

1 **Escaping Saddle Points in Successive Convex Approximation***

2 Amrit Singh Bedi, Ketan Rajawat, Vaneet Aggarwal, and Alec Koppel †

3

4 **Abstract.** Optimizing non-convex functions is of primary importance in modern machine learning due to the fact that it underlies the training of deep networks and nonlinear dimensionality reduction. While the capability of first-order algorithms to converge to local extrema rather than saddle points has gained attention recently, it is well known that their convergence rates are inferior to those which exploit additional structure. In particular, empirically, successive convex approximation (SCA) converges faster than first-order methods. However, SCA in non-convex settings without additional structure in general converge to first-order stationary points, which could either be local extrema or saddle points whose performance is typically inferior. To mitigate this issue, we propose a perturbation of SCA at appropriate times, which inherits the same convergence rate results as the non-perturbed version, while enhancing the quality of its limit points. In particular, we establish that SCA with perturbations converges to a second order stationary points. Experiments on multi-dimensional scaling, a machine learning problem whose training objective is non-convex, substantiates the performance gains associated with employing random perturbations.

16 **Key words.** Stochastic optimization, non-convex optimization, first order methods.

17 **AMS subject classifications.** 46N10, 93E35, 68W27.

18 **1. Introduction.** Optimization of non-convex continuous functions is an increasingly salient
19 problem as it underlies the training of deep networks [11], clustering, matrix completion [34], phase
20 retrieval [33], tensor decomposition [9], tensor completion [16, 37], dimensionality reduction [38],
21 among others. These tools are employed to discern patterns in data sets whose scale grows larger by
22 the day. Despite the prevalence of non-convex optimization, methods for solving it efficiently is an on-
23 going challenge [12]. The most common approaches are gradient-based search [9], convex relaxation
24 [17], and characterizing which geometries exhibit “benign landscapes” that may be exploited through
25 an appropriately chosen metric.

26 In this work, we adopt the spirit of gradient search for general non-convex problems. Gradient
27 methods are popular due to their ease of implementation, simple analysis, and ability to serve as a
28 template to mitigate the challenges of non-convexity. These challenges include the fact that finding
29 the global extrema of a non-convex problem is NP hard in general, and that the best limit point one
30 may hope for in general (in an almost sure sense) is convergence to local extrema. We note that even
31 attaining acceptable local extrema is non-trivial, since gradient methods converge to stationary points,
32 which could either be saddle points or local extrema. Ensuring convergence to local extrema can
33 be ensured by perturbing the optimization iterates by random noise [10, 24] or augmented step-size
34 rules [7]. While such favorable limiting performance make perturbed gradient search attractive as a
35 template upon which to study finite-time performance [13], it is well-known that successive convex
36 approximation (SCA) experimentally converges far faster [29, 28, 15].

* Submitted to the editors October 2, 2019.

† A. S. Bedi and A. Koppel are with US Army Research Laboratory, Adelphi, MD, USA {Email: amrit0714@gmail.com, alec.e.koppel.civ@mail.mil}. K. Rajawat is with the Department of Electrical Engineering, IIT Kanpur, India {Email: ketan@iitk.ac.in}. V. Aggarwal is with the School of Industrial Engineering and the School of Electrical and Computer Engineering, Purdue University, West Lafayette IN 47907, USA {Email: vaneet@purdue.edu}.

SCA, rather than computing its local approximation of the non-convex objective function using the first-order Taylor expansion, as in gradient methods, instead uses any convex surrogate function whose gradient coincides with the original objective. Thus, one may encode any higher-order algorithm through an appropriate choice of convex surrogate. As a result, SCA has been widely used to (a) separate the function in its variables, achieving parallel and distributed implementation, (b) discern computationally cheap updates either because the surrogate admits a closed-form or computationally affordable minimizer [30, 28]. In practice, the chosen surrogate could be well-behaved and capable of navigating narrow valleys, i.e., exploiting higher-order structure such as curvature, unlike gradient descent. The empirical performance of convex surrogates has been well-investigated within the context of majorization-minimization [19] and successive convex approximation methods [1, 26, 30]. Unfortunately, existing uses of SCA for non-convex problems are only theoretically guaranteed to converge to stationary points, which may be degenerate saddle points [28]. Thus, one would like to employ SCA due to its favorable convergence in practice, while guaranteeing that it does not become stuck at undesirable critical points of the objective. Doing so is the goal of this work through the use of randomized perturbations.

