{(k)}) \in \mathbb{R}^m$ and compute $\nabla f(\alpha^{(k)})_i$ as

$$\nabla f(\alpha^{(k)})_i = -z_i^T R^{(k)} = -z_i^T y + z_i^T X \alpha^{(k)}, \quad (7)$$

where $z_i \in \mathbb{R}^m$ is the i -th column of the design matrix X , i.e., the vector formed by collecting the i -th attribute of all the training points. We refer to this approach as the “method of residuals”. The other way is to expand the second line in (7)

$$\nabla f(\alpha^{(k)})_i = -z_i^T y + \sum_{j \neq 0} \alpha_j^{(k)} z_i^T z_j,$$

and keep track of the inner products $z_i^T z_j$ between z_i and the predictors z_j corresponding to non-zero coefficients of the current iterate. We call this the “method of active covariates”. The discussion in the next subsections reveals that the first approach is more convenient if, at each iteration, we only need to access a small subset of the coordinates of $\nabla f(\alpha^{(k)})$. It is clear from (6) that after computing $\nabla f(\alpha^{(k)})_i$ for $i = 1, \dots, p$ the solution to the linear subproblem in Algorithm 1 follows easily. The other quantities required by the algorithm

are the objective value (in order to monitor convergence) and the line search stepsize in (5), which can be obtained as

$$\begin{aligned} f(\alpha^{(k)}) &= \frac{1}{2}y^T y + \frac{1}{2}S^{(k)} - F^{(k)}, \\ \lambda^{(k)} &= \lambda_* := \frac{S^{(k)} - \tilde{\delta}\nabla f(\alpha^{(k)})_{i_*} - F^{(k)}}{S^{(k)} - 2\tilde{\delta}G_{i_*} + \tilde{\delta}^2 z_{i_*}^T z_{i_*}}, \end{aligned} \quad (8)$$

where $i_* = i_*^{(k)}$, $G_{i_*} = \nabla f(\alpha^{(k)})_{i_*} + z_{i_*}^T y$, and the terms $S^{(k)}, F^{(k)}$ can be updated recursively as

$$\begin{aligned} S^{(k+1)} &= (1 - \lambda_*)^2 S^{(k)} + 2\tilde{\delta}\lambda_*(1 - \lambda_*)G_{i_*} + \tilde{\delta}^2 \lambda_*^2 z_{i_*}^T z_{i_*} \\ F^{(k+1)} &= (1 - \lambda_*)F^{(k)} + \tilde{\delta}\lambda_* z_{i_*}^T y, \end{aligned}$$

with starting values $S^{(0)} = 0$ and $F^{(0)} = 0$. If we store the products $z_i^T y$ before the execution of the algorithm, the only non-trivial computation required here is $\nabla f(\alpha^{(k)})_{i_*}$ which was already computed to solve the subproblem in (6).

4.1 Randomized Frank-Wolfe Iterations

Although the FW method is generally very efficient for structured problems with a sparse solution, it also has a number of practical drawbacks. For example, it is well known that the total number of iterations required by a FW algorithm can be large, thus making the optimization prohibitive on very large problems. Even when (4) has an analytical solution due to the problem structure, the resulting complexity depends on the problem size [35], and can thus be impractical in cases where handling large-scale datasets is required. For example, in the specialization of the algorithm to problem (1), the main bottleneck is the computation of the FW vertex $i_*^{(k)}$ in (6) which corresponds to examining all the p candidate predictors and choosing the one most correlated with the current residuals (assuming the design matrix has been standardized s.t. the predictors have unit norm). Coincidentally, this is the same strategy underlying well-known methods for variable selection such as LARS and Forward Stepwise Regression (see Section 1).

A simple and effective way to avoid this dependence on p is to compute the FW vertex approximately, by limiting the search to a fixed number of extreme points on the boundary of the feasible set Σ [39, 5]. Specialized to the Lasso problem (1), this technique can be formalized as extracting a random sample $\mathcal{S} \subseteq \{1, \dots, p\}$ and solving

$$u^{(k)} = -\delta \operatorname{sign} \left(\nabla f(\alpha^{(k)})_{i_{\mathcal{S}}^{(k)}} \right) \mathbf{e}_{i_{\mathcal{S}}^{(k)}}, \quad \text{where } i_{\mathcal{S}}^{(k)} = \operatorname{argmax}_{i \in \mathcal{S}} \left| \nabla f(\alpha^{(k)})_i \right|. \quad (9)$$

Formally, one can think of a randomized FW iteration as the optimization of an approximate linear model, built by considering the partial gradient $\nabla f(\alpha^{(k)})_{|\mathcal{S}^{(k)}}$, i.e. the projection of the gradient onto the subspace identified by the sampled coordinates [48]. The number of coordinates of the gradient that need to be estimated with this scheme is $|\mathcal{S}|$ instead of p . If $|\mathcal{S}| \ll p$, this leads to a significant reduction in terms of computational overhead. Our stochastic specialization of the FW iteration for the Lasso problem takes thus the form of Algorithm 2

After selecting the variable $z_{i_*} \in \mathcal{S}$ best correlated with the current vector of residuals, the algorithm updates the current solution along the segment connecting $z_{i_*} \in \mathcal{S}$ with $\alpha^{(k)}$. Note how this approach differs from a method like LARS, where the direction to move the last iterate is equiangular to all the active variables. ^[1] It also differs from CD, which can make active more than one variable at each “epoch” or cycle through the predictors. The algorithm computes the stepsize by looking explicitly to the value of the objective, which can be computed analytically without increasing the cost of the iteration. Finally, the method updates the vector of residuals and proceeds to the next iteration.

Algorithm 2 Randomized Frank-Wolfe step for the Lasso problem

- 1: Choose the sampling set \mathcal{S} (see Section [4.5](#)).
- 2: Search for the predictor best correlated with the vector of residuals $R^{(k)} = (y - X\alpha^{(k)})$:

$$i_*^{(k)} = \operatorname{argmax}_{i \in \mathcal{S}} \left| \nabla f(\alpha^{(k)})_i \right| \equiv \left| z_i^T R^{(k)} \right| .$$

- 3: Set $\tilde{\delta}^{(k)} = -\delta \operatorname{sign}(\nabla f(\alpha^{(k)})_{i_*})$.
- 4: Compute the step-size $\lambda^{(k)}$ using [\(8\)](#).
- 5: Update the vector of coefficients as

$$\alpha^{(k+1)} = (1 - \lambda^{(k)})\alpha^{(k)} + \tilde{\delta}\lambda^{(k)}\mathbf{e}_{i_*^{(k)}} .$$

- 6: Update the vector of residuals $R^{(k)}$

$$R^{(k+1)} = (1 - \lambda^{(k)})R^{(k)} + \lambda^{(k)} \left(y - \tilde{\delta}z_{i_*^{(k)}} \right) . \quad (10)$$

Note that, although in this work we only discuss the basic Lasso problem, extending the proposed implementation to the more general ElasticNet model of [\[53\]](#) is straightforward. The derivation of the necessary analytical formulae is analogous to the one shown above. Furthermore, an extension of the algorithm to solve ℓ_1 -regularized logistic regression problems, another relevant tool in high-dimensional data analysis, can be easily obtained following the guidelines in [\[12\]](#).

4.2 Complexity and Implementation Details

In Algorithm [2](#) we compute the coordinates of the gradient using the method of residuals given by equation [\(7\)](#). Due to the randomization, this method becomes very advantageous with respect to the use of the alternative method based on the active covariates, even for very large p . Indeed, if we denote by s the cost of performing a dot product between a predictor z_i and another vector in \mathbb{R}^m , the overall cost of picking out the FW vertex in step 1 of our algorithm is $\mathcal{O}(s|\mathcal{S}|)$. Using the method of the active covariates would instead give an overall cost of $\mathcal{O}(s|\mathcal{S}|\|\alpha^{(k)}\|_0)$, which is always worse. Note however that this method may

¹As shown in [\[18\]](#), the LARS direction is $d_k = (X_{A_k}^T X_{A_k})^{-1} X_{A_k}^T R^{(k)}$ where X_{A_k} is the restriction of the design matrix to the active variables. The FW direction is just $d_k = e_{i_*} - \alpha^{(k)}$.

be better than the method of the residuals in a deterministic implementation by using caching tricks as proposed in [11], [12]. For instance, caching the dot products between all the predictors and the active ones and keeping updated all the coordinates of the gradient would cost $\mathcal{O}(p)$ except when new active variables appear in the solution, in which case the cost becomes $\mathcal{O}(ps)$. However, this would allow to find the FW vertex in $\mathcal{O}(p)$ operations. In this scenario, the fixed $\mathcal{O}(sp)$ cost of the method of residuals may be worse if the Lasso solution is very sparse. It is worth noting that the dot product cost s is proportional to the number of nonzero components in the predictors, which in typical high-dimensional problems is significantly lower than m .

