

A Random Block-Coordinate Douglas-Rachford Splitting Method with Low Computational Complexity for Binary Logistic Regression

Luis Briceño-Arias, Giovanni Chierchia, Emilie Chouzenoux, Jean-Christophe Pesquet

► **To cite this version:**

Luis Briceño-Arias, Giovanni Chierchia, Emilie Chouzenoux, Jean-Christophe Pesquet. A Random Block-Coordinate Douglas-Rachford Splitting Method with Low Computational Complexity for Binary Logistic Regression. 2017. <hal-01672507>

HAL Id: hal-01672507

<https://hal.archives-ouvertes.fr/hal-01672507>

Submitted on 25 Dec 2017

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.

A Random Block-Coordinate Douglas-Rachford Splitting Method with Low Computational Complexity for Binary Logistic Regression

Luis M. Briceño-Arias ·
Giovanni Chierchia · Emilie Chouzenoux ·
Jean-Christophe Pesquet

the date of receipt and acceptance should be inserted later

Abstract In this paper, we propose a new optimization algorithm for sparse logistic regression based on a stochastic version of the Douglas-Rachford splitting method. Our algorithm sweeps the training set by randomly selecting a mini-batch of data at each iteration, and it allows us to update the variables in a block coordinate manner. Our approach leverages the proximity operator of the logistic loss, which is expressed with the generalized Lambert W function. Experiments carried out on standard datasets demonstrate the efficiency of our approach w.r.t. stochastic gradient-like methods.

Keywords Proximity operator · Douglas-Rachford splitting · Block-coordinate descent · Logistic regression

1 Introduction

Sparse classification algorithms have gained much popularity in the context of supervised learning, thanks to their ability to discard irrelevant features during the training stage. Such algorithms aim at learning a weighted linear combination of basis functions that fits the training data, while encouraging as many weights as possible to be equal to zero. This amounts to solving an optimization problem that involves a loss function plus a sparse regularization term. Different types of classifiers arise by varying the loss function, the most popular being the hinge and the logistic losses [1, 2].

This work was partly supported by the the CNRS MASTODONS project under grant 2016TABASCO.

L. M. Briceño-Arias
Departamento de Matemática, Universidad Técnica Federico Santa María, Av Espanã 1681, Valparaíso, Chile.

G. Chierchia (corresponding author) and E. Chouzenoux
Université Paris Est, LIGM, CNRS UMR 8049, ESIEE Paris, UPEM, Noisy-le-Grand, France.

E. Chouzenoux and J.-C. Pesquet
Center for Visual Computing, INRIA Saclay, CentraleSupélec, University Paris-Saclay, Gif sur Yvette, France.

In the context of supervised learning, sparse regularization traces back to the work of Bradley and Mangasarian [3], who showed that the ℓ_1 -norm can efficiently perform feature selection by shrinking small coefficients to zero. Other forms of regularization have also been studied, such as the ℓ_0 -norm [4], the ℓ_p -norm with $p > 0$ [5], the ℓ_∞ -norm [6], and other nonconvex terms [7]. Mixed-norms have been investigated as well, due to their ability to impose a more structured form of sparsity [8–13].

Many efficient learning algorithms exist in the case of quadratic regularization, by benefiting from the advantages brought by Lagrangian duality [14]. This is unfortunately not true for sparse regularization, because the dual formulation is as difficult to solve as the primal one. Consequently, sparse linear classifiers are usually trained through the direct resolution of the primal optimization problem. Among the possible approaches, one can resort to linear programming [15], gradient-like methods [8, 16], proximal algorithms [7, 17, 18], and other optimization techniques [19].

Nowadays, it is well known that random updates can significantly reduce the computational time when a quadratic regularization is used [20, 21]. Therefore, a great deal of attention has been paid recently to stochastic approaches capable of handling a sparse regularization [22]. The list of investigated techniques includes block-coordinate descent strategies [23–26], stochastic forward-backward iterations [27–31], random Douglas-Rachford splitting methods [32], random primal-dual proximal algorithms [33], and stochastic majorization-minimization methods [34, 35].

In this paper, we propose a random-sweeping block-coordinate Douglas-Rachford splitting method. In addition to the stochastic behavior, it presents three distinctive features with respect to related approaches [32, 36, 37]. Firstly, the matrix to be inverted at the initial step is block-diagonal, while in the concurrent approaches, it did not present any specific structure. The block diagonal property implies that the inversion step actually amounts to inverting a set of smaller size matrices. Secondly, the proposed algorithm can take advantage explicitly from a strong convexity property possibly fulfilled by some of the functions involved in the optimization problem. Finally, the dual variables appear explicitly in the proposed scheme, making it possible to use clever block-coordinate descent strategies [38].

Moreover, the proposed algorithm appears to be well tailored to binary logistic regression with sparse regularization. In contrast to gradient-like methods, our approach deals with the logistic loss through its proximity operator. This results in an algorithm that is not tied up to the Lipschitz constant of the loss function, possibly leading to larger updates per iteration. In this regard, our second contribution is to show that the proximity operator of the binary logistic loss can be expressed in closed form using the generalized Lambert W function [39, 40]. We also provide comparisons with state-of-the-art stochastic methods using benchmark datasets.

The paper is organized as follows. In Section 2, we derive the proposed Douglas-Rachford algorithm. In Section 3, we introduce sparse logistic regression, along with the proximal operator of the logistic loss. In Section 4, we evaluate our approach on standard datasets, and compare it to three sparse classification algorithms proposed in the literature [28, 30, 41]. Finally, conclusions are drawn in Section 5.

NOTATION: $\Gamma_0(\mathcal{H})$ denotes the set of proper, lower semicontinuous, convex functions from a real Hilbert space \mathcal{H} to $] -\infty, +\infty]$. Let $\psi \in \Gamma_0(\mathcal{H})$. For every $\nu \in \mathcal{H}$, the subdifferential of ψ at ν is $\partial\psi(\nu) = \{\xi \in \mathcal{H} \mid (\forall \zeta \in \mathcal{H}) \langle \zeta - \nu \mid \xi \rangle + \psi(\nu) \leq \psi(\zeta)\}$, the proximity operator of ψ at ν is $\text{prox}_\psi(\nu) = \text{argmin}_{\xi \in \mathcal{H}} \frac{1}{2}\|\xi - \nu\|^2 + \psi(\nu)$, and the conjugate of ψ is $\psi^* = \sup_{\xi \in \mathcal{H}} \langle \xi \mid \cdot \rangle - \psi(\xi)$ in $\Gamma_0(\mathcal{H})$. The adjoint of a bounded

linear operator A from \mathcal{H} to a real Hilbert space \mathcal{G} is denoted by A^* . Let $(\Omega, \mathcal{F}, \mathbb{P})$ be the underlying probability space, the σ -algebra generated by a family Φ of random variables is denoted by $\sigma(\Phi)$.

2 Optimization method

Throughout this section, $\mathcal{H}_1, \dots, \mathcal{H}_B, \mathcal{G}_1, \dots, \mathcal{G}_L$ are separable real Hilbert spaces. In addition, $\mathcal{H} = \mathcal{H}_1 \oplus \dots \oplus \mathcal{H}_B$ denotes the Hilbertian sum of $\mathcal{H}_1, \dots, \mathcal{H}_B$. Any vector $\mathbf{v} \in \mathcal{H}$ can thus be uniquely decomposed as $(v_b)_{1 \leq b \leq B}$ where, for every $b \in \{1, \dots, B\}$, $v_b \in \mathcal{H}_b$. In the following, a similar notation will be used to denote vectors in any product space (in bold) and their components.

We will now aim at solving the following problem.

Problem 1 For every $b \in \{1, \dots, B\}$ and for every $\ell \in \{1, \dots, L\}$, let $f_b \in \Gamma_0(\mathcal{H}_b)$, let $h_\ell: \mathcal{G}_\ell \rightarrow \mathbb{R}$ be a differentiable convex function with β_ℓ -Lipschitz gradient, for some $\beta_\ell \in]0, +\infty[$, and let $A_{\ell,b}$ be a linear bounded operator from \mathcal{H}_b to \mathcal{G}_ℓ . The problem is to

$$\underset{\mathbf{w} \in \mathcal{H}}{\text{minimize}} \quad \sum_{b=1}^B f_b(w_b) + \sum_{\ell=1}^L h_\ell \left(\sum_{b=1}^B A_{\ell,b} w_b \right),$$

under the assumption that the set of solutions \mathcal{E} is nonempty.

In order to address Problem 1, we propose to employ the random-sweeping block-coordinate version of the Douglas-Rachford splitting method with stochastic errors described in Algorithm 1. Let us define

$$(\forall b \in \{1, \dots, B\}) \quad C_b = \left(\text{Id} + \tau_b \sum_{\ell=1}^L \frac{\gamma_\ell}{1 + \gamma_\ell \rho_\ell} A_{\ell,b}^* A_{\ell,b} \right)^{-1} : \mathcal{H}_b \rightarrow \mathcal{H}_b, \quad (1)$$

where $(\tau_b)_{1 \leq b \leq B}$ and $(\gamma_\ell)_{1 \leq \ell \leq L}$ are the positive constants introduced in Algorithm 1. The next result establishes the convergence of the proposed algorithm.