2. Related Work and Contributions. First-order methods have become widespread in machine learning, but their performance on non-convex problems can be problematic, as they may become stuck at saddle points. Inspired from statistical mechanics, augmentations of stochastic search that add random perturbations date back to at least [10, 24]. [24] establishes global convergence in distribution to the set of minimizers. The corresponding nonasymptotic analysis in the context of non-convex learning problems is provided in [25]. We focus on convergence to local extrema in finite time with high probability, a weaker condition but for a more stringent time horizon. The reason for this focus is that finite-time performance may be translated into sample complexity, a salient issue in training learning models. Alternatively, when favorable geometry is present, as in matrix completion or phase retrieval [23, 5, 34], faster local convergence may be attained. However, these performance gains break down when moving from specialized to general non-convex optimization problems.

Beyond first-order methods, majorization-minimization (MM) [35, 19], SCA [29], and quasi-newton methods [39] have been proposed for non-convex non-convex optimization in order to reduce time to convergence. The first two approaches break the non-convex problem into a sequence of minimization problems involving surrogate functions. At the point of approximation, MM minimizes a convex upper-bound for the non-convex objective, while SCA employs a surrogate that is convex without any upper bound condition, but instead gradient consistency. Therefore, SCA better approximates the non-convex objective and results in better algorithm accuracy [15]. Quasi-Newton methods, by contrast, improve upon gradient methods by estimating the Hessian in order to implement approximate Newton steps.

Motivated by the favorable attributes and generality of SCA, in this work we propose augmentations that facilitate extending its favorable learning rates to favorable limiting behavior. To do so, we note that for non-convex objectives, the convergence to first-order stationary points (FOSP) is established in [28, 15]. By contrast, here we employ perturbation methods to escape saddle points and converge to the second order stationary point (SOSP), which under suitable conditions coincide with local extrema [13].

To escape the saddle points, two popular techniques, namely, random perturbations [13] and second-order based methods [22] have been used. A cubic regularization based algorithm is proposed

80 in [22] with convergence to ϵ SOSP in $\mathcal{O}(\frac{1}{\epsilon^{3/2}})$ number of iterations. In practice, the second order
 81 information based methods require Hessian inversion at each algorithm iteration, which is computa-
 82 tionally expensive. On the other hand, in perturbation based methods, noise from a given distribution
 83 (uniform in this paper) is added to the algorithm update (when in the neighborhood of a saddle point)
 84 to escape the saddle point which is computationally inexpensive. This work proposes SCA with ran-
 85 dom perturbations in order to inherit the favorable experimental convergence of SCA together with
 86 improved limiting performance.

87 To do so, we interpret SCA algorithms as an inexact version of gradient descent, which permits us
 88 to establish its convergence is comparable to [3, 31]. Then, building upon this foundation, we augment
 89 its limit points through random perturbation arguments pioneered by [13], which in particular. yields
 90 its convergence to SOSP. The result is a unique modification of SCA [28], which we call perturbed
 91 successive convex approximation (P-SCA), that converges to ϵ second order stationary point in $\mathcal{O}(\frac{1}{\epsilon^2})$
 92 with an extra multiplicative factor of $\text{polylog}(d)$ where d is the underlying parameter dimension. More
 93 broadly, our analysis extends convergence of random perturbation methods to SOSP to any arbitrary
 94 inexact gradient descent algorithm under an additional hypothesis that the norm of the error in the
 95 gradient is bounded. Thus, our main contributions are as follows:

- 96 • To date, convergence rates of successive convex approximation is missing, i.e., existing analy-
 97 ses focus on asymptotic convergence [28]. By contrast, we present the convergence rate analy-
 98 sis for SCA.
- 99 • By virtue of these rates, we are then able to employ random perturbation based arguments to
 100 establish conditions under which randomization yields convergence to second-order stationary
 101 points. Doing so requires extending the ideas in [13] methods which employ inexact gradient
 102 directions. The reason for this is that we establish that SCA may be interpreted as a member
 103 of the family of inexact gradient algorithms, and hence is no longer a first order method.
- 104 • With this foundation laid, the results of [13] are verified and proved for gradient descent with
 105 in errors.
- 106 • Experimentally, we demonstrate the advantages of the proposed algorithm via solving a multi-
 107 dimensional scaling (MDS) problem in Sec. 6.

