



**HAL**  
open science

# Provably Convergent Working Set Algorithm for Non-Convex Regularized Regression

Alain Rakotomamonjy, Rémi Flamary, Gilles Gasso, Joseph Salmon

► **To cite this version:**

Alain Rakotomamonjy, Rémi Flamary, Gilles Gasso, Joseph Salmon. Provably Convergent Working Set Algorithm for Non-Convex Regularized Regression. 2020. hal-02875560v3

**HAL Id: hal-02875560**

**<https://hal.science/hal-02875560v3>**

Preprint submitted on 15 Oct 2020 (v3), last revised 19 Oct 2021 (v4)

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

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

---

# Provably Convergent Working Set Algorithm for Non-Convex Regularized Regression

---

**Alain Rakotomamonjy**

Criteo AI Lab, Paris  
LITIS, Univ. de Rouen  
alain.rakoto@insa-rouen.fr

**Rémi Flamary**

Univ Côte d’Azur, CNRS, OCA Lagrange  
remi.flamary@unice.fr

**Gilles Gasso**

LITIS, INSA de Rouen  
gilles.gasso@insa-rouen.fr

**Joseph Salmon**

IMAG, Université de Montpellier, CNRS  
Montpellier, France  
joseph.salmon@umontpellier.fr

## Abstract

Non-convex sparse regularizers are common tools for learning with high-dimensional data. For accelerating convergence of optimization problems involving those regularizers, a working set strategy addresses the optimization problem through an iterative algorithm by gradually incrementing the number of variables to optimize until the identification of the solution support. While working set methods have been theoretically supported only for convex regularizers, this paper proposes a working set algorithm for non-convex sparse regularizers with convergence guarantees. The algorithm, named *FireWorks*, is based on a non-convex reformulation of a recent primal-dual approach and leverages on the geometry of the residuals. Our theoretical guarantees build upon a lower bound of the objective function decrease between two inner solver iterations and shows the convergence to a stationary point of the problem. More importantly, we also show that convergence is preserved even when the inner solver is inexact, under sufficient decay of the error across iterations. Our experimental results demonstrate strong computational gain when using our working set strategy compared to the full problem solver for both block-coordinate descent or a proximal gradient solver.

## 1 Introduction

Many real-world learning problems are of (very) high dimension. This is the case for natural language processing problems with very large vocabulary or recommendation problems involving million of items. In such cases, one way of addressing the learning problem is to consider sparsity-inducing penalties. Likewise, when the solution of a learning problem is known to be sparse, using those penalties yield to models that can leverage this prior knowledge. The *Lasso* [29] and the *Basis pursuit* [6, 5] where the first approaches that have employed  $\ell_1$ -norm penalty for inducing sparsity.

The *Lasso* model has enjoyed large practical successes in the machine learning and signal processing communities [27, 9, 20, 34]. Nonetheless, it suffers from theoretical drawbacks (*e.g.*, biased estimates for large coefficients of the model) which can be overcome by considering non-convex sparsity-inducing penalties. Those penalties provide continuous approximations of the  $\ell_0$ - (pseudo)norm which is the true measure of sparsity. There exists a flurry of different penalties like the *Smoothly Clipped Absolute Deviation* (SCAD) [10], the *Log Sum penalty* (LSP) [4], the *capped- $\ell_1$  penalty* [36], the *Minimax Concave Penalty* (MCP) [35]. We refer the interested reader to [28] for a discussion on the pros and cons of such non-convex formulations.

Penalty	$r_\lambda( w )$	$\partial r_\lambda( w )$
Log sum	$\lambda \log(1 +  w /\theta)$	$\begin{cases} [-\frac{\lambda}{\theta}, \frac{\lambda}{\theta}] & \text{if } w = 0 \\ \{\lambda \frac{\text{sign}(w)}{\theta +  w }\} & \text{if } w \neq 0 \end{cases}$
MCP	$\begin{cases} \lambda w  - \frac{w^2}{2\theta} & \text{if }  w  \leq \lambda\theta \\ \theta\lambda^2/2 & \text{if }  w  > \theta\lambda \end{cases}$	$\begin{cases} [-\lambda, \lambda] & \text{if } w = 0 \\ \{\lambda \text{sign}(w) - \frac{w}{\theta}\} & \text{if } 0 <  w  \leq \lambda\theta \\ \{0\} & \text{if }  w  > \theta\lambda \end{cases}$
SCAD	$\begin{cases} \lambda w  & \text{if }  w  \leq \lambda \\ \frac{-w^2 + 2\theta\lambda w  - \lambda^2}{2(\theta-1)} & \text{if } \lambda <  w  \leq \lambda\theta \\ \frac{\lambda^2(1+\theta)}{2} & \text{if }  w  > \theta\lambda \end{cases}$	$\begin{cases} [-\lambda, \lambda] & \text{if } w = 0 \\ \{\lambda \text{sign}(w)\} & \text{if } 0 <  w  \leq \lambda \\ \{\frac{-w + \theta\lambda \text{sign}(w)}{\theta-1}\} & \text{if } \lambda <  w  \leq \lambda\theta \\ \{0\} & \text{if }  w  > \theta\lambda \end{cases}$

Table 1: Common non-convex penalties with their sub-differentials. Here  $\lambda > 0$ ,  $\theta > 0$  ( $\theta > 1$  for MCP,  $\theta > 2$  for SCAD).

In addition to theoretical statistical analyses, efforts have also been made for developing computationally efficient algorithms for non-convex regularized optimization problems. This includes coordinate descent algorithms [3], proximal gradient descent [16] or Newton method [33, 24]. However, all these methods share one kind of inefficiency in the sense that they spend a similar computational effort for each variable, even when those variables will end up being irrelevant (zero weight) in the final learnt model. In the non-convex setting, few methods have tried to lift this issue. One approach mixes importance sampling and randomized coordinate descent [11], while another one seeks at safely screening features that are irrelevant [25]. Working set (also known as active set) strategy aims at focusing computational effort on a subset of relevant variables, making them highly efficient for optimization problem with sparse solutions, provided that the algorithm is able to quickly identify the ‘‘relevant’’ features. In the literature, several works on working set algorithms address this selection issue mostly for convex optimization problems such as the Support Vector Machine problem [32, 14] or the Lasso problem [12, 30, 18, 22]. Working set strategies have been extended to non-convex sparse optimization problems [1, 2] but lack of theoretical understandings.

In this work, inspired by the Blitz algorithm proposed by Johnson and Guestrin [18](see also [21, 22] for its connection with safe screening rules) we propose a theoretically supported method for selecting working set in non-convex regularized sparse optimization problems. While Blitz can only be implemented for convex problems, leveraging on primal-dual aspects of the  $\ell_1$ -regularized problem, we show that a similar algorithm can be designed by exploiting the key role of the residual in the sparse regression problem. Our algorithm proposes a method for selecting variables to integrate into a working set, and provides a theoretical guarantee on objective value decrease. Based on those results, we provide, as far as we know, the first convergence guarantee of working set algorithm in a non-convex Lasso setting and we show that this convergence property is preserved in a inexact setting.

In summary, our contributions are the following:

- We propose a novel working set algorithm for non-convex regularized regression that selects features to integrate in the model based on a so-called ‘‘feasible’’ residual;
- We show that the algorithm enjoys properties such as convergence to a stationary point, even when the inner solver is inexact, under sufficient decay of the error along the iterations;
- Our experimental results show that our *FireWorks* algorithm achieves substantial computational gain (that can reach two order of magnitude) compared to the baseline approaches and even to other working set approach.

**Notation** We denote as  $\mathbf{X} \in \mathbb{R}^{n \times d}$  the design matrix. We write vectors of size  $d$  or size  $n$  in bold e.g.,  $\mathbf{y} \in \mathbb{R}^n$  or  $\mathbf{w} \in \mathbb{R}^d$ . We will consider several sets and they are noted in calligraphic mode. We have set of indices, mostly noted as  $\mathcal{A}$ , with  $\mathcal{A}$  being a subset of indices extracted from  $\{1, \dots, d\}$  and with cardinality noted  $|\mathcal{A}|$ . Given a set  $\mathcal{A}$ ,  $\bar{\mathcal{A}}$  denotes its complement in  $\{1, \dots, d\}$ . Set defined by (union of) function level-set will be denoted as  $\mathcal{C}$ , with indices defining the function. Vectors noted as  $\mathbf{w}_{\mathcal{A}}$  are of size  $|\mathcal{A}|$  and we note  $\tilde{\mathbf{w}}_{\mathcal{A}} \in \mathbb{R}^d$  for the vector of component  $w_{j,\mathcal{A}}$  for all  $j \in \mathcal{A}$  and 0 elsewhere. We will note  $\text{res}(\mathbf{w}) \triangleq \mathbf{y} - \mathbf{X}\mathbf{w}$  and  $\text{res}(\mathbf{w}_{\mathcal{A}}) \triangleq \mathbf{y} - \mathbf{X}_{\mathcal{A}}\mathbf{w}_{\mathcal{A}} = \mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}_{\mathcal{A}}$ .

