
Beyond L1: Faster and Better Sparse Models with `skglm`

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose a new fast algorithm to estimate any sparse generalized linear model
2 with convex or non-convex separable penalties. Our algorithm is able to solve
3 problems with millions of samples and features in seconds, by relying on coordinate
4 descent, working sets and Anderson acceleration. It handles previously unaddressed
5 models, and is extensively shown to improve state-of-art algorithms. We provide
6 a flexible, `scikit-learn` compatible package, which easily handles customized
7 datafits and penalties.

8 1 Introduction

9 Sparse generalized linear models play a central role in modern machine learning and signal processing.
10 The Lasso (Tibshirani, 1996) and its derivatives (Zou and Hastie, 2005; Ng, 2004; Candes et al., 2008;
11 Simon et al., 2013) have found numerous successful applications to large scale tasks in genomics
12 (Ghosh and Chinnaiyan, 2005), vision (Mairal, 2010), or neurosciences (Strohmeier et al., 2016). This
13 impact was made possible by two key factors: efficient algorithms and software implementations.

14 State-of-the-art algorithms for “smooth + non-smooth separable” problems predominantly rely on
15 coordinate descent (CD, Tseng and S.Yun 2009; Nesterov 2012), which, when it can be applied,
16 is more efficient than full gradient methods (Richtárik and Takáč, 2014, Sec. 6.1). Coordinate
17 descent can even be improved with Nesterov-like acceleration, to obtain improved convergence rates
18 (Lin et al., 2014; Fercoq and Richtárik, 2015). However, these better rates may fail to reflect in
19 practical accelerations. On the contrary, Bertrand and Massias (2021) relied on Anderson acceleration
20 (Anderson, 1965) to provide both better rates and practical acceleration for coordinate descent.

21 Even with efficient algorithms such as coordinate descent, the practical use of sparsity hits a com-
22 putational barrier for problems with more than millions of features (Le Morvan and Vert, 2018).
23 Multiple techniques have been proposed to make coordinate descent scale to huge problems. Notably,
24 algorithms can be accelerated by reducing the number of variables to optimize over, using screening
25 rules or working sets. Screening rules discard features from the problem in advance (El Ghaoui
26 et al. 2010; Bonnefoy et al. 2015) or dynamically (Fercoq et al., 2015; Ndiaye et al., 2017). On the
27 other side, working sets (Johnson and Guestrin, 2015; Massias et al., 2018) iteratively solve larger
28 subproblems and progressively include variables identified as relevant.

29 For the Lasso and a few convex models, coordinate descent has been broadly disseminated to
30 practitioners in off-the-shelf packages such as `glmnet` (Friedman et al., 2007) or `scikit-learn`
31 (Pedregosa et al., 2011). More recently, `celer`, a state-of-the-art convex working set algorithm
32 (Massias et al., 2020) allowed for successful applications of the Lasso in large scale problems in
33 medicine (Reidenbach et al., 2021; Kim et al., 2021) or seismology (Muir and Zhan, 2021).

34 Yet the Lasso is limited: non-convex sparse models enjoy better theoretical and empirical properties
35 (Breheny and Huang, 2011; Soubies et al., 2015). As illustrated in Figure 1, they yield sparser

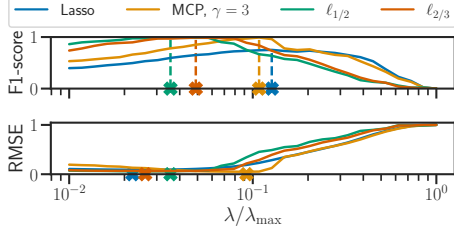


Figure 1: **Regularization paths computed with our algorithm.** Non-convex sparse penalties behave better than the L1 norm. Due to their lower bias, they achieve perfect support recovery, lower prediction error and their optimal regularization strength λ in estimation (top) and prediction (bottom) correspond.

36 solutions than convex penalties and mitigate the intrinsic Lasso bias. Yet, they have not so often
 37 been applied to huge scale applications. This is mostly an algorithmic barrier: while coordinate
 38 descent can be applied to non-convex penalties (Breheny and Huang, 2011; Mazumder et al., 2011;
 39 Bolte et al., 2014), screening rules and working sets are heavily dependent on convexity or quadratic
 40 datafits (Rakotomamonjy et al., 2019, 2022).

41 In this work, we solve this issue by designing a **state-of-the-art generic algorithm** to solve a wide
 42 range of sparse generalized linear models. The contributions are the following:

- 43 • We propose a non-convex converging working set algorithm relying on Anderson accelerated
 44 coordinate descent. For a specific class of non-convex penalties, we show:
 - 45 (a) Convergence of the proposed working set algorithm (Proposition 5).
 - 46 (b) Support identification of coordinate descent (Proposition 10).
 - 47 (c) Local convergence rates for the Anderson extrapolation (Proposition 13).
- 48 • We provide an extensive experimental comparison and we show state-of-the-art improve-
 49 ments on a wide range of convex and non-convex problems. In addition we release an
 50 efficient and modular python implementation, with a `scikit-learn` API, for practitioners
 51 to apply non-convex penalties to large scale problems.

52 2 Framework and proposed algorithm

53 2.1 Problem setting

In this paper, we consider problems of the form:

$$\hat{\beta} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \Phi(\beta) \triangleq \underbrace{F(X\beta)}_{\triangleq f(\beta)} + \sum_{j=1}^p g_j(\beta_j), \quad (1)$$

54 where F is smooth, and the functions g_j are proper and lower semicontinuous but not necessarily
 55 convex, whose proximal operator can be computed exactly. We write $g = \sum_j g_j$. Instances of
 56 Problem (1) include convex estimators: the Lasso, the elastic net, the sparse logistic regression,
 57 the dual of SVM with hinge loss. They also include non-convex penalties: $\ell_{0.5}$ and $\ell_{2/3}$ penalties
 58 (Foucart and Lai, 2009), the minimax concave penalty (MCP, Zhang 2010) or SCAD (Zhang, 2010),
 59 both with regression and classification losses. Formally, the assumptions are the following.

60 **Assumption 1.** $f : \mathbb{R}^p \rightarrow \mathbb{R}$ is convex and differentiable and for all $j \in [p]$, the restriction of $\nabla_j f$
 61 to the j -th coordinate is L_j -Lipschitz: for all $(x, h) \in \mathbb{R}^p \times \mathbb{R}$, $|\nabla_j f(x + h e_j) - \nabla_j f(x)| \leq L_j |h|$.

62 **Assumption 2.** For any $j \in [p]$, $g_j : \mathbb{R} \rightarrow \mathbb{R}$ is proper, closed, and lower bounded.

63 Following Attouch and Bolte 2009; Bolte et al. 2014 we focus on finding a critical point of Φ .

64 **Definition 3.** Using the Fréchet subdifferential (Kruger, 2003), a critical point $x \in \mathbb{R}^p$ is a point
 65 which satisfies $-\nabla f(x) \in \partial g(x)$.

66 Assumptions 1 and 2 are usual, and, under boundedness of the iterates, ensure convergence of
 67 forward-backward and coordinate descent algorithms to a critical point (Attouch et al. 2013, Thm
 68 5.1, Bolte et al. 2014, Thm. 3.1). In addition, our work focuses on the case where g_j 's present
 69 non-differentiability points, leading to the following extended notion of sparsity.

70 **Definition 4** (Generalized support). The generalized support of $\beta \in \mathbb{R}^p$ is the set of indices $j \in [p]$
 71 such that g_j is differentiable at β_j : $\operatorname{gsupp}(\beta) = \{j \in [p] : \partial g_j(\beta_j) \text{ is a singleton}\}$.