Proposition 1 For every $b \in \{1, \dots, B\}$, let $w_b^{[0]}, t_b^{[0]}$ and $(a_b^{[i]})_{i \in \mathbb{N}}$ be \mathcal{H}_b -valued random variables and, for every $\ell \in \{1, \dots, L\}$, let $\mathbf{v}_\ell^{[0]}$ and $\mathbf{s}_\ell^{[0]}$ be \mathcal{G}_ℓ^B -valued random variables and let $(d_\ell^{[i]})_{i \in \mathbb{N}}$ be \mathcal{G}_ℓ -valued random variables. In addition, let $(\varepsilon^{[i]})_{i \in \mathbb{N}}$ be identically distributed $\{0, 1\}^{B+L} \setminus \{\mathbf{0}\}$ -valued random variables and in Algorithm 1 assume that

- (i) $(\forall i \in \mathbb{N}) \quad \sigma(\varepsilon^{[i]})$ and $\boldsymbol{\chi}^{[i]} = \sigma(\mathbf{t}^{[0]}, \dots, \mathbf{t}^{[i]}, \mathbf{s}^{[0]}, \dots, \mathbf{s}^{[i]})$ are independent;
- (ii) $(\forall b \in \{1, \dots, B\}) \quad \sum_{i \in \mathbb{N}} \sqrt{\mathbb{E}(\|a_b^{[i]}\|^2 | \boldsymbol{\chi}^{[i]})} < +\infty$;
- (iii) $(\forall \ell \in \{1, \dots, L\}) \quad \sum_{i \in \mathbb{N}} \sqrt{\mathbb{E}(\|d_\ell^{[i]}\|^2 | \boldsymbol{\chi}^{[i]})} < +\infty$;
- (iv) $(\forall b \in \{1, \dots, B\}) \quad \mathbb{P}[\varepsilon_b^{[0]} = 1] > 0$ and $(\forall \ell \in \{1, \dots, L\}) \quad \mathbb{P}[\varepsilon_{B+\ell}^{[0]} = 1] > 0$.

Then, the sequence $(\mathbf{w}^{[i]})_{i \in \mathbb{N}}$ generated by Algorithm 1 converges weakly P-a.s. to an \mathcal{E} -valued random variable.

Algorithm 1 Random Douglas-Rachford splitting for solving Problem 1

INITIALIZATION

Set $(\tau_b)_{1 \leq b \leq B} \in]0, +\infty[^B$ and $\eta \in]0, 1]$.
 For every $\ell \in \{1, \dots, L\}$, set $\rho_\ell \geq 0$ such that $B\beta_\ell\rho_\ell \leq 1$.
 For every $\ell \in \{1, \dots, L\}$, set $\gamma_\ell > 0$ such that $\gamma_\ell\rho_\ell < 1$.

$$(\forall b \in \{1, \dots, B\}) \quad u_b^{[0]} = \sum_{\ell=1}^L \frac{1}{1 + \gamma_\ell\rho_\ell} A_{\ell,b}^* s_{\ell,b}^{[0]}$$

FOR $i = 0, 1, \dots$ Set $\mu^{[i]} \in]\eta, 2 - \eta[$ for $b = 1, \dots, B$

$$\begin{cases} w_b^{[i+1]} = w_b^{[i]} + \varepsilon_b^{[i]} \left(C_b \left(t_b^{[i]} - \tau_b u_b^{[i]} \right) - w_b^{[i]} \right) \\ t_b^{[i+1]} = t_b^{[i]} + \varepsilon_b^{[i]} \mu^{[i]} \left(\text{prox}_{\tau_b f_b} (2w_b^{[i+1]} - t_b^{[i]}) + a_b^{[i]} - w_b^{[i+1]} \right) \end{cases}$$

for $\ell = 1, \dots, L$

$$v_\ell^{[i+1]} = v_\ell^{[i]} + \varepsilon_{B+\ell}^{[i]} \left(\frac{s_\ell^{[i]} + \gamma_\ell (A_{\ell,b} w_b^{[i]})_{1 \leq b \leq B}}{1 + \gamma_\ell\rho_\ell} - v_\ell^{[i]} \right)$$

$$p_\ell^{[i]} = 2 \sum_{b=1}^B v_{\ell,b}^{[i+1]} - \sum_{b=1}^B s_{\ell,b}^{[i]}$$

$$q_\ell^{[i]} = \text{prox}_{\frac{B(1-\gamma_\ell\rho_\ell)}{\gamma_\ell} h_\ell} (p_\ell^{[i]}/\gamma_\ell) + d_\ell^{[i]}$$

for $b = 1, \dots, B$

$$s_{\ell,b}^{[i+1]} = s_{\ell,b}^{[i]} + \varepsilon_{B+\ell}^{[i]} \mu^{[i]} \left(\frac{p_\ell^{[i]} - \gamma_\ell q_\ell^{[i]}}{B(1 - \gamma_\ell\rho_\ell)} - v_{\ell,b}^{[i+1]} \right)$$

for $b = 1, \dots, B$

$$u_b^{[i+1]} = u_b^{[i]} + \sum_{\ell=1}^L \frac{\varepsilon_{B+\ell}^{[i]}}{1 + \gamma_\ell\rho_\ell} A_{\ell,b}^* (s_{\ell,b}^{[i+1]} - s_{\ell,b}^{[i]}).$$

Proof Problem 1 can be reformulated as minimizing $\mathbf{f} + \mathbf{h} \circ \mathbf{A}$ where

$$\mathbf{f}: \mathcal{H} \rightarrow]-\infty, +\infty]: \mathbf{w} \mapsto \sum_{b=1}^B f_b(w_b) \quad (2)$$

$$\mathbf{A}: \mathcal{H} \rightarrow \mathcal{G}: \mathbf{w} \mapsto (A_{\ell,1} w_1, \dots, A_{\ell,B} w_B)_{1 \leq \ell \leq L} \quad (3)$$

$$\mathbf{h}: \mathcal{G} \rightarrow \mathbb{R}: \mathbf{v} \mapsto \sum_{\ell=1}^L h_\ell(\Lambda_\ell \mathbf{v}_\ell) \quad (4)$$

$$(\forall \ell \in \{1, \dots, L\}) \quad \Lambda_\ell: \mathcal{G}_\ell^B \rightarrow \mathcal{G}_\ell: \mathbf{v}_\ell \mapsto \sum_{b=1}^B v_{\ell,b} \quad (5)$$

and $\mathbf{v} = (\mathbf{v}_\ell)_{1 \leq \ell \leq L}$ denotes a generic element of $\mathcal{G} = \mathcal{G}_1^B \oplus \cdots \oplus \mathcal{G}_L^B$ with $\mathbf{v}_\ell = (v_{\ell,b})_{1 \leq b \leq B} \in \mathcal{G}_\ell^B$ for every $\ell \in \{1, \dots, L\}$. Since $\text{dom}(\mathbf{h}) = \mathcal{G}$, from [42, Theorem 16.47(i)], Problem 1 is equivalent to

$$\text{find } \mathbf{w} \in \mathcal{H} \text{ such that } \mathbf{0} \in \partial \mathbf{f}(\mathbf{w}) + \mathbf{A}^* \nabla \mathbf{h}(\mathbf{A}\mathbf{w}), \quad (6)$$

which, from [36, Proposition 2.8] is also equivalent to

$$\text{find } (\mathbf{w}, \mathbf{v}) \in \mathcal{H} \times \mathcal{G} \text{ such that } (\mathbf{0}, \mathbf{0}) \in \mathbf{N}(\mathbf{w}, \mathbf{v}) + \mathbf{S}(\mathbf{w}, \mathbf{v}), \quad (7)$$

where $\mathbf{N}: (\mathbf{w}, \mathbf{v}) \mapsto \partial \mathbf{f}(\mathbf{w}) \times \partial \mathbf{h}^*(\mathbf{v})$ is maximally monotone and $\mathbf{S}: (\mathbf{w}, \mathbf{v}) \mapsto (\mathbf{A}^* \mathbf{v}, -\mathbf{A}\mathbf{w})$ is a skewed linear operator. Note that $\mathbf{A}^*: \mathbf{v} \mapsto (\sum_{\ell=1}^L A_{\ell,b}^* v_{\ell,b})_{1 \leq b \leq B}$ and, from (2), (4) and [42, Proposition 13.30 and Proposition 16.9], $\partial \mathbf{f}: \mathbf{w} \mapsto \times_{b=1}^B \partial f_b(w_b)$ and $\partial \mathbf{h}^*: \mathbf{v} \mapsto \times_{\ell=1}^L \partial(h_\ell \circ \mathbf{A}_\ell)^*(\mathbf{v}_\ell)$. Since, for every $\ell \in \{1, \dots, L\}$, $h_\ell \circ \mathbf{A}_\ell$ is convex differentiable with a $B\beta_\ell$ -Lipschitzian gradient $\nabla(h_\ell \circ \mathbf{A}_\ell) = \mathbf{A}_\ell^* \circ \nabla h_\ell \circ \mathbf{A}_\ell$, it follows from Baillon-Haddad theorem [42, Corollary 18.17] that $\nabla(h_\ell \circ \mathbf{A}_\ell)$ is $(B\beta_\ell)^{-1}$ -cocoercive and, hence, $\partial(h_\ell \circ \mathbf{A}_\ell)^* = (\nabla(h_\ell \circ \mathbf{A}_\ell))^{-1}$ is $(B\beta_\ell)^{-1}$ -strongly monotone. Therefore, for every $\rho_\ell \in [0, (B\beta_\ell)^{-1}]$, $(h_\ell \circ \mathbf{A}_\ell)^*$ is ρ_ℓ -strongly convex. By defining

$$(\forall \ell \in \{1, \dots, L\}) \quad \varphi_\ell = (h_\ell \circ \mathbf{A}_\ell)^* - \rho_\ell \|\cdot\|^2/2, \quad (8)$$

