

1 We thank reviewers for useful comments. **All reviewers:** We look at related questions: **(i) How is LCPP different**
2 **from proximal point and what role is played by g_k ,** **(ii) Other experiments beyond MCP and on inverse problem.**
3 **(i)** LCPP is different from proximal point as it uses proximal point in the objective and convexification in constraint.
4 Since, we convexify the constraint in g_k , we can show that subproblem can be solved efficiently and approximate final
5 solution will be feasible. Feasibility allows to characterize a simple MFCQ condition under which \bar{y}^k is bounded.
6 As demonstrated in Fig 1 below, (a)-(c) shows the original nonconvex set and convex subsets dynamically generated
7 in the run of algorithm. It is unclear whether these subsets will be algorithmically well behaved, i.e., the optimal
8 Lagrange multiplier for such constraints will be small. Indeed, as shown in Fig 1(d)-(e), as prox center gets close to
9 the point which violates MFCQ, the convex polyhedron flattens and bound on Lagrange multiplier blows up. Our
10 analysis precisely characterizes such bad conditions and can also provide a priori bounds, giving us full control over
11 complications that may arise in dynamically generated convex subproblems (see Appendix B). This need of bounded
12 \bar{y}^k is crucial for convergence to KKT of constrained optimization and is not a concern for unconstrained or simple
13 set constrained problems in which usual proximal point operates. **(ii)** Due to the space limit in the submitted version,
14 we only report preliminary results on MCP with logistic regression objective. We will compare with other nonconvex
15 penalties and inverse problem objective in a later version.

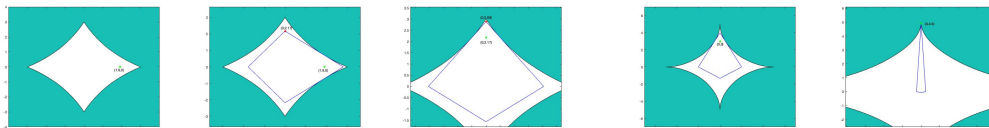


Figure 1: (a)-(e) SCAD constrained optimization. Prox center denoted by green.

16 **Reviewer 1: Weakness** The SCAD paper [12] showed the advantage of nonconvex surrogate penalty, and pointed out
17 the potential problem of biased estimation of large coefficients in Lasso since ℓ_1 penalty does not flatten. Empirical
18 merit of using nonconvex constraint to relax ℓ_0 is also demonstrated by examples in [25]. By suboptimality, we meant
19 to say that using nonconvex formulation is advantageous over ℓ_1 due to the above discussions in the literature and used
20 it to merely motivate the main topic of this paper. We can always improve the language as pointed out by the reviewer.
21 Figure 3 does show that LCPP outperforms Lasso across a wide range of nonzero patterns. See below in other comment.
22 We agree that disclosing form of the constraint beforehand will make presentation better, succinct and allow to have
23 fairer comparison with [5,20] in Section 1. At first glance, it may seem that constraint form is too restrictive, however,
24 it covers variety of nonconvex sparsity inducing functions considered in the literature, e.g., [16,29]. They assume
25 $g(x) = g_1(x) - h(x)$ and we identify more structure that can be exploited to get efficient method for different problems.
26 **Related work** Thanks for pointing out ‘CLASH’ paper which we will add in the related work. While there exists more
27 advanced models for specific type of problems such as CLASH, they often require strong assumption on data (i.e. RIP)
28 and specific function type. In comparison, LCPP aims to deal with general loss and a variety of nonconvex penalties.
29 The proof technique and convergence results are substantially different from CLASH.

30 **Other Comment: Fig 3:** We observe that it is better to impose sparsity constraint. In gisette rcv1.binary and real-sim,
31 using too many features results in overfitting. In mnist, it is possible to obtain nearly the best performance with only
32 using nearly half the features. **Why (5) is hard:** Hardness is not merely in (5) but the question that we want a feasible
33 solution to (5) with faster convergence and with less restrictive constraint qualification that works for nonconvex
34 constraints of interest. Variety of algorithms in classic nonlinear programming [2] showed asymptotic convergence with
35 or without feasibility and [5,20] addressed the complexity part partially but the constraint qualifications requirements
36 were strong to ensure feasibility. We believe addressing these questions is hard for general problems but rates can be
37 further improved for particular class of problems. **Why choose (5) over (4):** While (5) controls the target penalty level
38 directly, (4) controls the level implicitly by tuning weight λ . Unlike the convex case, first-order optimal point of (4)
39 does not guarantee a small value of penalty for large λ . **Matrix rank:** Problem of rank constraint on matrices instead
40 of ℓ_0 -norm of vector is a natural extension. See [36] for more information. Very interesting comment.

41 **Reviewer 2: W1.** We compare with Logistic regression which uses liblinear (dual coordinate descent) through sklearn
42 interface. We also compare with GIST, an efficient nonconvex proximal gradient method for DC penalized problem.
43 These are quite standard solvers for convex and nonconvex sparse models. W2. In our problem, η is a fixed target
44 constraint level and we don’t change its value during the algorithm. By adding a constant to η , we are dealing with a
45 new problem with a more relaxed constraint. If same constant is added on both side of (5) then algorithm does not
46 change at all. W3. To the best of our knowledge, operator splitting and nonconvex ADMM methods are provably
47 convergent only when the constraint function is linear. They are not applicable for (5).

48 **Reviewer 3** Thank you for pointing out Catalyst. As per our knowledge, Frank-Wolfe (FW) method has complexity
49 guarantee if the constraint set is convex which isn’t the case in (5). It seems possible to apply FW to solve the convex
50 subproblem in LCPP. We believe that adaptive strategy of Catalyst can be applied in our paper, since the subproblem is
51 convex and admits efficient proximal map. Table 1 describes the complexity to obtain an ϵ -KKT points for different
52 types of objective function $\psi(x)$. The constraint $g(x)$ is nonconvex and nonsmooth according to Assumption 2.1.

53 **Reviewer 4:** We will change the language in the abstract. Other concerns addressed in common response.