72 **Example 1.** With $\lambda, \gamma > 0$, MCP is defined as: $\text{MCP}_{\lambda, \gamma}(x) \mapsto \begin{cases} \lambda|x| - \frac{x^2}{2\gamma}, & \text{if } |x| \leq \gamma\lambda \\ \frac{1}{2}\gamma\lambda^2, & \text{if } |x| > \gamma\lambda \end{cases}$.

73 Penalties such as ℓ_1, ℓ_q ($0 < q < 1$), MCP or SCAD are only not differentiable at 0, and this
 74 corresponds to the usual notion of sparsity. But [Definition 4](#) goes beyond sparsity and extends to
 75 estimators such as SVM, where $g_j = \iota_{[0, C]}$ and the generalized support corresponds to support
 76 vectors $\{j \in [p] : \beta_j = 0 \text{ or } \beta_j = C\}$. The generalized support of a critical point is usually of
 77 cardinality much smaller than p , and its knowledge makes the problem easier and faster to solve. Our
 78 working set algorithm exploit this structure in order to converge faster.

79 2.2 Proposed algorithm

80 The proposed algorithm exploits two main ideas:

- 81 • A working set strategy, able to handle a large class of convex and non-convex penalties
 82 ([Algorithm 1](#)).
- 83 • An Anderson accelerated coordinate descent for non-convex problems ([Algorithm 2](#)). The
 84 building blocks of [Algorithm 2](#), coordinate descent (CD, [Algorithm 3](#)) and Anderson extrap-
 85 olation (Anderson, [Algorithm 4](#)), can be found in [Appendix A](#).

86 To avoid wasting computation on features outside the generalized support, working set algorithms
 87 iteratively select a subset of coordinates deemed important (the *working set*), and solve [Prob-
 88 lem \(1\)](#) restricted to them. The key question is thus the notion of *important* features. Stem-
 89 ming from [Definition 3](#), we rank features by their violation of the optimality condition: $\text{score}_j^\partial =$
 90 $\text{dist}(-\nabla_j f(\beta), \partial g_j(\beta))$. For example, the MCP Fréchet subdifferential at 0 is $\partial g_j(0) = [-\lambda, \lambda]$,
 91 and the proposed score reads

$$\text{score}_j^\partial = \begin{cases} \max\{0, |\nabla_j f(\beta)| - \lambda\} & \text{if } \beta_j = 0, \\ |\nabla_j f(\beta) + \nabla g_j(\beta_j)| & \text{otherwise.} \end{cases} \quad (2)$$

92 To control the working set growth, we use score_j^∂ to rank the features. Then, with $n_k =$
 93 $\max(n_{k-1}, 2 \lfloor \text{gsupp}(\beta^{(t)}) \rfloor)$ we take the n_k largest of them in the working set, while retaining
 94 features currently in the working set. This growth quickly rises to the unknown size of the generalized
 95 support while avoiding overshooting, as backed up by recent theory in [Ndiaye and Takeuchi \(2021\)](#).

96

Algorithm 1 skglm (proposed)

input: $X, \beta \in \mathbb{R}^p, n_{\text{out}} \in \mathbb{N},$
 $n_{\text{in}} \in \mathbb{N}, \text{ws_size} \in \mathbb{N}, \epsilon > 0$
 1 **for** $t = 1, \dots, n_{\text{out}}$ **do**
 2 $\text{score} = (\text{dist}(-\nabla_j f(\beta), \partial g_j(\beta_j)))_{j \in [p]}$
 3 $\text{ws_size} = \max(\text{ws_size}, 2 \times \lfloor \text{gsupp}(\beta) \rfloor)$
 // ws_size features with largest
 scores
 4 $\text{ws} = \text{arg_topK}(\text{score}, K = \text{ws_size})$
 5 **if** $\max_{j \in [p]} \text{dist}(-\nabla_j f(\beta), \partial g_j(\beta_j)) \leq \epsilon$
 then stop
 6 **else** // accelerated CD on working set
 7 $\beta \leftarrow \text{inner_solver}(X, \beta, \text{ws}, n_{\text{in}}, \epsilon)$
 8 **return** β

Algorithm 2 inner_solver

input: $X, \beta^{(0)} \in \mathbb{R}^p, \text{ws} \subset [p], n_{\text{in}}, \epsilon, M = 5$
 1 **for** $k = 1, \dots, n_{\text{in}}$ **do**
 2 $\beta^{(k)} \leftarrow \text{CD}(X, \beta^{(k-1)}, X\beta, \text{ws})$ // [Algo. \(3\)](#)
 3 **if** $k \bmod M = 0$ **then**
 // [Algo. \(4\)](#), $\mathcal{O}(M^2|\text{ws}| + M^3)$
 $\beta_{\text{ws}}^{\text{extr}} \leftarrow \text{Anderson}(\beta_{\text{ws}}^{(k-M)}, \dots, \beta_{\text{ws}}^{(M)})$
 // test objective $\mathcal{O}(n|\text{ws}|)$
 if $\Phi(\beta_{\text{ws}}^{\text{extr}}) < \Phi(\beta_{\text{ws}}^{(k)})$ **then**
 | $\beta_{\text{ws}}^{(k)} \leftarrow \beta_{\text{ws}}^{\text{extr}}, X\beta \leftarrow X_{\text{ws}}\beta_{\text{ws}}^{\text{extr}}$
 4 **if** $\max_{j \in \text{ws}} \text{dist}(-\nabla_j f(\beta), \partial g_j(\beta_j)) \leq \epsilon$
 then stop
 5 **if** $\max_{j \in \text{ws}} \text{dist}(-\nabla_j f(\beta), \partial g_j(\beta_j)) \leq \epsilon$
 then stop
 6 **if** $\max_{j \in \text{ws}} \text{dist}(-\nabla_j f(\beta), \partial g_j(\beta_j)) \leq \epsilon$
 then stop
 7 **if** $\max_{j \in \text{ws}} \text{dist}(-\nabla_j f(\beta), \partial g_j(\beta_j)) \leq \epsilon$
 then stop
 8 **return** $\beta^{(k)}$

97 **Proposition 5.** Let \mathcal{W}_t be the t -th working set. Suppose that [Algorithm 2](#) converges toward a critical
 98 point, and for all $t \geq 0, \mathcal{W}_t \subset \mathcal{W}_{t+1}$, then the iterates of [Algorithm 1](#) converges towards a critical
 99 point of [Problem \(1\)](#).

100 Proof of [Proposition 5](#) can be found in [Appendix B.1](#). The second key ingredient to our algorithm is to
 101 use state-of-the-art Anderson accelerated coordinate descent for non-convex problems. In [Section 2.3](#)
 102 we show that coordinate descent yields finite time support identification for a large class of non-
 103 convex problems ([Proposition 10](#)), which leads to acceleration ([Proposition 13](#)). As experiments
 104 demonstrate in [Section 3](#), this rate allows our algorithm to surpass state-of-the-art solvers.

105 **2.3 Anderson accelerated coordinate descent analysis for α -semi-convex penalties**

106 We now turn to our main technical contributions: we show that [Algorithm 2](#) achieves finite time
 107 support identification ([Proposition 10](#)) of the generalized support ([Definition 4](#)) for specific class of
 108 non-smooth non-convex penalties ([Assumption 6](#)), which includes the MCP ([Proposition 7](#)). Based
 109 on [Proposition 10](#), we are able to derive convergence rates for Anderson acceleration [Proposition 13](#).

110 We study our inner solver ([Algorithm 2](#)); for convenience we still refer to β and X for their counter-
 111 parts restricted to the working set. The following assumptions are required.