In the current implementation, the term $\sigma_i := z_i^T y$ will be pre-computed for any $i = 1, 2, \dots, p$ before starting the iterations of the algorithm. This allows to write (7) as $-z_i^T R^{(k)} = -\sigma_i + z_i^T X \alpha^{(k)}$. Equation (10) for updating residuals can therefore be replaced by an equation to update $p^{(k)} = X \alpha^{(k)}$, eliminating the dependency on m .

4.3 Relation to SVM Algorithms and Sparsity Certificates

The previous implementation suggests that the FW algorithm will be particularly suited to the case $p \gg m$ where a regression problem has a very large number of features but not so many training points. It is interesting to compare this scenario to the situation in SVM problems: in the SVM case, the FW vertices correspond to training points, and the standard FW algorithm is able to quickly discover the relevant “atoms” (the Support Vectors), but has no particular advantage when handling lots of features. In contrast, in Lasso applications, where we are using the z_i ’s as training points, the situation is somewhat inverted: the algorithm should discover the critical features in at most $\mathcal{O}(1/\epsilon)$ iterations and guarantee that at most $\mathcal{O}(1/\epsilon)$ attributes will be used to perform predictions. This is, indeed, the scenario in which Lasso is used for several applications of practical interest, as problems with a very large number of features arise frequently in specific domains like bio-informatics, web and text analysis and sensor networks.

In the context of SVMs, the randomized FW algorithm has been already discussed in [5]. However, the results in the mentioned paper were experimental in nature, and did not contain a proof of convergence, which is instead provided in this work. Note that, although we have presented the randomized search for the specific case of problem (1), the technique applies more generally to the case where Σ is a polytope (or has a separable structure with every block being a polytope, as in [27]). We do not feel this hypothesis to be restrictive, as basically every practical application of the FW algorithm proposed in the literature falls indeed into this setting.

4.4 Convergence Analysis

We show that the stochastic FW converges (in the sense of expected value) with the same rate as the standard algorithm. First, we need the following technical result.

Lemma 1. *Let \mathcal{S} be picked at random from the set of all equiprobable κ -subsets of $\{1, \dots, p\}$, $1 \leq \kappa \leq p$, and let v be any vector in \mathbb{R}^p . Then*

$$\mathbb{E} \left[\left(\sum_{i \in \mathcal{S}} \mathbf{e}_i \mathbf{e}_i^T \right) v \right] = \frac{\kappa}{p} v.$$

Proof. Let $A_{\mathcal{S}} = \sum_{i \in \mathcal{S}} \mathbf{e}_i \mathbf{e}_i^T$ and $(A_{\mathcal{S}})_{ij}$ an element of $A_{\mathcal{S}}$. For $i \neq j$, $(A_{\mathcal{S}})_{ij} = 0$ and $\mathbb{E}[(A_{\mathcal{S}})_{ij}] = 0$. For $i = j$, $(A_{\mathcal{S}})_{ii}$ is a Bernoulli random variable with expectation $\frac{\kappa}{p}$. In fact, $(A_{\mathcal{S}})_{ii} = 1$ iff $i \in \mathcal{S}$; as there are $\binom{p-1}{\kappa-1}$ κ -subsets of $\{1, \dots, p\}$ containing i ,

$$\mathbb{P}(i \in \mathcal{S}) = \binom{p-1}{\kappa-1} \binom{p}{\kappa}^{-1} = \kappa/p = \mathbb{P}((A_{\mathcal{S}})_{ii} = 1) = \mathbb{E}[(A_{\mathcal{S}})_{ii}].$$

Therefore, for $i \in \{1, \dots, p\}$,

$$\mathbb{E}[(A_{\mathcal{S}})_{i*} v] = \mathbb{E} \left[\sum_{j=1}^p (A_{\mathcal{S}})_{ij} v_j \right] = \sum_{j=1}^p v_j \mathbb{E}[(A_{\mathcal{S}})_{ij}] = \frac{\kappa}{p} v_i.$$

□

Note that selecting a random subset \mathcal{S} of size κ to solve (9) is equivalent to (i) building a random matrix $A_{\mathcal{S}}$ as in Lemma 1 (ii) computing the restricted gradient $\tilde{\nabla}_{\mathcal{S}} f = \frac{p}{\kappa} A_{\mathcal{S}} \nabla f(\alpha^{(k)})$ and then (iii) solving the linear sub-problem (6) substituting $\nabla f(\alpha^{(k)})$ by $\tilde{\nabla}_{\mathcal{S}} f$. In other words, the proposed randomization can be viewed as approximating $\nabla f(\alpha^{(k)})$ by $\tilde{\nabla}_{\mathcal{S}} f$. Lemma 1 implies that $\mathbb{E}[\tilde{\nabla}_{\mathcal{S}} f] = \nabla f(\alpha^{(k)})$, which is the key to prove our main result.

Proposition 2. *Let α^* be an optimal solution of problem (3). Then, for any $k \geq 1$, the iterates of Algorithm 1 with the randomized search rule satisfy*

$$\mathbb{E}_{\mathcal{S}^{(k)}}[f(\alpha^{(k)})] - f(\alpha^*) \leq \frac{4\tilde{C}_f}{k+2},$$

where $\mathbb{E}_{\mathcal{S}^{(k)}}$ denotes the expectation with respect to the k -th random sampling.

This result has a similar flavor to that in [27], and the analysis is similar to the one presented in [48]. However, in contrast to the above works, we do not assume any structure in the optimization problem or in the sampling. A detailed proof can be found in the Appendix. As in the deterministic case, Proposition 2 implies a complexity bound of $\mathcal{O}(1/\varepsilon)$ iterations to reach an approximate solution $\alpha^{(k)}$ such that $\mathbb{E}_{\mathcal{S}^{(k)}}[f(\alpha^{(k)})] - f(\alpha^*) \leq \varepsilon$.

4.5 Choosing the Sampling Size

When using a randomized FW iteration it is important to choose the sampling size in a sensible way. Indeed, some recent works showed how this choice entails a tradeoff between accuracy (in the sense of premature stopping) and complexity, and henceforth CPU time [5]. This kind of approximation is motivated by the following result, which suggests that it is reasonable to pick $|\mathcal{S}| \ll p$.

Theorem 1 ([39], Theorem 6.33). *Let $\mathcal{D} \subset \mathbb{R}$ s.t. $|\mathcal{D}| = p$ and let $\mathcal{D}' \subset \mathcal{D}$ be a random subset of size κ . Then, the probability that the largest element in \mathcal{D}' is greater than or equal to \tilde{p} elements of \mathcal{D} is at least $1 - (\frac{\tilde{p}}{p})^\kappa$.*

The value of this result lies in the ability to obtain probabilistic bounds for the quality of the sampling *independently of the problem size p* . For example, in the case of the Lasso problem, where $\mathcal{D} = \{|\nabla f(\alpha^{(k)})_1|, \dots, |\nabla f(\alpha^{(k)})_p|\}$ and $\mathcal{D}' = \{|\nabla f(\alpha^{(k)})_i| \text{ s.t. } i \in \mathcal{S}\}$, it is easy to see that it suffices to take $|\mathcal{S}| \approx 194$ to guarantee that, with probability at least 0.98, $|\nabla f(\alpha^{(k)})_{i_{\mathcal{S}}}|$ lies between the 2% largest gradient components (in absolute value), independently of p . This kind of sampling has been discussed for example in [6].

The issue with this strategy is that, for problems with very sparse solutions (which is the case for strong levels of regularization), even a large confidence interval does not guarantee that the algorithm can sample a good vertex in most of the iterations. Intuitively, the sampling strategy should allow the algorithm to detect the set of vertices active at the optimum, which correspond, at various stages of the optimization process, to descent directions for the objective function. In sparse approximation problems, extracting a sampling set without relevant features may imply adding “spurious” components to the solution, reducing the sparseness of the model we want to find.