## 2 Linear regression with non-convex regularizers

We first introduce the non-convex Lasso problem we are interested in as well as its first order optimality conditions. We emphasize on the form of the optimality conditions which will be key for designing our working set algorithm.

### 2.1 The optimization problem

We consider solving the problem of least-squares regression with a generic penalty of the form

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \triangleq \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \sum_{j=1}^d r_\lambda(|w_j|) , \quad (1)$$

where  $\mathbf{y} \in \mathbb{R}^n$  is a target vector,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_d] \in \mathbb{R}^{n \times d}$  is the design matrix with column-wise features  $\mathbf{x}_j \in \mathbb{R}^n$ ,  $\mathbf{w}$  is the coefficient vector of the model and the map  $r_\lambda : \mathbb{R}_+ \mapsto \mathbb{R}_+$  is concave and differentiable on  $[0, +\infty)$  with a regularization parameter  $\lambda > 0$ . In addition, we assume that  $r_\lambda(|w|)$  is a lower semi-continuous function. Note that most penalty functions such as SCAD, MCP or log sum (see their definitions in Table 1) satisfy such a property and that for those penalties,  $f(\cdot)$  is lower bounded.

We consider tools such as Fréchet subdifferentials and limiting-subdifferentials [19, 26, 23] well suited for non-smooth and non-convex optimization, so that a vector  $\mathbf{w}^*$  belongs to the set of minimizers (not necessarily global) of Problem (1) if the following Fermat's condition holds [7, 19]:

$$\forall j, \mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}\mathbf{w}^*) \in \partial r_\lambda(|w_j^*|) , \quad (2)$$

with  $\partial r_\lambda(|\cdot|)$  being the Fréchet subdifferential of  $r_\lambda(|\cdot|)$ , assuming it exists at  $\mathbf{w}^*$ . In particular, this is the case for the MCP, log sum and SCAD penalties presented in Table 1. For the sake of clarity, we present the optimality conditions for MCP and log sum, in the next examples.

**Example 1.** For the MCP penalty (see Table 1 for its definition and its subdifferential), it is easy to show that the  $\partial r_\lambda(|0|) = [-\lambda, \lambda]$ . Hence, the Fermat's condition becomes with the residual  $\text{res}(\mathbf{w}^*)$

$$\begin{cases} -\mathbf{x}_j^\top \text{res}(\mathbf{w}^*) = 0, & \text{if } |w_j^*| > \lambda\theta \\ -\mathbf{x}_j^\top \text{res}(\mathbf{w}^*) + \lambda \text{sign}(w_j^*) = \frac{w_j^*}{\theta}, & \text{if } 0 < |w_j^*| \leq \lambda\theta \\ |\mathbf{x}_j^\top \text{res}(\mathbf{w}^*)| \leq \lambda, & \text{if } w_j^* = 0 \end{cases} \quad (3)$$

**Example 2.** For the log sum penalty, one can explicitly compute  $\partial r_\lambda(|0|) = [-\frac{\lambda}{\theta}, \frac{\lambda}{\theta}]$  and leverage on smoothness of  $r_\lambda(|w|)$  when  $|w| > 0$  for computing  $\partial r_\lambda(|w|)$ . Then, the condition in Equation (2) can be written as:

$$\begin{cases} -\mathbf{x}_j^\top \text{res}(\mathbf{w}^*) + \lambda \frac{\text{sign}(w_j^*)}{\theta + |w_j^*|} = 0, & \text{if } w_j^* \neq 0 , \\ |\mathbf{x}_j^\top \text{res}(\mathbf{w}^*)| \leq \frac{\lambda}{\theta}, & \text{if } w_j^* = 0 . \end{cases} \quad (4)$$

As we can see, first order optimality conditions lead to simple equations and inclusions. More interestingly, one can note that regardless of the regularizer, the structure of optimality condition for a weight  $w_j^* = 0$  depends on the correlation of the feature  $\mathbf{x}_j$  with the optimal residual  $\mathbf{y} - \mathbf{X}\mathbf{w}^*$ . Hence, these conditions can be used for defining a region in which the optimal residual has to live.

## 3 Working Set Algorithm and Analysis

Before presenting the *FireWorks* algorithm, we first introduce in what follows, all necessary concepts underlying our working set algorithm and needed for analyzing its properties.

### 3.1 Restricted problem and optimality

Given a set  $\mathcal{A}$  of  $m$  indices belonging to  $\{1, \dots, d\}$ , we denote as restricted problem, the problem defined in Equation (1) but restricted to columns of  $\mathbf{X}$  defined by  $\mathcal{A}$ , leading thus to

$$\min_{\mathbf{w}_{\mathcal{A}} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{y} - \mathbf{X}_{\mathcal{A}}\mathbf{w}_{\mathcal{A}}\|_2^2 + \sum_{j=1}^m r_\lambda(|w_{j,\mathcal{A}}|) . \quad (5)$$

Naturally, a vector  $\mathbf{w}_{\mathcal{A}}^*$  minimizing this problem has to satisfy its own Fermat's condition. However, the next proposition derives another necessary condition for achieving Fermat's condition that will be useful in the sequel for characterizing optimality of  $\tilde{\mathbf{w}}_{\mathcal{A}}^*$  on the full problem.

**Proposition 1.** *If  $\mathbf{w}_{\mathcal{A}}^*$  satisfies Fermat's condition of problem (5), then for all  $j \in \mathcal{A}$ , we have*

$$|\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{\mathcal{A}} \mathbf{w}_{\mathcal{A}}^*)| \leq r'_\lambda(0) \quad (6)$$

*Proof.* At first, note that since the function  $r_\lambda$  is concave, then its gradient is positive and decreasing hence  $\forall w, r'_\lambda(w) \leq r'_\lambda(0)$ . Now, when  $j \in \{i \in \mathcal{A} : w_{i,\mathcal{A}}^* = 0\}$  then the inequality in Equation 6 naturally comes from the Fermat's condition. When  $j \in \{i \in \mathcal{A} : w_{i,\mathcal{A}}^* \neq 0\}$ , then for a stationary point, we must have  $\mathbf{x}_j^\top \text{res}(\mathbf{w}_{\mathcal{A}}^*) = r'_\lambda(|w_{j,\mathcal{A}}^*|)$ . Taking the absolute value of this equation and plugging in the inequality of the gradient concludes the proof.  $\square$

Based on this latter condition, we define the function  $h_j : \mathbb{R}^n \rightarrow \mathbb{R}$ , for  $j \in \{1, \dots, d\}$  as  $h_j(\mathbf{a}) = |\mathbf{x}_j^\top \mathbf{a}| - r'_\lambda(0)$ , the convex set  $\mathcal{C}_j$  expressed as the slab

$$\mathcal{C}_j \triangleq \{\mathbf{a} \in \mathbb{R}^n : h_j(\mathbf{a}) \leq 0\}$$

and

$$\mathcal{C}_j^\equiv \triangleq \{\mathbf{a} \in \mathbb{R}^n : h_j(\mathbf{a}) = 0\}.$$

Now, we are going to define a set useful for characterizing candidate stationary points of either Equations (1) or (5). We note as  $\mathcal{C} = \bigcap_{j=1}^d \mathcal{C}_j$  and  $\mathcal{C}_{\mathcal{A}} = \bigcap_{j \in \mathcal{A}} \mathcal{C}_j$  and from this, the necessary optimality condition defined in Proposition 1 can be written as  $\mathbf{y} - \mathbf{X}_{\mathcal{A}} \mathbf{w}_{\mathcal{A}}^* \in \mathcal{C}_{\mathcal{A}}$ . Note that  $\mathcal{C}$  is the dual feasible set for  $\ell_1$ -type convex regularizers. Now, assume  $\mathbf{w}_{\mathcal{A}}^*$  is a minimizer of its restricted problem then  $\tilde{\mathbf{w}}_{\mathcal{A}}^*$  (in  $\mathbb{R}^d$ ) satisfies the Fermat's condition of the full problem if and only if we also have

$$\mathbf{y} - \mathbf{X} \tilde{\mathbf{w}}_{\mathcal{A}}^* \in \mathcal{C}_{\bar{\mathcal{A}}} \quad , \quad (7)$$

where  $\bar{\mathcal{A}}$  is the complement of  $\mathcal{A}$  in  $\{1, \dots, d\}$ . Indeed, since  $\mathbf{w}_{\mathcal{A}}^*$  is optimal for the restricted problem, Fermat's condition is already satisfied for all  $j \in \mathcal{A}$ . Then, the above condition ensures that  $\forall j \in \bar{\mathcal{A}}$ , we have  $|\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X} \tilde{\mathbf{w}}_{\mathcal{A}}^*)| \leq r'_\lambda(0)$  since, as by definition,  $\tilde{w}_j^* = 0, \forall j \in \bar{\mathcal{A}}$ . As they play a key role in our algorithm and for monitoring optimality, we define the distance of a vector  $\mathbf{r}$  to the convex set  $\mathcal{C}_j$  and  $\mathcal{C}_j^\equiv$  as

$$\begin{aligned} \text{dist}(\mathbf{r}, \mathcal{C}_j) &\triangleq \min_{\mathbf{z} \in \mathbb{R}^n} \|\mathbf{z} - \mathbf{r}\|_2 \quad , \quad \text{s.t. } h_j(\mathbf{z}) \leq 0; \\ \text{dist}_S(\mathbf{s}, \mathcal{C}_j^\equiv) &\triangleq \min_{\mathbf{z} \in \mathbb{R}^n} \|\mathbf{z} - \mathbf{s}\|_2 \quad , \quad \text{s.t. } h_j(\mathbf{z}) = 0 \quad . \end{aligned}$$

