
Inexact Augmented Lagrangian Framework for Non-Convex Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We propose a practical inexact augmented Lagrangian method (iALM) for noncon-
2 vex problems with nonlinear constraints. We characterize the total computational
3 complexity of our method subject to a verifiable geometric condition, which is
4 closely related to the Polyak-Lojсиеwicz and Mangasarian-Fromowitz conditions.

5 In particular, when a first-order solver is used for the inner iterates, we prove that
6 iALM finds a first-order stationary point with $\tilde{O}(1/\epsilon^3)$ calls to the first-order oracle.
7 If, in addition, the problem is smooth and a second-order solver is used for the
8 inner iterates, iALM finds a second-order stationary point with $\tilde{O}(1/\epsilon^5)$ calls to
9 the second-order oracle. These complexity results match the known theoretical
10 results in the literature with a simple, implementable and versatile algorithm.

11 We provide numerical evidence on large-scale machine learning problems, in-
12 cluding the Burer-Monteiro factorization of semidefinite programs, and a novel
13 nonconvex relaxation of the standard basis pursuit template. We verify our geomet-
14 ric condition in all these examples.

15 1 Introduction

16 We study the following nonconvex optimization problem

$$\begin{cases} \min_{x \in \mathbb{R}^d} f(x) + g(x) \\ A(x) = 0, \end{cases} \quad (1)$$

17 where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a continuously-differentiable nonconvex function, and $A : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is
18 a nonlinear operator and $b \in \mathbb{R}^m$. We assume that $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is a proximal-friendly (possibly
19 nonsmooth) convex function.

20 A host of problems in computer science [? ? ?], machine learning [? ?], and signal processing [?
21 ?] naturally fall under the template (1), including max-cut, clustering, generalized eigenvalue
22 decomposition, as well as the quadratic assignment problem (QAP). [please add a reference for qap.](#)

23 To solve (1), this paper proposes an intuitive and easy-to-implement augmented Lagrangian algorithm,
24 and provides its total iteration complexity under an interpretable geometric condition. Before we
25 elaborate on the results, let us first motivate (1) with an application to semidefinite programming
26 (SDP):

27 **Vignette: Burer-Monteiro splitting.** A powerful convex relaxation for max-cut, clustering, and
28 several other problems described above is provided by the SDP

$$\begin{cases} \min_{X \in \mathbb{S}^{d \times d}} \langle C, X \rangle \\ B(X) = b, X \succeq 0, \end{cases} \quad (2)$$

29 where $C \in \mathbb{R}^{d \times d}$, X is a positive semidefinite $d \times d$ matrix, and $B : \mathbb{S}^{d \times d} \rightarrow \mathbb{R}^m$ is a linear operator.
 30 If the unique-games conjecture is true, SDPs achieve the best approximation for the underlying
 31 discrete problem [?].

32 Since d is often large, many first- and second-order methods for solving such SDP’s are immedi-
 33 ately ruled out, not only due to their high computational complexity, but also due to their storage
 34 requirements, which are $\mathcal{O}(d^2)$.

35 A contemporary challenge in optimization is therefore to solve SDPs using little space and in a
 36 scalable fashion. A recent algorithm, namely, homotopy conditional gradient method (HCGM)—
 37 based on Linear Minimization Oracles (LMO)—can solve (2) in a small space via sketching [?];
 38 however, such LMO-based methods are extremely slow in obtaining accurate solutions.

39 A key approach for solving (1), dating back to [? ?], is the so-called Burer-Monteiro (BM)
 40 factorization $X = UU^\top$, where $U \in \mathbb{R}^{d \times r}$ and r is selected according to the guidelines in [? ?
 41], which are shown to be optimal [?]. This factorization does not introduce any extraneous local
 42 minima [?]. Moreover, [?] established the connection between local minima of factorized
 43 problem 3 and global optimum for 2. [might have to double check this](#)

44 This factorization leads to the nonconvex problem

$$\begin{cases} \min_{U \in \mathbb{R}^{d \times r}} \langle C, UU^\top \rangle \\ B(UU^\top) = b, \end{cases} \quad (3)$$

45 which can be easily written in the form of (1). To solve (3), the inexact Augmented Lagrangian
 46 Method (iALM) is widely used [? ? ?], due to its cheap per iteration cost and its empirical success.
 47 Every (outer) iteration of iALM calls a solver to solve an intermediate augmented Lagrangian
 48 subproblem to near stationarity, and the user is free in the choice of this solver, which could use
 49 first-order, such as the proximal gradient descent [?], or second-order information, such as an
 50 BFGS [?].

51 Unlike its convex counterpart [? ? ?], the convergence rate and the complexity of iALM for (3)
 52 are not well-understood, see Section 5 for a review of the related literature. Indeed, addressing this
 53 important theoretical gap is one of the contributions of our work.

54 **Summary of contributions:**

- 55 ◦ Our framework is future-proof in the sense that we obtain the convergence rate of iALM for (1)
 56 with an arbitrary solver for finding first- and second-order stationary points.
- 57 ◦ We investigate the effect of using different solvers for augmented Lagrangian subproblems and
 58 provide overall iteration complexity bounds for finding first- and second-order stationary points of (1).
 59 Our complexity bounds match the best theoretical results in optimization, see Section 5.
- 60 ◦ We propose a geometric condition that simplifies the algorithmic analysis for iALM, and clarify its
 61 connection to well-known Polyak-Lojasiewicz and Mangasarian-Fromovitz conditions. [please add
 62 citations for these conditions here](#). We also verify this condition for key problems in Section 6. [how
 63 about the KL condition spars’19 reviewers pointed out?](#)

64 **Roadmap.** Section 2 collects the main tools and our notation. We present the iALM in Section 3
 65 and obtain its convergence rate to first- and second-order stationary points in Section 4, alongside
 66 their iteration complexities. We provide a comprehensive review of the literature and highlight our
 67 key differences in Section 5. Section 6 presents the numerical evidence and comparisons with the
 68 state-of-the-art techniques.

69 **2 Preliminaries**

70 **Notation.** We use the notation $\langle \cdot, \cdot \rangle$ and $\| \cdot \|$ for the standard inner product and the norm on \mathbb{R}^d . For
 71 matrices, $\| \cdot \|$ and $\| \cdot \|_F$ denote the spectral and the Frobenius norms, respectively. For the convex
 72 function $g : \mathbb{R}^d \rightarrow \mathbb{R}$, the subdifferential set at $x \in \mathbb{R}^d$ is denoted by $\partial g(x)$ and we will occasionally
 73 use the notation $\partial g(x)/\beta = \{z/\beta : z \in \partial g(x)\}$. When presenting iteration complexity results, we
 74 often use $\tilde{O}(\cdot)$ which suppresses the logarithmic dependencies.

75 We use the indicator function $\delta_{\mathcal{X}} : \mathbb{R}^d \rightarrow \mathbb{R}$ of a set $\mathcal{X} \subset \mathbb{R}^d$, which takes x to

$$\delta_{\mathcal{X}}(x) = \begin{cases} 0 & x \in \mathcal{X} \\ \infty & x \notin \mathcal{X}. \end{cases} \quad (4)$$

76 The distance function from a point x to \mathcal{X} is denoted by $\text{dist}(x, \mathcal{X}) = \min_{z \in \mathcal{X}} \|x - z\|$. For integers
77 $k_0 \leq k_1$, we denote $[k_0 : k_1] = \{k_0, \dots, k_1\}$.

78 For an operator $A : \mathbb{R}^d \rightarrow \mathbb{R}^m$ with components $\{A_i\}_{i=1}^m$, we let $DA(x) \in \mathbb{R}^{m \times d}$ denote the
79 Jacobian of A , where the i th row of $DA(x)$ is the gradient vector $\nabla A_i(x) \in \mathbb{R}^d$.

80 **Smoothness.** We require $f : \mathbb{R}^d \rightarrow \mathbb{R}$ and $A : \mathbb{R}^d \rightarrow \mathbb{R}^m$ to be smooth, namely, there exist
81 $\lambda_f, \lambda_A \geq 0$ such that

$$\begin{aligned} \|\nabla f(x) - \nabla f(x')\| &\leq \lambda_f \|x - x'\|, \\ \|DA(x) - DA(x')\| &\leq \lambda_A \|x - x'\|, \end{aligned} \quad (5)$$

82 for every $x, x' \in \mathbb{R}^d$.