108 **3. Preliminaries.** In this paper, we are solving a non-convex optimization problem with the
 109 help of gradient based methods. To follow the analysis, we describe some useful definitions which
 110 will be referred throughout the paper.

111 **Definition 1** For a differentiable function $U(\mathbf{x})$, a point \mathbf{x} is ϵ first order stationary point (FOSP)
 112 if $\|\nabla U(\mathbf{x})\| \leq \epsilon$.

113 **Definition 2** (Strict saddle point) For a twice differentiable function, a point \mathbf{x} is a strict saddle
 114 point if \mathbf{x} is a FOSP and $\lambda_{\min}(\nabla^2 U(\mathbf{x})) < 0$. For a strict saddle point, the second order stationary
 115 point (SOCP) is a local minima of the non-convex objective as defined next.

116 **Definition 3** (ϵ second order order stationarity) For a non-convex function $U(\mathbf{x})$, a point \mathbf{x} is ϵ
 117 second order stationary if $\|\nabla U(\mathbf{x})\| \leq \epsilon$ and $\lambda_{\min}(\nabla^2 U(\mathbf{x})) \geq -\sqrt{L_2}\epsilon$, where L_2 is Hessian
 118 Lipschitz constant.

119 **4. Problem Formulation.** In this work, we focus on the minimization of a non-convex objec-
 120 tive function $U(\mathbf{x})$ given by

121 (4.1)
$$\min_{\mathbf{x}} U(\mathbf{x})$$

123 where $\mathbf{x} \in \mathbb{R}^n$ is the optimization variable, $U(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice differentiable non-convex
 124 function. For a general non-convex optimization problem, it is well known that an algorithm cannot
 125 achieve global minima using first order algorithms [9, 13]. In literature, there are various algorithms
 126 proposed which converges to the approximate FOSP of the problem (4.1). A FOSP could be a local
 127 minima, local maxima, or a saddle point. It is sufficient to achieve convergence to FOSP in convex
 128 settings because it coincides with the global minima of the function. However, for the non-convex
 129 functions, it may be highly suboptimal to achieve FOSP if it denotes the local maxima and we are
 130 minimizing the function and vice versa. In addition, the converged FOSP can be a saddle point for the
 131 unconstrained optimization problem and algorithm might get stuck there. Hence, it becomes important
 132 to look for the convergence to a second order stationary point. In literature, various algorithms are
 133 proposed for problem (4.1) with the guarantee of convergence to a second order stationary point [6,
 134 13, 14]. However, they are based on the standard gradient descent algorithm.

135 Recently, an SCA based algorithm is proposed by [28] to solve the problem of (4.1) in an iterative
 136 manner. We utilize the similar ideas and use SCA to solve the unconstrained version of the prob-
 137 lem. The SCA algorithm is shown to outperform the other first order methods in literature [28]. In
 138 this method, the non-convex objective function $U(\mathbf{x})$ is approximated by a surrogate convex function
 139 $\tilde{U}(\mathbf{x}; \mathbf{y})$ approximated at \mathbf{y} . Utilizing this idea, the standard SCA algorithm to solve the problem in
 140 (4.1) is summarized in Algorithm 4.1.

Algorithm 4.1 Successive convex approximation (SCA)

- 1: Set $t = 0$, initialize $\eta = (0, 1]$,
- 2: **STOP** if \mathbf{x}_t is a FOSP
- 3: Solve the following optimization problem to get $\hat{\mathbf{x}}(\mathbf{x}_t)$

$$(4.2) \quad \hat{\mathbf{x}}(\mathbf{x}_t) := \arg \min_{\mathbf{x} \in \mathbb{R}^n} \tilde{U}(\mathbf{x}; \mathbf{x}_t),$$

- 4: Set $\mathbf{x}_{t+1} = \mathbf{x}_t + \eta(\hat{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t)$
 - 5: Set $t = t + 1$ and go to step 2
-

141 In Algorithm 4.1, at each step t , a convex approximation $\tilde{U}(\mathbf{x}; \mathbf{x}_t)$ of the non-convex objective
 142 function $U(\mathbf{x})$ is used to perform the update in step 3. The main idea to construct an approximation
 143 function $\tilde{U}(\cdot)$ is to utilize a simple convex function whose gradient calculation is easier. A detailed
 144 discussion about the selection of $\tilde{U}(\cdot)$ is provided in [28]. The approximated convex function \tilde{U}
 145 should preserve the first order condition of the original non-convex function such that the gradient of
 146 the both functions is same; see Assumption 2. It is important to note that unlike the standard successive
 147 approximation techniques, the proposed algorithm does not require the $\tilde{U}(\cdot)$ to upper bound the non-
 148 convex function U . This provides a lot of freedom in the selection of $\tilde{U}(\cdot)$ which depends upon the
 149 application of interest. Some specific examples are provided next.