112 **Assumption 6** (α -semi-convex). *For all $j \in [p]$ g_j/L_j is α -semi-convex, i.e., $g_j/L_j + \alpha\|\cdot\|^2/2$ is
 113 convex, with $\alpha < 1$.*

114 Note that in statistics, the admissible value range of hyperparameters for MCP and SCAD are datafit-
 115 dependent, (see [Breheny and Huang 2011](#), Sec. 2.1, normalized columns and $\gamma > 1 = 1/\|X_{:j}\| =$
 116 $1/L_j$ or [Soubies et al. 2015](#), Eq. 4.2) and yields α -semi-convexity for MCP and SCAD¹.

117 **Proposition 7** (α -semi-convexity of MCP). *Let $\text{MCP}_{\lambda,\gamma}(x) \triangleq \begin{cases} \lambda|x| - \frac{x^2}{2\gamma}, & \text{if } |x| \leq \gamma\lambda \\ \frac{1}{2}\gamma\lambda^2, & \text{if } |x| > \gamma\lambda \end{cases}$.*

118 *If $\gamma > 1/L_j$, then $\text{MCP}_{\lambda,\gamma}/L_j$ is α -semi-convex (i.e., [Assumption 6](#) holds).*

119 Note that [Assumption 6](#) does not hold for the ℓ_q -penalties ($0 < q < 1$), for which we propose an
 120 alternative in [Appendix C](#).

121 **Assumption 8** (Existence). *[Problem \(1\)](#) admits at least one critical point.*

122 In [Proposition 10](#), convergence of [Algorithm 2](#) toward a critical point $\hat{\beta}$ is assumed, and the following
 123 assumption is made on this critical point.

124 **Assumption 9** (Non degeneracy). *The considered critical point $\hat{\beta} \in \mathbb{R}^p$ is non-degenerated: for all
 125 $j \notin \text{gsupp}(\hat{\beta})$, $-\nabla f_j(\hat{\beta}) \in \text{interior}(\partial g_j(\hat{\beta}_j))$.*

126 [Assumption 9](#) is a generalization of qualification constraints ([Hare and Lewis, 2007](#), Sec. 1), and is
 127 usual in the machine learning literature ([Zhao and Yu, 2006](#); [Bach, 2008](#); [Vaier et al., 2015](#)). For
 128 the ℓ_1 -norm, if the entries of the design matrix X are drawn from a i.i.d normal distribution, then
 129 [Assumption 9](#) holds with high probability ([Candes and Tao, 2005](#); [Rudelson and Vershynin, 2008](#)).

130 Equipped with the previous assumptions we show that coordinate descent achieves model identifica-
 131 tion for this class of non-convex problems.

132 **Proposition 10** (Model identification of CD). *Suppose*

- 133 1. *Assumptions 1, 2, 6 and 8 hold.*
- 134 2. *The sequence $(\beta^{(k)})_{k \geq 0}$ generated by coordinate descent ([Algorithm 2](#) without extrapola-
 135 tion) converges toward a critical point $\hat{\beta}$.*
- 136 3. *Assumption 9 holds for $\hat{\beta}$.*

137 *Then, [Algorithm 2](#) (without extrapolation) identifies the model in finitely many iterations: there exists
 138 $K > 0$ such that for all $k \geq K$, $\beta_{\mathcal{S}^c}^{(k)} = \hat{\beta}_{\mathcal{S}^c}$.*

139 In other words, for k large enough, $\beta^{(k)}$ shares the generalized support of $\hat{\beta}$. The identification
 140 property was proved for a proximal gradient descent algorithm in the non-convex case ([Liang et al.,
 141 2016](#)) under the assumption that the non-smooth function g is partly smooth ([Lewis, 2002](#)). For
 142 ourselves, [Proposition 10](#) not rely on the partly smooth assumption to ensure identification property.
 143 Authors are not aware of previous identification results for coordinate descent in the non-convex case.

144 In addition, if f and g are locally regular on the generalized support at the considered critical point,
 145 our algorithm enjoys local acceleration when combined with Anderson extrapolation ([Proposition 13](#)).
 146

147 **Assumption 11** (Locally \mathcal{C}^3). *For all $j \in \mathcal{S} \triangleq \text{gsupp}(\hat{\beta})$, g_j is locally \mathcal{C}^3 around $\hat{\beta}_j$, and f is locally
 148 \mathcal{C}^3 around $\hat{\beta}$.*

¹However MCP and SCAD are not α -semiconvex for all hyperparameter values.

Table 1: Most popular packages for sparse generalized linear models.

Name	Acceleration	Huge scale	Ncvx	Modular
glmnet (Friedman et al., 2010)	✗	✗	✗	✗ (Fortran)
scikit-learn (Pedregosa et al., 2011)	✗	✗	✗	✗ (Cython)
lightning (Blondel and Pedregosa, 2016)	✗	✗	✗	✓ (Cython)
celer (Massias et al., 2018)	✓	✓	✗	✗ (Cython)
picasso (Ge et al., 2019)	✗	✗	✓	✗ (C++)
pyGLMnet (Jas et al., 2020)	✗	✗✗	✗	✓ (Python)
fireworks (Rakotomamonjy et al., 2022)	✗	✓	✓	N.A. (Python)
skgglm (ours)	✓	✓	✓✓	✓ (Python)

149 Assumption 11 on the function f is mild and holds for usual machine learning datafitting terms.
 150 Assumption 11 on the functions g_j , $j \in \mathcal{S}$, is stronger: for instance, for the MCP, it implies $\hat{\beta}_j \neq \gamma\lambda$
 151 for all $j \in \mathcal{S}$. However this assumption is standard in the literature, see Liang et al. 2016, Sec. 3.3

152 **Assumption 12.** (Local strong convexity) The Hessian of f at the considered critical point $\hat{\beta} \in \mathbb{R}^p$,
 153 restricted to its generalized support \mathcal{S} , is positive definite, i.e., $\nabla_{\mathcal{S},\mathcal{S}}^2 f(\hat{\beta}) + \nabla_{\mathcal{S},\mathcal{S}}^2 g(\hat{\beta}) \succ 0$.

154 Assumption 12 requires local strong convexity restricted to the generalized support \mathcal{S} , which is
 155 standard in the MCP / SCAD literature (Breheny and Huang 2011, Section 4.1) and is usual to derive
 156 local linear rates of convergence (Liang et al., 2016, Section 3.3). For instance, for the Lasso, if the
 157 entries of the design matrix X are drawn from a continuous distribution, then Assumption 12 holds
 158 with probability one (Tibshirani, 2013, Lemma 4).

159 **Proposition 13.** Consider a critical point $\hat{\beta}$ and suppose

- 160 1. Assumptions 1, 2 and 8 hold.
- 161 2. The functions f and g_j , $j \in [p]$ are piecewise quadratic (which is the case for the MCP
 162 regression).
- 163 3. The sequence $(\beta^{(k)})_{k \geq 0}$ generated by anderson accelerated coordinate descent with updates
 164 from 1 to p and p to 1 (Algorithm 2 with extrapolation) converges to a critical point $\hat{\beta}$.
- 165 4. Assumptions 9, 11 and 12 hold for $\hat{\beta}$.