A more effective strategy in this context would be ask for a certain probability that the sampling will include at least one of the “optimal” features. Letting \mathcal{S}^* be the index set of the active vertices at the optimum, and denoting $s = |\mathcal{S}^*|$ and $\kappa = |\mathcal{S}|$, we have

$$P(\mathcal{S}^* \cap \mathcal{S} = \emptyset) = \prod_{j=0}^{\kappa-1} \left(1 - \frac{s}{p-j}\right) \leq \left(1 - \frac{s}{p}\right)^\kappa, \quad (11)$$

with the latter inequality being a reasonable approximation if $\kappa \ll p$. From (11), we can guarantee that $\mathcal{S}^* \cap \mathcal{S} \neq \emptyset$ with probability at least ρ by imposing:

$$\left(1 - \frac{s}{p}\right)^\kappa \leq (1 - \rho) \Leftrightarrow \kappa \geq \frac{\ln(1 - \rho)}{\ln\left(1 - \frac{s}{p}\right)} = \frac{\ln(\text{confidence})}{\ln(\text{sparseness})}. \quad (12)$$

On the practical side, this sampling strategy often implies taking a larger κ . Assuming that the fraction of relevant features (s/p) is constant, we essentially get the bound for κ provided by Theorem 1, which is independent of p . However, for the worst case $s/p \rightarrow 0$, we get

$$\frac{\ln(1 - \rho)}{\ln\left(1 - \frac{s}{p}\right)} \approx \left(\frac{-\ln(1 - \rho)}{s}\right) p, \quad (13)$$

which suggests to use a sampling size proportional to p . A more involved strategy, which exploits the incremental structure of the FW algorithm, is using a large κ at early iterations and smaller values of κ as the solution gets more dense. If the optimal solution is very sparse, the algorithm requires few expensive iterations to converge. In contrast, if the solution is dense, the algorithm requires more, but faster, iterations (e.g. for a confidence $1 - \rho = 0.98$ and $s/p = 0.02$ the already mentioned $\kappa = 194$ suffices). For the problems considered in this paper, setting κ to a small fraction of p works well in practice, as shown by the experiments in the next Section.

Dataset	m	t	p
Synthetic-10000	200	200	10,000
Synthetic-50000	200	200	50,000
Pyrim	74	--	201,376
Triazines	186	--	635,376
E2006-tfidf	16,087	3,308	150,360
E2006-log1p	16,087	3,308	4,272,227

Table 1: List of the benchmark datasets used in the experiments.

5 Numerical Experiments

In this section, we assess the practical effectiveness of the randomized FW algorithm by performing experiments on both synthetic datasets and real-world benchmark problems with hundreds of thousands or even millions of features. The characteristics of the datasets are summarized in Table 1, where we denote by m the number of training examples, by t the number of test examples, and by p the number of features. The synthetic datasets were generated with the Scikit-learn function `make_regression`. Each of them comes in two versions corresponding to a different number of relevant features in the true linear model used to generate the data (32 and 100 features for the problem of size $p = 10000$, and 158 and 500 features for that of size $p = 50000$). The real large-scale regression problems **E2006-tfidf** and **E2006-log1p**, and the datasets **Pyrim** and **Triazines**, are available from [2].

In assessing the performance of our method, we used as a baseline the following algorithms, which in our view, and according to the properties summarized in Table 2, can be considered among the most competitive solvers for Lasso problems:

- The well-known CD algorithm by Friedman *et al.*, as implemented in the `Glmnet` package [12]. This method is highly optimized for the Lasso and is considered a state of the art solver in this context.
- The SCD algorithm as described in [41], which is significant both for being a stochastic method and for having better theoretical guarantees than the standard cyclic CD.
- The Accelerated Gradient Descent with projections for both the regularized and the constrained Lasso, as this algorithm guarantees an optimal complexity bound. We choose as a reference the implementation in the `SLEP` package by Liu *et al.* [32].

Among other possible first-order methods, the classical SGD suffers from a worse convergence rate, and its variant SMIDAS has a complexity bound which depends on p , thus we did not include them in our experiments. Indeed, the authors of [40] conclude that SCD is both faster and produces significantly sparser models compared to SGD. Finally, although the recent GeoLasso algorithm of [51] is interesting because of its close relationship with FW, its running time and memory requirements are clearly higher compared to the above solvers.

Since an appropriate level of regularization needs to be automatically selected in practice, the algorithms are compared by computing the entire regu-

Approach	Form	Number of Iterations	Complexity per Iteration	Sparse Its.
Accelerated Gradient + Proj. [33]	(1)	$\mathcal{O}(1/\sqrt{\epsilon})$	$\mathcal{O}(mp + p)\dagger_1$	No
Accelerated Gradient + Reg. Proj. [34]	(2)	$\mathcal{O}(1/\sqrt{\epsilon})$	$\mathcal{O}(mp + p)\dagger_1$	No
Cyclic Coordinate Descent (CD) [11] [12]	(2)	Unknown	$\mathcal{O}(mp)\dagger_2$	Yes
Stochastic Gradient Descent (SGD) [29]	(2)	$\mathcal{O}(1/\epsilon^2)$	$\mathcal{O}(p)$	No
Stochastic Mirror Descent [41]	(2)	$\mathcal{O}(\log(p)/\epsilon^2)$	$\mathcal{O}(p)$	No
GeoLasso [51]	(1)	$\mathcal{O}(1/\epsilon)$	$\mathcal{O}(mp + a^2)$	Yes
Frank-Wolfe (FW) [22]	(1)	$\mathcal{O}(1/\epsilon)$	$\mathcal{O}(mp)$	Yes
Stochastic Coord. Descent (SCD) [38]	(2)	$\mathcal{O}(p/\epsilon)$	$\mathcal{O}(m)\dagger_3$	Yes
Stochastic Frank-Wolfe	(1)	$\mathcal{O}(1/\epsilon)$	$\mathcal{O}(m S)$	Yes

Table 2: Methods proposed for scaling the Lasso and their complexities. Here, a denotes the number of active features at a given iteration, which in the worst case is $a = \text{rank}(X) \leq \min(m, p)$. A method is said to have *sparse iterations* if a non trivial bound for the number of non-zero entries of each iterate holds at any moment. \dagger_1 $\mathcal{O}(p)$ is required for the projections. \dagger_2 An iteration of cyclic coordinate descent corresponds to a complete cycle through the features. \dagger_3 An iteration of SCD corresponds to the optimization on a single feature.

larization path on a range of values of the regularization parameters λ and δ (depending on whether the method solves the penalized or the constrained formulation). Specifically, we first estimate two intervals $[\lambda_{\min}, \lambda_{\max}]$ and $[\delta_{\min}, \delta_{\max}]$, and then solve problems [\(2\)](#) and [\(1\)](#) on a 100-point parameter grid in logarithmic scale. For the penalized Lasso, we use $\lambda_{\min} = \lambda_{\max}/100$, where λ_{\max} is determined as in the Glmnet code. Then, to make the comparison fair (i.e. to ensure that all the methods solve the same problems according to the equivalence in Section [2](#)), we choose for the constrained Lasso $\delta_{\max} = \|\alpha_{\min}\|_1$ and $\delta_{\min} = \delta_{\max}/100$, where α_{\min} is the solution obtained by Glmnet with the smallest regularization parameter λ_{\min} and a high precision ($\epsilon = 10^{-8}$). The idea is to give the solvers the same “sparsity budget” to find a solution of the regression problem.