150 • **Block wise convex:** If the objective function $U(\mathbf{x}; \mathbf{y})$ is block wise convex with respect to each
 151 block i , the following approximation can be used

$$152 \quad (4.3) \quad \tilde{U}_i(\mathbf{x}_i; \mathbf{y}) := U(\mathbf{x}_i, \mathbf{y}_{-i}) + \frac{\tau_i}{2} (\mathbf{x}_i - \mathbf{y}_i)^T \mathbf{H}_i(\mathbf{y}) (\mathbf{x}_i - \mathbf{y}_i),$$

153 where $\mathbf{y} := \{\mathbf{y}_i\}_{i=1}^I$, $\mathbf{y}_{-i} := \{\mathbf{y}_j\}_{j \neq i}$.

154 • **Product of functions:** If the objective function is of the type $U(\mathbf{x}) = f_1(\mathbf{x})f_2(\mathbf{x})$ where $f_1(\mathbf{x})$
 155 and $f_2(\mathbf{x})$ are convex and positive, we can use

156
$$\tilde{U}(\mathbf{x}; \mathbf{y}) = f_1(\mathbf{x})f_2(\mathbf{y}) + f_1(\mathbf{y})f_2(\mathbf{x}) + \frac{\tau_i}{2}(\mathbf{x} - \mathbf{y})^T \mathbf{H}(\mathbf{y})(\mathbf{x} - \mathbf{y}).$$

 157

158 Algorithm 4.1 converges to the first order stationary point as derived in [28, Theorem 2], which
 159 states that for a constant step size η , the sequence $\{\mathbf{x}_t\}$ is bounded and each of its limit points is
 160 stationary point of the optimization problem in (4.1). But the results presented in [28] are asymptotic
 161 in nature and does not provide the number of iterations required to converge to an ϵ optimal solution.
 162 This results is important because it characterizes the convergence behavior of the proposed algorithm.
 163 Moreover, as stated earlier, the convergence to FOSP is not a sufficient condition of convergence to
 164 local minima. This is because a FOSP could be a saddle point and the proposed algorithm might get
 165 stuck there.

166 In this paper, we are interested in deriving the number of iterations required to achieve an ϵ FOSP.
 167 In addition to that, we want to characterize the number of iterations required to achieve an ϵ second
 168 order stationary point. In literature there are different techniques proposed to escape the saddle points
 169 and converge to SOSP. For instant, some of the approached include methods utilizing second order
 170 information [20], and random perturbation based methods [13]. Motivated from the advantages of
 171 perturbation based methods as discussed in [13] and advantages of SCA methods [28], we use the
 172 perturbed method to make the SCA algorithm escape saddle points effectively. A perturbed successive
 173 convex approximation (P-SCA) algorithm is summarized in Algorithm 4.2.

Algorithm 4.2 :P-SCA Algorithm

- 1: Set $t = 0$, initialize $\chi \leftarrow 3 \max\{\log\left(\frac{dL_1\Delta_U}{c\epsilon^2\delta}, 4\right)\}$, $\eta \leftarrow \frac{c}{L_1}$, $r \leftarrow \frac{\epsilon\sqrt{c}}{L_1\chi^2}$, $g_{\text{th}} \leftarrow \frac{\epsilon\sqrt{c}}{\chi^2}$, $f_{\text{th}} \leftarrow \frac{c}{\chi^3}\sqrt{\frac{\epsilon^3}{L_2}}$, $t_{\text{th}} = \frac{(1-s)\chi L_1}{c^2\sqrt{\epsilon L_2}}$, $t_{\text{noise}} = -t_{\text{th}} - 1$, $s \in (0, 1)$, and \mathbf{x}_0
- 2: If $\|\nabla U(\mathbf{x}_t)\| \leq g_{\text{th}}$ and $t - t_{\text{noise}} > t_{\text{th}}$ then
 - $\tilde{\mathbf{x}}_t \leftarrow \mathbf{x}_t$, and $t_{\text{noise}} \leftarrow t$
 - $\mathbf{x}_t \leftarrow \tilde{\mathbf{x}}_t + \xi_t$, where ξ is uniformly $\sim \mathbb{B}_0(r)$ with radius r
- 3: If $t - t_{\text{noise}} = t_{\text{th}}$ and $U(\mathbf{x}_t) - U(\tilde{\mathbf{x}}_{t_{\text{noise}}}) > -(1-s)f_{\text{th}}$ then
 - return** $\tilde{\mathbf{x}}_{t_{\text{noise}}}$
- 4: Solve the following optimization problem to get $\hat{\mathbf{x}}(\mathbf{x}_t)$