83 **Augmented Lagrangian method (ALM).** ALM is a classical algorithm, which first appeared in [?
84 ?] and extensively studied afterwards in [? ?]. For solving (1), ALM suggests solving the problem

$$\min_x \max_y \mathcal{L}_\beta(x, y) + g(x), \quad (6)$$

85 where, for penalty weight $\beta > 0$, \mathcal{L}_β is the corresponding augmented Lagrangian, defined as

$$\mathcal{L}_\beta(x, y) := f(x) + \langle A(x), y \rangle + \frac{\beta}{2} \|A(x)\|^2. \quad (7)$$

86 The minimax formulation in (6) naturally suggests the following algorithm for solving (1). For dual
87 step sizes $\{\sigma_k\}_k$, consider the iterations

$$x_{k+1} \in \underset{x}{\text{argmin}} \mathcal{L}_\beta(x, y_k) + g(x), \quad (8)$$

88

$$y_{k+1} = y_k + \sigma_k A(x_{k+1}).$$

89 However, computing x_{k+1} above requires solving the nonconvex problem (8) to optimality, which is
90 typically intractable. Instead, it is often easier to find an approximate first- or second-order stationary
91 point of (8).

92 Hence, we argue that by gradually improving the stationarity precision and increasing the penalty
93 weight β above, we can reach a stationary point of the main problem in (1), as detailed in Section 3.

94 **Optimality conditions.** First-order necessary optimality conditions for (1) are well-studied. Indeed,
95 $x \in \mathbb{R}^d$ is a first-order stationary point of (1) if there exists $y \in \mathbb{R}^m$ such that

$$\begin{cases} -\nabla f(x) - DA(x)^\top y \in \partial g(x) \\ A(x) = 0, \end{cases} \quad (9)$$

96 where $DA(x)$ is the Jacobian of A at x . Recalling (7), we observe that (9) is equivalent to

$$\begin{cases} -\nabla_x \mathcal{L}_\beta(x, y) \in \partial g(x) \\ A(x) = 0, \end{cases} \quad (10)$$

97 which is in turn the necessary optimality condition for (6). **mfs: check approx. optimality conditions,**
98 **how they apply in this setting.** Inspired by this, we say that x is an (ϵ_f, β) first-order stationary point
99 of (6) if there exists a $y \in \mathbb{R}^m$ such that

$$\begin{cases} \text{dist}(-\nabla_x \mathcal{L}_\beta(x, y), \partial g(x)) \leq \epsilon_f \\ \|A(x)\| \leq \epsilon_f, \end{cases} \quad (11)$$

100 for $\epsilon_f \geq 0$. In light of (11), a suitable metric for evaluating the stationarity of a pair $(x, y) \in \mathbb{R}^d \times \mathbb{R}^m$
101 is

$$\text{dist}(-\nabla_x \mathcal{L}_\beta(x, y), \partial g(x)) + \|A(x)\|, \quad (12)$$

102 which we use as the first-order stopping criterion. As an example, for a convex set $\mathcal{X} \subset \mathbb{R}^d$, suppose
 103 that $g = \delta_{\mathcal{X}}$ is the indicator function on \mathcal{X} . Let also $T_{\mathcal{X}}(x) \subseteq \mathbb{R}^d$ denote the tangent cone to \mathcal{X} at x ,
 104 and with $P_{T_{\mathcal{X}}(x)} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we denote the orthogonal projection onto this tangent cone. Then, for
 105 $u \in \mathbb{R}^d$, it is not difficult to verify that

$$\text{dist}(u, \partial g(x)) = \|P_{T_{\mathcal{X}}(x)}(u)\|. \quad (13)$$

106 When $g = 0$, a first-order stationary point $x \in \mathbb{R}^d$ of (1) is also second-order stationary if

$$\lambda_{\min}(\nabla_{xx}\mathcal{L}_{\beta}(x, y)) \geq 0, \quad (14)$$

107 where $\nabla_{xx}\mathcal{L}_{\beta}$ is the Hessian of \mathcal{L}_{β} with respect to x , and $\lambda_{\min}(\cdot)$ returns the smallest eigenvalue of
 108 its argument. Analogously, x is an $(\epsilon_f, \epsilon_s, \beta)$ second-order stationary point if, in addition to (11), it
 109 holds that

$$\lambda_{\min}(\nabla_{xx}\mathcal{L}_{\beta}(x, y)) \geq -\epsilon_s, \quad (15)$$

110 for $\epsilon_s \geq 0$. Naturally, for second-order stationarity, we use $\lambda_{\min}(\nabla_{xx}\mathcal{L}_{\beta}(x, y))$ as the stopping
 111 criterion.

112 **Smoothness lemma.** This next result controls the smoothness of $\mathcal{L}_{\beta}(\cdot, y)$ for a fixed y . The proof
 113 is standard but nevertheless is included in Appendix C for completeness.

114 **Lemma 2.1 (smoothness)** *For fixed $y \in \mathbb{R}^m$ and $\rho, \rho' \geq 0$, it holds that*

$$\|\nabla_x \mathcal{L}_{\beta}(x, y) - \nabla_x \mathcal{L}_{\beta}(x', y)\| \leq \lambda_{\beta} \|x - x'\|, \quad (16)$$

115 *for every $x, x' \in \{x'' : \|x''\| \leq \rho, \|A(x'')\| \leq \rho'\}$, where*

$$\begin{aligned} \lambda_{\beta} &\leq \lambda_f + \sqrt{m}\lambda_A \|y\| + (\sqrt{m}\lambda_A \rho' + d\lambda_A^2)\beta \\ &=: \lambda_f + \sqrt{m}\lambda_A \|y\| + \lambda''(A, \rho, \rho')\beta. \end{aligned} \quad (17)$$

116 *Above, λ_f, λ_A were defined in (5) and*

$$\lambda'_A := \max_{\|x\| \leq \rho} \|DA(x)\|. \quad (18)$$

117 3 Algorithm

118 To solve the equivalent formulation of (1) presented in (6), we propose the inexact ALM (iALM),
 119 detailed in Algorithm 1.

120 At the k^{th} iteration, Step 2 of Algorithm 1 calls a solver that finds an approximate stationary point
 121 of the augmented Lagrangian $\mathcal{L}_{\beta_k}(\cdot, y_k)$ with the accuracy of ϵ_{k+1} , and this accuracy gradually
 122 increases in a controlled fashion.

123 The increasing sequence of penalty weights $\{\beta_k\}_k$ and the dual update (Steps 4 and 5) are responsible
 124 for continuously enforcing the constraints in (1). The appropriate choice for $\{\beta_k\}_k$ will be specified
 125 in Sections 4.1 and 4.2.

126 As we will see in the convergence analysis, the particular choice of the dual step sizes $\{\sigma_k\}_k$ in
 127 Algorithm 1 ensures that the dual variable y_k remains bounded, see [?] for a precedent in the ALM
 128 literature where a similar choice for σ_k is considered.

129 4 Convergence Rate

130 In this section, we detail the convergence rate of Algorithm 1 for finding first-order and second-
 131 order stationary points, along with total iteration complexity results. All the proofs are deferred
 132 to Appendix A. Theorem 4.1 below characterizes the convergence rate of Algorithm 1 for finding
 133 stationary points in terms of the number of outer iterations.

134

Algorithm 1 Inexact ALM for solving (1)

Input: Non-decreasing, positive, unbounded sequence $\{\beta_k\}_{k \geq 1}$, stopping thresholds $\tau_f > 0$ and $\tau_s > 0$.

Initialization: Initial primal variable $x_1 \in \mathbb{R}^d$, initial dual variable $y_0 \in \mathbb{R}^m$, initial dual step size $\sigma_1 > 0$.

for $k = 1, 2, \dots$ **do**

1. **(Update tolerance)** $\epsilon_{k+1} = 1/\beta_k$.
2. **(Inexact primal solution)** Obtain $x_{k+1} \in \mathbb{R}^d$ such that

$$\text{dist}(-\nabla_x \mathcal{L}_{\beta_k}(x_{k+1}, y_k), \partial g(x_{k+1})) \leq \epsilon_{k+1}$$

for first-order stationarity

$$\lambda_{\min}(\nabla_{xx} \mathcal{L}_{\beta_k}(x_{k+1}, y_k)) \geq -\epsilon_{k+1}$$

for second-order-stationarity, if $g = 0$ in (1).