Warm-start strategy

As usually done in these cases, and for all the methods, we compute the path using a warm-start strategy where each solution is used as an initial guess for the optimization with the next value of the parameter. Note that we always start from the most sparse solution. This means that in the cases of CD, SCD and regularized SLEP we start from λ_{\max} towards λ_{\min} , while for FW and constrained SLEP we go from δ_{\min} towards δ_{\max} . Furthermore, since $\delta < \|\alpha^R\|_1$, where α^R is the unconstrained solution, we know that the solution will lie on the boundary, therefore we adopt a heuristic strategy whereby the previous solution is scaled so that its ℓ_1 -norm is δ . Both algorithms are initialized with the null solution as the initial guess. Regarding the stopping criterion for each

problem in the path, we stop the current run when $\|\alpha^{(k+1)} - \alpha^{(k)}\|_\infty \leq \varepsilon$ for all the algorithms. Other choices are possible (for example, FW methods are usually stopped using a duality gap-based criterion [22]), but this is the condition used by default to stop CD in Glmnet. A value of $\varepsilon = 10^{-3}$ is used in the following experiments. All the considered algorithms have been coded in C++. The code and datasets used in this paper are available from public repositories on Github (<https://github.com/efrandi/FW-Lasso>) and Dataverse (<https://goo.gl/PTQ05R>), respectively. The SLEP package has a Matlab interface, but the key functions are all coded in C. Overall, we believe our comparison can be considered very fair. We executed the experiments on a 3.40GHz Intel i7 machine with 16GB of main memory running CentOS. For the randomized experiments, results were averaged over 10 runs.

5.1 “Sanity Check” on the Synthetic Datasets

The aim of these experiments is not to measure the performance of the algorithms (which will be assessed below on four larger, real-life datasets), but rather to compare their ability to capture the evolution of the most relevant features of the models, and discuss how this relates to their behaviour from an optimization point of view. To do this, we monitor the values of the 10 most relevant features along the path, as computed by both the CD and FW algorithms, and plot the results in Figures 1 and 2. To determine the relevant variables, we use as a reference the regularization path obtained by Glmnet with $\varepsilon = 10^{-8}$ (which is assumed to be a reasonable approximation of the exact solution), and compute the 10 variables having, on average, the highest absolute value along the path. As this experiment is intended mainly as a sanity check to verify that our solver reconstructs the solution correctly, we do not include other algorithms in the comparison. In order to implement the random sampling strategy in the FW algorithm, we first calculated the average number of active features along the path, rounded up to the nearest integer, as an empirical estimate of the sparsity level. Then, we chose $|\mathcal{S}|$ based on the probability ρ of capturing at least one of the relevant features at each sampling, according to (13). A confidence level of 99% was used for this experiment, leading to sampling sizes of 372 and 324 points for the two problems of size 10000, and of 1616 and 1572 points for those of size 50000.

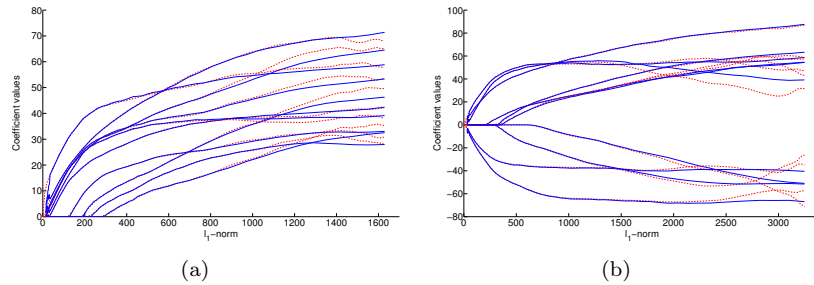


Figure 1: Growth of the 10 most significant features on the synthetic problem of size 10000, with 32 (a) and 100 (b) relevant features. Results for CD are in red dashed lines, and in blue continuous lines for FW.

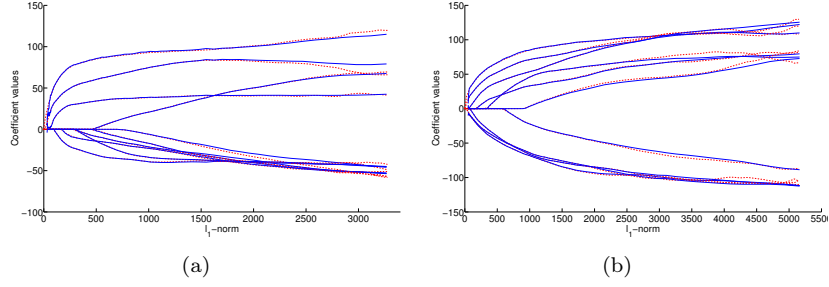


Figure 2: Growth of the 10 most significant features on the synthetic problems of size 50000, with 158 (a) and 500 (b) relevant features. Results for CD are in red dashed lines, and in blue continuous lines for FW.

Figure 3 depicts, for two of the considered problems, the prediction error on the test set along the path found by both algorithms. It can be seen that both methods are able to find the same value of the best prediction error (i.e. to correctly identify the best model). The FW algorithm also seems to obtain slightly better results towards the end of the path, consistently with the fact that the coefficients of the most relevant variables tend to be more stable, compared with CD, when the regularization is weaker.

5.2 Results on Large-scale Datasets

In this section, we report the results on the problems **Pyrim**, **Triazines**, **E2006-tfidf** and **E2006-log1p**. These datasets represent actual large-scale problems of practical interest. The **Pyrim** and **Triazines** datasets stem from two quantitative structure-activity relationship models (QSAR) problems, where the biological responses of a set of molecules are measured and statistically related to the molecular structure on their surface. We expanded the original feature set by means of product features, i.e. modeling the response variable y as a linear combination of polynomial basis functions, as suggested in [20]. For

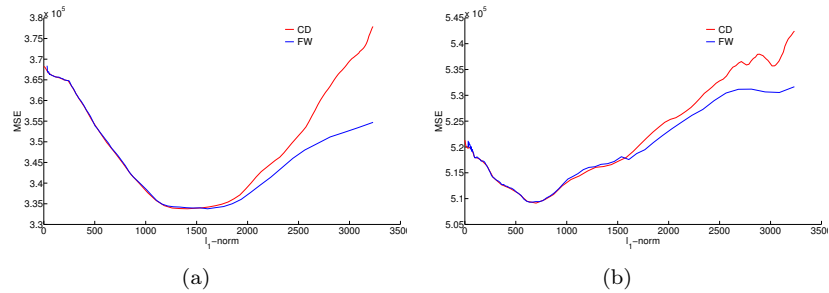


Figure 3: Test error (ℓ_1 -norm vs. MSE) for problems **Synthetic-10000** (100 relevant features) **Synthetic-50000** (158 relevant features). Results for CD are in red, and in blue for FW.

this experiment, we used product features of order 5 and 4 respectively, which leads to large-scale regression problems with $p = 201,376$ and $p = 635,376$. Problems **E2006-tfidf** and **E2006-log1p** stem instead from the real-life NLP task of predicting the risk of investment (i.e. the stock return volatility) in a company based on available financial reports [25].

To implement the random sampling for the FW algorithm, we cannot use the technique described above for the experiments on the synthetic datasets, as in real problems we do not have an *a-priori* estimate of the sparsity level. We therefore adopt a simpler strategy, and set $|\mathcal{S}|$ to a fixed, small fraction of the total number of features. Our choices are summarized in Table 3.

% of p	Pyrim	Triazines	E2006-tfidf	E2006-log1p
1%	2,014	6,354	1,504	42,723
2%	4,028	12,708	3,008	85,445
3%	6,042	19,062	4,511	128,167

Table 3: Sizes of the sampling set $|\mathcal{S}|$ for the large-scale datasets.

As a measure of the performance of the considered algorithms, we report the CPU time in seconds, the total number of iterations and the number of requested dot products (which account for most of the required running time for all the algorithms²) along the entire regularization path. Note that, in assessing the number of iterations, we consider one complete cycle of CD to be roughly equivalent to a full deterministic iteration of FW (since in both cases the complexity is determined by a full pass through all the coordinates) and to p random coordinate explorations in SCD. Finally, in order to evaluate the sparsity of the solutions, we report the average number of active features along the path. Results are displayed in Tables 4 and 5. In the second table, the speed-ups with respect to the CD algorithm are also reported. It can be seen how for all the choices of the sampling size the FW algorithm allows for a substantial improvement in computational performance, as confirmed by both the CPU times and the machine-independent number of requested dot products (which are roughly proportional to the running times). The plain SCD algorithm performs somewhat worse than CD, something we attribute mainly to the fact that the Glmnet implementation of CD is a highly optimized one, using a number of *ad-hoc* tricks tailored to the Lasso problem that we decided to preserve in our comparison. If we used a plain implementation of CD, we would expect to obtain results very close to those exhibited by SCD.