166 Then there exists $K \in \mathbb{N}$, and a \mathcal{C}^1 function $\psi : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$ such that, for all $k \in \mathbb{N}$, $k \geq K$:

$$\beta_j^{(k)} = \hat{\beta}_j, \text{ for all } j \in \mathcal{S}^c, \quad (3)$$

167 Let $T \triangleq \mathcal{J}\psi(\hat{\beta})$, $H \triangleq \nabla_{\mathcal{S},\mathcal{S}}^2 f(\hat{\beta}) + \nabla_{\mathcal{S},\mathcal{S}}^2 g(\hat{\beta})$, $\zeta \triangleq (1 - \sqrt{1 - \rho(T)}) / (1 + \sqrt{1 - \rho(T)})$ and
 168 $B \triangleq (T - \text{Id})^\top (T - \text{Id})$. Then $\rho(T) < 1$ and the iterates of Anderson extrapolation enjoy local
 169 accelerated convergence rate:

$$\|\beta_{\mathcal{S}}^{(k-K)} - \hat{\beta}_{\mathcal{S}}\|_B \leq \left(\sqrt{\kappa(H)} \frac{2\zeta^{M-1}}{1+\zeta^{2(M-1)}} \right)^{(k-K)/M} \|\beta_{\mathcal{S}}^{(K)} - \hat{\beta}_{\mathcal{S}}\|_B. \quad (4)$$

170 The proof can be found in Appendix B.5.

171 **Related work.** Most Anderson acceleration convergence results are shown for quadratic objectives for
 172 specific algorithms: gradient descent (Golub and Varga, 1961; Anderson, 1965), ADMM (Poon and
 173 Liang, 2019), coordinate descent (Bertrand et al., 2020). Outside of the quadratic case, convergence
 174 results are usually significantly weaker (Scieur et al., 2016; Sidi, 2017; Brezinski et al., 2018; Mai
 175 and Johansson, 2019; Ouyang et al., 2020). Regarding the smooth non-convex case, Wei et al.
 176 (2021) proposed a stochastic Anderson acceleration and proved convergence towards a critical point.
 177 Proposition 13 generalizes Scieur et al. (2020, Prop 2.1) and Bertrand and Massias (2021, Prop. 4) to
 178 the proximal convex and α -semi-convex cases. To our knowledge this is one of the first quantitative
 179 results for Anderson acceleration in a non-convex setting.

180 2.4 Comparison with existing work

181 In this section we compare our contribution to existing algorithms and implementations, which are
 182 summarized in Table 1. *Huge scale* refers to the fact that the algorithm can run on problems with

183 millions of variables. *Non-convex* tells if the algorithm handles non-convex penalties. *Modular*
184 indicates that it is easy to add a new model, through a different datafitting term or penalty.

185 The packages `glmnet` (Friedman et al., 2010), `scikit-learn` (Pedregosa et al., 2011) and
186 `lightning` (Blondel and Pedregosa, 2016) implement coordinate descent (cyclic or random). They
187 rely on compiled code such as Fortran or Cython, making it very difficult to implement new models²
188 or faster algorithms like working set³. They do not handle non-convex penalties.

189 More recent algorithms such as `blitz` (Johnson and Guestrin, 2015), `celer` (Massias et al., 2018),
190 `picasso` (Ge et al., 2019) or `fireworks` (Rakotomamonjy et al., 2022) use working set strategies.
191 `celer` and `blitz` are state-of-the-art algorithms for the Lasso, but their score to prioritize features
192 relies on duality. `fireworks` extends `blitz` to some non-convex penalties (writing as difference of
193 convex functions), with $\text{score}_j^{\text{fireworks}} = \text{dist}(-\nabla_j f(\beta), \partial g_j(0))$. Yet this rule does not consider the
194 subdifferential of g at the current point, but at 0, which is a coarse information. Finally, `fireworks`
195 does not provide accelerated convergence rates and does not come with a public implementation.
196 `picasso` (Ge et al., 2019) lacks modularity (penalties are hardcoded), and the solver is not suited for
197 huge scale (it does not support large sparse matrices). Deng and Lan (2019) proposed an algorithm
198 based on inertially accelerated coordinate descent, which fails to provide practical speedups according
199 to Bertrand and Massias (2021).

200 Contrary to these algorithms, ours is generic and relies only on the knowledge of ∇f and prox_g . For
201 any new penalty, this information can be written in a few lines of Python code, compiled with `numba`
202 (Lam et al., 2015) for speed efficiency. We therefore improve of state-of-the-art algorithms in the
203 convex case, and generalize to virtually any datafit and penalty, even nonconvex.

204 3 Experiments

205 Our package relying on `numpy` and `numba` (Lam et al., 2015; Harris et al., 2020) is attached in the
206 supplementary material and will be open-sourced upon publication under BSD 3-Clause License. We
207 use datasets from `libsvm`⁴ (Fan et al. 2008, see table 2).

208 We compare multiple algorithms to solve popular Machine Learning and inverse problems: Lasso,
209 Elastic net, multitask sparse regression, MCP regression. The compared algorithms are the following:

- 210 • `scikit-learn` (Pedregosa et al., 2011), which implements coordinate descent in Cython,
- 211 • `celer` (Massias et al., 2020), which combines working sets, screening rules, coordinate
212 descent, and Anderson acceleration in the dual, in Cython,
- 213 • `blitz` (Johnson and Guestrin, 2015), which combines working sets with prox-Newton
214 iterations (Lee et al., 2012) in C++,
- 215 • coordinate descent (CD, Tseng and S.Yun 2009),
- 216 • `skglm` (Algorithm 1, ours), using $M = 5$ iterates for the Anderson extrapolation.

217 **Other solvers.** Experiment per experiment, there exist niche solvers (such as aggressive Gap Safe
218 Rules, Ndiaye et al. 2020). Since our goal is a *general purpose* algorithm able to deal with many
219 models, we do not include them in the comparison. In addition, we focus on solving a single
220 instance of Problem (1), rather than a regularization path (*i.e.*, a sequence of problems for multiple
221 regularization strengths). As `glmnet` is designed to compute regularization paths, we could not
222 include it in the comparison. The reader can refer to Johnson and Guestrin (2015, Fig. 4) or Figure 8
223 in Appendix E for comparisons on single optimization problems with `glmnet`. `glmnet` and additional
224 algorithms are discussed in Appendix E.

225 **How to do a fair comparison between solvers?** To plot the convergence curves, we use the
226 `benchopt`⁵ benchmarking package (Moreau et al., 2022). In order to automate and reproduce
227 optimization benchmarks it treats solvers as black boxes. It launches them several times with
228 increasing maximum number of iterations, and stores the resulting objective values and times to reach
229 it. As each point on a solver curve is obtained in a different run, the curves are not monotonic, and
230 there may be several points corresponding to the same time. This merely reflects the variability in

²<https://github.com/scikit-learn/scikit-learn/pull/10745> (4 years old)

³<https://github.com/scikit-learn/scikit-learn/pull/7853> (5 years old)

⁴<https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

⁵<https://github.com/benchopt/benchopt>

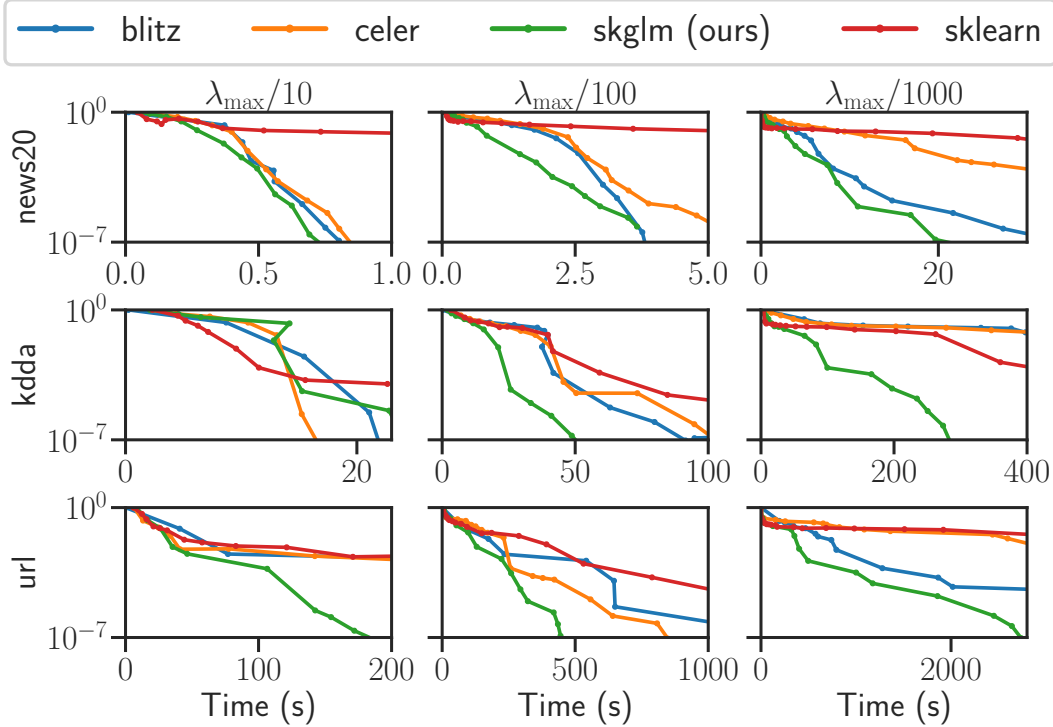


Figure 2: **Lasso, duality gap.** Normalized duality gap as a function of time for the Lasso on multiple datasets, for multiple values of λ .