it follows from [42, Proposition 10.8] that, for every $\ell \in \{1, \dots, L\}$, $\varphi_\ell \in \Gamma_0(\mathcal{G}_\ell^B)$ and, hence, $\partial \varphi_\ell := \partial(h_\ell \circ \mathbf{A}_\ell)^* - \rho_\ell \mathbf{Id}$ is maximally monotone. Consequently, Problem 1 can be rewritten equivalently as

$$\text{find } (\mathbf{w}, \mathbf{v}) \in \mathcal{H} \times \mathcal{G} \text{ such that } \begin{cases} (\forall b \in \{1, \dots, B\}) \quad 0 \in \partial f_b(w_b) + B_b(\mathbf{w}, \mathbf{v}) \\ (\forall \ell \in \{1, \dots, L\}) \quad 0 \in \partial \varphi_\ell(\mathbf{v}_\ell) + B_\ell(\mathbf{w}, \mathbf{v}), \end{cases} \quad (9)$$

which, for strictly positive constants $(\tau_b)_{1 \leq b \leq B}$ and $(\gamma_\ell)_{1 \leq \ell \leq L}$, is equivalent to

$$\text{find } (\mathbf{w}, \mathbf{v}) \in \mathcal{H} \times \mathcal{G} \text{ such that } \begin{cases} (\forall b \in \{1, \dots, B\}) \quad 0 \in \tau_b \partial f_b(w_b) + \tau_b B_b(\mathbf{w}, \mathbf{v}) \\ (\forall \ell \in \{1, \dots, L\}) \quad 0 \in \gamma_\ell \partial \varphi_\ell(\mathbf{v}_\ell) + \gamma_\ell B_\ell(\mathbf{w}, \mathbf{v}), \end{cases} \quad (10)$$

where

$$\begin{cases} B_b: (\mathbf{w}, \mathbf{v}) \mapsto \sum_{\ell=1}^L A_{\ell,b}^* v_{\ell,b} \\ B_\ell: (\mathbf{w}, \mathbf{v}) \mapsto -(A_{\ell,1} w_1, \dots, A_{\ell,B} w_B) + \rho_\ell \mathbf{v}_\ell. \end{cases} \quad (11)$$

Since $\mathbf{S}: (\mathbf{w}, \mathbf{v}) \mapsto (\mathbf{A}^* \mathbf{v}, -\mathbf{A}\mathbf{w})$ and $\mathbf{D}: \mathcal{G} \rightarrow \mathcal{G}: \mathbf{v} \mapsto (\rho_\ell \mathbf{v}_\ell)_{1 \leq \ell \leq L}$ are linear and monotone operators in $\mathcal{H} \times \mathcal{G}$ and \mathcal{G} , respectively, the operator

$$\mathbf{B}: (\mathbf{w}, \mathbf{v}) \mapsto (\mathbf{A}^* \mathbf{v}, -\mathbf{A}\mathbf{w} + \mathbf{D}\mathbf{v}) = ((B_b(\mathbf{w}, \mathbf{v}))_{1 \leq b \leq B}, (B_\ell(\mathbf{w}, \mathbf{v}))_{1 \leq \ell \leq L})$$

is maximally monotone in $\mathcal{H} \times \mathcal{G}$. Therefore, by defining the strongly positive diagonal linear operator

$$\begin{aligned} \mathbf{U}: \mathcal{H} \times \mathcal{G} &\rightarrow \mathcal{H} \times \mathcal{G} \\ (\mathbf{w}, \mathbf{v}) &\mapsto (\mathbf{T}\mathbf{w}, \mathbf{\Gamma}\mathbf{v}), \end{aligned} \quad (12)$$

where $\mathbf{T}: \mathbf{w} \mapsto (\tau_b w_b)_{1 \leq b \leq B}$ and $\mathbf{\Gamma}: \mathbf{v} \mapsto (\gamma_\ell \mathbf{v}_\ell)_{1 \leq \ell \leq L}$, the operator

$$\mathbf{UB}: (\mathbf{w}, \mathbf{v}) \mapsto (\mathbf{TA}^* \mathbf{v}, -\mathbf{\Gamma A w} + \mathbf{\Gamma D v}) = ((\tau_b B_b(\mathbf{w}, \mathbf{v}))_{1 \leq b \leq B}, (\gamma_\ell B_\ell(\mathbf{w}, \mathbf{v}))_{1 \leq \ell \leq L}) \quad (13)$$

is maximally monotone in $(\mathcal{H} \times \mathcal{G}, \|\cdot\|_{U^{-1}})$, where

$$(\forall (\mathbf{w}, \mathbf{v}) \in \mathcal{H} \times \mathcal{G}) \quad \|(\mathbf{w}, \mathbf{v})\|_{U^{-1}} = \sqrt{\sum_{b=1}^B \tau_b^{-1} \|w_b\|^2 + \sum_{\ell=1}^L \gamma_\ell^{-1} \|\mathbf{v}_\ell\|^2}. \quad (14)$$

Note that the renormed product space $(\mathcal{H} \times \mathcal{G}, \|\cdot\|_{U^{-1}})$ is the Hilbert sum $\mathcal{H}_1 \oplus \cdots \oplus \mathcal{H}_B \oplus \mathcal{G}_1^B \oplus \cdots \oplus \mathcal{G}_L^B$ where, for every $b \in \{1, \dots, B\}$ and $\ell \in \{1, \dots, L\}$, \mathcal{H}_b and \mathcal{G}_ℓ^B are endowed by the norm $\|\cdot\|_{\tau_b}: w_b \mapsto \|w_b\|/\sqrt{\tau_b}$ and $\|\cdot\|_{\gamma_\ell}: \mathbf{v}_\ell \mapsto \|\mathbf{v}_\ell\|/\sqrt{\gamma_\ell}$, respectively. Therefore, since $\tau_b \partial f_b$ and $\gamma_\ell \partial \varphi_\ell$ are maximally monotone in $(\mathcal{H}_b, \|\cdot\|_{\tau_b})$ and $(\mathcal{G}_\ell^B, \|\cdot\|_{\gamma_\ell})$, respectively, we conclude that (10) is a particular case of the primal inclusion in [32, Proposition 5.1].

Now we write Algorithm 1 as a particular case of the random block-coordinate Douglas-Rachford splitting proposed in [32, Proposition 5.1] applied to (10) in $(\mathcal{H} \times \mathcal{G}, \|\cdot\|_{U^{-1}})$. Given $(\mathbf{t}, \mathbf{s}) \in \mathcal{H} \times \mathcal{G}$, let $(\mathbf{w}, \mathbf{v}) = J_{\mathbf{UB}}(\mathbf{t}, \mathbf{s}) = (\mathbf{Id} + \mathbf{UB})^{-1}(\mathbf{t}, \mathbf{s})$. It follows from (13) that

$$\begin{cases} \mathbf{w} = \mathbf{t} - \mathbf{TA}^* \mathbf{v} \\ \mathbf{v} = (\mathbf{Id} + \mathbf{\Gamma D})^{-1}(\mathbf{s} + \mathbf{\Gamma A w}), \end{cases} \quad (15)$$

which leads to

$$\mathbf{w} = (\mathbf{Id} + \mathbf{TA}^*(\mathbf{Id} + \mathbf{\Gamma D})^{-1} \mathbf{\Gamma A})^{-1} (\mathbf{t} - \mathbf{TA}^*(\mathbf{Id} + \mathbf{\Gamma D})^{-1} \mathbf{s}). \quad (16)$$

In order to derive an explicit formula for the matrix inversion in (16), set $\mathbf{z} = \mathbf{t} - \mathbf{TA}^*(\mathbf{Id} + \mathbf{\Gamma D})^{-1} \mathbf{s}$. We have $\mathbf{z} = \mathbf{w} + \mathbf{TA}^*(\mathbf{Id} + \mathbf{\Gamma D})^{-1} \mathbf{\Gamma A w}$ and, since (3) and \mathbf{D} is diagonal, we obtain

$$\mathbf{TA}^*(\mathbf{Id} + \mathbf{\Gamma D})^{-1} \mathbf{\Gamma A}: \mathbf{w} \mapsto \left(\tau_b \sum_{\ell=1}^L \frac{\gamma_\ell A_{\ell,b}^* A_{\ell,b} w_b}{1 + \gamma_\ell \rho_\ell} \right)_{1 \leq b \leq B},$$

and, hence,

$$(\forall b \in \{1, \dots, B\}) \quad w_b = \left(\mathbf{Id} + \tau_b \sum_{\ell=1}^L \frac{\gamma_\ell A_{\ell,b}^* A_{\ell,b}}{1 + \gamma_\ell \rho_\ell} \right)^{-1} z_b = C_b z_b. \quad (17)$$

Therefore, (16) can be written equivalently as

$$(\forall b \in \{1, \dots, B\}) \quad w_b = C_b \left(t_b^{[i]} - \tau_b \sum_{\ell=1}^L \frac{A_{\ell,b}^* s_{\ell,b}}{1 + \gamma_\ell \rho_\ell} \right) = C_b (t_b - \tau_b u_b), \quad (18)$$

where

$$(\forall b \in \{1, \dots, B\}) \quad u_b = \sum_{\ell=1}^L \frac{A_{\ell,b}^* s_{\ell,b}}{1 + \gamma_\ell \rho_\ell}. \quad (19)$$