Furthermore, FW is always able to find the sparsest solution among the considered methods. The extremely large gap in average sparsity between FW and CD on one side, and the SLEP solvers on the other, is due to the fact that the latter compute in general dense iterates. Although the Accelerated Gradient Descent solver is fast and competitive from an optimization point of view, providing always the lower number of iterations as predicted by the theory, it is not able to keep the solutions sparse along the path. This behavior clearly highlights the advantage of using incremental approximations in the context of sparse recovery and feature selection. Importantly, note that the small number

²Note that the SLEP code uses highly optimized libraries for matrix multiplication, therefore matrix-vector computations can be faster than naive C++ implementations.

of features found by FW is not a result of the randomization technique: it is robust with respect to the sampling size, and additional experiments performed using a deterministic FW solver revealed that the average number of nonzero entries in the solution does not change even if the randomization is completely removed.

	CD	SCD	SLEP Reg.	SLEP Const.
Pyrim				
Time (s)	6.22e + 00	1.59e + 01	5.43e + 00	6.86e + 00
Iterations	2.54e + 02	1.44e + 02	1.00e + 02	1.12e + 02
Dot products	2.08e + 07	2.90e + 07	8.56e + 07	1.29e + 08
Active features	68.4	116.6	3,349	13,030
Triazines				
Time (s)	2.75e + 01	8.42e + 01	4.27e + 01	5.93e + 01
Iterations	2.62e + 02	1.59e + 02	1.01e + 02	1.11e + 02
Dot products	6.80e + 07	1.01e + 08	2.87e + 08	4.67e + 08
Active features	150.0	330.8	29,104	130,303
E2006-tfidf				
Time (s)	9.10e + 00	2.19e + 01	1.24e + 01	2.27e + 01
Iterations	3.48e + 02	2.01e + 02	1.06e + 02	2.50e + 02
Dot products	2.04e + 07	3.03e + 07	5.97e + 07	1.37e + 08
Active features	149.5	275.3	444.8	724.3
E2006-log1p				
Time (s)	1.60e + 02	4.92e + 02	1.00e + 02	1.42e + 02
Iterations	3.55e + 02	1.99e + 02	1.11e + 02	1.18e + 02
Dot products	5.73e + 08	8.50e + 08	1.78e + 09	2.85e + 09
Active features	281.3	1158.2	12,806	54,704

Table 4: Results for the baseline solvers on the large-scale problems **Pyrim**, **Triazines**, **E2006-tfidf** and **E2006-log1p**.

To better assess the effect of using an incremental algorithm in obtaining a sparse model, we plot in Figure 4 the evolution of the number of active features along the path on problems **E2006-tfidf** and **E2006-log1p**. It can be clearly seen how CD and FW (with the latter performing the best overall) are able to recover sparser solutions, and can do so without losing accuracy in the model, as we show below.

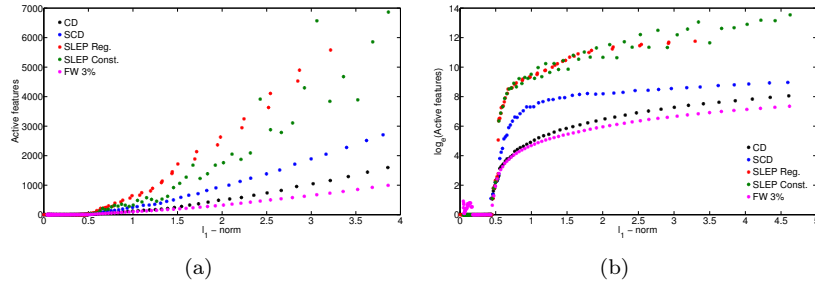


Figure 4: Sparsity patterns (ℓ_1 -norm vs. active coordinates) for problems **E2006-tfidf** (a) and **E2006-log1p** (b). The latter is plotted in a natural logarithmic scale due to the high number of features found by the SLEP solvers.

	FW 1%	FW 2%	FW 3%
Pyrim			
Time (s)	2.28e − 01	4.47e − 01	6.60e − 01
Speed-up	27.3×	13.9×	9.4×
Iterations	2.77e + 02	2.80e + 02	2.77e + 02
DotProd	7.61e + 05	1.53e + 06	2.28e + 06
Active features	27.6	28.1	27.9
Triazines			
Time (s)	2.61e + 00	5.31e + 00	8.19e + 00
Speed-up	10.5×	5.2×	3.4×
Iterations	7.15e + 02	7.29e + 02	7.43e + 02
DotProd	5.18e + 06	1.06e + 07	1.61e + 07
Active features	120.6	117.5	118.7
E2006-tfidf			
Time (s)	8.83e − 01	1.76e + 00	2.74e + 00
Speed-up	10.3×	5.2×	3.3×
Iterations	1.27e + 03	1.35e + 03	1.41e + 03
DotProd	1.97e + 06	4.35e + 06	6.84e + 06
Active features	123.7	125.8	127.1
E2006-log1p			
Time (s)	1.93e + 01	4.14e + 01	6.59e + 01
Speed-up	8.3×	3.9×	2.4×
Iterations	1.75e + 03	1.91e + 03	1.99e + 03
DotProd	7.90e + 07	1.71e + 08	2.68e + 08
Active features	196.9	199.8	203.7

Table 5: Performance metrics for stochastic FW on the large-scale problems **Pyrim**, **Triazines**, **E2006-tfidf** and **E2006-log1p**.

In order to evaluate the accuracy of the obtained models, we plot in Figures 5 and 6 the mean square error (MSE) against the ℓ_1 -norm of the solution along the regularization path, computed both on the original training set (curves 5(a-c) and 6(a-c)) and on the test set (curves 5(b-d) and 6(b-d)). Note that the value of the objective function in Problem (1) coincides with the mean squared error (MSE) on the training set. First of all, we can see how the decrease in the objective value is basically identical in all cases, which indicates that with our sampling choices the use of a randomized algorithm does not affect the optimization accuracy. Second, the test error curves show that the predictive capability of all the FW models is competitive with that of the models found by the CD algorithm (particularly in the case of the larger problem **E2006-log1p**). It is also important to note that in all cases the best model, corresponding to the minimum of the test error curves, is found for a relatively low value of the constraint parameter, indicating that sparse solutions are preferable and that solutions involving more variables tend to cause overfitting, which is yet another incentive to use algorithms that can naturally induce sparsity. Again, it can be seen how the minima of all the curves coincide, indicating that all the algorithms are able to correctly identify the best compromise between sparsity and training error. The fact that we are able to attain the same models obtained by a state of the art algorithm such as Glmnet using a sampling size as small as 3% of the total number of features is particularly noteworthy. Combined with the consistent advantages in CPU time over other competing solvers and its attractive sparsity properties, it shows how the randomized FW represents a solid, high-performance option for solving high-dimensional Lasso problems.

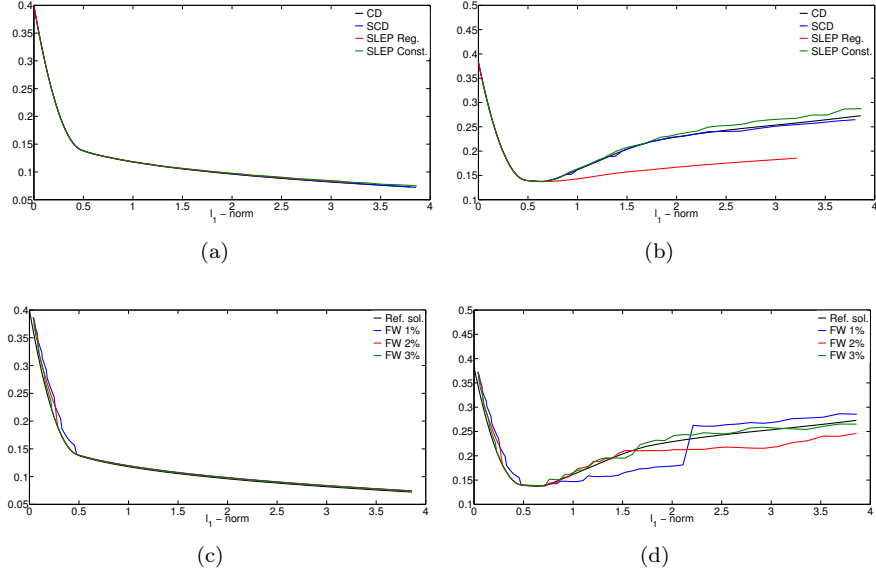


Figure 5: Error curves (ℓ_1 -norm vs. MSE) for problem **E2006-tfd**: on top, training error (a) and test error (b) for CD, SCD and SLEP; on bottom, training error (c) and test error (d) for stochastic FW.