3. **(Update dual step size)**

$$\sigma_{k+1} = \sigma_1 \min \left(\frac{\|A(x_1)\| \log^2 2}{\|A(x_{k+1})\| (k+1) \log^2(k+2)}, 1 \right).$$

4. **(Dual ascent)** $y_{k+1} = y_k + \sigma_{k+1} A(x_{k+1})$.
5. **(Stopping criterion)** If

$$\text{dist}(-\nabla_x \mathcal{L}_{\beta_k}(x_{k+1}, y_{k+1}), \partial g(x_{k+1})) + \|A(x_{k+1})\| \leq \tau_f,$$

for first-order stationarity and if also $\lambda_{\min}(\nabla_{xx} \mathcal{L}_{\beta_k}(x_{k+1}, y_{k+1})) \geq -\tau_s$ for second-order stationarity, then quit and return x_{k+1} as an (approximate) stationary point of (1).

end for

135 **Theorem 4.1 (convergence rate)** Let $\rho := \sup_{k \in [K]} \|x\|$.¹ Suppose that f and A satisfy (5) and let

$$\lambda'_f = \max_{\|x\| \leq \rho} \|\nabla f(x)\|, \quad \lambda'_A = \max_{\|x\| \leq \rho} \|DA(x)\|, \quad (19)$$

136 be the (restricted) Lipschitz constants of f and A , respectively. For integers $2 \leq k_0 \leq k_1$, consider
137 the interval $K = [k_0 : k_1]$, and let $\{x_k\}_{k \in K}$ be the output sequence of Algorithm 1 on the interval
138 K .² For $\nu > 0$, assume that

$$\nu \|A(x_k)\| \leq \text{dist} \left(-DA(x_k)^\top A(x_k), \frac{\partial g(x_k)}{\beta_{k-1}} \right), \quad (20)$$

139 for every $k \in K$. We consider two cases:

- 140 • If a first-order solver is used in Step 2, then x_k is an $(\epsilon_{k,f}, \beta_k)$ first-order stationary point of (1)
141 with

$$\begin{aligned} \epsilon_{k,f} &= \frac{1}{\beta_{k-1}} \left(\frac{2(\lambda'_f + \lambda'_A y_{\max})(1 + \lambda'_A \sigma_k)}{\nu} + 1 \right) \\ &=: \frac{Q(f, g, A, \sigma_1)}{\beta_{k-1}}, \end{aligned} \quad (21)$$

142 for every $k \in K$, where the expression for $y_{\max}(x_1, y_0, \sigma_1)$ is given in (41) due to the limited
143 space.

¹If necessary, to ensure that $\rho < \infty$, one can add a small factor of $\|x\|^2$ to \mathcal{L}_β in (7). Then it is easy to verify that the iterates of Algorithm 1 remain bounded, provided that the penalty weight β is large enough, $\sup_x \|\nabla f(x)\|/\|x\| < \infty$ and $\sup_x \|A(x)\| < \infty$.

²The choice of $k_1 = \infty$ is valid here too.

- 144 • If a second-order solver is used in Step 2, then x_k is an $(\epsilon_{k,f}, \epsilon_{k,s}, \beta_k)$ second-order stationary
 145 point of (1) with $\epsilon_{k,s}$ specified above and with

$$\begin{aligned}\epsilon_{k,s} &= \epsilon_{k-1} + \sigma_k \sqrt{m} \lambda_A \frac{2\lambda'_f + 2\lambda'_A y_{\max}}{\nu \beta_{k-1}} \\ &= \frac{\nu + \sigma_k \sqrt{m} \lambda_A 2\lambda'_f + 2\lambda'_A y_{\max}}{\nu \beta_{k-1}} =: \frac{Q'(f, g, A, \sigma_1)}{\beta_{k-1}}.\end{aligned}\quad (22)$$

146 Loosely speaking, Theorem 4.1 states that Algorithm 1 converges to a (first- or second-) order
 147 stationary point of (1) at the rate of $1/\beta_k$. A few remarks are in order.

148 **Regularity.** The key geometric condition in Theorem 4.1 is (20) which, broadly speaking, ensures
 149 that the primal updates of Algorithm 1 reduce the feasibility gap as the penalty weight β_k grows. We
 150 will verify this condition for several examples in Section 6.

151 This condition is closely related to those in the existing literature. In the special case where $g = 0$
 152 in (1), it is easy to verify that (20) reduces to the Polyak-Lojasiewicz (PL) condition for minimizing
 153 $\|A(x)\|^2$ [?]. PL condition itself is a special case of Kurdyka-Lojasiewicz with $\theta = 1/2$, see [? ,
 154 Definition 1.1]. When $g = 0$, it is also easy to see that (20) is weaker than the Mangasarian-Fromovitz
 155 condition in nonlinear optimization [?].

156 By its definition, we may think of (20) as a local condition, which should hold within a neighborhood
 157 of the constraint set $\{x : A(x) = 0\}$ rather than everywhere in \mathbb{R}^d . There is a constant complexity
 158 algorithm in [?] to reach this so-called ‘‘information zone’’, which supplements Theorem 4.1.

159 Moreover, in contrast to most conditions in the nonconvex optimization literature, such as [?], the
 160 condition in (20) appears to be easier to verify, as we see in Section 6.

161 **AE: Fatih, I think you had added the following two references to our response for icml. Could**
 162 **you discuss their relevance here? [2] Rockafellar, Lagrange Multipliers and Optimality, 1993**
 163 **[3] Bertsekas, On penalty and multiplier methods for constrained minimization. 1996**

164 **AE: the spars review talks about the "Pong-Li" work. Fatih, do you know what is that?**

165 **Penalty method.** A classical algorithm to solve (1) is the penalty method, which is characterized by
 166 the absence of the dual variable ($y = 0$) in (7). Indeed, ALM can be interpreted as an adaptive penalty
 167 or smoothing method with a variable center as determined by the dual variable. It is worth noting
 168 that, with same proof technique, one can establish the same convergence rate of Theorem 4.1 for
 169 the penalty method. However, while both methods have the same convergence rate in theory, iALM
 170 outperforms the penalty method in practice by virtue of its variable center and has been excluded
 171 from this presentation.

172 **Computational complexity.** Theorem 4.1 allows us to specify the number of (outer) iterations that
 173 Algorithm 1 requires to reach a near-stationary point of problem (1) with a prescribed precision and,
 174 in particular, specifies the number of calls made to the solver in Step 2. In this sense, Theorem 4.1
 175 does not fully capture the computational complexity of Algorithm 1, as it does not take into account
 176 the computational cost of the solver in Step 2.

177 To better understand the total iteration complexity of Algorithm 1, we consider two scenarios in the
 178 following. In the first scenario, we take the solver in Step 2 to be the Accelerated Proximal Gradient
 179 Method (APGM), a well-known first-order algorithm [?]. In the second scenario, we will use the
 180 second-order trust region method developed in [?].

181 4.1 First-Order Optimality

182 Let us first consider the case where the solver in Step 2 is the first-order algorithm APGM, described
 183 in detail in [?]. At a high level, APGM makes use of $\nabla_x \mathcal{L}_\beta(x, y)$, the proximal operator prox_g and
 184 classical Nesterov acceleration for the iterates [?] to reach first-order stationarity for the first update
 185 in (8). Suppose that $g = \delta_{\mathcal{X}}$ is the indicator function on a bounded convex set $\mathcal{X} \subset \mathbb{R}^d$ and let

$$\rho' = \max_{x \in \mathcal{X}} \|x\|, \quad (23)$$

186 be the radius of a ball centered at the origin that includes \mathcal{X} . Then, adapting the results in [?] to our
 187 setup, APGM reaches x_k in Step 2 of Algorithm 1 after

$$\mathcal{O}\left(\frac{\lambda_{\beta_k}^2 \rho'^2}{\epsilon_{k+1}}\right) \quad (24)$$

188 (inner) iterations, where λ_{β_k} denotes the Lipschitz constant of $\nabla_x \mathcal{L}_{\beta_k}(x, y)$, bounded in (17). For
 189 the clarity of the presentation, we have used a looser bound in (24) compared to [?]. Using (24), we
 190 derive the following corollary, describing the total iteration complexity of Algorithm 1 in terms of the
 191 number calls made to the first-order oracle in APGM.

192 **Corollary 4.2** For $b > 1$, let $\beta_k = b^k$ for every k . If we use APGM from [?] for Step 2 of
 193 Algorithm 1, the algorithm finds an (ϵ_f, β_k) first-order stationary point, after T calls to the first-order
 194 oracle, where

$$T = \mathcal{O}\left(\frac{Q^3 \rho'^2}{\epsilon^3} \log_b\left(\frac{Q}{\epsilon}\right)\right) = \tilde{\mathcal{O}}\left(\frac{Q^3 \rho'^2}{\epsilon^3}\right). \quad (25)$$

195 For Algorithm 1 to reach a near-stationary point with an accuracy of ϵ_f in the sense of (11) and
 196 with the lowest computational cost, we therefore need to perform only one iteration of Algorithm 1,
 197 with β_1 specified as a function of ϵ_f by (21) in Theorem 4.1. In general, however, the constants in
 198 (21) are unknown and this approach is thus not feasible. Instead, the homotopy approach taken by
 199 Algorithm 1 ensures achieving the desired accuracy by gradually increasing the penalty weight.³ This
 200 homotopy approach increases the computational cost of Algorithm 1 only by a factor logarithmic in
 201 the ϵ_f , as detailed in the proof of Corollary 4.2.

202 4.2 Second-Order Optimality

203 Let us now consider the second-order optimality case where the solver in Step 2 is the the trust region
 204 method developed in [?]. Trust region method minimizes quadratic approximation of the function
 205 within a dynamically updated trust-region radius. Second-order trust region method that we consider
 206 in this section makes use of Hessian (or an approximation of Hessian) of the augmented Lagrangian
 207 in addition to first order oracles.