Moreover, from (15), we deduce that

$$(\forall \ell \in \{1, \dots, L\}) \quad \mathbf{v}_\ell = \frac{\mathbf{s}_\ell + \gamma_\ell (A_{\ell,b} w_b)_{1 \leq b \leq B}}{1 + \gamma_\ell \rho_\ell}, \quad (20)$$

and, hence, we have $J_{UB}: (\mathbf{t}, \mathbf{s}) \mapsto ((Q_b(\mathbf{t}, \mathbf{s}))_{1 \leq b \leq B}, (Q_\ell(\mathbf{t}, \mathbf{s}))_{1 \leq \ell \leq L})$, where

$$\begin{cases} (\forall b \in \{1, \dots, B\}) & Q_b: (\mathbf{t}, \mathbf{s}) \mapsto C_b(t_b - \tau_b u_b) \\ (\forall \ell \in \{1, \dots, L\}) & Q_\ell: (\mathbf{t}, \mathbf{s}) \mapsto \frac{\mathbf{s}_\ell + \gamma_\ell (A_{\ell,b} Q_b(\mathbf{t}, \mathbf{s}))_{1 \leq b \leq B}}{1 + \gamma_\ell \rho_\ell}. \end{cases} \quad (21)$$

Now, it follows from [42, Proposition 16.44] that, for every $b \in \{1, \dots, B\}$ and $\ell \in \{1, \dots, L\}$, $J_{\tau_b \partial f_b} = \text{prox}_{\tau_b f_b}$ and $J_{\gamma_\ell \partial \varphi_\ell} = \text{prox}_{\gamma_\ell \varphi_\ell}$ and, for every $\ell \in \{1, \dots, L\}$ and $(\mathbf{r}_\ell, \mathbf{z}_\ell) \in \mathcal{G}_\ell^B \times \mathcal{G}_\ell^B$, we have

$$\begin{aligned} \mathbf{r}_\ell = \text{prox}_{\gamma_\ell \varphi_\ell} \mathbf{z}_\ell &\Leftrightarrow \frac{\mathbf{z}_\ell - \mathbf{r}_\ell}{\gamma_\ell} \in \partial \varphi_\ell(\mathbf{r}_\ell) \\ &\Leftrightarrow \frac{\mathbf{z}_\ell - \mathbf{r}_\ell}{\gamma_\ell} \in \partial(h_\ell \circ \mathbf{A}_\ell)^* \mathbf{r}_\ell - \rho_\ell \mathbf{r}_\ell \\ &\Leftrightarrow \mathbf{z}_\ell - (1 - \gamma_\ell \rho_\ell) \mathbf{r}_\ell \in \gamma_\ell \partial(h_\ell \circ \mathbf{A}_\ell)^* \mathbf{r}_\ell. \end{aligned} \quad (22)$$

Therefore, if $\gamma_\ell \rho_\ell < 1$ we have that (22) is equivalent to

$$\mathbf{r}_\ell = \text{prox}_{\frac{\gamma_\ell}{1 - \gamma_\ell \rho_\ell} (h_\ell \circ \mathbf{A}_\ell)^*} \left(\frac{\mathbf{z}_\ell}{1 - \gamma_\ell \rho_\ell} \right) \quad (23)$$

and, from Moreau's decomposition formula [42, Theorem 14.13(ii)], we obtain

$$\mathbf{r}_\ell = \frac{1}{1 - \gamma_\ell \rho_\ell} \left(\mathbf{z}_\ell - \gamma_\ell \text{prox}_{\frac{1 - \gamma_\ell \rho_\ell}{\gamma_\ell} (h_\ell \circ \mathbf{A}_\ell)} \left(\frac{\mathbf{z}_\ell}{\gamma_\ell} \right) \right). \quad (24)$$

Noting that, for every $\ell \in \{1, \dots, L\}$, $\mathbf{A}_\ell \circ \mathbf{A}_\ell^* = B \text{Id}$, from [42, Proposition 24.14], (22) and (24) we deduce that

$$(\forall b \in \{1, \dots, B\}) \quad r_{\ell,b} = \frac{1}{B(1 - \gamma_\ell \rho_\ell)} \left(\sum_{d=1}^B z_{\ell,d} - \gamma_\ell \text{prox}_{\frac{B(1 - \gamma_\ell \rho_\ell)}{\gamma_\ell} h_\ell} \left(\frac{1}{\gamma_\ell} \sum_{d=1}^B z_{\ell,d} \right) \right) \quad (25)$$

and, hence,

$$\text{prox}_{\gamma_\ell \varphi_\ell} \mathbf{z}_\ell = \frac{1}{B(1 - \gamma_\ell \rho_\ell)} \left(\sum_{d=1}^B z_{\ell,d} - \gamma_\ell \text{prox}_{\frac{B(1 - \gamma_\ell \rho_\ell)}{\gamma_\ell} h_\ell} \left(\frac{1}{\gamma_\ell} \sum_{d=1}^B z_{\ell,d} \right) \right)_{1 \leq b \leq B}. \quad (26)$$

Therefore, by defining, for every $i \in \mathbb{N}$ and $\ell \in \{1, \dots, L\}$, $\mathbf{e}_\ell^{[i]} \in \mathcal{G}_\ell^B$ via

$$(\forall \ell \in \{1, \dots, L\}) \quad \mathbf{e}_\ell^{[i]} = \left(-\frac{\gamma_\ell}{B(1 - \gamma_\ell \rho_\ell)} d_\ell^{[i]} \right)_{1 \leq b \leq B}, \quad (27)$$

we deduce that Algorithm 1 can be written equivalently as

For $i = 0, 1, \dots$

$$\left[\begin{array}{l} \text{For } b = 1, \dots, B \\ \left[\begin{array}{l} w_b^{[i+1]} = w_b^{[i]} + \varepsilon_b^{[i]} \left(Q_b(\mathbf{t}^{[i]}, \mathbf{s}^{[i]}) - w_b^{[i]} \right) \\ t_b^{[i+1]} = t_b^{[i]} + \varepsilon_b^{[i]} \mu^{[i]} \left(\text{prox}_{\tau_b f_b}(2w_b^{[i+1]} - t_b^{[i]}) + a_b^{[i]} - w_b^{[i+1]} \right) \end{array} \right. \\ \text{For } \ell = 1, \dots, L \\ \left[\begin{array}{l} \mathbf{v}_\ell^{[i+1]} = \mathbf{v}_\ell^{[i]} + \varepsilon_{B+\ell}^{[i]} \left(Q_\ell(\mathbf{t}^{[i]}, \mathbf{s}^{[i]}) - \mathbf{v}_\ell^{[i]} \right) \\ \mathbf{s}_\ell^{[i+1]} = \mathbf{s}_\ell^{[i]} + \varepsilon_{B+\ell}^{[i]} \mu^{[i]} \left(\text{prox}_{\gamma_\ell \varphi_\ell}(2\mathbf{v}_\ell^{[i+1]} - \mathbf{s}_\ell^{[i]}) + \mathbf{e}_\ell^{[i]} - \mathbf{v}_\ell^{[i+1]} \right). \end{array} \right. \end{array} \right.$$

Defining, for every $i \in \mathbb{N}$, $\mathbf{a}^{[i]} = ((a_b^{[i]})_{1 \leq b \leq B}, (\mathbf{e}_\ell^{[i]})_{1 \leq \ell \leq L}) \in \mathcal{H} \times \mathcal{G}$, we have

$$\begin{aligned} \sum_{i \in \mathbb{N}} \sqrt{\mathbb{E}(\|\mathbf{a}^{[i]}\|_{\mathcal{U}^{-1}}^2 | \mathcal{X}^{[i]})} &= \sum_{i \in \mathbb{N}} \sqrt{\sum_{b=1}^B \tau_b^{-1} \mathbb{E}(\|a_b^{[i]}\|^2 | \mathcal{X}^{[i]}) + \sum_{\ell=1}^L \gamma_\ell^{-1} \mathbb{E}(\|\mathbf{e}_\ell^{[i]}\|^2 | \mathcal{X}^{[i]})} \\ &\leq \sum_{b=1}^B \tau_b^{-1/2} \sum_{i \in \mathbb{N}} \sqrt{\mathbb{E}(\|a_b^{[i]}\|^2 | \mathcal{X}^{[i]})} + \sum_{\ell=1}^L \gamma_\ell^{-1/2} \sum_{i \in \mathbb{N}} \sqrt{\mathbb{E}(\|\mathbf{e}_\ell^{[i]}\|^2 | \mathcal{X}^{[i]})} \\ &< +\infty, \end{aligned} \quad (28)$$

where the last inequality follows from (ii), (iii), (27) and

$$\sum_{i \in \mathbb{N}} \sqrt{\mathbb{E}(\|\mathbf{e}_\ell^{[i]}\|^2 | \mathcal{X}^{[i]})} = \frac{\gamma_\ell}{\sqrt{B}(1 - \gamma_\ell \rho_\ell)} \sum_{i \in \mathbb{N}} \sqrt{\mathbb{E}(\|d_\ell^{[i]}\|^2 | \mathcal{X}^{[i]})} < +\infty. \quad (29)$$

Altogether, since operator $J_{\mathcal{U}B}$ is weakly sequentially continuous because it is continuous and linear, the result follows from [32, Proposition 5.1] when the error term in the computation of $J_{\mathcal{U}B}$ is zero. \square

Remark 1

- (i) In Proposition 1, the binary variables $\varepsilon_b^{[i]}$ and $\varepsilon_{B+\ell}^{[i]}$ signal whether the variables $t_b^{[i]}$ and $\mathbf{s}_\ell^{[i]}$ are activated or not at iteration i . Assumption (iv) guarantees that each of the latter variables is activated with a nonzero probability at each iteration. In particular, it must be pointed out that the variables $p_\ell^{[i]}$ and $q_\ell^{[i]}$ only need to be computed when $\varepsilon_{B+\ell}^{[i]} = 1$.
- (ii) Note that Algorithm 1 may look similar to the stochastic approach proposed in [32, Corollary 5.5] (see also [36, Remark 2.9], and [37] for deterministic variants). It exhibits however three key differences. Most importantly, the operator inversions performed at the initial step amount to inverting a set of positive definite self-adjoint operators defined on the spaces $(\mathcal{H}_b)_{1 \leq b \leq B}$. We will see in our application that this reduces to invert a set of small size symmetric positive definite matrices. Another advantage is that the smoothness of the functions $(h_\ell)_{1 \leq \ell \leq L}$ is taken into account, and a last one is that the dual variables appear explicitly in the iterations.