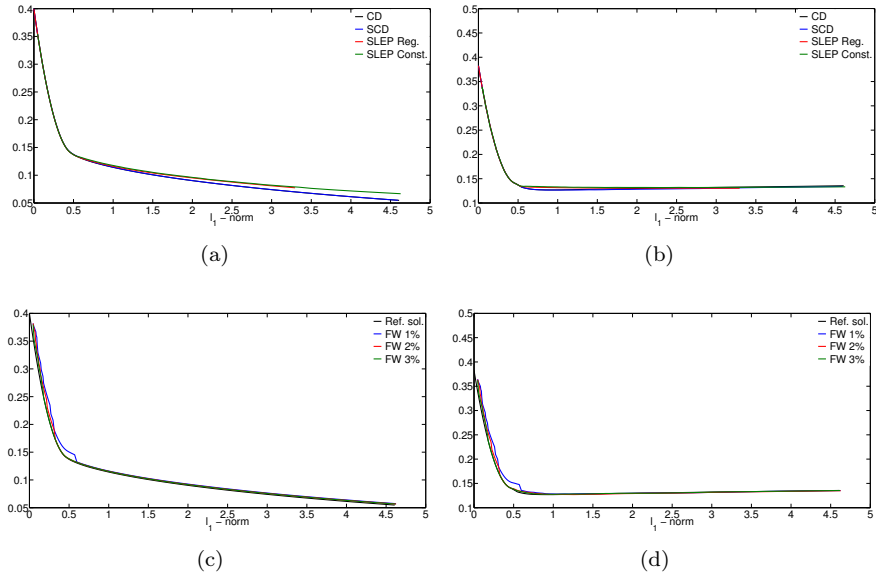


Figure 6: Error curves (ℓ_1 -norm vs. MSE) for problem **E2006-log1p**: on top, training error (a) and test error (b) for CD, SCD and SLEP; on bottom, training error (c) and test error (d) for stochastic FW.

6 Conclusions and Perspectives

In this paper, we have studied the practical advantages of using a randomized Frank-Wolfe algorithm to solve the constrained formulation of the Lasso regression problem on high-dimensional datasets involving a number of variables ranging from the hundred thousands to a few millions. We have presented a theoretical proof of convergence based on the expected value of the objective function. Our experiments show that we are able to obtain results that outperform those of other state-of-the-art solvers such as the Glmnet algorithm, a standard among practitioners, without sacrificing the accuracy of the model in a significant way. Importantly, our solutions are consistently more sparse than those found using several popular first-order methods, demonstrating the advantage of using an incremental, greedy optimization scheme in this context.

In a future work, we intend to address the issue of whether it is possible to find suitable sampling conditions which can lead to a stronger stochastic convergence result, i.e. to certifiable probability bounds for approximate solutions. Finally, we remark that the proposed approach can be readily extended to other similar problems such as ElasticNet or more general l_1 -regularized problems such as logistic regression, or to related applications such as the sparsification of SVM models. Another possibility to tackle various Lasso formulations is to exploit an equivalent formulation in terms of SVMs, an area where FW methods have already shown promising results. Together, all these elements strengthen the conclusions of our previous research work, showing that FW algorithms can provide a complete and flexible framework to efficiently solve a wide variety of large-scale Machine Learning and Data Mining problems.

Acknowledgments

The research leading to these results has received funding from the European Research Council under the European Union's Seventh Framework Programme (FP7/2007-2013) / ERC AdG A-DATADRIVE-B (290923). This paper reflects only the authors' views and the Union is not liable for any use that may be made of the contained information. Research Council KUL: GOA/10/09 MaNet, CoE PFV/10/002 (OPTEC), BIL12/11T; Flemish Government: FWO: projects: G.0377.12 (Structured systems), G.088114N (Tensor based data similarity); PhD/Postdoc grants; iMinds Medical Information Technologies SBO 2014; IWT: POM II SBO 100031; Belgian Federal Science Policy Office: IUAP P7/19 (DYSCO, Dynamical systems, control and optimization, 2012-2017). The second author received funding from CONICYT Chile through FONDECYT Project 11130122. The first author thanks the colleagues from the Department of Computer Science and Engineering, University of Bologna, for their hospitality during the period in which this research was conceived.

References

- [1] Andreas Argyriou, Marco Signoretto, and Johan A. K. Suykens. Hybrid algorithms with applications to sparse and low rank regularization. In Johan A. K. Suykens, Marco Signoretto, and Andreas Argyriou, editors, *Regular-*

- ization, Optimization, Kernels, and Support Vector Machines, chapter 3, pages 53–82. Chapman & Hall/CRC, 2014.
- [2] Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: a library for support vector machines*, <http://www.csie.ntu.edu.tw/~cjlin/libsvm>, 2011.
 - [3] Kenneth Clarkson. Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm. *ACM Transactions on Algorithms*, 6(4):63:1–63:30, 2010.
 - [4] Bradley Efron, Trevor Hastie, Iain Johnstone, and Robert Tibshirani. Least angle regression. *Ann. Stat.*, 32(2):407–499, 2004.
 - [5] E. Frandi, R. Nanculef, and J. A. K. Suykens. Complexity issues and randomization strategies in Frank-Wolfe algorithms for machine learning. In *7th NIPS Workshop on Optimization for Machine Learning*, 2014.
 - [6] E. Frandi, R. Nanculef, and J. A. K. Suykens. A PARTAN-accelerated Frank-Wolfe algorithm for large-scale SVM classification. In *Proceedings of the IJCNN 2015*, 2015.
 - [7] Emanuele Frandi, Maria Grazia Gasparo, Stefano Lodi, Ricardo Nanculef, and Claudio Sartori. A new algorithm for training SVMs using approximate minimal enclosing balls. In *Proceedings of the 15th Iberoamerican Congress on Pattern Recognition*, pages 87–95. Springer, 2010.
 - [8] Emanuele Frandi, Maria Grazia Gasparo, Stefano Lodi, Ricardo Nanculef, and Claudio Sartori. Training support vector machines using Frank-Wolfe methods. *IJPRAI*, 27(3), 2011.
 - [9] Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 1:95–110, 1956.
 - [10] Jerome Friedman. Fast sparse regression and classification. *Int. J. Forecast.*, 28:722–738, 2012.
 - [11] Jerome Friedman, Trevor Hastie, Holger Höfling, and Robert Tibshirani. Pathwise coordinate optimization. *Ann. Appl. Stat.*, 1(2):302–332, 2007.
 - [12] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.*, 33(1):1–22, 2010.
 - [13] Dan Garber and Elan Hazan. Faster rates for the Frank-Wolfe method over strongly-convex sets. In *Proceedings of the 32nd ICML*, 2015.
 - [14] Joachim Giesen, Martin Jaggi, and Søren Laue. Optimizing over the growing spectrahedron. In *20th Annual European Symposium on Algorithms*, LNCS, pages 503–514. Springer, 2012.
 - [15] Jacques Guélat and Patrice Marcotte. Some comments on Wolfe’s “away step”. *Math. Prog.*, 35:110–119, 1986.
 - [16] Jiang Gui and Hongzhe Li. Penalized Cox regression analysis in the high-dimensional and low-sample size settings, with applications to microarray gene expression data. *Bioinformatics*, 21(13):3001–3008, 2005.