These definitions can also be used for defining the most violated optimality condition which is a key component of methods proposed by [1, 11]. Indeed, given a set  $\mathcal{A}$ , the solution  $\mathbf{w}_{\mathcal{A}}$  of Equation (5) and the residual  $\text{res}(\mathbf{w}_{\mathcal{A}})$ , we have the index  $j^* = \arg \max_{j \in \bar{\mathcal{A}}} \text{dist}(\text{res}(\mathbf{w}_{\mathcal{A}}), \mathcal{C}_j)$  which is the index of the most violated constraint among non-active variables for the residual  $\text{res}(\mathbf{w}_{\mathcal{A}})$ .

### 3.2 Feasible Residual Working Set Algorithm for non-convex Lasso

A working set algorithm for solving problem (1) consists in sequentially solving a series of restricted problem as defined in Equation (5) with a sequence of working sets  $\mathcal{A}_0, \mathcal{A}_1, \dots, \mathcal{A}_k$ . The main differences among working set algorithms lie on the way the set is being updated. For instance, the approach of [1], denoted in the experiment as MaxVC, selects the single most violated constraint (as defined above) in the non-active set to be included in the new working set leading to the algorithm presented in the supplementary material. Flamary et al. [11] followed similar approach but considered a randomized selection in which the probability of selection is related to  $\text{dist}(\text{res}(\mathbf{w}_{\mathcal{A}}^*), \mathcal{C}_j)$ .

Our algorithm is inspired by Blitz [18] which is a working set algorithm dedicated to convex constrained optimization problem. But as the problem we address is a non-convex one, we manipulate different mathematical objects that need to be redefined. The procedure is presented in Algorithm 1.

It starts by selecting a small subset of indices (either randomly or cleverly like the indices with largest  $|\mathbf{x}_j^\top \mathbf{y}|$ , for instance) as initial working set and by choosing a vector  $\mathbf{s}_1$  such that  $\mathbf{s}_1 \in \mathcal{C} = \bigcap_{j=1}^d \mathcal{C}_j$ ,

---

**Algorithm 1:** FireWorks: Feasible Residual Working Set Algorithm

---

**Input:**  $\{\mathbf{X}, \mathbf{y}\}$ ,  $\mathcal{A}_1$  active set,  $\mathbf{s}_1 \in \mathcal{C}$ , a sequence of  $\tau_k$  / mechanism for defining  $\tau_k$ , initial vector  $\tilde{\mathbf{w}}_{\mathcal{A}_0}$

**Output:**  $\tilde{\mathbf{w}}_{\mathcal{A}_k}$

- 1: **for**  $k = 1, 2, \dots$  **do**
- 2:  $\mathbf{w}_{\mathcal{A}_k} = \arg \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{y} - \mathbf{X}_{\mathcal{A}_k} \mathbf{w}\|_2^2 + \sum_{j \in \mathcal{A}_k} r_\lambda(|w_j|)$ ; // warm-start solver with  $\mathbf{w}_{\mathcal{A}_{k-1}}$
- 3:  $\mathbf{r}_k = \mathbf{y} - \mathbf{X}_{\mathcal{A}_k} \mathbf{w}_{\mathcal{A}_k}$ ; // get residual
- 4:  $\alpha_k = \max\{\alpha \in [0, 1] : \alpha \mathbf{r}_k + (1 - \alpha) \mathbf{s}_k \in \mathcal{C}\}$
- 5:  $\mathbf{s}_{k+1} = \alpha_k \mathbf{r}_k + (1 - \alpha_k) \mathbf{s}_k$ ; // define the most "feasible" residual
- 6:  $\mathcal{A}_k = \mathcal{A}_k / \{j \in \mathcal{A}_k : w_{j, \mathcal{A}_k} = 0\}$ ; // we prune the set from inactive features
- 7: compute  $\tau_k$ ; // e.g., sort  $\text{dist}_S(\mathbf{s}_{k+1}, \mathcal{C}_j^-)$  so as to keep constant number of features to add
- 8:  $\mathcal{A}_{k+1} = \{j : \text{dist}_S(\mathbf{s}_{k+1}, \mathcal{C}_j^-) \leq \tau_k\} \cup \mathcal{A}_k$
- 9: **end for**
- 10: Build  $\tilde{\mathbf{w}}_{\mathcal{A}_k}$

---

for instance setting  $\mathbf{s}_1 = \mathbf{0}$ . From this vector  $\mathbf{s}_1$ , we will generate a sequence of  $\{\mathbf{s}_k\}$  that plays a key role in the selection of the constraints to be integrated in the next restricted model. The approach, at iteration  $k$ , starts by exactly solving the restricted problem with the set  $\mathcal{A}_k$  and then by computing the residual  $\mathbf{r}_k = \text{res}(\mathbf{w}_{\mathcal{A}_k}^*)$ . As noted in Equation (7), if  $\mathbf{r}_k \in \mathcal{C}_{\bar{\mathcal{A}}}$  then the vector  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  is a stationary point of the full problem. If  $\mathbf{r}_k \notin \mathcal{C}_{\bar{\mathcal{A}}}$ , we need to update the working set  $\mathcal{A}_k$ . For doing so, we proceed by defining  $\mathbf{s}_{k+1}$  as the vector on the segment  $[\mathbf{s}_k, \mathbf{r}_k]$ , nearest to  $\mathbf{r}_k$  that belongs to  $\mathcal{C}$ . Then, we update the working set by prioritizing the  $j$ -th coordinate *w.r.t.* the distance of  $\mathbf{s}_{k+1}$  to  $\mathcal{C}_j^-$  and the constraint associated to  $j$  is included in the new working set if  $\text{dist}_S(\mathbf{s}_{k+1}, \mathcal{C}_j^-) \leq \tau_k$ .  $\tau_k$  is a strictly positive term that defines the number of constraints to be added to the current working set. In practice, we have chosen  $\tau_k$  so that a fixed number  $n_k$  of constraints is added to  $\mathcal{A}_k$  for each  $k$ .

We provide the following intuition on why this algorithm works in practice. At first, note that by construction  $\mathbf{s}_{k+1}$  is a convex combination of two vectors one of which is the residual hence justifies its interpretation as a residual. However, the main difference between the  $\mathbf{s}_k$ 's and  $\mathbf{r}_k$ 's is that the former belongs to  $\mathcal{C}$  and thus to any  $\mathcal{C}_{\bar{\mathcal{A}}}$  while  $\mathbf{r}_k$  belongs to  $\mathcal{C}$  only for a potential  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  optimal for the full problem. Then, when  $\mathbf{w}_{\mathcal{A}_k}^*$  is a stationary point for the restricted problem but not for the full problem,  $\mathbf{r}_k \in \mathcal{C}_{\bar{\mathcal{A}}}$  but  $\mathbf{r}_k \notin \mathcal{C}$ . Hence,  $\mathbf{s}_{k+1}$  represents a residual candidate for optimality and constraints near this residual can be interpreted as the constraints that need to be relaxed by integrating related features  $j$  in the working set (allowing thus  $w_j$  to be potentially non-zero at the next iteration). The indices  $j$  for which distances of  $\mathbf{s}_{k+1}$  to  $\mathcal{C}_j^-$  are below a given threshold are then integrated into the working set. The mechanism for constraint selection is illustrated in Figure 1.

### 3.3 Some properties of the algorithm

In this subsection, we provided some analyzes of the proposed algorithm. At first, let us introduce an optimality condition of a vector solving the restricted problem.

**Proposition 2.** *Given a working set  $\mathcal{A}_k$  and  $\mathbf{w}_{\mathcal{A}_k}^*$  solving the related restricted problem,  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  is also optimal for the full problem if and only if  $\alpha = 1$  in Algorithm 1, step 4 (which also means  $\mathbf{s}_{k+1} = \mathbf{r}_k$ ).*

*Proof.* (sketch) apply the fact that at optimality, we have  $\mathbf{r}_k \in \mathcal{C}_{\bar{\mathcal{A}}}$ . □

Now, we are going to characterize the decrease in objective value obtained between two updates of working sets.