$$(4.4) \quad \hat{\mathbf{x}}(\mathbf{x}_t) := \arg \min_{\mathbf{x}} \tilde{U}(\mathbf{x}; \mathbf{x}_t)$$

- 5: Set $\mathbf{x}_{t+1} = \mathbf{x}_t + \eta(\hat{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t)$
 - 6: Set $t = t + 1$ and go to Step 2
-

174 In Algorithm 4.2, the standard SCA algorithm updates are used to reach FOSP. If the condition
 175 for FOSP are satisfied (which is nothing but the gradient norm is smaller than a threshold value g_{th}),
 176 a random perturbation from uniform ball $\mathbb{B}_0(r)$ of radius r around $\tilde{\mathbf{x}}_t$ is added. This happens at most
 177 once every t_{th} number of iterations. After adding perturbation, if the function value is not decreased
 178 by the threshold $(1-s)f_{\text{th}}$ (where $s \in (0, 1)$ is a constant and defined later in the paper) then $\tilde{\mathbf{x}}_{t_{\text{noise}}}$ is
 179 returned as the output of Algorithm 4.2. Here t_{noise} denotes the last iteration at which the optimization
 180 variable was perturbed with respect to current iteration of the algorithm. It is proved in [28] that after

181 the addition of this random perturbations, the standard gradient descent iterates escapes the saddle
 182 points efficiently. We extend the analysis in [28] and [13] papers to prove that the proposed P-SCA
 183 algorithm also escapes the saddle points in a similar manner and $\hat{\mathbf{x}}_{\text{noise}}$ will be in the neighborhood of
 184 ϵ SOSP. We note that the steps in Algorithm 4.2 do not exhibit a gradient descent like update. Hence,
 185 it is difficult to analyze the dynamics of the proposed algorithm. To perform the analysis, it becomes
 186 important to look at the proposed algorithm from the perspective of inexact gradient descent algorithms
 187 [32, 2, 8, 31]. The meaning of term ‘inexact’ means that the gradient value used for the algorithm
 188 update is in error. The error may arise due to stochastic nature of the gradient or partial availability
 189 of the gradient [3] etc. In the above mentioned references, all the inexact gradient algorithms are
 190 considered for convex objective function settings. The analysis for asymptotic convergence to FOSP
 191 under non-convex objective is performed in [3] but with diminishing step size. The work in this paper
 192 focuses on more general scenario of non-convex objective and inexactness of the gradient [40].

193 Hence, the main idea for the inexact version is to interpret the algorithm as a gradient descent
 194 algorithm with error. It is proved in literature that under some regularity conditions on the error
 195 term in the gradient, the algorithm converges to (or provide close approximation to) the optimal value
 196 under convexity assumption . In order to proceed with analysis for the proposed P-SCA algorithm, we
 197 first show that the SCA algorithm is nothing but an inexact version of the standard gradient descent
 198 algorithm. From step 5 of the Algorithm 4.2, we have

$$200 \quad (4.5) \quad \mathbf{x}_{t+1} = \mathbf{x}_t + \eta(\hat{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t),$$

201 Next, add and subtract the original gradient $\nabla U(\mathbf{x}_t)$ as follows

$$203 \quad (4.6) \quad \mathbf{x}_{t+1} = \mathbf{x}_t + \eta(\nabla U(\mathbf{x}_t) - \nabla U(\mathbf{x}_t) + \hat{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t).$$

204 Finally, we can write the P-SCA algorithm update in the following alternative form

$$205 \quad (4.7) \quad \mathbf{x}_{t+1} = \underbrace{\mathbf{x}_t - \eta \nabla U(\mathbf{x}_t)}_{\text{GD update}} + \eta \underbrace{[\nabla U(\mathbf{x}_t) + \mathbf{x}_t - \hat{\mathbf{x}}(\mathbf{x}_t)]}_{:= \mathbf{e}_t}.$$