208 As shown in [?], finding approximate second-order stationary points of convex-constrained problems
 209 is in general NP-hard. For this reason, we focus in this section on the special case of (1) with $g = 0$.

210 Let us compute the total computational complexity of Algorithm 1 with the trust region method in
 211 Step 2, in terms of the number of calls made to the second-order oracle. By adapting the result in [?
 212] to our setup, we find that the number of (inner) iterations required in Step 2 of Algorithm 1 to
 213 produce x_{k+1} is

$$\mathcal{O}\left(\frac{\lambda_{\beta_k, H}^2 (\mathcal{L}_{\beta_k}(x_1, y) - \min_x \mathcal{L}_{\beta_k}(x, y))}{\epsilon_k^3}\right), \quad (26)$$

214 where $\lambda_{\beta, H}$ is the Lipschitz constant of the Hessian of the augmented Lagrangian, which is of the
 215 order of β , as can be proven similar to Lemma 2.1 and x_1 is the initial iterate of the given outer
 216 loop. In [?], the term $\mathcal{L}_{\beta}(x_1, y) - \min_x \mathcal{L}_{\beta}(x, y)$ is bounded by a constant independent of ϵ . We
 217 assume a uniform bound for this quantity $\forall \beta_k$, instead of for one value of β_k as in [?]. Using (26)
 218 and Theorem 4.1, we arrive at the following:

219 **Corollary 4.3** For $b > 1$, let $\beta_k = b^k$ for every k . We assume that

$$\mathcal{L}_{\beta}(x_1, y) - \min_x \mathcal{L}_{\beta}(x, y) \leq L_u, \quad \forall \beta. \quad (27)$$

220 If we use the trust region method from [?] for Step 2 of Algorithm 1, the algorithm finds an
 221 ϵ -second-order stationary point of (1) in T calls to the second-order oracle where

$$T \leq \mathcal{O}\left(\frac{L_u Q'^5}{\epsilon^5} \log_b\left(\frac{Q'}{\epsilon}\right)\right) = \tilde{\mathcal{O}}\left(\frac{L_u Q'^5}{\epsilon^5}\right). \quad (28)$$

222 Before closing this section, we note that the remark after Corollary 4.2 applies here as well.

³In this context, homotopy loosely corresponds to the gradual enforcement of the constraints by increasing the penalty weight.

223 **5 Related Work**

224 ALM has a long history in the optimization literature, dating back to [? ?]. In the special case of (1)
 225 with a convex function f and a linear operator A , standard, inexact and linearized versions of ALM
 226 have been extensively studied [? ? ? ?].

227 Classical works on ALM focused on the general template of (1) with nonconvex f and nonlinear A ,
 228 with arguably stronger assumptions and required exact solutions to the subproblems of the form (8),
 229 which appear in Step 2 of Algorithm 1, see for instance [?].

230 A similar analysis was conducted in [?] for the general template of (1). The authors considered
 231 inexact ALM and proved convergence rates for the outer iterates, under specific assumptions on
 232 the initialization of the dual variable. However, unlike our results, the authors did not analyze how
 233 to solve the subproblems inexactly and they did not provide total complexity results and verifiable
 234 conditions.

235 Problem (1) with similar assumptions to us is also studied in [?] and [?] for first-order and second-
 236 order stationarity, respectively, with explicit iteration complexity analysis. As we have mentioned
 237 in Section 4, our iteration complexity results matches these theoretical algorithms with a simpler
 238 algorithm and a simpler analysis. In addition, these algorithms require setting final accuracies since
 239 they utilize this information in the algorithm. In contrast to [? ?], Algorithm 1 does not set accuracies
 240 a priori.

241 [?] also considers the same template (1) for first-order stationarity with a penalty-type method instead
 242 of ALM. Even though the authors show $\mathcal{O}(1/\epsilon^2)$ complexity, this result is obtained by assuming that
 243 the penalty parameter remains bounded. We note that such an assumption can also be used to match
 244 our complexity results.

245 [?] studies the general template (1) with specific assumptions involving local error bound conditions
 246 for the (1). These conditions are studied in detail in [?], but their validity for general SDPs (2) has
 247 never been established. This work also lacks the total iteration complexity analysis presented here.

248 Another work [?] focused on solving (1) by adapting the primal-dual method of Chambolle and
 249 Pock [?]. The authors proved the convergence of the method and provided convergence rate by
 250 imposing error bound conditions on the objective function that do not hold for standard SDPs.

251 [? ?] is the first work that proposes the splitting $X = UU^\top$ for solving SDPs of the form (2).
 252 Following these works, the literature on Burer-Monteiro (BM) splitting for the large part focused on
 253 using ALM for solving the reformulated problem (3).

254 However, this approach has a few drawbacks: First, it requires exact solutions in Step 2 of Algo-
 255 rithm 1 in theory, which in practice is replaced with inexact solutions. Second, their results only
 256 establish convergence without providing the rates. In this sense, our work provides a theoretical
 257 understanding of the BM splitting with inexact solutions to Step 2 of Algorithm 1 and complete
 258 iteration complexities.

259 [? ?] are among the earliest efforts to show convergence rates for BM splitting, focusing on
 260 the special case of SDPs without any linear constraints. For these specific problems, they prove
 261 the convergence of gradient descent to global optima with convergence rates, assuming favorable
 262 initialization. These results, however, do not apply to general SDPs of the form (2) where the difficulty
 263 arises due to the linear constraints.

264 [?] focused on the quadratic penalty formulation of (1), namely,

$$\min_{X \succeq 0} \langle C, X \rangle + \frac{\mu}{2} \|B(x) - b\|^2, \quad (29)$$

265 which after BM splitting becomes

$$\min_{U \in \mathbb{R}^{d \times r}} \langle C, UU^\top \rangle + \frac{\mu}{2} \|B(UU^\top) - b\|^2, \quad (30)$$

266 for which they study the optimality of the second-order stationary points. These results are for
 267 establishing a connection between the stationary points of (30) and global optima of (29). In contrast,
 268 we focus on the relation of the stationary points of (6) to the constrained problem (1).

269 Another popular method for solving SDPs are due to [? ? ?], focusing on the case where the
 270 constraints in (1) can be written as a Riemannian manifold after BM splitting. In this case, the authors

271 apply the Riemannian gradient descent and Riemannian trust region methods for obtaining first- and
 272 second-order stationary points, respectively. They obtain $\mathcal{O}(1/\epsilon^2)$ complexity for finding first-order
 273 stationary points and $\mathcal{O}(1/\epsilon^3)$ complexity for finding second-order stationary points.

274 While these complexities appear better than ours, the smooth manifold requirement in these works
 275 is indeed restrictive. In particular, this requirement holds for max-cut and generalized eigenvalue
 276 problems, but it is not satisfied for other important SDPs such as quadratic programming (QAP),
 277 optimal power flow and clustering with general affine constraints. In addition, as noted in [?], per
 278 iteration cost of their method for max-cut problem is an astronomical $\mathcal{O}(d^6)$.

279 Lastly, there also exists a line of work for solving SDPs in their original convex formulation, in a
 280 storage efficient way [? ? ?]. These works have global optimality guarantees by their virtue of
 281 directly solving the convex formulation. On the downside, these works require the use of eigenvalue
 282 routines and exhibit significantly slower convergence as compared to nonconvex approaches [?].

283 6 Numerical Evidence

284 We first begin with a caveat: It is known that quasi-Newton methods, such as BFGS and IBFGS, might
 285 not converge for nonconvex problems [? ?]. For this reason, we have used the trust region method as
 286 the second-order solver in our analysis in Section 4, which is well-studied for nonconvex problems [?
 287]. Empirically, however, BFGS and IBFGS are extremely successful and we have therefore opted for
 288 those solvers in this section since the subroutine does not affect Theorem 4.1 as long as the subsolver
 289 performs well in practice.