**Proposition 3.** *Assume that  $\mathbf{w}_{\mathcal{A}_k}$  and  $\mathbf{w}_{\mathcal{A}_{k+1}}$  are respectively the solutions of the restricted problem with the working set  $\mathcal{A}_k$  and  $\mathcal{A}_{k+1}$ , with  $\mathcal{A}_{k+1} = \{j\} \cup \mathcal{A}_k$ . Denote as  $\mathbf{r}_k \triangleq \text{res}(\mathbf{w}_{\mathcal{A}_k})$  then, we*

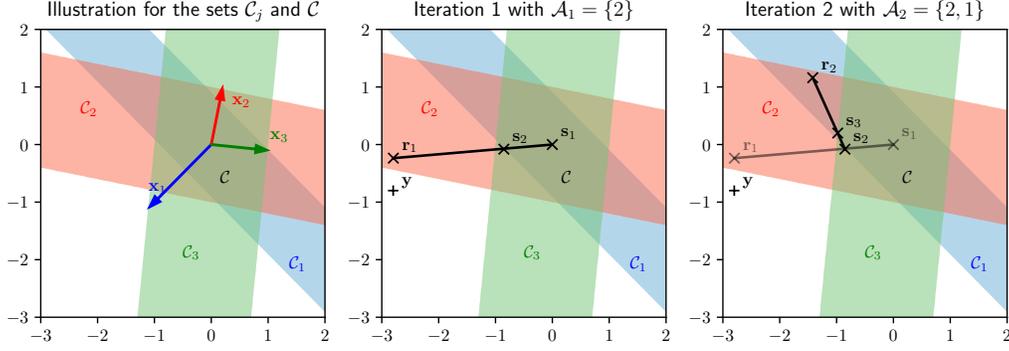


Figure 1: Illustrating the constraint selection. (left) Given three variables, we plot their associate slabs  $\{\mathcal{C}_j\}_{j=1}^3$ .  $\mathcal{C}$  is the intersection of the 3 slabs. We assume that the working set is  $\{2\}$ . (middle) After the first iteration, the residual  $\mathbf{r}_1$  satisfies constraint  $h_2(\mathbf{a}) \leq 0$  and thus lies in the  $h_2$ 's feasible region  $\mathcal{C}_2$ . Then, the segment  $[s_1, \mathbf{r}_1]$  gives us the most feasible point  $s_2 \in \mathcal{C}$ . If  $\tau_1$  is chosen so as to select only one feature, it is then  $j = 1$ . The new working set is  $\{2, 1\}$ . (right) After optimizing over this working set, the residual  $\mathbf{r}_2$  satisfies constraints  $h_2$  and  $h_1$  and thus lies in the  $\mathcal{C}_1 \cap \mathcal{C}_2$  region.

have the following inequality

$$\|\tilde{\mathbf{w}}_{\mathcal{A}_{k+1}} - \tilde{\mathbf{w}}_{\mathcal{A}_k}\|_2 \geq \frac{1}{\|\mathbf{X}\|_2} \text{dist}(\mathbf{r}_k, \mathcal{C}_j) .$$

*Proof.* We have the following inequalities

$$\begin{aligned} \|\mathbf{r}_{k+1} - \mathbf{r}_k\|_2 &= \|\mathbf{X}(\tilde{\mathbf{w}}_{\mathcal{A}_{k+1}} - \tilde{\mathbf{w}}_{\mathcal{A}_k})\|_2 \\ &\leq \|\mathbf{X}\|_2 \|\tilde{\mathbf{w}}_{\mathcal{A}_{k+1}} - \tilde{\mathbf{w}}_{\mathcal{A}_k}\|_2 . \end{aligned}$$

Now recall that  $\mathbf{r}_k \notin \mathcal{C}_{\mathcal{A}_{k+1}}$  since it violates the  $h_j$  constraint while  $\mathbf{r}_{k+1} \in \mathcal{C}_{\mathcal{A}_{k+1}}$  as  $\mathbf{w}_{\mathcal{A}_{k+1}}$  has been optimized over  $\mathcal{A}_{k+1}$ . As such, we also have  $h_j(\mathbf{r}_{k+1}) \leq 0$ . Now by definition of  $\text{dist}(\mathbf{r}_k, \mathcal{C}_j)$  either  $\mathbf{r}_{k+1}$  is the minimizer of the distance optimization problem, hence  $\text{dist}(\mathbf{r}_k, \mathcal{C}_j) = \|\mathbf{r}_{k+1} - \mathbf{r}_k\|_2$  or  $\text{dist}(\mathbf{r}_k, \mathcal{C}_j) \leq \|\mathbf{r}_{k+1} - \mathbf{r}_k\|_2$ . Plugging this latter inequality in Equation (8) concludes the proof.  $\square$

Given the right hand side of the equation in Proposition 3, we now show that the distance of the residual  $\mathbf{r}_k$  at step  $k$  to a constraint  $\mathcal{C}_j$ , defined by a feature  $j$  that is not yet in the active set, is lower bounded by a term depending on the parameter  $\tau_{k-1}$  which governs the number of features that has been added to the active set at step  $k-1$ .

**Lemma 1.** *At step  $k \geq 2$ , consider a constraint  $\mathcal{C}_j$  such that  $h_j(\mathbf{r}_k) > 0$  and  $h_j(\mathbf{s}_k) < 0$ , then*

$$\text{dist}(\mathbf{r}_k, \mathcal{C}_j) \geq \frac{1 - \alpha_k}{\alpha_k} \tau_{k-1} . \quad (8)$$

The proof of this lemma has been deported to the supplementary. From the above Proposition 3 and Lemma 1, we can ensure that the sequence of  $\{\tilde{\mathbf{w}}_{\mathcal{A}_k}\}$  produced by Algorithm 1 converges towards a stationary point under mild conditions on the inner solver.

**Proposition 4.** *Suppose that for each step  $k$ , the algorithm solving the inner problem ensures a decrease in the objective value in the form*

$$f(\tilde{\mathbf{w}}_{\mathcal{A}_{k+1}}) - f(\tilde{\mathbf{w}}_{\mathcal{A}_k}) \leq -\gamma_k \|\tilde{\mathbf{w}}_{\mathcal{A}_{k+1}} - \tilde{\mathbf{w}}_{\mathcal{A}_k}\|_2^2 .$$

with  $\forall k, \gamma_k > 0$ . For the inner solver, we also impose that when solving the problem with set  $\mathcal{A}_{k+1}$ , the inner solver is warm-started with  $\tilde{\mathbf{w}}_{\mathcal{A}_k}$ . Assume also that  $\|\mathbf{X}\|_2 > 0$ ,  $\tau_k > 0$  and  $h_j$  satisfies assumption in Lemma 1, then the sequence of  $\alpha_k$  produced by Algorithm 1 converges towards 1 and  $\forall j, \lim_{k \rightarrow \infty} |\mathbf{x}_j^\top \mathbf{r}_k| \leq r'_\lambda(0)$ .

*Proof.* Using result in Proposition 3 and Lemma 1 and the above assumption, we have, for  $k \geq 2$ ,

$$\begin{aligned} f(\tilde{\mathbf{w}}_{\mathcal{A}_{k+1}}) &\leq f(\tilde{\mathbf{w}}_{\mathcal{A}_k}) - \frac{\gamma_k}{\|\mathbf{X}\|_2^2} \left( \frac{1 - \alpha_k}{\alpha_k} \right)^2 \tau_{k-1}^2 \\ &\leq f(\tilde{\mathbf{w}}_{\mathcal{A}_2}) - \frac{1}{\|\mathbf{X}\|_2^2} \sum_{\ell=2}^k \gamma_\ell \left( \frac{1 - \alpha_\ell}{\alpha_\ell} \right)^2 \tau_{\ell-1}^2. \end{aligned}$$

This means that  $\frac{1}{\|\mathbf{X}\|_2^2} \sum_{\ell=2}^k \gamma_\ell \left( \frac{1 - \alpha_\ell}{\alpha_\ell} \right)^2 \tau_{\ell-1}^2 \leq f(\tilde{\mathbf{w}}_{\mathcal{A}_2}) - f(\tilde{\mathbf{w}}_{\mathcal{A}_{k+1}})$ . Since  $f$  is bounded from below, the right hand side is less than some positive constant, hence  $\sum_{\ell=2}^{\infty} \gamma_\ell \left( \frac{1 - \alpha_\ell}{\alpha_\ell} \right)^2 \tau_{\ell-1}^2 < \infty$ .