231 solvers running time across runs; we refer to [Figure 11](#) in [Appendix E.7](#) for the inevitability of this
 232 phenomenon with black box solvers.

233 3.1 Convex problems

234 **Lasso.** In [Figure 2](#) we compare solvers for the Lasso: ($f = \frac{1}{2n} \|y - X \cdot\|^2$, $g_j = \lambda |\cdot|$). We
 235 parametrize λ as a fraction of $\lambda_{\max} = \|X^\top y\|_\infty / n$, smallest regularization strength for which $\hat{\beta} = 0$.
 236 For large scale datasets (*rcv1*, *news20*), *skglm* yields performances better or similar to the state-of-
 237 the-art algorithms *blitz* and *celer*. For huge scale datasets (*kdda* and *url*), *skglm* yields significant
 238 speedups over them. The improvement over the popular *scikit-learn* can be of two orders of
 239 magnitude. Thus, *while dealing with many more models, our algorithm still yields state-of-the-art*
 240 *speed for basic ones.*

241 **Elastic net.** Our approach easily generalizes to other problems, such as the elastic net ($f =$
 242 $\frac{1}{2n} \|y - X \cdot\|^2$, $g_j = \lambda(\rho |\cdot| + \frac{1-\rho}{2} (\cdot)^2)$). [Figure 3](#) shows the duality gap as a function of time
 243 for *skglm* (ours), *sklearn*, and our numba implementation of coordinate descent. The proposed
 244 algorithm is orders of magnitude faster than *scikit-learn* and vanilla coordinate descent, in
 245 particular for large datasets and low regularization parameter values (*finance*, $\lambda_{\max}/1000$). Note that
 246 *blitz* does not implement a solver for the elastic net. Many Lasso solvers would easily handle the
 247 elastic net, but relying on Cython/C++ code makes the implementation time-consuming. By contrast,
 248 it takes 40 lines of code to define an $\ell_1 + \ell_2$ -squared penalty with our implementation. An additional
 249 experiment on the dual of SVM with hinge loss is in [Appendix E.4](#).

250 3.2 Non-convex problems

251 In this subsection we propose a comparison on two non convex problems.

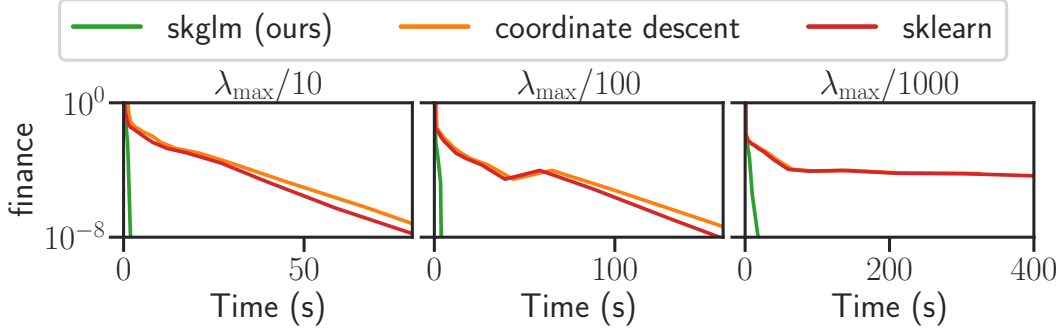


Figure 3: **Elastic net, duality gap.** Normalized duality gap as a function of time for the elastic net for multiple values of λ , $\rho = 0.5$.

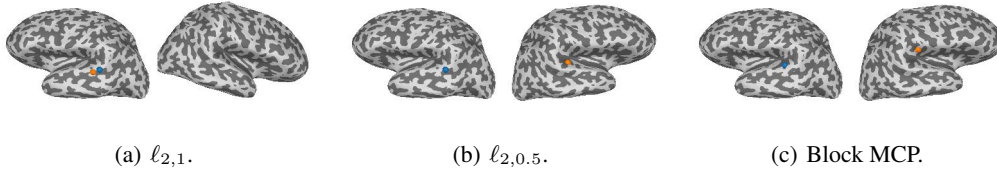


Figure 4: **Real data, brain source locations recovered by convex and non-convex penalties after a right auditory stimulation.** 4(a) shows that a convex penalty fails at identifying one source in each hemisphere, while 4(b) and 4(c) demonstrates the capability of non-convex penalties to recover the correct solution.

252 **MCP regression.** MCP regression is Problem (1) with $f = \frac{1}{2n} \|y - X \cdot\|^2$, $g_j = \text{MCP}_{\lambda, \gamma}$ for
 253 $\gamma > 1$. As usual for this problem, we scale the columns of X to have norm \sqrt{n} . On Figure 5,
 254 we compare our algorithm to `picasso` on a dense dataset ($n = 1000, p = 5000$); as this package
 255 does not support large sparse design matrices, for the `rcv1` dataset we use an iterative reweighted
 256 L1 algorithm (Candes et al., 2008). Since the derivative of the MCP vanishes for values bigger than
 257 $\lambda\gamma$, this approach requires solving weighted Lasso with some 0 weights. Up to our knowledge, our
 258 algorithm is the only efficient one with such a property. Our algorithm handles problems of large
 259 size, converges to a critical point, and, due to its progressive inclusion of features, is able to reach a
 260 sparser critical point than it competitors.

261 **Application to neuroscience** To demonstrate the usefulness of our algorithm for practitioners,
 262 we apply it to the magneto-/electroencephalographic (M/EEG) inverse problem. It consists in
 263 reconstructing the spatial cortical current density at the origin of M/EEG measurements made at
 264 the surface of the scalp. Non-convex penalties (Strohmeier et al., 2015) exhibit several advantages
 265 over convex ones (Gramfort et al., 2013): they yield sparser physiologically-plausible solutions and
 266 mitigate the ℓ_1 amplitude bias. Here the setting is multitask: $Y \in \mathbb{R}^{n \times T}$ and thus we use block
 267 penalties (details in Appendix D). We use real data from the `mne` software (Gramfort et al., 2014);
 268 the experiment is a right auditory stimulation, with two expected neural sources to recover in each
 269 auditory cortex. In Figure 4, while the $\ell_{2,1}$ penalty fails at localizing one source in each hemisphere,
 270 the non-convex penalties recover the correct locations. This emphasizes on the critical need for fast
 271 solvers for non-convex sparse penalties as well as our algorithm’s ability to handle the latter.