- [17] Zaid Harchaoui, Anatoli Juditski, and Arkadi Nemirovski. Conditional gradient algorithms for norm-regularized smooth convex optimization. *Math. Prog. Ser. A*, 13(1):1–38, 2014.
- [18] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer New York Inc., 2001.
- [19] Arthur E Hoerl and Robert W Kennard. Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, 1970.
- [20] L. Huang, J. Jia, B. Yu, B. G. Chun, P. Maniatis, and M. Naik. Predicting execution time of computer programs using sparse polynomial regression. In *Advances in Neural Information Processing Systems*, pages 883–891, 2010.
- [21] Georgiana Ifrim, Gökhan Bakir, and Gerhard Weikum. Fast logistic regression for text categorization with variable-length n-grams. In *Proceedings of the 14th ACM SIGKDD*, pages 354–362, 2008.
- [22] Martin Jaggi. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- [23] Martin Jaggi. An equivalence between the lasso and support vector machines. In Johan A. K. Suykens, Marco Signoretto, and Andreas Argyriou, editors, *Regularization, Optimization, Kernels, and Support Vector Machines*, chapter 1, pages 1–26. Chapman & Hall/CRC, 2014.
- [24] Seung-Jean Kim, Kwangmoo Koh, Michael Lustig, Stephen Boyd, and Dmitry Gorinevsky. An interior-point method for large-scale l_1 -regularized least squares. *IEEE J. Sel. Top. Sign. Proces.*, 1(4):606–617, 2007.
- [25] Shimon Kogan, Dimitry Levin, Bryan R. Routledge, Jacob S. Sagi, and Noah A. Smith. Predicting risk from financial reports with regression. In *Proceedings of the NAACL '09*, pages 272–280, 2009.
- [26] Simon Lacoste-Julien and Martin Jaggi. An affine invariant linear convergence analysis for Frank-Wolfe algorithms. *arXiv:1312.7864v2*, 2014.
- [27] Simon Lacoste-Julien, Martin Jaggi, Mark Schmidt, and Patrick Pletscher. Block-coordinate Frank-Wolfe optimization for structural SVMs. In *Proceedings of the 30th International Conference on Machine Learning*, 2013.
- [28] Guanghui Lan. The complexity of large-scale convex programming under a linear optimization oracle. *arXiv:1309.5550v2*, 2014.
- [29] John Langford, Lihong Li, and Tong Zhang. Sparse online learning via truncated gradient. In *Advances in Neural Information Processing Systems*, pages 905–912, 2009.
- [30] Mihee Lee, Haipeng Shen, Jianhua Z. Huang, and J. S. Marron. Biclustering via sparse singular value decomposition. *Biometrics*, 66(4):1087–1095, Dec 2010.

- [31] Yuanqing Li, Praneeth Namburi, Zhuliang Yu, Cuntai Guan, Jianfeng Feng, and Zhenghui Gu. Voxel selection in fmri data analysis based on sparse representation. *IEEE Trans. Biomed. Eng.*, 56(10):2439–2451, 2009.
- [32] J. Liu, S. Ji, and J. Ye. *SLEP: Sparse Learning with Efficient Projections*, <http://www.yelab.net/software/SLEP/>. Arizona State University, 2009.
- [33] Jun Liu and Jieping Ye. Efficient euclidean projections in linear time. In *Proceedings of the 26th ICML*, pages 657–664. ACM, 2009.
- [34] Jun Liu and Jieping Ye. Efficient ℓ_1/ℓ_q norm regularization. *arXiv:1009.4766*, 2010.
- [35] Ricardo Nanculef, Emanuele Frandi, Claudio Sartori, and Héctor Allende. A novel Frank-Wolfe algorithm. Analysis and applications to large-scale SVM training. *Inf. Sci.*, 285:66–99, 2014.
- [36] Yu. Nesterov. Gradient methods for minimizing composite functions. *Math. Prog. Ser. B*, 140(1):125–161, 2013.
- [37] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Pasos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.*, 12:2825–2830, 2011.
- [38] Peter Richtárik and Martin Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Math. Prog. Ser. A*, 144(1):1–38, 2014.
- [39] Bernard Schölkopf and Alexander Smola. *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. MIT Press, 2001.
- [40] Shai Shalev-Shwartz, Nathan Srebro, and Tong Zhang. Trading accuracy for sparsity in optimization problems with sparsity constraints. *SIAM J. on Optimization*, 20(6):2807–2832, 2010.
- [41] Shai Shalev-Shwartz and Ambuj Tewari. Stochastic methods for ℓ_1 -regularized loss minimization. *J. Mach. Learn. Res.*, 12:1865–1892, 2011.
- [42] M. Signoretto, E. Frandi, Z. Karevan, and J. A. K. Suykens. High level high performance computing for multitask learning of time-varying models. In *IEEE CIBD 2014*, 2014.
- [43] Robert Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc., Series B*, 58(1):267–288, 1996.
- [44] Robert Tibshirani. Regression shrinkage and selection via the lasso: a retrospective. *J. R. Stat. Soc., Series B*, 73(3):273–282, 2011.
- [45] Robert Tibshirani, Michael Saunders, Saharon Rosset, Ji Zhu, and Keith Knight. Sparsity and smoothness via the fused lasso. *J. R. Stat. Soc., Series B*, 67(1):91–108, 2005.

- [46] Joel A. Tropp. Greed is good: Algorithmic results for sparse approximation. *IEEE Trans. Inf. Theory*, 50(10):2231–2242, 2004.
- [47] B. A. Turlach. On algorithms for solving least squares problems under an l_1 penalty or an l_1 constraint. In *Proc. Am. Stat. Assoc., Stat. Comp. Sec.*, pages 2572–2577, 2005.
- [48] Yijie Wang and Xiaoning Qian. Stochastic coordinate descent Frank-Wolfe algorithm for large-scale biological network alignment. In *GlobalSIP14 - Workshop on Genomic Signal Processing and Statistics*, 2014.
- [49] Jason Weston, André Elisseeff, Bernhard Schölkopf, and Mike Tipping. Use of the zero norm with linear models and kernel methods. *J. Mach. Learn. Res.*, 3:1439–1461, 2003.
- [50] Philip Wolfe. Convergence theory in nonlinear programming. In *Integer and Nonlinear Programming*, pages 1–36. North-Holland, 1970.
- [51] Quan Zhou, Shiji Song, Gao Huang, and Cheng Wu. Efficient Lasso training from a geometrical perspective. *Neurocomputing (in press)*, 2015.
- [52] Hui Zou. The adaptive lasso and its oracle properties. *JASA*, 101(476):1418–1429, 2006.
- [53] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *J. R. Stat. Soc., Series B*, 67:301–320, 2005.
- [54] Hui Zou and Hao Helen Zhang. On the adaptive elastic-net with a diverging number of parameters. *Ann. Stat.*, 37(4):1733, 2009.

Appendix: proof of Proposition 2

Let f be a convex differentiable real function on \mathbb{R}^p . Given $\mathcal{S} \subseteq \{1, \dots, p\}$, we define the *restricted gradient* of f with respect to \mathcal{S} as its scaled projection i.e.

$$\tilde{\nabla}_{\mathcal{S}} f(\cdot) = \frac{p}{\kappa} \left(\sum_{i \in \mathcal{S}} \mathbf{e}_i \mathbf{e}_i^T \right) \nabla f(\cdot), \quad (14)$$

where $\kappa = |\mathcal{S}|$. We define the *curvature* constant of f over a compact set Σ as

$$C_f = \sup_{x \in \Sigma, y \in \Sigma_x} \frac{1}{\lambda^2} (f(y) - f(x) - (y - x)^T \nabla f(x)) , \quad (15)$$

where $\Sigma_x = \{y \in \Sigma : y = x + \lambda(s - x), s \in \Sigma, \lambda \in (0, 1]\}$.