- (iii) If, for every $\ell \in \{1, \dots, L\}$, $\rho_\ell = 0$ and $B = 1$, Algorithm 1 simplifies to Algorithm 2, where unnecessary indices have been dropped and we have set

$$(\forall i \in \mathbb{N}) \quad \begin{cases} \tilde{u}^{[i]} = -\tau u^{[i]} \\ (\forall \ell \in \{1, \dots, L\}) \quad \tilde{s}_\ell^{[i]} = -\tau s_\ell^{[i]}. \end{cases} \quad (30)$$

In this case,

$$C = \left(\text{Id} + \tau \sum_{\ell=1}^L \gamma_\ell A_\ell^* A_\ell \right)^{-1}. \quad (31)$$

Algorithm 2 Random Douglas-Rachford splitting for solving Problem 1 when $\rho_\ell = 0$ and $B = 1$

INITIALIZATION

$$\left[\begin{array}{l} \text{Set } \tau \in]0, +\infty[\text{ and } \eta \in]0, 1]. \\ \text{For every } \ell \in \{1, \dots, L\}, \text{ set } \gamma_\ell > 0. \\ \tilde{u}^{[0]} = \sum_{\ell=1}^L A_\ell^* \tilde{s}_\ell^{[0]} \end{array} \right.$$

FOR $i = 0, 1, \dots$

$$\left[\begin{array}{l} \text{Set } \mu^{[i]} \in]\eta, 2 - \eta[\\ w^{[i+1]} = w^{[i]} + \varepsilon^{[i]} \left(C \left(t^{[i]} + \tilde{u}^{[i]} \right) - w^{[i]} \right) \\ t^{[i+1]} = t^{[i]} + \varepsilon^{[i]} \mu^{[i]} \left(\text{prox}_{\tau f} (2w^{[i+1]} - t^{[i]}) + a^{[i]} - w^{[i+1]} \right) \\ \text{for } \ell = 1, \dots, L \\ \left[\begin{array}{l} q_\ell^{[i]} = \text{prox}_{\frac{h_\ell}{\gamma_\ell}} \left(2A_\ell w^{[i]} - \frac{\tilde{s}_\ell^{[i]}}{\tau \gamma_\ell} \right) + d_\ell^{[i]} \\ \tilde{s}_\ell^{[i+1]} = \tilde{s}_\ell^{[i]} + \varepsilon_{\ell+1}^{[i]} \mu^{[i]} \tau \gamma_\ell \left(q_\ell^{[i]} - A_\ell w^{[i]} \right) \end{array} \right. \\ \tilde{u}^{[i+1]} = \tilde{u}^{[i]} + \sum_{\ell=1}^L \varepsilon_{\ell+1}^{[i]} A_\ell^* \left(\tilde{s}_\ell^{[i+1]} - \tilde{s}_\ell^{[i]} \right). \end{array} \right.$$

When $(\forall \ell \in \{1, \dots, L\}) \gamma_\ell = 1/\tau$, it turns out this algorithm is exactly the same as the one resulting from a direct application of [32, Corollary 5.5] [43].

- (iv) The situation when, for a given ℓ , $h_\ell \in \Gamma_0(\mathcal{G}_\ell)$ is not Lipschitz-differentiable can be seen as the limit case when $\beta_\ell \rightarrow +\infty$. It can then be shown that Algorithm 1 remains valid by setting $\rho_\ell = 0$.

3 Sparse logistic regression

The proposed algorithm can be applied in the context of binary linear classification. A binary linear classifier can be modeled as a function that predicts the output

$y \in \{-1, +1\}$ associated to a given input $\mathbf{x} \in \mathbb{R}^N$. This prediction is defined through a linear combination of the input components, yielding the decision variable

$$d_{\mathbf{w}}(x) = \text{sign}(\mathbf{x}^\top \mathbf{w}), \quad (32)$$

where $\mathbf{w} \in \mathbb{R}^N$ is the weight vector to be estimated. In supervised learning, this weight vector is determined from a set of input-output pairs

$$\mathcal{S} = \{(\mathbf{x}_\ell, y_\ell) \in \mathbb{R}^N \times \{-1, +1\} \mid \ell \in \{1, \dots, L\}\}, \quad (33)$$

which is called *training set*. More precisely, the learning task can be defined as the trade-off between fitting the training data and reducing the model complexity, leading to an optimization problem expressed as

$$\underset{\mathbf{w} \in \mathbb{R}^N}{\text{minimize}} \mathbf{f}(\mathbf{w}) + \sum_{\ell=1}^L h(y_\ell \mathbf{x}_\ell^\top \mathbf{w}), \quad (34)$$

where $\mathbf{f} \in \Gamma_0(\mathbb{R}^N)$ is a regularization function and $h \in \Gamma_0(\mathbb{R})$ stands for the loss function. In the context of sparse learning, a popular choice for the regularization is the ℓ_1 -norm. Although many choices for the loss function are possible, we are primarily interested in the logistic loss, which is detailed in the next section.

3.1 Logistic regression

Logistic regression aims at maximizing the posterior probability density function of the weights given the training data, here assumed to be a realization of statistically independent input-output random variables. This leads us to

$$\underset{\mathbf{w} \in \mathbb{R}^N}{\text{maximize}} \varphi(\mathbf{w}) \prod_{\ell=1}^L \pi(y_\ell \mid \mathbf{x}_\ell, \mathbf{w}) \theta_\ell(\mathbf{x}_\ell \mid \mathbf{w}), \quad (35)$$

where φ is the weight prior probability density function and, for every $\ell \in \{1, \dots, L\}$, θ_ℓ is the conditional data likelihood of the ℓ -th input knowing the weight values, while $\pi(y_\ell \mid \mathbf{x}_\ell, \mathbf{w})$ is the conditional probability of the ℓ -th output knowing the ℓ -th input and the weights. Let us model this conditional probability with the sigmoid function defined as

$$\pi(y_\ell \mid \mathbf{x}_\ell, \mathbf{w}) = \frac{1}{1 + \exp(-y_\ell \mathbf{x}_\ell^\top \mathbf{w})}, \quad (36)$$

and assume that the inputs and the weights are statistically independent and that φ is log-concave. Then, the negative-logarithm of the energy in (35) yields an instance of Problem (34) in which

$$(\forall v \in \mathbb{R}) \quad h(v) = \log(1 + \exp(-v)) \quad (37)$$

and, for every $\mathbf{w} \in \mathbb{R}^N$, $\mathbf{f}(\mathbf{w}) = -\log \varphi(\mathbf{w})$. (The term $\prod_{\ell=1}^L \theta_\ell(\mathbf{x}_\ell \mid \mathbf{w})$ can be discarded since the inputs and the weights are assumed statistically independent.) The function in (37) is called *logistic loss*. For completeness, note that other loss functions, leading to different kinds of classifiers, are the *hinge loss* [44]

$$(\forall v \in \mathbb{R}) \quad h^{\text{hinge}}(v) = (\max\{0, 1 - v\})^q \quad (38)$$

with $q \in \{1, 2\}$, and the *Huber loss* [45]

$$(\forall v \in \mathbb{R}) \quad h^{\text{huber}}(v) = \begin{cases} 0 & \text{if } v \geq 1 \\ -v & \text{if } v \leq -1 \\ \frac{1}{4}(v-1)^2 & \text{otherwise.} \end{cases} \quad (39)$$

These functions can be also handled by the proposed algorithm.

3.2 Optimization algorithm

Let us blockwise decompose the weight variable $\mathbf{w} \in \mathbb{R}^N$ as

$$\mathbf{w}^\top = [w_1^\top \dots w_B^\top], \quad (40)$$

where, for every $b \in \{1, \dots, B\}$, $w_b \in \mathbb{R}^{N_b}$ and N_1, \dots, N_B are strictly positive integers such that $N = N_1 + \dots + N_B$. Let us also decompose the input vector as $\mathbf{x}^\top = [x_1^\top \dots x_B^\top]$. Finally, let us assume that the regularization function is block-separable, i.e. $\mathbf{f} = \bigoplus_{b=1}^B f_b$, where, for every $b \in \{1, \dots, B\}$, $f_b \in \Gamma_0(\mathbb{R}^{N_b})$. A typical example of such functions is given by

$$(\forall b \in \{1, \dots, B\}) \quad f_b = \lambda \|\cdot\|_{\kappa_b}, \quad (41)$$

where $\lambda \in [0, +\infty[$ and $\|\cdot\|_{\kappa_b}$, $\kappa_b \in [1, +\infty]$, denotes the ℓ_{κ_b} -norm of \mathbb{R}^{N_b} . In particular, when for every $b \in \{1, \dots, B\}$ $\kappa_b = 1$, \mathbf{f} reduces to the standard ℓ_1 regularizer, whereas setting $\kappa_b \equiv 2$ results in a potential promoting group sparsity [46].