Since the latter sum is bounded, it implies that  $\gamma_\ell \left( \frac{1 - \alpha_\ell}{\alpha_\ell} \right)^2 \tau_{\ell-1}^2 \rightarrow 0$  as  $\ell \rightarrow \infty$ , and as  $\gamma_\ell > 0$ ,  $\tau_\ell > 0$ , we have  $\lim_{\ell \rightarrow \infty} \alpha_\ell = 1$ . Now using the definition of  $\mathbf{s}_{k+1}$ , we have  $\forall j$ ,  $\mathbf{x}_j^\top \mathbf{r}_k = \frac{1}{\alpha_k} \mathbf{x}_j^\top \mathbf{s}_{k+1} - \frac{1 - \alpha_k}{\alpha_k} \mathbf{x}_j^\top \mathbf{s}_k$ . Then, taking the absolute value, triangle inequality, using the fact that  $\forall k$ ,  $\mathbf{s}_k \in \mathcal{C}$  and taking the limit concludes the proof.  $\square$

The above proposition ensures convergence to a stationary point under some conditions on the inner solver. Several algorithms may satisfy this assumption. For instance, any first-order iterative algorithm which selects its step size as  $1/t_k$  based on line search criterion of the form  $\forall k$ ,  $f(\mathbf{w}_{k+1}) \leq f(\mathbf{w}_k) - \frac{\sigma}{2} t_k \|\mathbf{w}_{k+1} - \mathbf{w}_k\|_2^2$ , where  $\sigma$  is a constant in the interval  $(0, 1)$ , provides such a guarantee. This is the case of the generalized proximal algorithm of Gong et al. [16] or proximal Newton approaches [24], assuming that  $f$  is differentiable with gradient Lipschitz and  $r_\lambda(\cdot)$  admits a proximal operator. As non-convex block coordinate descent algorithms [3] can also be interpreted as proximal algorithm, they also satisfy this sufficient decrease condition under the same assumptions than proximal approaches.

**Inexact inner solver** One key point when considering a meta-solver like Blitz [18] or a working set algorithm is that for some approaches, theoretical properties hold only when the inner solver is solved exactly. This is for instance the case for the SimpleSVM algorithm of Vishwanathan et al. [32] or the active set algorithm proposed by Boisbunon et al. [1] which convergence is based on non-cyclicity of the working set selection (prohibiting pruning) and thus on the ability of solving exactly the inner problem. For the approach we propose, we show next that the distance between two consecutive inexact solutions of the inner problem is still lower bounded.

**Proposition 5.** *Let  $\mathbf{w}_{\mathcal{A}_k}$  and  $\mathbf{w}_{\mathcal{A}_{k+1}}$  the approximate solutions of the inner problem with respectively the working sets  $\mathcal{A}_k$  and  $\mathcal{A}_{k+1}$ . Assume that  $\mathbf{w}_{\mathcal{A}_{k+1}}$  as being obtained through a tolerance of  $\xi_{k+1} \leq \tau_k$  of its Fermat's condition (e.g., for the log sum penalty, Equation (4) are satisfied up to  $\xi_{k+1}$ ), then the following inequality holds :*

$$\|\mathbf{w}_{\mathcal{A}_{k+1}} - \mathbf{w}_{\mathcal{A}_k}\|_2^2 \geq \frac{1}{\|\mathbf{X}\|_2^2} (\text{dist}(\mathbf{r}_k, h_j) - \xi_{k+1}).$$

*Proof.* First note that if  $\mathbf{w}_{\mathcal{A}_{k+1}}$  is such that  $\mathbf{r}_{k+1} \in \mathcal{C}_{\mathcal{A}_{k+1}}$  then we are in the same condition than in Proposition 1 and the same proof applies. Let us assume then that  $\mathbf{r}_{k+1} \notin \mathcal{C}_{\mathcal{A}_{k+1}}$  and  $\text{dist}(\mathbf{r}_{k+1}, \mathcal{C}_j) \leq \xi_{k+1}$ . Define as  $\mathbf{u}$  the point in  $\mathcal{C}_j$  that defines the distance of  $\mathbf{r}_k$  to  $\mathcal{C}_j$  and as  $\mathbf{p}$  the point that minimizes the distance between  $\mathbf{r}_{k+1}$  and the segment  $[\mathbf{u}, \mathbf{r}_k]$ . Then, owing to simple geometrical arguments and orthogonality we have :  $\|\mathbf{r}_{k+1} - \mathbf{r}_k\|^2 = \|\mathbf{r}_{k+1} - \mathbf{p}\|^2 + \|\mathbf{p} - \mathbf{r}_k\|^2$  and thus  $\|\mathbf{r}_{k+1} - \mathbf{r}_k\| \geq \|\mathbf{r}_k - \mathbf{p}\|$ . Now, because  $\mathbf{p}$  belongs to the segment defined by  $\mathbf{u}$  and  $\mathbf{r}_k$ , we have

$$\|\mathbf{r}_{k+1} - \mathbf{r}_k\| \geq \|\mathbf{r}_k - \mathbf{u}\| - \|\mathbf{u} - \mathbf{p}\| \geq \text{dist}(\mathbf{r}_k, \mathcal{C}_j) - \xi_{k+1}$$

where the last inequality comes from the fact that  $\|\mathbf{u} - \mathbf{p}\| = \text{dist}(\mathbf{r}_{k+1}, \mathcal{C}_j) \leq \xi_{k+1}$ . Plugging this inequality into Equation (8) completes the proof.  $\square$

Note that the above lower bound is meaningful only if the tolerance  $\xi_{k+1}$  is smaller than the distance of the residual to the set  $\mathcal{C}_j$ . This is a reasonable assumption to be made since we expect  $\mathbf{r}_k$  to violate  $\mathcal{C}_j$ . Now, we can derive condition of convergence towards a stationary point of the full problem.

**Corollary 1.** *Under the assumption of Proposition 4 and assuming that the sequence of tolerance is such that  $\sum_k \xi_k < \infty$ , then Algorithm 1 produces a sequence of iterates that converges towards a stationary point.*

The proof follows the same steps as for Proposition 4, with the addition that sequence  $\{\xi_k\}$  is convergent and thus has been omitted. Note that the assumption of convergent sum of errors is a common assumption, notably in the proximal algorithm literature [8, 31] and it helps guaranteeing convergence towards exact stationary point instead of an approximate convergence (up to a tolerance  $\tau$ ).

**Relation with maximum violated constraint.** The mechanism we have proposed for updating the working set is based on the current residual  $\mathbf{r}_k$  and a feasible residual  $\mathbf{s}_k$ . By changing how  $\mathbf{s}_{k+1}$  is defined, we can retrieve the classical maximal violated constraint approach. Indeed, if we set at Line 5 of Algorithm 1,  $\forall k, \mathbf{s}_k = 0$  and  $\mathbf{s}_{k+1} = \alpha_k \mathbf{r}_k$ , with  $\alpha_k \in [0, 1]$  then  $\mathbf{s}_{k+1}$  is a rescaling of the current residual and the scale is chosen so that  $\mathbf{s}_{k+1} \in \mathcal{C}$ . Using simple inequality argument, it is straightforward to show that  $\alpha_k = \min(\min_{j \in \bar{\mathcal{A}}} \frac{\lambda}{|\mathbf{x}_j^\top \mathbf{r}_k|}, 1)$  and the minimum in  $j$  occurs for the largest value of  $|\mathbf{x}_j^\top \mathbf{r}_k|$ . Recall again that the polynomial convergence of this algorithm is guaranteed for exact inner solver and when no working set pruning (removing from the set  $\mathcal{A}_k$  variables which weights are 0) occurs.

Table 2: Running time in seconds of different algorithms on different problems. In the first column, we reported data, the tolerance on the stopping criterion and the constant  $K$  such that  $\lambda = K \max_j |\mathbf{x}_j^\top \mathbf{y}|$  (the larger the  $K$ , the sparser  $\mathbf{w}^*$  is). The small *Toy* dataset has  $n = 100$ ,  $d = 1000$  and  $p = 30$ ; the large one has  $n = 1000$ ,  $d = 5000$ ,  $p = 500$ . For each inner solver, we bold the most efficient algorithm. The symbol ”-” denotes that the algorithm did not finish one iteration in 24 hours. The number in parenthesis is the number of non-zero weights in  $\mathbf{w}^*$ . All experiments have been run on one single core of an Intel Xeon CPU E5-2680 clocked at 2,4Ghz.