272 **Ablation study.** To evaluate the influence of the two components of Algorithm 1, an ablation study
 273 (Figure 6 and Figure 10 in Appendix E) is performed. Four algorithms are compared: with/without
 274 working sets and with/without Anderson acceleration. Figure 6 represents the duality gap of the Lasso
 275 as a function of time for multiple datasets and values of the regularization parameters λ (parametrized
 276 as a fraction of λ_{\max}). First, Figure 6 shows that working sets always bring significant speedups. Then,
 277 when combined with working set, Anderson acceleration bring significant speed-ups, especially for
 278 hard problems with low regularization parameters. An interesting observation is that on large scale

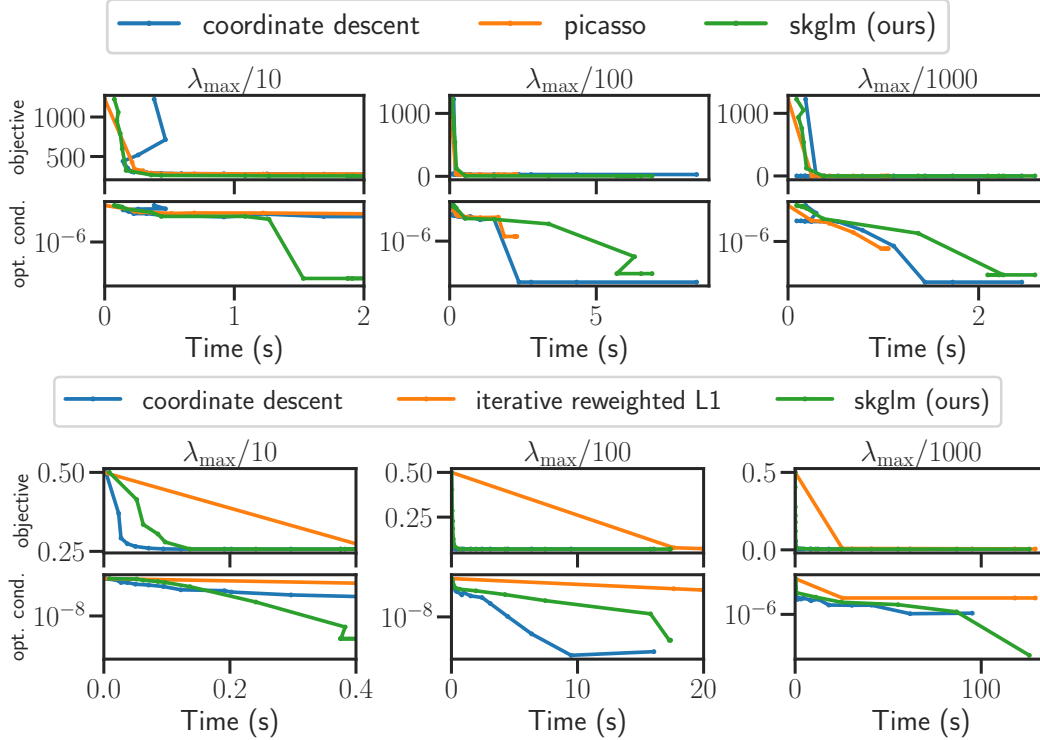


Figure 5: **MCP, objective value and violation of first order condition.** Objective value and violation of optimality condition of the iterates, $\text{dist}(-\nabla f(\beta^{(k)}), \partial g(\beta^{(k)}))$, as a function of time for the MCP for multiple values of λ ($\gamma = 3$) on a simulated dense dataset (top) and the rcv1 dataset (normalized columns).

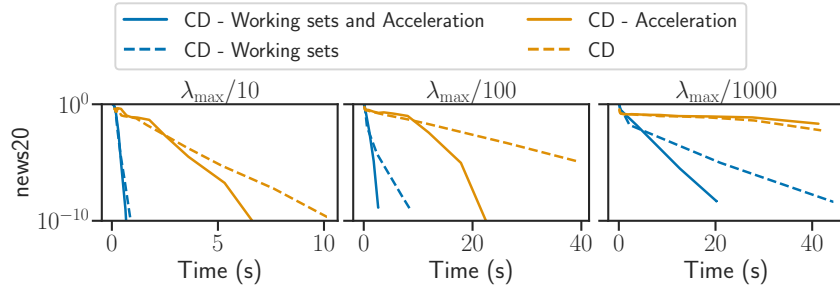


Figure 6: **Lasso, duality gap.** Normalized duality gap as a function of time for the Lasso.

279 datasets (*news20* and *finance* in Appendix E.6) and for low regularization parameters ($\lambda_{\max}/100$ and
 280 $\lambda_{\max}/1000$) Anderson acceleration *without* working set does not bring acceleration. This highlights
 281 the importance of combining Anderson acceleration with working sets.

282 **Conclusion and broader impact.** In this paper, we have proposed an accelerated versatile algo-
 283 rithm for a specific class of non-smooth non-convex problems. Based on working sets, coordinate
 284 descent and Anderson acceleration, we have improved state-of-the-art on convex problems, and han-
 285 dled previously out-of-reach problems. Thorough experiments demonstrated the speed and interest
 286 of our approach. A limitation of this work is the considered function class (α -semi-convex), which
 287 can be seen as restrictive. One possible extension would be weakly convex functions (Davis and
 288 Drusvyatskiy, 2019, Sec. 1). We deeply believe that the high quality code provided will benefit to
 289 practitioners, and ease the use of non-convex penalties for real world problems, from neuroimaging to
 290 genomics. We proposed an optimization algorithm and do not see potential negative societal impacts.

291 **References**

- 292 D. G. Anderson. Iterative procedures for nonlinear integral equations. *Journal of the ACM*, 12(4):
293 547–560, 1965.
- 294 H. Attouch and J. Bolte. On the convergence of the proximal algorithm for nonsmooth functions
295 involving analytic features. *Mathematical Programming*, 116(1):5–16, 2009.
- 296 H. Attouch, J. Bolte, and B. F. Svaiter. Convergence of descent methods for semi-algebraic and
297 tame problems: proximal algorithms, forward–backward splitting, and regularized Gauss–Seidel
298 methods. *Mathematical Programming*, 137(1):91–129, 2013.
- 299 F. Bach. Consistency of the group Lasso and multiple kernel learning. *J. Mach. Learn. Res.*, 9:
300 1179–1225, 2008.
- 301 Q. Bertrand and M. Massias. Anderson acceleration of coordinate descent. In *AISTATS*, 2021.
- 302 Q. Bertrand, Q. Kloppenstein, M. Blondel, S. Vaiter, A. Gramfort, and J. Salmon. Implicit differentia-
303 tion of lasso-type models for hyperparameter optimization. *ICML*, 2020.
- 304 M. Blondel and F. Pedregosa. Lightning: large-scale linear classification, regression and ranking in
305 python, 2016.
- 306 J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex
307 and nonsmooth problems. *Mathematical Programming*, 146(1):459–494, 2014.
- 308 A. Bonnefoy, V. Emiya, L. Ralaivola, and R. Gribonval. Dynamic screening: accelerating first-order
309 algorithms for the Lasso and Group-Lasso. *IEEE Trans. Signal Process.*, 63(19):20, 2015.
- 310 S. Boyd, N. Parikh, and E. Chu. *Distributed optimization and statistical learning via the alternating
311 direction method of multipliers*. Now Publishers Inc, 2011.
- 312 P. Breheny and J. Huang. Coordinate descent algorithms for nonconvex penalized regression, with
313 applications to biological feature selection. *The annals of applied statistics*, 5(1):232, 2011.
- 314 C. Brezinski, M. Redivo-Zaglia, and Y. Saad. Shanks sequence transformations and anderson
315 acceleration. *SIAM Review*, 60(3):646–669, 2018.
- 316 E. J. Candes and T. Tao. Decoding by linear programming. *IEEE transactions on information theory*,
317 51(12):4203–4215, 2005.
- 318 E. J. Candes, M. B. Wakin, and S. P. Boyd. Enhancing sparsity by reweighted ℓ_1 minimization.
319 *Journal of Fourier analysis and applications*, 14(5):877–905, 2008.
- 320 D. Davis and D. Drusvyatskiy. Stochastic model-based minimization of weakly convex functions.
321 *SIAM Journal on Optimization*, 29(1):207–239, 2019.
- 322 Q. Deng and C. Lan. Efficiency of coordinate descent methods for structured nonconvex optimization.
323 *arXiv preprint arXiv:1909.00918*, 2019.
- 324 L. El Ghaoui, V. Viallon, and T. Rabbani. Safe feature elimination for the lasso and sparse supervised
325 learning problems. *arXiv preprint arXiv:1009.4219*, 2010.
- 326 R. E. Fan, K. W. Chang, C. J. Hsieh, X. R. Wang, and C. J. Lin. Liblinear: A library for large linear
327 classification. *JMLR*, 9:1871–1874, 2008.
- 328 O. Fercoq and P. Richtárik. Accelerated, parallel, and proximal coordinate descent. *SIAM Journal on
329 Optimization*, 25(4):1997–2023, 2015.
- 330 O. Fercoq, A. Gramfort, and J. Salmon. Mind the duality gap: safer rules for the lasso. In *ICML*,
331 pages 333–342. PMLR, 2015.
- 332 S. Foucart and M.-J. Lai. Sparsest solutions of underdetermined linear systems via ℓ_q -minimization
333 for $0 < q \leq 1$. *Applied and Computational Harmonic Analysis*, 26(3):395–407, 2009.