For any $\alpha \in \Sigma$ we define its primal gap and duality gap as

$$h(\alpha) = f(\alpha) - f(\alpha^*) \quad (16)$$

$$g(\alpha) = \max_{u \in \Sigma} (\alpha - u)^T \nabla f(\alpha) , \quad (17)$$

respectively. Convexity of the function f implies that $f(\alpha) + (u - \alpha)^T \nabla f(\alpha)$ is nowhere greater than $f(\alpha)$. Therefore,

$$g(\alpha) \geq h(\alpha) \quad \forall \alpha \in \Sigma . \quad (18)$$

For any iterate $\alpha^{(k)}$ generated by our algorithm, we define its *expected primal gap and duality gap* as

$$h_{k+1} = \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha^{(k+1)})] \quad (19)$$

$$g_{k+1} = \mathbb{E}_{\mathcal{S}^{(k)}} [g(\alpha^{(k+1)})] , \quad (20)$$

respectively. Here we denote by $\mathcal{S}^{(k)}$ the random subset of $\{1, \dots, p\}$ used to approximate the gradient at each iteration. Clearly,

$$g_k \geq h_k \quad \forall k. \quad (21)$$

Lemma 2. Let $\alpha_\lambda^{(k+1)} = \alpha^{(k)} + \lambda(u^{(k)} - \alpha^{(k)})$ be a step in the direction of

$$u^{(k)} \in \underset{u \in \Sigma}{\operatorname{argmin}} (u - \alpha^{(k)})^T \tilde{\nabla}_{\mathcal{S}^{(k)}} f(\alpha^{(k)}), \quad (22)$$

with step-size $\lambda \in (0, 1]$. Then

$$h_{k+1}(\lambda) = \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha_\lambda^{(k+1)})] \leq h_k - \lambda g_k + \lambda^2 C_f \quad (23)$$

Proof. Since $\alpha^{(k)}, u^{(k)} \in \Sigma$ and $\alpha_\lambda^{(k+1)} \in \Sigma_{\alpha^{(k)}}$, it follows from (15) that

$$f(\alpha_\lambda^{(k+1)}) \leq f(\alpha^{(k)}) + \lambda(u^{(k)} - \alpha^{(k)})^T \nabla f(\alpha^{(k)}) + \lambda^2 C_f .$$

Expectation on both sides with respect to $\mathcal{S}^{(k)}$ yields

$$\mathbb{E}_{\mathcal{S}^{(k)}} [f(\alpha_\lambda^{(k+1)})] \leq f(\alpha^{(k)}) + \lambda \mathbb{E}_{\mathcal{S}^{(k)}} [(u^{(k)} - \alpha^{(k)})^T \nabla f(\alpha^{(k)})] + \lambda^2 C_f. \quad (24)$$

Since $\mathbb{E}_{\mathcal{S}^{(k)}} [\tilde{\nabla}_{\mathcal{S}^{(k)}} f(\alpha^{(k)})] = \nabla f(\alpha^{(k)})$, by the definition of $u^{(k)}$ and by the order preserving and linearity properties of expectation, we obtain

$$\begin{aligned} \mathbb{E}_{\mathcal{S}^{(k)}} [(u^{(k)} - \alpha^{(k)})^T \nabla f(\alpha^{(k)})] &= \mathbb{E}_{\mathcal{S}^{(k)}} [\min_{u \in \Sigma} (u - \alpha^{(k)})^T \tilde{\nabla}_{\mathcal{S}^{(k)}} f(\alpha^{(k)})] \\ &\leq \min_{u \in \Sigma} \mathbb{E}_{\mathcal{S}^{(k)}} [(u - \alpha^{(k)})^T \tilde{\nabla}_{\mathcal{S}^{(k)}} f(\alpha^{(k)})] \\ &= \min_{u \in \Sigma} (u - \alpha^{(k)})^T \nabla f(\alpha^{(k)}) \\ &= -g(\alpha^{(k)}). \end{aligned} \quad (25)$$

Substitution into (24) and expectation with respect to $\mathcal{S}^{(k-1)}$ finally yield

$$\mathbb{E}_{\mathcal{S}^{(k)}} [f(\alpha_\lambda^{(k+1)})] \leq \mathbb{E}_{\mathcal{S}^{(k-1)}} [f(\alpha^{(k)})] - \lambda \mathbb{E}_{\mathcal{S}^{(k-1)}} [g(\alpha^{(k)})] + \lambda^2 C_f .$$

Subtracting $f(\alpha^*)$ from both sides, (19) and (20) yield the result. \square

Lemma 3. The initialization $\alpha^{(1)} = \mathbf{u}^*$ with $\mathbf{u}^* \in \arg \min_{\mathbf{u} \in \mathcal{V}(\Sigma)} f(\mathbf{u})$ guarantees $h_k \leq C_f \quad \forall k > 1$.

Proof. First, note that $h_{k+1} \leq h_k \quad \forall k > 1$. Indeed,

$$\min_{\lambda \in (0,1]} h(\alpha_\lambda^{(k+1)}) = \min_{\lambda \in (0,1]} h(\alpha^{(k)} + \lambda(u^{(k)} - \alpha^{(k)})) \leq h(\alpha^{(k)}) .$$

Thus

$$\begin{aligned} h_{k+1} &= \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha^{(k+1)})] = \mathbb{E}_{\mathcal{S}^{(k)}} \left[\min_{\lambda \in (0,1]} h(\alpha_\lambda^{(k+1)}) \right] \\ &\leq \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha^{(k)})] = h_k . \end{aligned}$$

Now, from Lemma 2, any step in the direction of (22) with step size $\lambda \in (0, 1]$ satisfies

$$h_{k+1}(\lambda) = \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha_\lambda^{(k+1)})] \leq h_k - \lambda g_k + \lambda^2 C_f \leq h_k - \lambda h_k + \lambda^2 C_f .$$

Suppose $h_k > C_f$. In this case, as $-\lambda h_k + \lambda^2 C_f < 0$, we can choose $\lambda = 1$ to obtain

$$h_{k+1}(\lambda)|_{\lambda=1} < h_k .$$

But

$$h_{k+1}(\lambda)|_{\lambda=1} = \mathbb{E}_{\mathcal{S}^{(k)}} [h(\mathbf{u}^{(k)})] \leq \mathbb{E}_{\mathcal{S}^{(k)}} [h(\mathbf{u}^*)] = \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha^{(1)})] = h_1 .$$

Thus, $h_1 < h_{k+1}$. This is a contradiction, since $h_{k+1} \leq h_k \forall k > 1$. \square

Lemma 4. *At each iteration k of Algorithm 2,*

$$h_{k+1} \leq h_k - \frac{h_k^2}{4C_f} . \quad (26)$$

Proof. At iteration k , Algorithm 2 updates $\alpha^{(k)}$ by a line search in the direction of (22). Hence

$$h_{k+1} = \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha^{(k+1)})] = \mathbb{E}_{\mathcal{S}^{(k)}} \left[\min_{\lambda \in (0,1]} h(\alpha_\lambda^{(k+1)}) \right] . \quad (27)$$

By the order preserving and linearity properties of expectation

$$\mathbb{E}_{\mathcal{S}^{(k)}} \left[\min_{\lambda \in (0,1]} h(\alpha_\lambda^{(k+1)}) \right] \leq \min_{\lambda \in (0,1]} \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha_\lambda^{(k+1)})] . \quad (28)$$

From lemma 2, we have that any step in the direction of (22) with step size λ satisfies

$$h_{k+1}(\lambda) = \mathbb{E}_{\mathcal{S}^{(k)}} [h(\alpha_\lambda^{(k+1)})] \leq h_k - \lambda g_k + \lambda^2 C_f \leq h_k - \lambda h_k + \lambda^2 C_f . \quad (29)$$

Combining (29) and (28) produces

$$h_{k+1} = \min_{\lambda \in (0,1]} h_{k+1}(\lambda) \leq \min_{\lambda \in (0,1]} (h_k - \lambda h_k + \lambda^2 C_f) . \quad (30)$$

From lemma 3, $h_k < 2C_f$. Thus, the minimum at the right hand side is obtained for $\bar{\lambda} = h_k/2C_f$ (take derivative, equal to 0, solve and check that $\bar{\lambda} < 1$). Substituting this value of λ yields

$$h_{k+1} \leq h_k - \frac{h_k^2}{2C_f} + \frac{h_k^2}{4C_f} = h_k - \frac{h_k^2}{4C_f} . \quad (31)$$

\square

Proof of Proposition 2. With the above results in hand, we can now prove the convergence result in the main paper i.e.

$$h_k = \mathbb{E}_{\mathcal{S}^{(k)}} \left[f(\alpha^{(k+1)}) \right] - f(\alpha^*) \leq \frac{4C_f}{k+2}.$$

Proof. We prove the claim by induction on k . The base case $k = 0$ is trivial to verify from lemma 3 (as $4/3 > 1$). Now, from Lemma 4 and the inductive hypothesis $h_k \leq \frac{4C_f}{k+2}$, we obtain

$$h_{k+1} \leq h_k - \frac{h_k^2}{4C_f} \leq \frac{h_k}{1 + \frac{h_k}{4C_f}} = \frac{1}{\frac{1}{h_k} + \frac{1}{4C_f}} \leq \frac{1}{\frac{k+2}{4C_f} + \frac{1}{4C_f}} = \frac{4C_f}{(k+1) + 2}.$$

which completes the inductive step and yields the claimed bound. \square