Data and Setting	MM prox	GIST	MaxVC Gist	FireWorks Gist	MM BCD	BCD	MaxVC BCD	FireWorks BCD
Toy small - 1.00e-03 - 0.07	1.4±0.4 (34)	0.8±0.2 (34)	0.3±0.2 (34)	<b>0.2±0.1 (34)</b>	3.4±0.9 (34)	14.2±4.9 (34)	1.9±0.8 (34)	<b>1.5±0.9 (34)</b>
Toy small - 1.00e-05 - 0.07	1.5±0.4 (34)	1.4±0.6 (34)	0.7±0.8 (34)	<b>0.4±0.1 (34)</b>	3.3±0.8 (34)	22.9±11.0 (34)	8.3±9.7 (34)	<b>2.7±1.2 (34)</b>
Toy small - 1.00e-03 - 0.01	11.2±1.2 (71)	6.3±2.2 (71)	1.6±0.6 (71)	<b>1.3±0.6 (71)</b>	83.7±18.6 (71)	73.7±21.7 (71)	15.6±4.5 (71)	<b>8.2±2.0 (71)</b>
Toy small - 1.00e-05 - 0.01	17.6±6.0 (66)	14.1±9.8 (66)	7.1±5.3 (66)	<b>4.6±2.8 (66)</b>	88.2±23.3 (66)	154.6±93.6 (66)	67.0±44.5 (66)	<b>40.8±24.1 (66)</b>
Toy large - 1.00e-03 - 0.07	41.1±15.3 (365)	26.2±13.0 (365)	<b>5.8±1.3 (365)</b>	8.2±3.3 (365)	1040.8±0.0 (365)	355.9±83.8 (365)	82.7±19.3 (365)	<b>73.5±9.7 (365)</b>
Toy large - 1.00e-05 - 0.07	-	50.5±7.6 (371)	36.8±13.3 (371)	<b>31.7±7.4 (371)</b>	1356.7±178 (371)	1030.5±471.7 (371)	561.7±208.8 (371)	<b>465.6±111.4 (371)</b>
Toy large - 1.00e-03 - 0.01	589.5±185.4 (758)	91.6±22.9 (758)	65.4±14.5 (758)	<b>34.9±4.1 (758)</b>	52848.8±0.0 (758)	1192.1±340.1 (758)	777.5±181.5 (758)	<b>337.0±46.3 (758)</b>
Toy large - 1.00e-05 - 0.01	-	<b>583.8±140.7 (759)</b>	1020.6±250.6 (759)	609.4±177.6 (759)	60897±5990 (759)	7847±2774 (759)	12720±2520 (759)	<b>6699±1686 (759)</b>

Data and Setting	MM prox	GIST	MaxVC Gist	FireWorks Gist	MM BCD	BCD	MaxVC BCD	FireWorks BCD
Leukemia - 1.00e-03 - 0.07	6.3±2.0 (7)	17.9±0.4 (7)	<b>0.2±0.0 (7)</b>	0.4±0.0 (7)	3.8±0.7 (7)	144.4±1.1 (7)	<b>0.8±0.0 (7)</b>	<b>0.8±0.0 (7)</b>
Leukemia - 1.00e-05 - 0.07	8.0±2.7 (9)	26.1±0.6 (9)	<b>0.3±0.0 (9)</b>	0.5±0.0 (9)	4.6±1.1 (9)	218.8±1.1 (9)	1.2±0.0 (9)	<b>1.1±0.0 (9)</b>
Leukemia - 1.00e-03 - 0.01	31.4±6.2 (41)	186.1±1.7 (41)	<b>5.4±0.0 (41)</b>	5.5±0.0 (41)	53.6±9.6 (41)	1168.3±0.2 (41)	19.9±0.0 (41)	<b>17.4±0.0 (41)</b>
Leukemia - 1.00e-05 - 0.07	71.4±7.5 (46)	525.2±8.5 (46)	20.3±0.0 (46)	<b>14.6±0.0 (46)</b>	65.5±4.9 (46)	1412.8±0.3 (46)	71.5±0.0 (46)	<b>42.7±0.0 (46)</b>
NewsGroup-3 - 1.00e-02 - 0.01	955.8±389.1	6041.1±7.2	<b>6.5±0.0</b>	8.3±0.0	7926.6±3183.6	3792.4±6.2	<b>4.9±0.0</b>	5.6±0.0
NewsGroup-3 - 1.00e-03 - 0.01	1200.6±402.7	5790.6±8.0	49.8±0.1	<b>36.6±0.0</b>	12078.0±3879.1	24070.5±18	53.2±0.1	<b>36.8±0.0</b>
NewsGroup-3 - 1.00e-04 - 0.01	1237.9±415.5	5734.0±3.9	1439.3±2.4	<b>326.1±0.2</b>	12130.8±3849.7	37639.8±19	279.2±0.2	<b>167.7±0.1</b>
NewsGroup-7 - 1.00e-02 - 0.01	-	26711.1±44	1001.2±2.7	<b>343.6±0.9</b>	-	77378.7±74	421.7±0.8	<b>172.5±0.1</b>
NewsGroup-7 - 1.00e-03 - 0.01	-	26685.6±14	2163.6±4.4	<b>876.9±0.6</b>	-	91603.9±0.0	728.9±2.9	<b>312.3±0.6</b>
NewsGroup-7 - 1.00e-04 - 0.01	-	26752.5±15	4285.2±6.1	<b>1632.5±3.2</b>	-	117749.0±0.0	1093.7±3.7	<b>554.2±1.0</b>
Criteo - 1.00e-02 - 0.005	-	-	-	-	-	-	41095.3±2200	<b>31052.7±1200</b>
Criteo - 1.00e-03 - 0.005	-	-	-	-	-	-	49006.7±1400	<b>37534.6±1600</b>
Criteo - 1.00e-04 - 0.005	-	-	-	-	-	-	59303.8±1300	<b>42773.9±1000</b>

## 4 Numerical Experiments

**Set-up** We now present some numerical studies showing the computational gain achieved by our approach. As an inner solver and baseline algorithms, we have considered a proximal algorithm [16] and a block-coordinate descent approach [3]; they are respectively denoted as GIST and BCD. They have been implemented in Python/Numpy and the code will be shared online upon publication. We have integrated those solvers into the maximum-violating constraint (MaxVC) working set approach (algorithm in the appendix) and our approach denoted as FireWorks for FeasIble RESidual WORKing Set). Note that for MaxVC, we add the same number of constraints in the working set as in our algorithm. This is already a better baseline than the genuine one proposed in [1] As another baseline, we have considered a solver based on majorization-minimization (MM) approach, which consists in iteratively minimizing a majorization of the non-convex objective function as in [17, 13, 25]. Each iteration results in a weighted Lasso problem that we solve with a Blitz-based convex proximal Lasso or BCD Lasso (up to precision of  $10^{-5}$  for its optimality conditions). For these approaches, we leverage on the closed-form proximal operator available for several (non-convex)

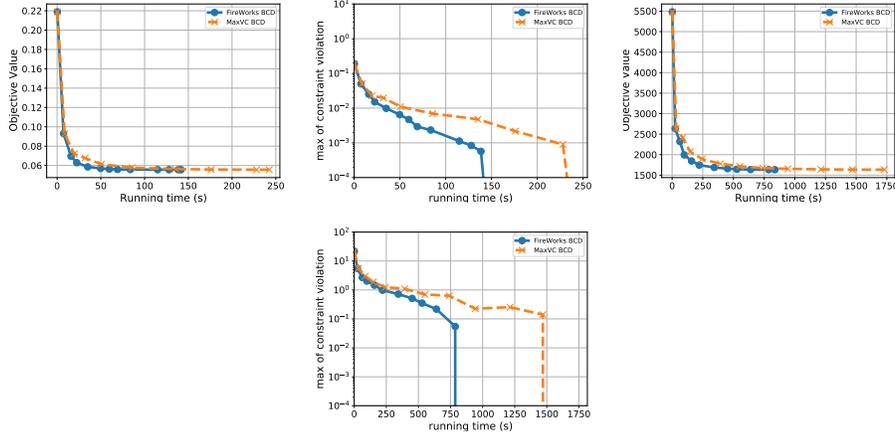


Figure 2: Example of evolution of the objective value and the maximum violation constraint on the 0-valued weights. The tolerance on the inner problem is set to  $10^{-6}$ . (most-left) performance on *NewsGroup-3*. (most-right) performance on *NewsGroup-7*.

regularizers. For our experiments, we have used the log-sum penalty which has an hyperparameter  $\theta$  that has been set to 1. For all algorithms, the stopping criterion is based on the tolerance over the Fermat’s optimality condition 2. The used performance measure for comparing all algorithms is the CPU running time. For all problems, we have set  $\tau_k$  so as to add the same fixed number  $n_k$  of features into the working set of MaxVC and our FireWorks algorithm.

**Toy problem** Here, the regression matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  is drawn uniformly from a standard Gaussian distribution (zero-mean unit variance). For given  $n, d$  and a number  $p$  of active variables, the true coefficient vector  $\mathbf{w}^{\text{true}}$  is obtained as follows. The  $p$  non-zero positions are chosen randomly, and their values are drawn from a zero-mean unit variance Gaussian distribution, to which we added  $\pm 0.1$  according to  $\text{sign}(w_j^{\text{true}})$ . Finally, the target vector is obtained as  $\mathbf{y} = \mathbf{X}\mathbf{w}^{\text{true}} + \mathbf{e}$  where  $\mathbf{e}$  is a zero-mean Gaussian noise with standard deviation  $\sigma = 0.01$ . For these problems,  $n_k = 30$  features are added to the working set at each iteration. Table 2 presents the running time for different algorithms to reach convergence under various settings. We note that our FireWorks algorithm is faster than the genuine inner solver and (at least on par) with the MaxVC approach especially in setting where  $\lambda$  is properly tuned with respect to the number of variables, *ie* when the solution is not too sparse. Note that the MM+Blitz approaches is performing worse than all other methods in almost all settings.