207 The first term in (4.7) is similar to the standard gradient descent algorithm. The second term in (4.7)
 208 corresponds to the error in the gradient and denoted by \mathbf{e}_t . Hence the final gradient used for the update
 209 is $\nabla U(\mathbf{x}_t) + \mathbf{e}_t$ as follows

$$210 \quad (4.8) \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \eta(\nabla U(\mathbf{x}_t) + \mathbf{e}_t).$$

212 For the further analysis in this paper, we will utilize the algorithm update form as presented in (4.7).

213 **5. Convergence analysis.** This section details about the convergence rate analysis of the pro-
 214 posed P-SCA algorithm and show that the number of iterations required to converge to a second order
 215 stationary point are of the order of $\mathcal{O}\left((\text{polylog}(d))\frac{L_1(U(\mathbf{x}_0) - U(\mathbf{x}^*))}{\epsilon^2}\right)$. This section also establishes
 216 the result that the proposed algorithm indeed escapes the saddle point with probability $1 - \delta$ for any
 217 $\delta > 0$. Note that the updates in the proposed algorithm utilize a convex approximations $\tilde{U}(\mathbf{x}; \mathbf{x}_t)$ at
 218 each step t for the non-convex function $U(\mathbf{x})$ [28]. Before discussing the mathematical results, some
 219 assumptions are required to hold for the objective function and its convex approximation. All the
 220 required assumptions are provided next.

221 **Assumption 1.** *The continuously differentiable objective function $U(\mathbf{x})$ and its gradient $\nabla U(\mathbf{x})$
222 are Lipschitz continuous with parameter $L_0 > 0$ and $L_1 > 0$, respectively. This implies that*

223 (5.1) $\|U(\mathbf{x}) - U(\mathbf{y})\| \leq L_0 \|\mathbf{x} - \mathbf{y}\|,$

224 (5.2) $\|\nabla U(\mathbf{x}) - \nabla U(\mathbf{y})\| \leq L_1 \|\mathbf{x} - \mathbf{y}\|$

226 for all \mathbf{x} and \mathbf{y} .

227 Apart from the objective function, the convex surrogate function needs to satisfy the following as-
228 sumption.

229 **Assumption 2.** *The convex approximation function $\tilde{U}(\mathbf{x}; \mathbf{y})$ to non-convex objective function
230 $U(\mathbf{x})$ at \mathbf{y} is continuously differentiable with respect to its first argument such that
231 B1) $\tilde{U}(\cdot; \mathbf{y})$ is uniformly strongly convex with parameter C , which implies that*

232 (5.3) $(\nabla \tilde{U}(\mathbf{x}; \mathbf{y}) - \nabla \tilde{U}(\mathbf{z}; \mathbf{y}))^T (\mathbf{x} - \mathbf{z}) \geq C \|\mathbf{x} - \mathbf{z}\|^2$

234 for all \mathbf{x} and \mathbf{z} .

235 B2) *The approximation function gradient $\nabla \tilde{U}(\cdot; \mathbf{y})$ is equal to the original function gradient at
236 the approximation point \mathbf{y}*

237 (5.4) $\nabla \tilde{U}(\mathbf{y}; \mathbf{y}) = \nabla U(\mathbf{y}).$

239 This property is the key to represent the P-SCA algorithm as an inexact gradient descent
240 algorithm.

241 Another assumption required for the objective function which upper bounds the rate of change of the
242 Hessian.

243 **Assumption 3.** *The objective function $U(\mathbf{x})$ is Hessian Lipschitz with parameter L_2 , which im-
244 plies that*

245 (5.5) $\|\nabla^2 U(\mathbf{x}) - \nabla^2 U(\mathbf{y})\| \leq L_2 \|\mathbf{x} - \mathbf{y}\|$

247 for all \mathbf{x} and \mathbf{y} .