290 6.1 Clustering

291 Given data points $\{z_i\}_{i=1}^n$, the entries of the corresponding Euclidean distance matrix $D \in \mathbb{R}^{n \times n}$
 292 are $D_{i,j} = \|z_i - z_j\|^2$. Clustering is then the problem of finding a co-association matrix $Y \in \mathbb{R}^{n \times n}$
 293 such that $Y_{ij} = 1$ if points z_i and z_j are within the same cluster and $Y_{ij} = 0$ otherwise. In [?], the
 294 authors provide a SDP relaxation of the clustering problem, specified as

$$\begin{cases} \min_{Y \in \mathbb{R}^{n \times n}} \text{tr}(DY) \\ Y\mathbf{1} = \mathbf{1}, \text{tr}(Y) = s, Y \succeq 0, Y \geq 0, \end{cases} \quad (31)$$

295 where s is the number of clusters and Y is both positive semidefinite and has nonnegative entries.
 296 Standard SDP solvers do not scale well with the number of data points n , since they often require
 297 projection onto the semidefinite cone with the complexity of $\mathcal{O}(n^3)$. We instead use the Burer-
 298 Monteiro factorization to solve (31), sacrificing convexity to reduce the computational complexity.
 299 More specifically, we solve the program

$$\begin{cases} \min_{V \in \mathbb{R}^{n \times r}} \text{tr}(DVV^\top) \\ VV^\top \mathbf{1} = \mathbf{1}, \|V\|_F^2 \leq s, V \geq 0, \end{cases} \quad (32)$$

300 where $\mathbf{1} \in \mathbb{R}^n$ is the vector of all ones. Note that $Y \geq 0$ in (31) is replaced above by the much
 301 stronger but easier-to-enforce constraint $V \geq 0$ in (32), see [?] for the reasoning behind this
 302 relaxation. Now, we can cast (32) as an instance of (1). Indeed, for every $i \leq n$, let $x_i \in \mathbb{R}^r$ denote
 303 the i th row of V . We next form $x \in \mathbb{R}^d$ with $d = nr$ by expanding the factorized variable V , namely,

$$x = [x_1^\top, \dots, x_n^\top]^\top \in \mathbb{R}^d,$$

304 and then set

$$f(x) = \sum_{i,j=1}^n D_{i,j} \langle x_i, x_j \rangle, \quad g = \delta_C,$$

305

$$A(x) = [x_1^\top \sum_{j=1}^n x_j - 1, \dots, x_n^\top \sum_{j=1}^n x_j - 1]^\top, \quad (33)$$

306 where C is the intersection of the positive orthant in \mathbb{R}^d with the Euclidean ball of radius \sqrt{s} . In
 307 Appendix D, we somewhat informally verify that Theorem 4.1 applies to (1) with f, g, A specified
 308 above.

309 In our simulations, we use two different solvers for Step 2 of Algorithm 1, namely, APGM and
 310 IBFGS. APGM is a solver for nonconvex problems of the form (8) with convergence guarantees
 311 to first-order stationarity, as discussed in Section 4. IBFGS is a limited-memory version of BFGS
 312 algorithm in [?] that approximately leverages the second-order information of the problem. We
 313 compare our approach against the following convex methods:

- 314 • HCGM: Homotopy-based Conditional Gradient Method in [?] which directly solves (31).
- 315 • SDPNAL+: A second-order augmented Lagrangian method for solving SDP's with nonneg-
- 316 ativity constraints [?].

317 As for the dataset, our experimental setup is similar to that described by [?]. We use the publicly-
 318 available fashion-MNIST data in [?], which is released as a possible replacement for the MNIST
 319 handwritten digits. Each data point is a 28×28 gray-scale image, associated with a label from ten
 320 classes, labeled from 0 to 9. First, we extract the meaningful features from this dataset using a simple
 321 two-layer neural network with a sigmoid activation function. Then, we apply this neural network to
 322 1000 test samples from the same dataset, which gives us a vector of length 10 for each data point,
 323 where each entry represents the posterior probability for each class. Then, we form the ℓ_2 distance
 324 matrix D from these probability vectors. The results are depicted in Figure 1. We implemented 3
 325 algorithms on MATLAB and used the software package for SDPNAL+ which contains mex files. It is
 326 predictable that the performance of our nonconvex approach would also improve by using mex files.

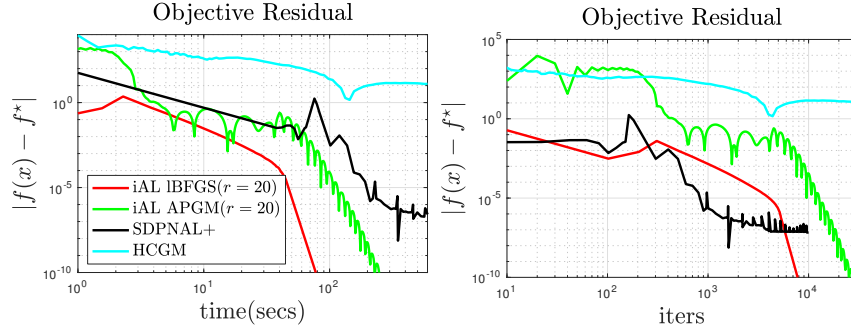


Figure 1: Convergence of different algorithms for clustering with fashion-MNIST dataset. Here, we set the rank as $r = 20$ for the nonconvex approaches. The solution rank for the template (31) is the number of clusters s [?, Theorem 1]. However, as discussed in [?], setting rank $r > s$ leads more accurate reconstruction at the expense of speed, hence our choice of $r = 20$.

327 6.2 Basis Pursuit

328 Basis Pursuit (BP) finds sparsest solutions of an under-determined system of linear equations, namely,

$$\begin{cases} \min_z \|z\|_1 \\ Bz = b, \end{cases} \quad (34)$$

329 where $B \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$. BP has found many applications in machine learning, statistics and
 330 signal processing [? ? ?]. A plethora of primal-dual convex optimization algorithms are available in
 331 the literature to solve BP, including [? ?]. There also exists a line of work [?] that handles sparse
 332 regression via regularization with ℓ_1 norm.

333 Here, we take a different approach and cast (34) as an instance of (1). Note that any $z \in \mathbb{R}^d$
 334 can be decomposed as $z = z^+ - z^-$, where $z^+, z^- \in \mathbb{R}^d$ are the positive and negative parts of
 335 z , respectively. Then consider the change of variables $z^+ = u_1^{\circ 2}$ and $z^- = u_2^{\circ 2} \in \mathbb{R}^d$, where \circ
 336 denotes element-wise power. Next, we concatenate u_1 and u_2 as $x := [u_1^\top, u_2^\top]^\top \in \mathbb{R}^{2d}$ and define
 337 $\bar{B} := [B, -B] \in \mathbb{R}^{n \times 2d}$. Then, (34) is equivalent to (1) with

$$\begin{aligned} f(x) &= \|x\|^2, & g(x) &= 0 \\ A(x) &= \bar{B}x^{\circ 2} - b. \end{aligned} \quad (35)$$

338 In Appendix E, we verify with minimal detail that Theorem 4.1 indeed applies to (1) with the above
 339 f, A .

340 We draw the entries of B independently from a zero-mean and unit-variance Gaussian distribution.
 341 For a fixed sparsity level k , the support of $z_* \in \mathbb{R}^d$ and its nonzero amplitudes are also drawn from
 342 the standard Gaussian distribution. Then the measurement vector is created as $b = Bz + \epsilon$, where ϵ
 343 is the noise vector with entries drawn independently from the zero-mean Gaussian distribution with
 344 variance $\sigma^2 = 10^{-6}$.

345 **AE: the image sizes throughout the paper are inconsistent which is not nice. The font size in
 Fig 3 is also very different from the rest of the paper which is not nice. please change.**

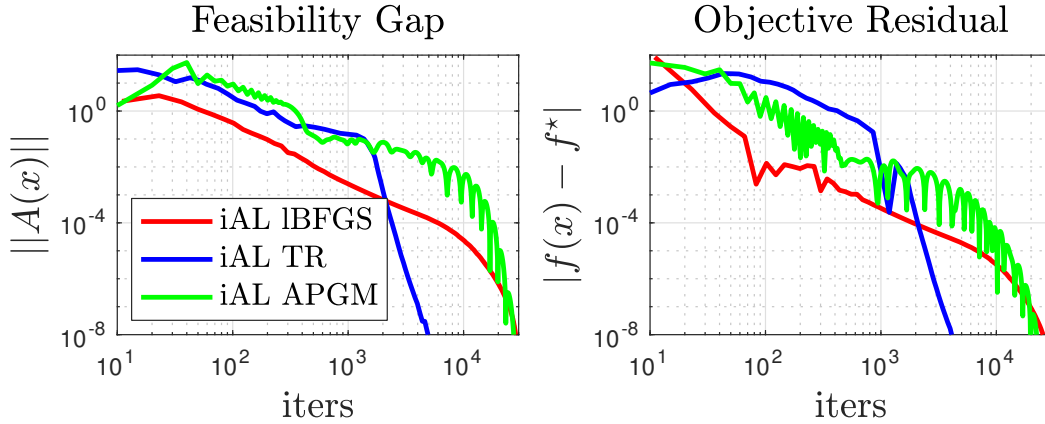


Figure 2: Convergence with different subsolvers for the aforementioned nonconvex relaxation.

346 Figure 2 compiles our results for the proposed relaxation. It is, indeed, interesting to see that these
 347 type of nonconvex relaxations gives the solution of convex one and first order methods succeed.
 348

349 **Discussion:** The true potential of our reformulation is in dealing with more structured norms rather
 350 than ℓ_1 , where computing the proximal operator is often intractable. One such case is the latent group
 351 lasso norm [?], defined as

$$\|z\|_{\Omega} = \sum_{i=1}^I \|z_{\Omega_i}\|,$$

352 where $\{\Omega_i\}_{i=1}^I$ are (not necessarily disjoint) index sets of $\{1, \dots, d\}$. Although not studied here, we
 353 believe that the nonconvex framework presented in this paper can serve to solve more complicated
 354 problems, such as the latent group lasso. We leave this research direction for future work.