**Real data** We have reported comparisons on three real datasets. The first one is the *Leukemia* dataset [15] which has a dense regression matrix with  $n = 72$  and  $d = 7129$ . We have also considered sparse problem such as *newsgroups* dataset in which we have kept only 3 categories (*religion* and *graphics*) resulting in  $n = 1441$ ,  $d = 26488$  and 7 categories *comp* leading to  $n = 4891$ ,  $d = 94414$ . For these two problems, we have respectively 223173 and 676247 non-zeros elements in the related design matrix  $\mathbf{X}$ . We have also used a large-scale dataset which is a subset of the Criteo Kaggle dataset composed of 2M samples and 1M features, with about 78M non-zero elements in  $\mathbf{X}$ . For *Leukemia*, we have  $n_k = 30$  at each iteration, whereas we have added 300 and 1000 respectively for the *newsgroup* and *Criteo* problem.

Figure 2 presents a example of how objective value and maximum constraint violation (measured as  $\max_j (|\mathbf{x}_j^\top \mathbf{r}_k| - r'_\lambda(0))$ ) evolves during the optimization process for the two *NewsGroup* datasets. We see in those examples that both MaxVC and FireWorks algorithms achieve approximatively the same objective value whereas our FireWorks approach converges faster.

Quantitative results are reported in the bottom part of Table 2. At first, we can note that the convex relaxation approach using MM and Blitz is always more efficient than the baseline non-convex methods using either BCD or GIST. More globally, the table also shows that using FireWorks leads to a speedup of at least one order of magnitude compared to the baseline algorithm and the MM approach. For large  $\lambda$  leading to sparse solutions, MaxVC is the most efficient approach on *Leukemia*,

while for large-scale datasets *newsgroup-3* and *newsgroup-7*, FireWorks benefits from pruning and it is substantially faster than all competitors. For *Criteo*, only the BCD working set algorithms are able to terminate in reasonable time and FireWorks is more efficient than MaxVC. Again the MM+Blitz approach is performing worse than the two non-convex active set algorithms and fails to converge in a reasonable time for large datasets.

## 5 Conclusions

We have introduced in this paper a working set based meta-algorithm for non-convex regularized regression. By generalizing the concept of primal-dual approach, but in a non-convex setting, we were able to derive a novel rule for updating the variables optimized by an iterative incremental algorithm. From a theoretical point of view, we showed convergence of the algorithm, even when the inner problem is not solved exactly but up to a certain tolerance. This is in contrast with the classical maximal violating constraint approach which convergence requires the exact resolution of each inner problem. Our experimental results show the computational gain achieved for a given solver when applied directly on the full variables or within our working set algorithm.

As a future work, we plan to investigate other non-convex learning problems for which this relation between residual and dual variable can have theoretical or algorithmic impacts.

## References

- [1] A. Boisbunon, R. Flamary, and A. Rakotomamonjy. Active set strategy for high-dimensional non-convex sparse optimization problems. In *ICASSP*, pages 1517–1521. IEEE, 2014.
- [2] A. Boisbunon, R. Flamary, A. Rakotomamonjy, A. Giros, and J. Zerubia. Large scale sparse optimization for object detection in high resolution images. In *IEEE Workshop in Machine Learning for Signal Processing (MLSP)*, 2014.
- [3] P. Breheny and J. Huang. Coordinate descent algorithms for nonconvex penalized regression, with applications to biological feature selection. *Ann. Appl. Stat.*, 5(1):232, 2011.
- [4] E. J. Candès, M. B. Wakin, and S. P. Boyd. Enhancing sparsity by reweighted  $l_1$  minimization. *J. Fourier Anal. Applicat.*, 14(5-6):877–905, 2008.
- [5] S. Chen and D. Donoho. Basis pursuit. In IEEE, editor, *Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers*, 1994.
- [6] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001.
- [7] F. H. Clarke. *Method of Dynamic and Nonsmooth Optimization*. SIAM, 1989.
- [8] P. L. Combettes and V. R. Wajs. Signal recovery by proximal forward-backward splitting. *Multiscale Modeling & Simulation*, 4(4):1168–1200, 2005.
- [9] D. L. Donoho. Compressed sensing. *IEEE Trans. Inf. Theory*, 52(4):1289–1306, 2006.
- [10] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Amer. Statist. Assoc.*, 96(456):1348–1360, 2001.
- [11] R. Flamary, A. Rakotomamonjy, and G. Gasso. Importance sampling strategy for non-convex randomized block-coordinate descent. In *2015 IEEE 6th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, pages 301–304, 2015.
- [12] J. Friedman, T. J. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.*, 33(1):1–22, 2010.
- [13] Gilles Gasso, Alain Rakotomamonjy, and Stéphane Canu. Recovering sparse signals with a certain family of nonconvex penalties and dc programming. *IEEE Trans. Signal Process.*, 57(12):4686–4698, 2009.
- [14] T. Glasmachers and C. Igel. Maximum-gain working set selection for SVMs. *Journal of Machine Learning Research*, 7(Jul):1437–1466, 2006.
- [15] Todd R Golub, Donna K Slonim, Pablo Tamayo, Christine Huard, Michelle Gaasenbeek, Jill P Mesirov, Hilary Coller, Mignon L Loh, James R Downing, Mark A Caligiuri, et al. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *science*, 286(5439):531–537, 1999.

- [16] P. Gong, C. Zhang, Z. Lu, J. Huang, and J. Ye. A general iterative shrinkage and thresholding algorithm for non-convex regularized optimization problems. In *ICML*, pages 37–45, 2013.
- [17] David R Hunter and Kenneth Lange. A tutorial on MM algorithms. *The American Statistician*, 58(1):30–37, 2004.
- [18] T. B. Johnson and C. Guestrin. Blitz: A principled meta-algorithm for scaling sparse optimization. In *ICML*, volume 37, pages 1171–1179, 2015.
- [19] A Ya Kruger. On Fréchet subdifferentials. *Journal of Mathematical Sciences*, 116(3):3325–3358, 2003.
- [20] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly. Compressed sensing MRI. *IEEE Signal Processing Magazine*, 25(2):72–82, 2008.
- [21] M. Massias, A. Gramfort, and J. Salmon. From safe screening rules to working sets for faster lasso-type solvers. In *NIPS-OPT*, 2017.
- [22] M. Massias, A. Gramfort, and J. Salmon. Celer: a Fast Solver for the Lasso with Dual Extrapolation. In *ICML*, volume 80, pages 3315–3324, 2018.
- [23] B.S. Mordukhovich, N. M. Nam, and N. D. Yen. Fréchet subdifferential calculus and optimality conditions in nondifferentiable programming. *Optimization*, 55(5-6):685–708, 2006.
- [24] A. Rakotomamonjy, R. Flamary, and G. Gasso. DC proximal Newton for nonconvex optimization problems. *IEEE transactions on neural networks and learning systems*, 27(3):636–647, 2015.
- [25] A. Rakotomamonjy, G. Gasso, and J. Salmon. Screening rules for lasso with non-convex sparse regularizers. In *ICML*, volume 97, pages 5341–5350, 2019.
- [26] R. T. Rockafellar and R. J.-B. Wets. *Variational analysis*, volume 317. Springer Science & Business Media, 2009.
- [27] S. K. Shevade and S. S. Keerthi. A simple and efficient algorithm for gene selection using sparse logistic regression. *Bioinformatics*, 19(17):2246–2253, 2003.
- [28] E. Soubies, L. Blanc-Féraud, and G. Aubert. A unified view of exact continuous penalties for  $\ell_2$ - $\ell_0$  minimization. *SIAM J. Optim.*, 27(3):2034–2060, 2017.
- [29] R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 58(1):267–288, 1996.
- [30] R. Tibshirani, J. Bien, J. Friedman, T. J. Hastie, N. Simon, J. Taylor, and R. J. Tibshirani. Strong rules for discarding predictors in lasso-type problems. *J. R. Stat. Soc. Ser. B Stat. Methodol.*, 74(2):245–266, 2012.
- [31] S. Villa, S. Salzo, L. Baldassarre, and A. Verri. Accelerated and inexact forward-backward algorithms. *SIAM Journal on Optimization*, 23(3):1607–1633, 2013.
- [32] S. V. N. Vishwanathan, A. J. Smola, and M. N. Murty. Simplesvm. In *ICML*, pages 760–767, 2003.
- [33] R. Wang, N. Xiu, and S. Zhou. Fast Newton method for sparse logistic regression. *arXiv preprint arXiv:1901.02768*, 2019.
- [34] J. Ye and J. Liu. Sparse methods for biomedical data. *ACM Sigkdd Explorations Newsletter*, 14(1):4–15, 2012.
- [35] C.-H. Zhang. Nearly unbiased variable selection under minimax concave penalty. *Ann. Statist.*, 38(2):894–942, 2010.
- [36] T. Zhang. Analysis of multi-stage convex relaxation for sparse regularization. *Journal of Machine Learning Research*, 11(Mar):1081–1107, 2010.