In the context described above, (34) is a particular case of Problem 1 where, for every $b \in \{1, \dots, B\}$, $\mathcal{H}_b = \mathbb{R}^{N_b}$, for every $\ell \in \{1, \dots, L\}$, $\mathcal{G}_\ell = \mathbb{R}$, $h_\ell = h$ and $A_{\ell,b} = y_\ell x_{\ell,b}^\top$. Note that since, for every $\ell \in \{1, \dots, L\}$, $y_\ell^2 = 1$, $A_{\ell,b}^* A_{\ell,b} = x_{\ell,b} x_{\ell,b}^\top$. Moreover, h defined in (37) is twice differentiable with

$$(\forall v \in \mathbb{R}) \quad h'(v) = -\frac{\exp(-v)}{1 + \exp(-v)}, \quad (42)$$

$$h''(v) = \frac{\exp(-v)}{(1 + \exp(-v))^2}. \quad (43)$$

Since h'' is maximized in $v = 0$, we have $\sup_{v \in \mathbb{R}} |h''(v)| = 1/4$, which implies that h' is $1/4$ -Lipschitz continuous and we have thus, for every $\ell \in \{1, \dots, L\}$, $\beta_\ell = 1/4$.

The problem can thus be solved with Algorithm 1, the convergence of which is guaranteed almost surely under the assumptions of Proposition 1.

3.3 Proximity operator

An efficient computation of the proximity operators of functions $(f_b)_{1 \leq b \leq B}$ and h plays a crucial role in the implementation of Algorithm 1. There exists an extensive literature on the computation of the proximity operators of functions like (41) [47]. In particular, when $\kappa_b = 1$ (resp. $\kappa_b = 2$), this proximity operator reduces to a component-wise (resp. blockwise) soft-thresholding [48]. Regarding the logistic loss in (37), although some numerical methods exist [49, 50], to the best of our knowledge,

no thorough investigation of the form of its proximity operator has been made. The next proposition will contribute to fill such a void. The result relies on the generalized W-Lambert function recently introduced in [39, 40], defined via

$$(\forall \bar{v} \in \mathbb{R})(\forall v \in \mathbb{R})(\forall r \in]0, +\infty[) \quad \bar{v}(\exp(\bar{v}) + r) = v \quad \Leftrightarrow \quad \bar{v} = W_r(v). \quad (44)$$

When $r \in [\exp(-2), +\infty[$, W_r is uniquely defined and strictly increasing, but when $r \in]0, \exp(-2)[$, there exist three branches for W_r . We will retain the only one which can take nonnegative values (denoted by $W_{r,0}$ in [39, Theorem 4]) and is also strictly increasing. This function can be efficiently evaluated through a Newton-based method devised by Mező *et al.* [39] and available on line.¹

Proposition 2 *Let $\gamma \in]0, +\infty[$ and let $h: v \mapsto \log(1 + \exp(-v))$. We have*

$$(\forall v \in \mathbb{R}) \quad \text{prox}_{\gamma h}(v) = v + W_{\exp(-v)}(\gamma \exp(-v)). \quad (45)$$

Proof Let $v \in \mathbb{R}$ and $\gamma \in]0, +\infty[$. For every $p \in \mathbb{R}$, it follows from the definition of $\text{prox}_{\gamma h}$, (42) and (44) that

$$p = \text{prox}_{\gamma h}(v) \quad \Leftrightarrow \quad v - p = -\frac{\gamma \exp(-p)}{1 + \exp(-p)} = -\frac{\gamma}{\exp(p) + 1} \quad (46)$$

$$\begin{aligned} &\Leftrightarrow (p - v)(\exp(p) + 1) = \gamma \\ &\Leftrightarrow (p - v)(\exp(p - v) + \exp(-v)) = \gamma \exp(-v) \\ &\Leftrightarrow p - v = W_{\exp(-v)}(\gamma \exp(-v)) \end{aligned} \quad (47)$$

and the result follows. \square

From a numerical standpoint, it must be emphasized that the exponentiation in (45) may be problematic, as it yields an arithmetic overflow when v tends to $-\infty$. To overcome this issue, one can use the asymptotic equivalence² between the proximity operator of the logistic function and other more tractable functions.

Proposition 3 *Let $\gamma \in]0, +\infty[$ and let $h: v \mapsto \log(1 + \exp(-v))$. Then, as $v \rightarrow -\infty$,*

$$\text{prox}_{\gamma h}(v) = v + \gamma(1 - \exp(\gamma + v) + (1 + \gamma)\exp(2(\gamma + v))) + \vartheta(\exp(2v)). \quad (48)$$

Proof Define

$$(\forall v \in \mathbb{R}) \quad \varphi(v) = W_{\exp(-v)}(\gamma \exp(-v)). \quad (49)$$

According to Proposition 2,

$$\varphi = \text{prox}_{\gamma h} - \text{Id}. \quad (50)$$

It follows from (46) that $\text{ran } \varphi = \text{ran}(\text{prox}_{\gamma h} - \text{Id}) \subset]0, \gamma[$. Moreover, from [42, Section 24.2] we deduce that $\varphi = -\gamma \text{prox}_{h^*/\gamma}(\cdot/\gamma)$, which is decreasing and continuous by virtue of [42, Proposition 24.31]. Therefore, $\lim_{v \rightarrow -\infty} \varphi(v)$ exists and from (44) we have

$$(\forall v \in \mathbb{R}) \quad \varphi(v) \exp(\varphi(v)) = (\gamma - \varphi(v)) \exp(-v). \quad (51)$$

¹<https://sites.google.com/site/istvanmezo81/others>

²Hereafter, following Landau's notation, we will write that $F(v) = \vartheta(G(v))$, where $F: \mathbb{R} \rightarrow \mathbb{R}$ and $G: \mathbb{R} \rightarrow \mathbb{R}$, if $F(v)/G(v) \rightarrow 0$ as $v \rightarrow +\infty$ (or $v \rightarrow -\infty$).

Since the left side of the equality above is bounded, we deduce that $\lim_{v \rightarrow -\infty} \varphi(v) = \gamma$. Subsequently, we define

$$(\forall v \in \mathbb{R}) \quad u(v) = \varphi(v) - \gamma \quad \text{satisfying} \quad \lim_{v \rightarrow -\infty} u(v) = 0. \quad (52)$$

Hence, (51) can be rewritten as

$$(\gamma + u(v)) \exp(\gamma) \exp(u(v)) = -u(v) \exp(-v) \quad (53)$$

and by using the first order Taylor expansion around $\xi = 0$, $\exp(\xi) = 1 + \xi + \vartheta(\xi)$ and the fact that $u(v)^2 + (\gamma + u(v))\vartheta(u(v)) = \vartheta(u(v))$, we obtain

$$\begin{aligned} u(v) &= -\exp(\gamma + v)(\gamma + u(v))(1 + u(v) + \vartheta(u(v))) \\ &= -\exp(\gamma + v)(\gamma + (\gamma + 1)u(v) + \vartheta(u(v))) \\ &= -\gamma \exp(\gamma + v) - (\gamma + 1) \exp(\gamma + v)u(v) - \exp(\gamma + v)\vartheta(u(v)). \end{aligned} \quad (54)$$

We deduce from this relation that

$$u(v) = -\frac{\gamma \exp(\gamma + v) + \exp(\gamma + v)\vartheta(u(v))}{1 + (\gamma + 1) \exp(\gamma + v)}. \quad (55)$$

It follows that

$$\lim_{v \rightarrow -\infty} u(v) \exp(-v) = -\gamma \exp(\gamma), \quad (56)$$

which implies that $\exp(\gamma + v)\vartheta(u(v)) = \vartheta(\exp(2v))$ and, from (55) we obtain

$$u(v) = -\frac{\gamma \exp(\gamma + v)}{1 + (\gamma + 1) \exp(\gamma + v)} + \vartheta(\exp(2v)). \quad (57)$$

Combining (50), (52), and (57) yields

$$\begin{aligned} \text{prox}_{\gamma h}(v) &= v + \gamma \left(1 - \frac{\exp(\gamma + v)}{1 + (\gamma + 1) \exp(\gamma + v)} \right) + \vartheta(\exp(2v)) \\ &= v + \gamma \left(\frac{1 + \gamma \exp(\gamma + v)}{1 + (\gamma + 1) \exp(\gamma + v)} \right) + \vartheta(\exp(2v)) \\ &= v + \gamma (1 - \exp(\gamma + v) + (1 + \gamma) \exp(2(\gamma + v))) + \vartheta(\exp(2v)), \end{aligned} \quad (58)$$

where the last equality follows from the second order Taylor expansion around $\xi = 0$. \square

4 Experimental results

In order to assess the performance of Algorithm 1, we performed the training on standard datasets^{3,4} (see Table 1), and we compared it with the following approaches.

³<http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html>

⁴<http://yann.lecun.com/exdb/mnist>

- Stochastic Forward-Backward splitting (SFB) [29–31, 51]

$$\begin{aligned}
& w^{[0]} \in \mathbb{R}^N \\
& \text{For } i = 0, 1, \dots \\
& \left[\begin{array}{l} \text{Select } \mathbb{L}^{[i]} \subset \{1, \dots, L\} \\ w^{[i+1]} = \text{prox}_{\gamma_i f} \left(w^{[i]} - \gamma_i \sum_{\ell \in \mathbb{L}^{[i]}} y_\ell x_\ell h'(y_\ell x_\ell^\top w^{[i]}) \right) \end{array} \right]
\end{aligned}$$

where $(\gamma_i)_{i \in \mathbb{N}}$ is a decreasing sequence of positive values.