355 References

356 **A Proof of Theorem 4.1**

357 For every $k \geq 2$, recall from (7) and Step 2 of Algorithm 1 that x_k satisfies

$$\begin{aligned} & \text{dist}(-\nabla f(x_k) - DA(x_k)^\top y_{k-1} \\ & \quad - \beta_{k-1} DA(x_k)^\top A(x_k), \partial g(x_k)) \\ & = \text{dist}(-\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_{k-1}), \partial g(x_k)) \leq \epsilon_k. \end{aligned} \quad (36)$$

358 With an application of the triangle inequality, it follows that

$$\begin{aligned} & \text{dist}(-\beta_{k-1} DA(x_k)^\top A(x_k), \partial g(x_k)) \\ & \leq \|\nabla f(x_k)\| + \|DA(x_k)^\top y_{k-1}\| + \epsilon_k, \end{aligned} \quad (37)$$

359 which in turn implies that

$$\begin{aligned} & \text{dist}(-DA(x_k)^\top A(x_k), \partial g(x_k)/\beta_{k-1}) \\ & \leq \frac{\|\nabla f(x_k)\|}{\beta_{k-1}} + \frac{\|DA(x_k)^\top y_{k-1}\|}{\beta_{k-1}} + \frac{\epsilon_k}{\beta_{k-1}} \\ & \leq \frac{\lambda'_f + \lambda'_A \|y_{k-1}\| + \epsilon_k}{\beta_{k-1}}, \end{aligned} \quad (38)$$

360 where λ'_f, λ'_A were defined in (19). We next translate (38) into a bound on the feasibility gap $\|A(x_k)\|$.

361 Using the regularity condition (20), the left-hand side of (38) can be bounded below as

$$\text{dist}(-DA(x_k)^\top A(x_k), \partial g(x_k)/\beta_{k-1}) \geq \nu \|A(x_k)\|. \quad (\text{see (20)}) \quad (39)$$

362 By substituting (39) back into (38), we find that

$$\|A(x_k)\| \leq \frac{\lambda'_f + \lambda'_A \|y_{k-1}\| + \epsilon_k}{\nu \beta_{k-1}}. \quad (40)$$

363 In words, the feasibility gap is directly controlled by the dual sequence $\{y_k\}_k$. We next establish that
364 the dual sequence is bounded. Indeed, for every $k \in K$, note that

$$\begin{aligned} \|y_k\| & = \|y_0 + \sum_{i=1}^k \sigma_i A(x_i)\| \quad (\text{Step 5 of Algorithm 1}) \\ & \leq \|y_0\| + \sum_{i=1}^k \sigma_i \|A(x_i)\| \quad (\text{triangle inequality}) \\ & \leq \|y_0\| + \sum_{i=1}^k \frac{\|A(x_1)\| \log^2 2}{k \log^2(k+1)} \quad (\text{Step 4}) \\ & \leq \|y_0\| + c \|A(x_1)\| \log^2 2 =: y_{\max}, \end{aligned} \quad (41)$$

365 where

$$c \geq \sum_{i=1}^{\infty} \frac{1}{k \log^2(k+1)}. \quad (42)$$

366 Substituting (41) back into (40), we reach

$$\begin{aligned} \|A(x_k)\| & \leq \frac{\lambda'_f + \lambda'_A y_{\max} + \epsilon_k}{\nu \beta_{k-1}} \\ & \leq \frac{2\lambda'_f + 2\lambda'_A y_{\max}}{\nu \beta_{k-1}}, \end{aligned} \quad (43)$$

367 where the second line above holds if k_0 is large enough, which would in turn guarantees that
368 $\epsilon_k = 1/\beta_{k-1}$ is sufficiently small since $\{\beta_k\}_k$ is increasing and unbounded. It remains to control

369 the first term in (12). To that end, after recalling Step 2 of Algorithm 1 and applying the triangle
 370 inequality, we can write that

$$\begin{aligned} & \text{dist}(-\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_k), \partial g(x_k)) \\ & \leq \text{dist}(-\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_{k-1}), \partial g(x_k)) \\ & \quad + \|\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_k) - \nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_{k-1})\|. \end{aligned} \quad (44)$$

371 The first term on the right-hand side above is bounded by ϵ_k , by Step 5 of Algorithm 1. For the
 372 second term on the right-hand side of (44), we write that

$$\begin{aligned} & \|\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_k) - \nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_{k-1})\| \\ & = \|DA(x_k)^\top (y_k - y_{k-1})\| \quad (\text{see (7)}) \\ & \leq \lambda'_A \|y_k - y_{k-1}\| \quad (\text{see (19)}) \\ & = \lambda'_A \sigma_k \|A(x_k)\| \quad (\text{see Step 5 of Algorithm 1}) \\ & \leq \frac{2\lambda'_A \sigma_k}{\nu \beta_{k-1}} (\lambda'_f + \lambda'_A y_{\max}). \quad (\text{see (43)}) \end{aligned} \quad (45)$$

373 By combining (44,45), we find that

$$\begin{aligned} & \text{dist}(\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_k), \partial g(x_k)) \\ & \leq \frac{2\lambda'_A \sigma_k}{\nu \beta_{k-1}} (\lambda'_f + \lambda'_A y_{\max}) + \epsilon_k. \end{aligned} \quad (46)$$

374 By combining (43,46), we find that

$$\begin{aligned} & \text{dist}(-\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_k), \partial g(x_k)) + \|A(x_k)\| \\ & \leq \left(\frac{2\lambda'_A \sigma_k}{\nu \beta_{k-1}} (\lambda'_f + \lambda'_A y_{\max}) + \epsilon_k \right) \\ & \quad + 2 \left(\frac{\lambda'_f + \lambda'_A y_{\max}}{\nu \beta_{k-1}} \right). \end{aligned} \quad (47)$$

375 Applying $\sigma_k \leq \sigma_1$, we find that

$$\begin{aligned} & \text{dist}(-\nabla_x \mathcal{L}_{\beta_{k-1}}(x_k, y_k), \partial g(x_k)) + \|A(x_k)\| \\ & \leq \frac{2\lambda'_A \sigma_1 + 2}{\nu \beta_{k-1}} (\lambda'_f + \lambda'_A y_{\max}) + \epsilon_k. \end{aligned} \quad (48)$$

376 For the second part of the theorem, we use the Weyl's inequality and Step 5 of Algorithm 1 to write

$$\begin{aligned} & \lambda_{\min}(\nabla_{xx} \mathcal{L}_{\beta_{k-1}}(x_k, y_{k-1})) \geq \lambda_{\min}(\nabla_{xx} \mathcal{L}_{\beta_{k-1}}(x_k, y_k)) \\ & \quad - \sigma_k \left\| \sum_{i=1}^m A_i(x_k) \nabla^2 A_i(x_k) \right\|. \end{aligned} \quad (49)$$

377 The first term on the right-hand side is lower bounded by $-\epsilon_{k-1}$ by Step 2 of Algorithm 1. We next
 378 bound the second term on the right-hand side above as

$$\begin{aligned} & \sigma_k \left\| \sum_{i=1}^m A_i(x_k) \nabla^2 A_i(x_k) \right\| \\ & \leq \sigma_k \sqrt{m} \max_i \|A_i(x_k)\| \|\nabla^2 A_i(x_k)\| \\ & \leq \sigma_k \sqrt{m} \lambda_A \frac{2\lambda'_f + 2\lambda'_A y_{\max}}{\nu \beta_{k-1}}, \end{aligned}$$

379 where the last inequality is due to (5,43). Plugging into (49) gives

$$\begin{aligned} & \lambda_{\min}(\nabla_{xx} \mathcal{L}_{\beta_{k-1}}(x_k, y_{k-1})) \\ & \geq -\epsilon_{k-1} - \sigma_k \sqrt{m} \lambda_A \frac{2\lambda'_f + 2\lambda'_A y_{\max}}{\nu \beta_{k-1}}, \end{aligned}$$

380 which completes the proof of Theorem 4.1.