# Supplementary material

## Provably Convergent Working Set Algorithm for Non-Convex Regularized Regression

### 5.1 Maximum-Violating Constraint Working Set Algorithm

The maximum-violating constraint algorithm is a simple algorithm that solves at each iteration a sub-problem with a subset of variables and then add some others that violate the most the statement  $\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}_{\mathcal{A}_k}^* \in \mathcal{C}_{\bar{\mathcal{A}}_k}$ , where “the most” is evaluated in term of distance to each set  $\mathcal{C}_j$ , with  $j \in \bar{\mathcal{A}}_k$ . Hence, at each iteration, we compute all these distances, sort them in descending order and add to the current working set, the  $n_k$  variables that yield to the largest distances. The algorithm is presented below.

---

#### Algorithm 2: Maximum Violating Constraints Algorithm

---

**Input:**  $\{\mathbf{X}, \mathbf{y}\}$ ,  $\mathcal{A}_1$  active set,  $n_k$  number of variables to add at iteration  $k$ , initial vector  $\tilde{\mathbf{w}}_{\mathcal{A}_0}$   
**Output:**  $\tilde{\mathbf{w}}_{\mathcal{A}_k}$

- 1: **for**  $k = 1, 2, \dots$  **do**
- 2:    $\mathbf{w}_{\mathcal{A}_k} = \arg \min_{\mathbf{w} \in \mathcal{A}_k} \frac{1}{2} \|\mathbf{y} - \mathbf{X}_{\mathcal{A}_k} \mathbf{w}\|_2^2 + \sum_{j \in \bar{\mathcal{A}}_k} r_\lambda(|w_j|)$ ;     // warm-start solver with  $\mathbf{w}_{\mathcal{A}_{k-1}}$
- 3:    $\mathbf{r}_k = \mathbf{y} - \mathbf{X}_{\mathcal{A}_k} \mathbf{w}_{\mathcal{A}_k}$ ;     // current residual
- 4:    $\mathbf{v} = \text{argsort dist}(\mathbf{r}_k, \mathcal{C}_j)$  in descending order
- 5:    $\mathcal{A}_{k+1} = \mathbf{v}[1 : n_k] \cup \mathcal{A}_k$ ;     // update working set by adding the  $n_k$  most violating variables
- 6: **end for**
- 7: Build  $\tilde{\mathbf{w}}_{\mathcal{A}_k}$

---

### 5.2 Proof of Proposition 2

**Proposition 2.** *Given a working set  $\mathcal{A}_k$  and  $\mathbf{w}_{\mathcal{A}_k}^*$  solving the related restricted problem,  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  is also optimal for the full problem if and only if  $\alpha = 1$  (which also means  $\mathbf{s}_{k+1} = \mathbf{r}_k$ ).*

*Proof.* Assume that  $\mathbf{w}_{\mathcal{A}_k}^*$  and  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  are optimal respectively for the restricted and the full problem. Let us show that in this case  $\alpha_k = 1$ . Since  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  is optimal for the full problem, we thus have  $\forall j \in \bar{\mathcal{A}}_k$ ,  $|\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}_{\mathcal{A}_k}^*)| \leq r'(0)$ . And thus we have the following equivalent statement

$$\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}_{\mathcal{A}_k}^* \in \mathcal{C} \Leftrightarrow \mathbf{y} - \mathbf{X}_{\mathcal{A}_k} \mathbf{w}_{\mathcal{A}_k}^* \in \mathcal{C} \Leftrightarrow \mathbf{r}_k \in \mathcal{C}$$

and thus  $\alpha_k = 1$ .

Now assume that  $\alpha_k = 1$  and let us show that  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  is optimal for the full problem. Since  $\alpha_k = 1$ , we have  $\mathbf{s}_{k+1} = \mathbf{r}_k$  and thus  $\mathbf{r}_k \in \mathcal{C}$ . The latter means that  $\forall j \in \bar{\mathcal{A}}_k$ ,  $|\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}_{\mathcal{A}_k} \mathbf{w}_{\mathcal{A}_k}^*)| \leq r'(0)$  and thus  $\forall j \in \bar{\mathcal{A}}_k$ ,  $|\mathbf{x}_j^\top (\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}_{\mathcal{A}_k}^*)| \leq r'(0)$ . Given this last property and the definition of  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  based on  $\mathbf{w}_{\mathcal{A}_k}^*$ , we can conclude that  $\tilde{\mathbf{w}}_{\mathcal{A}_k}^*$  is optimal for the full problem.  $\square$

### 5.3 Proof of Lemma 1

The proof follows similar steps as those given by Johnson and Guestrin [18].

**Lemma 1.** *At step  $k \geq 2$ , consider a constraint  $\mathcal{C}_j$  such that  $h_j(\mathbf{r}_k) > 0$  and  $h_j(\mathbf{s}_k) < 0$  then*

$$\text{dist}(\mathbf{r}_k, \mathcal{C}_j) \geq \frac{1 - \alpha_k}{\alpha_k} \tau_{k-1} . \tag{9}$$

*Proof.* Denote as  $j$  the index of the function  $h_j$  such that  $h_j(\mathbf{r}_k) > 0$  and  $h_j(\mathbf{s}_k) < 0$ . Let's  $\mathbf{z}_k \in \{\mathbf{z} \in \mathbb{R}^n : h_j(\mathbf{z}) = 0\}$ . The following equality holds

$$\begin{aligned}
\text{dist}(\mathbf{r}_k, \mathcal{C}_j) &= \|\mathbf{z}_k - \mathbf{r}_k\|_2 \\
&= \left\| \mathbf{z}_k - \frac{1}{\alpha_k}(\mathbf{s}_{k+1} - (1 - \alpha_k)\mathbf{s}_k) \right\| \\
&= \left\| \mathbf{z}_k - \frac{1}{\alpha_k}\mathbf{s}_{k+1} + \frac{1 - \alpha_k}{\alpha_k}\mathbf{s}_k \right\| \\
&= \left\| -\mathbf{z}_k + \frac{1}{\alpha_k}\mathbf{s}_{k+1} - \frac{1 - \alpha_k}{\alpha_k}\mathbf{s}_k \right\| \\
&= \frac{1 - \alpha_k}{\alpha_k} \left\| -\frac{\alpha_k}{1 - \alpha_k}\mathbf{z}_k + \frac{1}{1 - \alpha_k}\mathbf{s}_{k+1} - \mathbf{s}_k \right\| \tag{10}
\end{aligned}$$

Note that because  $h_j(\mathbf{r}_k) > 0$  and  $h_j(\mathbf{s}_k) < 0$ ,  $\alpha_k \neq 0$  since  $h_j$  is a continuous function. By construction, we have  $h_j(\mathbf{z}_k) = 0$  as  $\mathbf{z}_k$  is a minimizer of the distance and  $h_j(\mathbf{s}_{k+1}) = 0$  as we have chosen  $j$  as the index of the set that makes  $\mathbf{s}_{k+1} \notin \mathcal{C}$ . Since  $h_j(\cdot) \leq 0$  is a convex set and the coefficients  $-\frac{\alpha_k}{1 - \alpha_k}$  and  $\frac{1}{1 - \alpha_k}$  do not lead to a convex combination of  $\mathbf{z}_k$  and  $\mathbf{s}_{k+1}$  and hence, we have  $h_j(-\frac{\alpha_k}{1 - \alpha_k}\mathbf{z}_k + \frac{1}{1 - \alpha_k}\mathbf{s}_{k+1}) \geq 0$ . On the other hand by construction, we have  $\mathbf{s}_k \in \mathcal{C}_j$ . Furthermore, we have  $\text{dist}_S(\mathbf{s}_k, \mathcal{C}_j^-) \geq \tau_{k-1}$ . Indeed, since  $h_j(\mathbf{r}_k) > 0$ , we have  $j \notin \mathcal{A}_k$  as by construction  $\mathbf{r}_k \in \mathcal{C}_{\mathcal{A}}$  ( $\mathbf{w}_{\mathcal{A}_k}$  has been optimized over  $\mathcal{A}_k$ ). Because  $j \notin \mathcal{A}_k$  means that  $\text{dist}_S(\mathbf{s}_k, \mathcal{C}_j^-) \geq \tau_{k-1}$ , by definition of the construction of  $\mathcal{A}_k$  in Algorithm 1.

Now as  $h_j(-\frac{\alpha_k}{1 - \alpha_k}\mathbf{z}_k + \frac{1}{1 - \alpha_k}\mathbf{s}_{k+1}) \geq 0$  and  $\text{dist}_S(\mathbf{s}_k, \mathcal{C}_j^-) \geq \tau_k$ , the norm in the above equation (10) is lower bounded by  $\tau_k$  and we have

$$\text{dist}(\mathbf{r}_k, \mathcal{C}_j) \geq \frac{1 - \alpha_k}{\alpha_k} \tau_k.$$

□