248 All of the above mentioned assumptions are standard and have been considered in literature [28].
249 Assumption 1 states that the non-convex objective function is Lipschitz which means that the gradient
250 of the objective function is smooth and does not changes arbitrarily. This is an important assumption
251 since it provides an upper bound on the gradient at two difference points in terms of the difference
252 between the points itself. This is a standard assumption in the analysis of inexact gradient descent
253 algorithms, for instance see [27]. Next, assumption 2 is for the convex approximation $\tilde{U}(\mathbf{x}; \mathbf{y})$ to the
254 non-convex objective function $U(\mathbf{x})$ at \mathbf{y} , and states that $\tilde{U}(\mathbf{x}; \mathbf{y})$ is strongly convex at given \mathbf{y} with
255 parameter C . This assumption assures that the convex optimization problem in (4.4) has a unique
256 solution. The last assumption 3 states that the gradient of $U(\mathbf{x})$ is smooth and hence the Hessian of
257 the function does not changes abruptly. All of the above mentioned assumptions are utilized to prove
258 the convergence results in this paper.

259 Before proceeding with the analysis, we remark that the analysis performed in [13] holds for the
260 gradient descent algorithm which is a first order method and is not applicable to SCA based methods.

261 This is because there is a minimization step included in the proposed Algorithm 4.2 which is not
 262 present in [13]. To proceed with the analysis, we first prove that the proposed algorithm is an inexact
 263 version of the gradient descent algorithm. Therefore, Assumption 1 is required to perform the analysis
 264 and handle the error in the inexact version of the gradient. This is a standard assumption in the analysis
 265 of inexact gradient descent algorithms [27]. Then, all the results of [13] are verified and proved for the
 266 gradient descent algorithm with gradient in error.

267 We now show that it is possible for the SCA algorithm to converge to the second order stationary
 268 point with a simple modification. In the SCA algorithm updates, once the gradient norm is less than a
 269 threshold value g_{th} , we add a small random perturbation to current iterate $\tilde{\mathbf{x}}_t$. In other words, a random
 270 perturbation is added to the algorithm updates at most once after every t_{th} iterations. To make the
 271 analysis simple, similar to [13] we assume that the random perturbation ξ_t is uniformly sampled from
 272 a ball $\mathbb{B}_0(r)$ of radius r centered around $\tilde{\mathbf{x}}_t$. After adding the perturbation at iteration t , if the function
 273 value does not decrease by $(1 - s)f_{\text{th}}$ after t_{th} iterations, then the algorithm stops and return $\tilde{\mathbf{x}}_{t_{\text{noise}}}$.
 274 This paper proves that the output $\tilde{\mathbf{x}}_{t_{\text{noise}}}$ is essentially the second order stationary point for the problem
 275 in (4.1). The main result is presented next in Theorem 5.1.

276 **Theorem 5.1.** *Let the assumptions 1-3 are satisfied, there exists a constant c_{\max} such that, for any
 277 $\delta > 0$, $\epsilon \leq \frac{L_1^2}{L_2}$, $\Delta_U \geq U(\mathbf{x}_0) - U(\mathbf{x}^*)$, and $c \leq c_{\max}$, the output of Algorithm 4.2 is a ϵ second order
 278 stationary point with probability $(1 - \delta)$ for problem (4.1). To achieve the ϵ second order stationary
 279 point, the number of iterations required for the algorithm is given by*

$$280 \quad (5.6) \quad \mathcal{O}\left(\frac{L_1 \Delta_U}{\epsilon^2} \log^4\left(\frac{d L_1 \Delta_U}{\epsilon^2 \delta}\right)\right).$$

282 It is interesting to note that the convergence rate of the proposed algorithm remains the same as
 283 that of the one in [13, Theorem 3]. The proof of the Theorem 5.1 and succeeding lemmas follows the
 284 similar ideas to [13]. Theorem 5.1 establishes the fact that the algorithm terminates in finite number
 285 of iterations. Before discussing the proof of Theorem 5.1, there are some other results which we need
 286 to discuss first. The algorithm operation can be divided in two scenarios 1) the iterate is not close to
 287 FOSP and the gradient $\nabla U(\mathbf{x}_t)$ is large and 2) when gradient is small but the Hessian $\nabla^2 U(\mathbf{x}_t)$ has
 288 a very small minimum eigenvalue. First, to prove the convergence to FOSP, we need to show that a
 289 single iteration of the proposed algorithm results in a descent direction for the objective function $U(\mathbf{x})$
 290 which is stated in Lemma 5.2, whose proof is provided in the Appendix A.

291 **Lemma 5.2.** *Under the Assumptions 1-2, and for a L_2 Hessian Lipschitz objective function $U(\mathbf{x})$,
 292 one iterate of Algorithm 4.2 is a descent direction, which implies that*

$$293 \quad (5.7) \quad U(\mathbf{x}_{t+1}) \leq U(\mathbf{x}_t) - \eta' \|(\tilde{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t)\|^2$$

295 where $\eta' := \eta C - \eta^2 \frac{L_1}{2}$ with $\eta \leq \frac{2C}{L_1}$.