- Regularized Dual Averaging (RDA) [28]

$$\begin{aligned}
& w^{[0]} \in \mathbb{R}^N, z^{[0]} = 0 \\
& \text{For } i = 0, 1, \dots \\
& \left[\begin{array}{l} \text{Select } \mathbb{L}^{[i]} \subset \{1, \dots, L\} \\ z^{[i+1]} = z^{[i]} + \sum_{\ell \in \mathbb{L}^{[i]}} y_\ell x_\ell h'(y_\ell x_\ell^\top w^{[i]}) \\ w^{[i+1]} = \text{prox}_{\gamma_i f} \left(-\gamma_i z^{[i+1]} \right) \end{array} \right]
\end{aligned}$$

where $(\gamma_i)_{i \in \mathbb{N}}$ is a decreasing sequence of positive values.

- Block-Coordinate Primal-Dual splitting (BCPD) [41]

$$\begin{aligned}
& w^{[0]} \in \mathbb{R}^N, v^{[0]} \in \mathbb{R}^L \\
& \text{For } i = 0, 1, \dots \\
& \left[\begin{array}{l} \text{Select } \mathbb{L}^{[i]} \subset \{1, \dots, L\} \\ w^{[i+1]} = \text{prox}_{\tau f} \left(w^{[i]} - \tau u^{[i]} \right) \\ (\forall \ell \in \mathbb{L}^{[i]}) \quad v_\ell^{[i+1]} = \text{prox}_{\sigma h^*} \left(v_\ell^{[i]} + \sigma y_\ell x_\ell^\top (2w^{[i+1]} - w^{[i]}) \right) \\ (\forall \ell \notin \mathbb{L}_i) \quad v_\ell^{[i+1]} = v_\ell^{[i]} \\ u^{[i+1]} = u^{[i]} + \sum_{\ell \in \mathbb{L}^{[i]}} (v_\ell^{[i+1]} - v_\ell^{[i]}) y_\ell x_\ell \end{array} \right]
\end{aligned}$$

where $\tau > 0$ and $\sigma > 0$ are such that $\tau \sigma \left\| \sum_{\ell=1}^L x_\ell x_\ell^\top \right\| \leq 1$.

The algorithmic parameters are reported in Table 2. For all the algorithms, mini-batches of size 1000 were randomly selected using a uniform distribution, and the initial vector $w^{[0]}$ was randomly drawn from the normal distribution with zero mean and unit variance. For the datasets with more than two classes (MNIST and RCV1), the “one-versus-all” approach is used [52]. All experiments were carried out with Matlab 2015a on an Intel i7 CPU at 3.40 GHz and 12 GB of RAM.

Table 3 reports the classification performance achieved by the compared algorithms, which includes the classification errors computed on a (disjoint) test set, as well as the sparsity degree of the solution. For all the considered datasets, the regularization parameter λ was selected with a cross-validation procedure. The results show that the proposed algorithm finds a solution that yields the same accuracy as state-of-the-art methods, while being sparser than the ones produced by gradient-like methods (SFB and RDA).

Table 1 Training sets used in the experiments (K is the number of classes).

Dataset	N	L	K
W8A	300	49749	2
MNIST	717	60000	10
RCV1	12560	30879	20

Table 2 Algorithmic parameters used in the experiments.

Dataset	SFB / RDA	Algo 1				BCPD	
	γ_i	γ, τ	$\mu^{[i]}$	ρ	B	τ	σ
W8A	$10^{-1}/\sqrt{i+1}$				1		
MNIST	$1/\sqrt{i+1}$	1	1.5	0.1	1	0.1	$\tau^{-1} \left\ \sum_{\ell=1}^L x_\ell x_\ell^\top \right\ ^{-1}$
RCV1	$10/\sqrt{i+1}$				9		

Table 3 Classification performance on test sets (after training for a fixed number of iterations).

DATASET	Algo 1	SFB	RDA	BCPD
	Errors – Zeros	Errors – Zeros	Errors – Zeros	Errors – Zeros
w8A	9.73% – 19.60%	9.92% – 5.65%	9.99% – 0.33%	10.44% – 50.50%
MNIST	8.49% – 41.57%	8.37% – 5.25%	8.60% – 11.13%	8.45% – 58.64%
RCV1	6.62% – 83.43%	6.67% – 35.22%	6.60% – 32.90%	6.25% – 98.57%

Figure 1 reports the training performance versus time of the considered algorithms, which includes the criterion in (34), and the distance to the solution $w^{[\infty]}$ obtained after many iterations for each compared method. The results indicate that the proposed approach converges faster to a smaller value of the objective criterion. This could be related to the implicit preconditioning present in Algorithm 1 through the matrix Q . Another interesting feature of our algorithm is the free choice of parameters γ and μ_i . Conversely, in both SFB and RDA, the parameter γ_i (also referred to as *learning rate*) needs to be carefully selected by hand, causing such algorithms to slow down or even diverge if the learning rate is chosen too small or too high.

5 Conclusion

In this paper, we have proposed a block-coordinate Douglas-Rachford algorithm for sparse logistic regression. In contrast to gradient-like methods, our approach relies on the proximity operator of the logistic loss, for which we derived a closed-form expression that can be efficiently implemented. Thanks to this feature, our approach removes restrictions on the choice of the algorithm parameters, unlike gradient-like methods, for which it is essential that the learning rate is carefully chosen. This is confirmed by our numerical results, which indicate that the training performance of the proposed algorithm compares favorably with state-of-the-art stochastic methods.

References

1. L. Rosasco, E. De Vito, A. Caponnetto, M. Piana, and A. Verri, “Are loss functions all the same?” *Neural Comput.*, vol. 16, no. 5, pp. 1063–1076, May 2004.

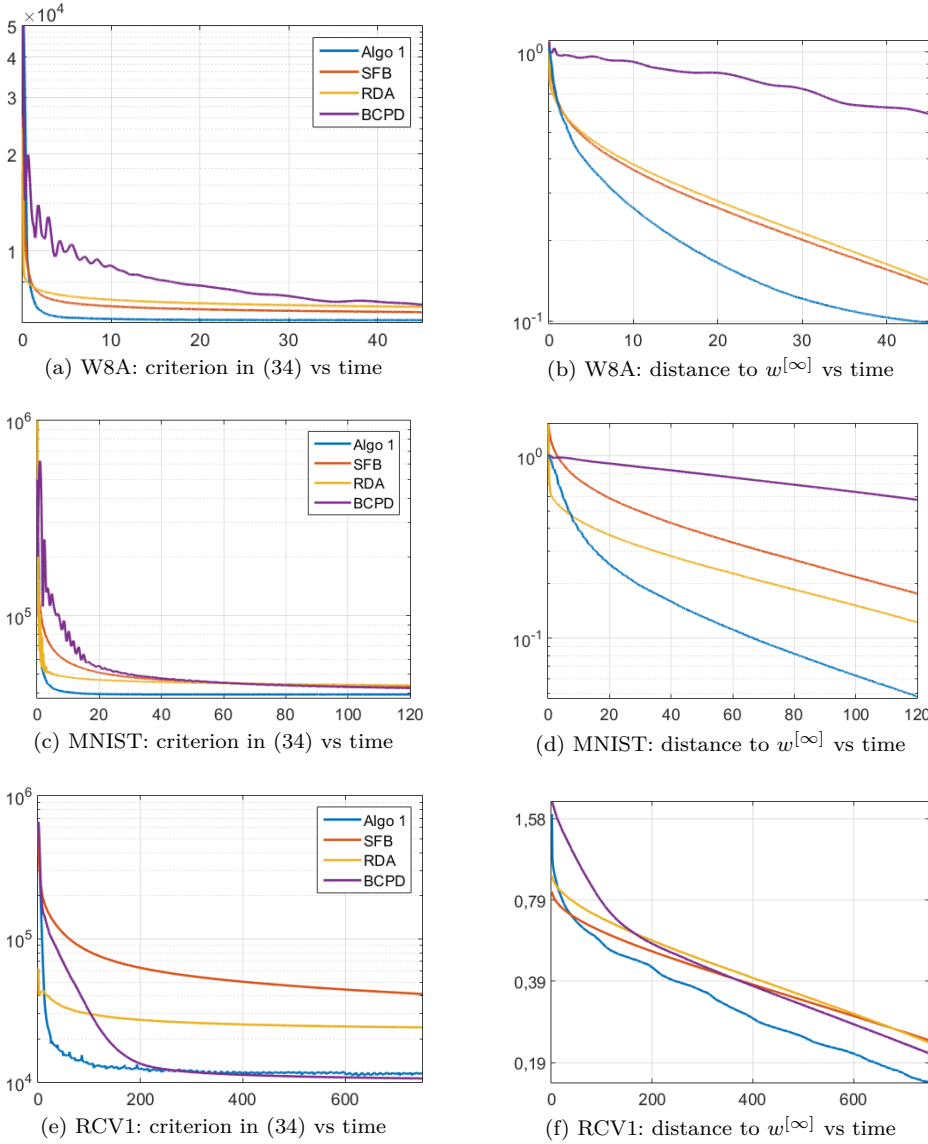


Fig. 1 Comparison of training performance (time is expressed in seconds).

2. P. L. Bartlett, M. I. Jordan, and J. D. McAuliffe, "Convexity, classification, and risk bounds," *J. Am. Stat. Assoc.*, vol. 101, no. 473, pp. 138–156, 2006.
3. P. S. Bradley and O. L. Mangasarian, "Feature selection via concave minimization and support vector machines," in *Proc. of ICML*, Madison, USA, 1998, pp. 82–90.
4. J. Weston, A. Elisseeff, B. Schölkopf, and M. Tipping, "Use of the zero-norm with linear models and kernel methods," *Mach. Learn.*, vol. 3, pp. 1439–1461, 2002.
5. Y. Liu, H. Helen Zhang, C. Park, and J. Ahn, "Support vector machines with adaptive L_q penalty," *Comput. Stat. Data Anal.*, vol. 51, no. 12, pp. 6380–6394, Aug. 2007.
6. H. Zou and M. Yuan, "The f_∞ -norm support vector machine," *Stat. Sin.*, vol. 18, pp. 379–398, 2008.