- 334 J. Friedman, T. J. Hastie, H. Höfling, and R. Tibshirani. Pathwise coordinate optimization. *Ann. Appl.*
335 *Stat.*, 1(2):302–332, 2007.
- 336 J. Friedman, T. J. Hastie, and R. Tibshirani. Regularization paths for generalized linear models via
337 coordinate descent. *J. Stat. Softw.*, 33(1):1, 2010.
- 338 J. Ge, X. Li, H. Jiang, H. Liu, T. Zhang, M. Wang, and T. Zhao. Picasso: A sparse learning library
339 for high dimensional data analysis in r and python. *The Journal of Machine Learning Research*, 20
340 (1):1692–1696, 2019.
- 341 D. Ghosh and A. M. Chinnaiyan. Classification and selection of biomarkers in genomic data using
342 lasso. *Journal of Biomedicine and Biotechnology*, 2005(2):147, 2005.
- 343 G. H. Golub and R. S. Varga. Chebyshev semi-iterative methods, successive overrelaxation iterative
344 methods, and second order richardson iterative methods. *Numerische Mathematik*, 3(1):147–156,
345 1961.
- 346 A. Gramfort, D. Strohmeier, J. Haueisen, M. S. Hämäläinen, and M. Kowalski. Time-frequency mixed-
347 norm estimates: Sparse M/EEG imaging with non-stationary source activations. *NeuroImage*, 70:
348 410–422, 2013.
- 349 A. Gramfort, M. Luessi, E. Larson, D. A. Engemann, D. Strohmeier, C. Brodbeck, L. Parkkonen, and
350 M. S. Hämäläinen. MNE software for processing MEG and EEG data. *NeuroImage*, 86:446 – 460,
351 2014.
- 352 W. L. Hare and A. S. Lewis. Identifying active manifolds. *Algorithmic Operations Research*, 2(2):
353 75–75, 2007.
- 354 C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau,
355 E. Wieser, J. Taylor, S. Berg, N. J. Smith, et al. Array programming with numpy. *arXiv preprint*
356 *arXiv:2006.10256*, 2020.
- 357 M. Jas, T. Achakulvisut, A. Idrizović, D. Acuna, M. Antalek, V. Marques, T. Odland, R. Garg,
358 M. Agrawal, Y. Umegaki, P. Foley, H. Fernandes, D. Harris, B. Li, O. Pieters, S. Otterson, G. De
359 Toni, C. Rodgers, E. Dyer, M. Hamalainen, K. Kording, and P. Ramkumar. Pyglmnet: Python
360 implementation of elastic-net regularized generalized linear models. *Journal of Open Source*
361 *Software*, 5(47):1959, 2020.
- 362 T. B. Johnson and C. Guestrin. Blitz: A principled meta-algorithm for scaling sparse optimization. In
363 *ICML*, volume 37, pages 1171–1179, 2015.
- 364 Y. J. Kim, N. Brackbill, E. Batty, J. Lee, C. Mitelut, W. Tong, EJ Chichilnisky, and L. Paninski.
365 Nonlinear decoding of natural images from large-scale primate retinal ganglion recordings. *Neural*
366 *Computation*, 33(7):1719–1750, 2021.
- 367 Q. Klopfenstein, Q. Bertrand, A. Gramfort, J. Salmon, and S. Vaiteer. Model identification and local
368 linear convergence of coordinate descent. *arXiv preprint arXiv:2010.11825*, 2020.
- 369 A. Y. Kruger. On Fréchet subdifferentials. *Journal of Mathematical Sciences*, 116(3):3325–3358,
370 2003.
- 371 S. K. Lam, A. Pitrou, and S. Seibert. Numba: A llvm-based python jit compiler. In *Proceedings of*
372 *the Second Workshop on the LLVM Compiler Infrastructure in HPC*, pages 1–6, 2015.
- 373 M. Le Morvan and J.-P. Vert. WHInter: A working set algorithm for high-dimensional sparse second
374 order interaction models. In *ICML*, pages 3635–3644. PMLR, 2018.
- 375 J. D. Lee, Y. Sun, and M. Saunders. Proximal Newton-type methods for convex optimization.
376 *Advances in Neural Information Processing Systems*, 25:827–835, 2012.
- 377 A. S. Lewis. Active sets, nonsmoothness, and sensitivity. *SIAM Journal on Optimization*, 13(3):
378 702–725, 2002.