381 **B Proof of Corollary 4.2**

382 Let K denote the number of (outer) iterations of Algorithm 1 and let ϵ_f denote the desired accuracy
 383 of Algorithm 1, see (11). Recalling Theorem 4.1, we can then write that

$$\epsilon_f = \frac{Q}{\beta_K}, \quad (50)$$

384 or, equivalently, $\beta_K = Q/\epsilon_f$. We now count the number of total (inner) iterations T of Algorithm 1
 385 to reach the accuracy ϵ_f . From (17) and for sufficiently large k , recall that $\lambda_{\beta_k} \leq \lambda''\beta_k$ is the
 386 smoothness parameter of the augmented Lagrangian. Then, from (24) ad by summing over the outer
 387 iterations, we bound the total number of (inner) iterations of Algorithm 1 as

$$\begin{aligned} T &= \sum_{k=1}^K \mathcal{O}\left(\frac{\lambda_{\beta_{k-1}}^2 \rho'^2}{\epsilon_k}\right) \\ &= \sum_{k=1}^K \mathcal{O}\left(\beta_{k-1}^3 \rho'^2\right) \quad (\text{Step 1 of Algorithm 1}) \\ &\leq \mathcal{O}\left(K\beta_{K-1}^3 \rho'^2\right) \quad (\{\beta_k\}_k \text{ is increasing}) \\ &\leq \mathcal{O}\left(\frac{KQ^3 \rho'^2}{\epsilon_f^3}\right). \quad (\text{see (50)}) \end{aligned} \quad (51)$$

388 In addition, if we specify $\beta_k = b^k$ for all k , we can further refine T . Indeed,

$$\beta_K = b^K \implies K = \log_b\left(\frac{Q}{\epsilon_f}\right), \quad (52)$$

389 which, after substituting into (51) gives the final bound in Corollary 4.2.

390 **C Proof of Lemma 2.1**

391 Note that

$$\mathcal{L}_\beta(x, y) = f(x) + \sum_{i=1}^m y_i A_i(x) + \frac{\beta}{2} \sum_{i=1}^m (A_i(x))^2, \quad (53)$$

392 which implies that

$$\begin{aligned} \nabla_x \mathcal{L}_\beta(x, y) &= \nabla f(x) + \sum_{i=1}^m y_i \nabla A_i(x) + \frac{\beta}{2} \sum_{i=1}^m A_i(x) \nabla A_i(x) \\ &= \nabla f(x) + DA(x)^\top y + \beta DA(x)^\top A(x), \end{aligned} \quad (54)$$

393 where $DA(x)$ is the Jacobian of A at x . By taking another derivative with respect to x , we reach

$$\begin{aligned} \nabla_x^2 \mathcal{L}_\beta(x, y) &= \nabla^2 f(x) + \sum_{i=1}^m (y_i + \beta A_i(x)) \nabla^2 A_i(x) \\ &\quad + \beta \sum_{i=1}^m \nabla A_i(x) \nabla A_i(x)^\top. \end{aligned} \quad (55)$$

394 It follows that

$$\begin{aligned} &\|\nabla_x^2 \mathcal{L}_\beta(x, y)\| \\ &\leq \|\nabla^2 f(x)\| + \max_i \|\nabla^2 A_i(x)\| (\|y\|_1 + \beta \|A(x)\|_1) \\ &\quad + \beta \sum_{i=1}^m \|\nabla A_i(x)\|^2 \\ &\leq \lambda_h + \sqrt{m} \lambda_A (\|y\| + \beta \|A(x)\|) + \beta \|DA(x)\|_F^2. \end{aligned} \quad (56)$$

395 For every x such that $\|x\| \leq \rho$ and $\|A(x)\| \leq \rho$, we conclude that

$$\|\nabla_x^2 \mathcal{L}_\beta(x, y)\| \leq \lambda_f + \sqrt{m} \lambda_A (\|y\| + \beta \rho') + \beta \max_{\|x\| \leq \rho} \|DA(x)\|_F^2, \quad (57)$$

396 which completes the proof of Lemma 2.1.

397 D Clustering

398 We only verify the condition in (20). Note that

$$A(x) = VV^\top \mathbf{1} - \mathbf{1}, \quad (58)$$

$$\begin{aligned} DA(x) &= \begin{bmatrix} w_{1,1}x_1^\top & \cdots & w_{1,n}x_1^\top \\ \vdots & & \\ w_{n,1}x_n^\top & \cdots & w_{n,n}x_n^\top \end{bmatrix} \\ &= [V \quad \cdots \quad V] + \begin{bmatrix} x_1^\top & & \\ & \ddots & \\ & & x_n^\top \end{bmatrix}, \end{aligned} \quad (59)$$

400 where $w_{i,i} = 2$ and $w_{i,j} = 1$ for $i \neq j$. In the last line above, n copies of V appear and the last
401 matrix above is block-diagonal. For x_k , define V_k as in the example and let $x_{k,i}$ be the i th row of V_k .
402 Consequently,

$$\begin{aligned} DA(x_k)^\top A(x_k) &= \begin{bmatrix} (V_k^\top V_k - I_n)V_k^\top \mathbf{1} \\ \vdots \\ (V_k^\top V_k - I_n)V_k^\top \mathbf{1} \end{bmatrix} \\ &\quad + \begin{bmatrix} x_{k,1}(V_k V_k^\top \mathbf{1} - \mathbf{1})_1 \\ \vdots \\ x_{k,n}(V_k V_k^\top \mathbf{1} - \mathbf{1})_n \end{bmatrix}, \end{aligned} \quad (60)$$

403 where $I_n \in \mathbb{R}^{n \times n}$ is the identity matrix. Let us make a number of simplifying assumptions. First,
404 we assume that $\|x_k\| < \sqrt{s}$, which can be easily enforced in the iterates. Under this assumption, it
405 follows that

$$(\partial g(x_k))_i = \begin{cases} 0 & (x_k)_i > 0 \\ \{a : a \leq 0\} & (x_k)_i = 0, \end{cases} \quad i \leq d. \quad (61)$$

406 Second, we assume that V_k has nearly orthonormal columns, namely, $V_k^\top V_k \approx I_n$. This can also be
407 easily enforced in each iterate of Algorithm 1 and naturally corresponds to well-separated clusters.
408 While a more fine-tuned argument can remove these assumptions, they will help us simplify the
409 presentation here. Under these assumptions, the (squared) right-hand side of (20) becomes

$$\begin{aligned} &\text{dist} \left(-DA(x_k)^\top A(x_k), \frac{\partial g(x_k)}{\beta_{k-1}} \right)^2 \\ &= \left\| (-DA(x_k)^\top A(x_k))_+ \right\|^2 \quad (a_+ = \max(a, 0)) \\ &= \left\| \begin{bmatrix} x_{k,1}(V_k V_k^\top \mathbf{1} - \mathbf{1})_1 \\ \vdots \\ x_{k,n}(V_k V_k^\top \mathbf{1} - \mathbf{1})_n \end{bmatrix} \right\|^2 \quad (x_k \in C \Rightarrow x_k \geq 0) \\ &= \sum_{i=1}^n \|x_{k,i}\|^2 (V_k V_k^\top \mathbf{1} - \mathbf{1})_i^2 \\ &\geq \min_i \|x_{k,i}\|^2 \cdot \sum_{i=1}^n (V_k V_k^\top \mathbf{1} - \mathbf{1})_i^2 \\ &= \min_i \|x_{k,i}\|^2 \cdot \|V_k V_k^\top \mathbf{1} - \mathbf{1}\|^2. \end{aligned} \quad (62)$$

410 Given a prescribed ν , ensuring $\|x_{k,i}\| \geq \nu$ guarantees (20). This requirement corresponds again to
 411 well-separated clusters. When the clusters are sufficiently separated and the algorithm is initialized
 412 close enough to the constraint set, there is indeed no need to separately enforce this condition. In
 413 practice, often n exceeds the number of true clusters and a more fine-tuned analysis is required to
 414 establish (20) by restricting the argument to a particular subspace of \mathbb{R}^n .

415 E Basis Pursuit

416 We only verify the regularity condition in (20) for (1) with f, A, g specified in (35). Note that

$$DA(x) = 2\bar{B}\text{diag}(x), \quad (63)$$

417 where $\text{diag}(x) \in \mathbb{R}^{2d \times 2d}$ is the diagonal matrix formed by x . The left-hand side of (20) then reads as

$$\begin{aligned} & \text{dist} \left(-DA(x_k)^\top A(x_k), \frac{\partial g(x_k)}{\beta_{k-1}} \right) \\ &= \text{dist}(-DA(x_k)^\top A(x_k), \{0\}) \quad (g \equiv 0) \\ &= \|DA(x_k)^\top A(x_k)\| \\ &= 2\|\text{diag}(x_k)\bar{B}^\top(\bar{B}x_k^{\circ 2} - b)\|. \quad (\text{see (63)}) \end{aligned} \quad (64)$$