7. L. Laporte, R. Flamary, S. Canu, S. Déjean, and J. Mothe, “Non-convex regularizations for feature selection in ranking with sparse SVM,” *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 6, pp. 1118–1130, Jun. 2014.
8. B. Krishnapuram, L. Carin, M. A. T. Figueiredo, and A. J. Hartemink, “Sparse multinomial logistic regression: Fast algorithms and generalization bounds,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 27, no. 6, Jun. 2005.
9. L. Meier, S. Van De Geer, and P. Bühlmann, “The group Lasso for logistic regression,” *J. R. Stat. Soc. B*, vol. 70, no. 1, pp. 53–71, 2008.
10. J. Duchi and Y. Singer, “Boosting with structural sparsity,” in *Proc. of ICML*, Montreal, Canada, Jun. 2009, pp. 297–304.
11. G. Obozinski, B. Taskar, and M. I. Jordan, “Joint covariate selection and joint subspace selection for multiple classification problems,” *Stat. Comput.*, vol. 20, no. 2, pp. 231–252, 2010.
12. G.-X. Yuan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin, “A comparison of optimization methods and software for large-scale L1-regularized linear classification,” *Mach. Learn.*, vol. 11, pp. 3183–3234, Dec. 2010.
13. L. Rosasco, S. Villa, S. Mosci, M. Santoro, and A. Verri, “Nonparametric sparsity and regularization,” *J. Mach. Learn. Res.*, vol. 14, pp. 1665–1714, Jul. 2013.
14. M. Blondel, A. Fujino, and N. Ueda, “Large-scale multiclass support vector machine training via euclidean projection onto the simplex,” in *Proc. of ICPR*, Stockholm, Sweden, 24–28 August 2014, pp. 1289–1294.
15. L. Wang and X. Shen, “On l_1 -norm multi-class support vector machines: methodology and theory,” *J. Am. Statist. Assoc.*, vol. 102, pp. 583–594, 2007.
16. J. Mairal, “Optimization with first-order surrogate functions,” in *Proc. of ICML*, 2013, pp. 783–791.
17. G. Chierchia, N. Pustelnik, J.-C. Pesquet, and B. Pesquet-Popescu, “A proximal approach for sparse multiclass SVM,” *Preprint arXiv:1501.03669*, Feb. 2015.
18. M. Barlaud, W. Belhajali, P. L. Combettes, and L. Fillatre, “Classification and regression using a constrained convex splitting method,” *IEEE Trans. Signal Process.*, 2017.
19. M. Tan, L. Wang, and I. W. Tsang, “Learning sparse SVM for feature selection on very high dimensional datasets,” in *Proc. of ICML*, Haifa, Israel, 21–24 June 2010, pp. 1047–1054.
20. C. J. Hsieh, K. W. Chang, C. J. Lin, S. S. Keerthi, and S. Sundararajan, “A dual coordinate descent method for large-scale linear SVM,” in *Proc. of ICML*, 2008, pp. 408–415.
21. S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher, “Block-coordinate Frank-Wolfe optimization for structural SVMs,” in *Proc. of ICML*, vol. 28, no. 1, 2013, pp. 53–61.
22. M. Pereyra, P. Schniter, E. Chouzenoux, J.-C. Pesquet, J.-Y. Tournéret, A. Hero, and S. McLaughlin, “A survey of stochastic simulation and optimization methods in signal processing,” *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 2, pp. 224–241, Mar. 2016.
23. S. Shalev-Shwartz and A. Tewari, “Stochastic methods for l1-regularized loss minimization,” *J. Mach. Learn. Res.*, vol. 12, pp. 1865–1892, Jun. 2011.
24. M. Blondel, K. Seki, and K. Uehara, “Block coordinate descent algorithms for large-scale sparse multiclass classification,” *Mach. Learn.*, vol. 93, no. 1, pp. 31–52, Oct. 2013.
25. O. Fercoq and P. Richtárik, “Accelerated, parallel, and proximal coordinate descent,” *SIAM J. Opt.*, vol. 25, no. 4, pp. 1997–2023, 2015.
26. Z. Lu and L. Xiao, “On the complexity analysis of randomized block-coordinate descent methods,” *Math. Program.*, vol. 152, no. 1–2, pp. 615–642, 2015.
27. J. Langford, L. Li, and T. Zhang, “Sparse online learning via truncated gradient,” *J. Mach. Learn. Res.*, vol. 10, pp. 777–801, Mar. 2009.
28. L. Xiao, “Dual averaging methods for regularized stochastic learning and online optimization,” *J. Mach. Learn. Res.*, vol. 11, pp. 2543–2596, Oct. 2010.
29. P. Richtárik and M. Takáč, “Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function,” *Math. Program.*, vol. 144, no. 1, pp. 1–38, Apr. 2014.
30. L. Rosasco, S. Villa, and B. C. Vũ, “Stochastic forward-backward splitting for monotone inclusions,” *J. Optim. Theory Appl.*, vol. 169, no. 2, pp. 388–406, May 2016.
31. P. Combettes and J.-C. Pesquet, “Stochastic approximations and perturbations in forward-backward splitting for monotone operators,” *Pure Appl. Func. Anal.*, vol. 1, no. 1, pp. 13–37, Jan. 2016.
32. P. L. Combettes and J.-C. Pesquet, “Stochastic quasi-Fejér block-coordinate fixed point iterations with random sweeping,” *SIAM J. Opt.*, vol. 25, no. 2, pp. 1221–1248, Jul. 2015.

33. J.-C. Pesquet and A. Repetti, "A class of randomized primal-dual algorithms for distributed optimization," *J. Nonlinear Convex Anal.*, vol. 16, no. 12, pp. 2453–2490, 2015.
34. J. Mairal, "Stochastic majorization-minimization algorithms for large-scale optimization," in *Proc. of NIPS*, 2013, pp. 2283–2291.
35. E. Chouzenoux and J.-C. Pesquet, "A stochastic majorize-minimize subspace algorithm for online penalized least squares estimation," *IEEE Trans. Signal Process.*, 2017.
36. L. M. Briceño-Arias and P. L. Combettes, "A monotone + skew splitting model for composite monotone inclusions in duality," *SIAM J. Opt.*, vol. 21, no. 4, pp. 1230–1250, 2011.
37. R. I. Boğ and C. Hendrich, "A Douglas–Rachford type primal-dual method for solving inclusions with mixtures of composite and parallel-sum type monotone operators," *SIAM J. Opt.*, vol. 23, no. 4, pp. 2541–2565, 2013.
38. D. Perekrestenko, V. Cevher, and M. Jaggi, "Faster coordinate descent via adaptive importance sampling," in *Proc. of AISTATS*, Fort Lauderdale, Florida, USA, 20–22 April 2017.
39. I. Mező and Á. Baricz, "On the generalization of the lambert W function," *Trans. Amer. Math. Soc.*, 2017.
40. A. Maignan and T. Scott, "Fleshing out the generalized lambert W function," *ACM Communications in Computer Algebra*, vol. 50, no. 2, pp. 45–60, Jun. 2016.
41. G. Chierchia, N. Pustelnik, and J.-C. Pesquet, "Random primal-dual proximal iterations for sparse multiclass SVM," in *Proc. of MLSP*, Salerno, Italy, Sep. 2016.
42. H. H. Bauschke and P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, 2nd ed. New York: Springer, 2017.
43. G. Chierchia, A. Cherni, E. Chouzenoux, and J.-C. Pesquet, "Approche de Douglas-Rachford aléatoire par blocs appliquée à la régression logistique parcimonieuse," in *Actes du GRETSI*, Juan-les-Pins, France, Sep. 2017.
44. C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, Sep. 1995.
45. A. Martins and R. Astudillo, "From softmax to sparsemax: A sparse model of attention and multi-label classification," in *Proc. of ICML*, New York, USA, Jun. 2016, pp. 1614–1623.
46. F. Bach, R. Jenatton, J. Mairal, and G. Obozinski, "Optimization with sparsity-inducing penalties," *Foundations and Trends in Machine Learning*, vol. 4, no. 1, pp. 1–106, Jan. 2012.
47. G. Chierchia, E. Chouzenoux, P. L. Combettes, and J.-C. Pesquet. (2017) The proximity operator repository (user's guide). [Online]. Available: <http://proximity-operator.net/>
48. P. L. Combettes and J.-C. Pesquet, "Proximal splitting methods in signal processing," in *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, H. H. Bauschke, R. S. Burachik, P. L. Combettes, V. Elser, D. R. Luke, and H. Wolkowicz, Eds. New York: Springer-Verlag, 2011, pp. 185–212.
49. N. Parikh and S. Boyd, "Proximal algorithms," *Foundations and Trends in Optimization*, vol. 1, no. 3, pp. 123–231, 2014.
50. P.-W. Wang, M. Wytoczek, and J. Z. Kolter, "Epigraph projections for fast general convex programming," in *Proc. of ICML*, 2016.
51. Y. F. Atchadé, G. Fort, and E. Moulines, "On perturbed proximal gradient algorithms," *J. Mach. Learn. Res.*, vol. 18, no. 1, pp. 310–342, Jan. 2017.
52. R. Rifkin and A. Klautau, "In defense of one-vs-all classification," *Mach. Learn.*, vol. 5, pp. 101–141, 2004.