- 379 J. Liang, J. Fadili, and G. Peyré. A multi-step inertial forward-backward splitting method for
380 non-convex optimization. *Advances in Neural Information Processing Systems*, 29:4035–4043,
381 2016.
- 382 Q. Lin, Z. Lu, and L. Xiao. An accelerated proximal coordinate gradient method. In *NeurIPS*, pages
383 3059–3067. 2014.
- 384 D. C. Liu and J. Nocedal. On the limited memory BFGS method for large scale optimization.
385 *Mathematical programming*, 45(1):503–528, 1989.
- 386 V. V. Mai and M. Johansson. Anderson acceleration of proximal gradient methods. In *ICML*. 2019.
- 387 J. Mairal. *Sparse coding for machine learning, image processing and computer vision*. PhD thesis,
388 École normale supérieure de Cachan, 2010.
- 389 M. Massias, A. Gramfort, and J. Salmon. Celer: a fast solver for the lasso with dual extrapolation.
390 2018.
- 391 M. Massias, S. Vaïter, A. Gramfort, and J. Salmon. Dual extrapolation for sparse generalized linear
392 models. *J. Mach. Learn. Res.*, 2020.
- 393 R. Mazumder, J. H. Friedman, and T. Hastie. Sparsenet: Coordinate descent with nonconvex penalties.
394 *Journal of the American Statistical Association*, 106(495):1125–1138, 2011.
- 395 T. Moreau, M. Massias, A. Gramfort, P. Ablin, B. Charlier, P.-A. Bannier, M. Dagréou, T. Dupré
396 la Tour, G. Durif, C. F. Dantas, Q. Klopfenstein, et al. Benchopt: Reproducible, efficient and
397 collaborative optimization benchmarks. *arXiv preprint arXiv:2206.13424*, 2022.
- 398 J. B. Muir and Z. Zhan. Seismic wavefield reconstruction using a pre-conditioned wavelet–curvelet
399 compressive sensing approach. *Geophysical Journal International*, 227(1):303–315, 2021.
- 400 E. Ndiaye and I. Takeuchi. Continuation path with linear convergence rate. *arXiv preprint*
401 *arXiv:2112.05104*, 2021.
- 402 E. Ndiaye, O. Fercoq, A. Gramfort, and J. Salmon. Gap safe screening rules for sparsity enforcing
403 penalties. *J. Mach. Learn. Res.*, 18(128):1–33, 2017.
- 404 E. Ndiaye, O. Fercoq, and J. Salmon. Screening rules and its complexity for active set identification.
405 *Journal of Convex Analysis*, 2020.
- 406 Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM*
407 *Journal on Optimization*, 22(2):341–362, 2012.
- 408 A. Y. Ng. Feature selection, l_1 vs. l_2 regularization, and rotational invariance. In *ICML*, page 78,
409 2004.
- 410 J. Nutini. *Greed is good: greedy optimization methods for large-scale structured problems*. PhD
411 thesis, University of British Columbia, 2018.
- 412 J. Nutini, M. W. Schmidt, I. H. Laradji, M. P. Friedlander, and H. A. Koepke. Coordinate descent
413 converges faster with the Gauss-Southwell rule than random selection. In *ICML*, pages 1632–1641,
414 2015.
- 415 W. Ouyang, Y. Peng, Y. Yao, J. Zhang, and B. Deng. Anderson acceleration for nonconvex ADMM
416 based on Douglas-Rachford splitting. In *Computer Graphics Forum*, volume 39, pages 221–239.
417 Wiley Online Library, 2020.
- 418 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
419 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
420 E. Duchesnay. Scikit-learn: Machine learning in Python. *JMLR*, 12:2825–2830, 2011.
- 421 C. Poon and J. Liang. Trajectory of alternating direction method of multipliers and adaptive accelera-
422 tion. In *NeurIPS*, pages 7357–7365, 2019.

- 423 A. Rakotomamonjy, G. Gasso, and J. Salmon. Screening rules for lasso with non-convex sparse
424 regularizers. In *ICML*, pages 5341–5350, 2019.
- 425 A. Rakotomamonjy, R. Flamary, G. Gasso, and J. Salmon. Provably convergent working set algorithm
426 for non-convex regularized regression. In *AISTATS*, 2022.
- 427 D. A. Reidenbach, A. Lal, L. Slim, O. Mosafi, and J. Israeli. Gepsi: A python library to simulate
428 gwas phenotype data. *bioRxiv*, 2021.
- 429 P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for
430 minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2014.
- 431 M. Rudelson and R. Vershynin. On sparse reconstruction from fourier and gaussian measurements.
432 *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of*
433 *Mathematical Sciences*, 61(8):1025–1045, 2008.
- 434 D. Scieur, A. d’Aspremont, and F. Bach. Regularized nonlinear acceleration. In *Advances In Neural*
435 *Information Processing Systems*, pages 712–720, 2016.
- 436 D. Scieur, A. d’Aspremont, and F. Bach. Regularized nonlinear acceleration. *Mathematical Program-*
437 *ming*, 179(1):47–83, 2020.
- 438 A. Sidi. *Vector extrapolation methods with applications*. SIAM, 2017.
- 439 N. Simon, J. Friedman, T. J. Hastie, and R. Tibshirani. A sparse-group lasso. *J. Comput. Graph.*
440 *Statist.*, 22(2):231–245, 2013.
- 441 E. Soubies, L. Blanc-Féraud, and G. Aubert. A continuous exact ℓ_0 penalty (cel0) for least squares
442 regularized problem. *SIAM Journal on Imaging Sciences*, 8(3):1607–1639, 2015.
- 443 D. Strohmeier, A. Gramfort, and J. Haueisen. MEG/EEG source imaging with a non-convex penalty in
444 the time-frequency domain. In *Pattern Recognition in Neuroimaging, 2015 International Workshop*
445 *on*, 2015.
- 446 D. Strohmeier, Y. Bekhti, J. Haueisen, and A. Gramfort. The iterative reweighted mixed-norm
447 estimate for spatio-temporal MEG/EEG source reconstruction. *IEEE transactions on medical*
448 *imaging*, 35(10):2218–2228, 2016.
- 449 R. Tibshirani. Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Ser. B Stat. Methodol.*,
450 58(1):267–288, 1996.
- 451 R. Tibshirani, J. Bien, J. Friedman, T. Hastie, N. Simon, J. Taylor, and R. J. Tibshirani. Strong rules
452 for discarding predictors in lasso-type problems. *Journal of the Royal Statistical Society: Series B*
453 *(Statistical Methodology)*, 74(2):245–266, 2012.
- 454 R. J. Tibshirani. The lasso problem and uniqueness. *Electronic Journal of statistics*, 7:1456–1490,
455 2013.
- 456 P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization.
457 *Mathematical Programming*, 117(1):387–423, 2009.
- 458 S. Vaiter, G. Peyré, and J. Fadili. Low complexity regularization of linear inverse problems. In
459 *Sampling Theory, a Renaissance*, pages 103–153. Springer, 2015.
- 460 F. Wei, C. Bao, and Y. Liu. Stochastic anderson mixing for nonconvex stochastic optimization.
461 *Advances in Neural Information Processing Systems*, 34, 2021.
- 462 F. Wen, L. Chu, P. Liu, and R. Qiu. A survey on nonconvex regularization-based sparse and low-rank
463 recovery in signal processing, statistics, and machine learning. *IEEE Access*, 6:69883–69906,
464 2018.
- 465 C.-H. Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of*
466 *statistics*, 38(2):894–942, 2010.

- 467 P. Zhao and B. Yu. On model selection consistency of lasso. *J. Mach. Learn. Res.*, 7:2541–2563,
468 2006.
- 469 H. Zou and T. J. Hastie. Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser.*
470 *B Stat. Methodol.*, 67(2):301–320, 2005.

471 **Checklist**

- 472 1. For all authors...
- 473 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
474 contributions and scope? [Yes]
- 475 (b) Did you describe the limitations of your work? [Yes] See limitations paragraph
- 476 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 477 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
478 them? [Yes]
- 479 2. If you are including theoretical results...
- 480 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See [Proposi-](#)
481 [tions 10, 13 and 14.](#)
- 482 (b) Did you include complete proofs of all theoretical results? [Yes] See [Appendix B.](#)
- 483 3. If you ran experiments...
- 484 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
485 mental results (either in the supplemental material or as a URL)? [Yes]
- 486 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
487 were chosen)? [Yes] See [Section 3.](#)
- 488 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
489 ments multiple times)? [N/A]
- 490 (d) Did you include the total amount of compute and the type of resources used (e.g., type
491 of GPUs, internal cluster, or cloud provider)? [N/A]
- 492 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 493 (a) If your work uses existing assets, did you cite the creators? [Yes] In particular we
494 acknowledge the python ecosystem.
- 495 (b) Did you mention the license of the assets? [N/A]
- 496 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 497 (d) Did you discuss whether and how consent was obtained from people whose data you're
498 using/curating? [N/A]
- 499 (e) Did you discuss whether the data you are using/curating contains personally identifiable
500 information or offensive content? [N/A]
- 501 5. If you used crowdsourcing or conducted research with human subjects...
- 502 (a) Did you include the full text of instructions given to participants and screenshots, if
503 applicable? [N/A]
- 504 (b) Did you describe any potential participant risks, with links to Institutional Review
505 Board (IRB) approvals, if applicable? [N/A]
- 506 (c) Did you include the estimated hourly wage paid to participants and the total amount
507 spent on participant compensation? [N/A]