418 To bound the last line above, let x_* be a solution of (1) and note that $\bar{B}x_*^{\circ 2} = b$ by definition. Let also
 419 $z_k, z_* \in \mathbb{R}^d$ denote the vectors corresponding to x_k, x_* . Corresponding to x_k , also define $u_{k,1}, u_{k,2}$
 420 naturally and let $|z_k| = u_{k,1}^{\circ 2} + u_{k,2}^{\circ 2} \in \mathbb{R}^d$ be the amplitudes of z_k . To simplify matters, let us assume
 421 also that B is full-rank. We then rewrite the last line of (64) as

$$\begin{aligned} & \|\text{diag}(x_k)\bar{B}^\top(\bar{B}x_k^{\circ 2} - b)\|^2 \\ &= \|\text{diag}(x_k)\bar{B}^\top\bar{B}(x_k^{\circ 2} - x_*^{\circ 2})\|^2 \quad (\bar{B}x_*^{\circ 2} = b) \\ &= \|\text{diag}(x_k)\bar{B}^\top B(x_k - x_*)\|^2 \\ &= \|\text{diag}(u_{k,1})B^\top B(z_k - z_*)\|^2 \\ & \quad + \|\text{diag}(u_{k,2})B^\top B(z_k - z_*)\|^2 \\ &= \|\text{diag}(u_{k,1}^{\circ 2} + u_{k,2}^{\circ 2})B^\top B(z_k - z_*)\|^2 \\ &= \|\text{diag}(|z_k|)B^\top B(z_k - z_*)\|^2 \\ &\geq \eta_n(B\text{diag}(|z_k|))^2 \|B(z_k - z_*)\|^2 \\ &= \eta_n(B\text{diag}(|z_k|))^2 \|Bz_k - b\|^2 \quad (Bz_* = \bar{B}x_*^{\circ 2} = b) \\ &\geq \min_{|T|=n} \eta_n(B_T) \cdot |z_{k,(n)}|^2 \|Bz_k - b\|^2, \end{aligned} \quad (65)$$

422 where $\eta_n(\cdot)$ returns the n th largest singular value of its argument. In the last line above, B_T is the
 423 restriction of B to the columns indexed by T of size n . Moreover, $z_{k,(n)}$ is the n th largest entry of z
 424 in magnitude. Given a prescribed ν , (20) therefore holds if

$$|z_{k,(n)}| \geq \sqrt{\frac{\nu}{\min_{|T|=n} \eta_n(B_T) \cdot \|Bz_k - b\|^2}}, \quad (66)$$

425 for every iteration k . If Algorithm 1 is initialized close enough to the solution z^* , there will be no
 426 need to directly enforce this condition.

427 **Discussion** The true potential of the reformulation of BP in (35) is in dealing with more structured
 428 norms than ℓ_1 , where computing the proximal operator is often intractable. One such case is the
 429 latent group lasso norm [?], defined as

$$\|z\|_\Omega = \sum_{i=1}^I \|z_{\Omega_i}\|,$$

430 where $\{\Omega_i\}_{i=1}^I$ are (not necessarily disjoint) index sets of $\{1, \dots, d\}$. Although not studied here, we
 431 believe that the non-convex framework presented in this paper can serve to solve more complicated
 432 problems, such as the latent group lasso. We leave this research direction for future work.

433 **E.1 ℓ_∞ Denoising with a Generative Prior**

434 The authors of [?] have proposed to project onto the range of a Generative Adversarial network
 435 (GAN) [?], as a way to defend against adversarial examples. For a given noisy observation $x^* + \eta$,
 436 they consider a projection in the ℓ_2 norm. We instead propose to use our augmented Lagrangian
 437 method to denoise in the ℓ_∞ norm, a much harder task:

$$\begin{aligned} \min_{x,z} \quad & \|x^* + \eta - x\|_\infty \\ \text{s.t.} \quad & x = G(z). \end{aligned} \tag{67}$$

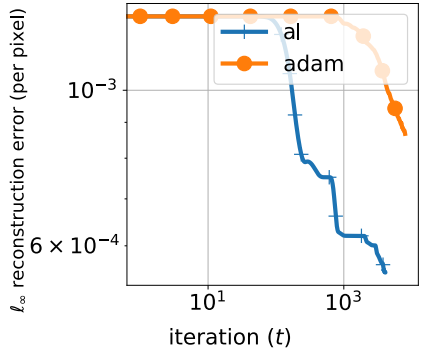


Figure 3: Augmented Lagrangian vs Adam for ℓ_∞ denoising (left). ℓ_2 vs ℓ_∞ denoising as defense against adversarial examples

438 We use a pretrained generator for the MNIST dataset, given by a standard deconvolutional neural
 439 network architecture. We compare the successful optimizer Adam against our method. Our algorithm
 440 involves two forward/backward passes through the network, as opposed to Adam that requires only
 441 one. For this reason we let our algorithm run for 4000 iterations, and Adam for 8000 iterations.
 442 For a particular example, we plot the objective value vs iteration count in figure E.1. Our method
 443 successfully minimizes the objective value, while Adam does not succeed.

444 **E.2 Generalized Eigenvalue Problem**

445 Generalized eigenvalue problem has extensive applications in machine learning, statistics and data
 446 analysis [?]. The well-known nonconvex formulation of the problem is [?].

$$\begin{cases} \min_{x \in \mathbb{R}^n} x^\top C x \\ x^\top B x = 1, \end{cases} \tag{68}$$

447 where $B, C \in \mathbb{R}^{n \times n}$ are symmetric matrices and B is positive definite, i.e. $B \succ 0$. The generalized
 448 eigenvector computation is equivalent to performing principal component analysis (PCA) of C in
 449 the norm B . Moreover, it is also equivalent to computing the top eigenvector of symmetric matrix
 450 $S = B^{-1/2} C B^{-1/2}$ and multiplying the resulting vector by $B^{-1/2}$. However, for sufficiently large n ,
 451 computing $B^{-1/2}$ is extremely expensive. The natural convex sdp relaxation for (68) involves lifting
 452 $Y = x x^\top$ and removes the non-convex $\text{rank}(Y) = 1$ constraint,

$$\begin{cases} \min_{Y \in \mathbb{R}^{n \times n}} \text{tr}(CY) \\ \text{tr}(BY) = 1, \quad X \succeq 0. \end{cases} \tag{69}$$

453 Here, we solve (68) because it directly fits into our template with,

$$\begin{aligned} f(x) &= x^\top C x, \quad g(x) = 0 \\ A(x) &= x^\top B x - 1. \end{aligned} \tag{70}$$

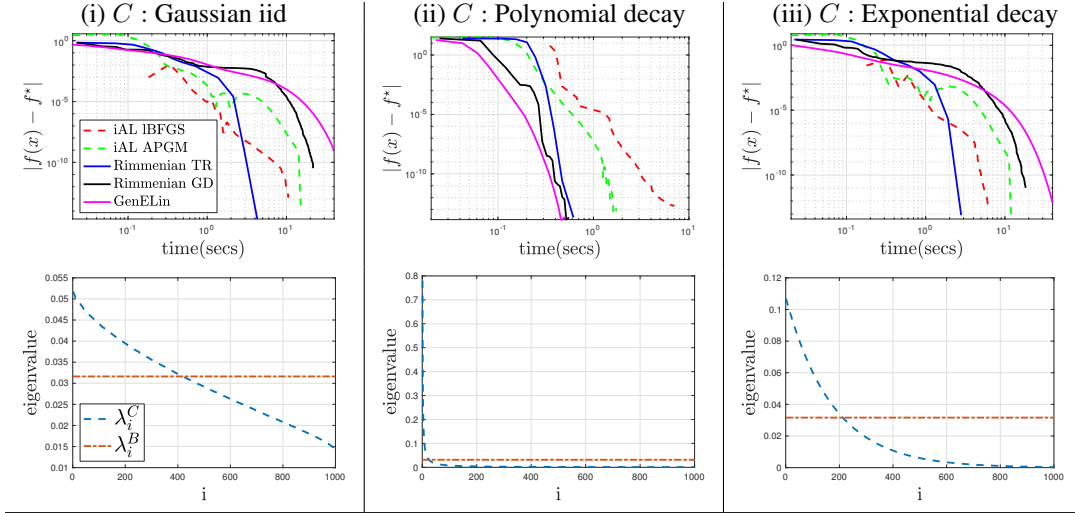


Figure 4: (Top) Objective convergence for calculating top generalized eigenvalue and eigenvector of B and C . (Bottom) Eigenvalue structure of the matrices. For (i),(ii) and (iii), C is positive semidefinite; for (iv), (v) and (vi), C contains negative eigenvalues. [(i): Generated by taking symmetric part of iid Gaussian matrix. (ii): Generated by randomly rotating $\text{diag}(1^{-p}, 2^{-p}, \dots, 1000^{-p})(p = 1)$. (iii): Generated by randomly rotating $\text{diag}(10^{-p}, 10^{-2p}, \dots, 10^{-1000p})(p = 0.0025)$.]

454 We compare our approach against 3 different methods. Manifold based Riemannian gradient descent
455 and Riemannian trust region methods in[?]] and generalized eigenvector via linear system solver
456 (abbreviated as. GenELin) in [?]. We have used Manopt software package in [?] for the manifold
457 based methods. For GenELin, we have utilized Matlab's backslash operator as the linear solver. The
458 results are compiled in Figure 4.