296 The above mentioned analysis helps us to derive an upper bound on the gradient error norm. Consider
 297 the additional term in (4.7), we have

$$298 \quad (5.8) \quad \|\mathbf{e}_t\| = \|\nabla U(\mathbf{x}_t) + \mathbf{x}_t - \hat{\mathbf{x}}(\mathbf{x}_t)\| \leq \|\nabla U(\mathbf{x}_t)\| + \|\hat{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t\|.$$

300 From the upper bound in (A.3) and Lipschitz property of the objective function, we have

$$301 \quad (5.9) \quad \|\mathbf{e}_t\| \leq L_0 \left(1 + \frac{1}{C}\right) =: \mathcal{D}.$$

303 Before proceeding, let us define the following parameters as specified in [28] which are used for
 304 deriving the mathematical results next.

$$305 \quad F := \frac{\eta L_1}{L_2^2} \gamma^3 \cdot \log^{-3} \left(\frac{d\kappa}{\delta} \right), \quad G := \frac{\sqrt{\eta L_1}}{L_2} \gamma^2 \cdot \log^{-2} \left(\frac{d\kappa}{\delta} \right)$$

$$306 \quad (5.10) \quad L := \sqrt{\eta L_1} \frac{\gamma}{L_2} \cdot \log^{-1} \left(\frac{d\kappa}{\delta} \right), \quad \mathcal{T} := \frac{\log \left(\frac{d\kappa}{\delta} \right)}{\eta \gamma}$$

$$307$$

308 where γ is the negative eigenvalue and the condition number κ is given by $\kappa = \frac{L_1}{\gamma} \geq 1$. Once the
 309 algorithm iterate \mathbf{x}_t is near a FOSP, it holds that $\|\nabla U(\mathbf{x}_t)\| \leq g_{\text{th}}$ (gradient is small) and the minimum
 310 eigen value of Hessian is largely negative $\lambda_{\min}(\nabla^2 U(\mathbf{x}_t)) \leq -\sqrt{L_2 \epsilon}$. When \mathbf{x}_t is near a FOSP, then
 311 we add perturbation to \mathbf{x}_t followed by SCA updates for t_{th} steps. We prove that after these t_{th} steps, the
 312 function value will decrease by at least f_{th} with high probability. This result is formalized in Lemma
 313 [5.3](#).

314 **Lemma 5.3.** *Under the Assumption 1-2, for any $\delta \in (0, d\kappa/e]$, if for given $\tilde{\mathbf{x}}$, $\|\nabla U(\tilde{\mathbf{x}})\| \leq G$
 315 (gradient is small) and $\lambda_{\min}(\nabla^2 U(\tilde{\mathbf{x}})) \leq -\gamma$ (sufficiently negative), then if we perform $\mathbf{x}_0 = \tilde{\mathbf{x}} + \boldsymbol{\xi}$
 316 where $\boldsymbol{\xi} \in \mathbb{B}_{\tilde{\mathbf{x}}}(r)$ with radius $r = \frac{L}{\kappa \log(\frac{d\kappa}{\delta})}$, it holds that*

$$317 \quad (5.11) \quad U(\mathbf{x}_T) - U(\tilde{\mathbf{x}}) \leq -F + 16L_2 \eta^3 \mathcal{D}^3$$

319 for any $T \geq \frac{1}{c_{\max} \mathcal{T}}$ with probability $1 - \delta$ when the step size is selected as $\eta \leq \frac{C_{\max}}{L_1}$ and $\delta =$
 320 $\frac{dL_1}{\sqrt{L_2 \epsilon}} \exp(-\chi)$.

321 The result in Lemma 5.3 states that adding the perturbation decreases the function values further. By
 322 selecting $\eta = \frac{c}{L_1}$, $\gamma = \sqrt{L_2 \epsilon}$, and $\delta = \frac{dL_1}{\sqrt{L_2 \epsilon}} \exp(-\chi)$. We can restate the result in Lemma 5.3
 323 follows. If $\tilde{\mathbf{x}}_t$ is a FOSP, after adding perturbation $\mathbf{x}_t = \tilde{\mathbf{x}}_t + \boldsymbol{\xi}_t$, it holds with probability $1 - \delta$ that

$$324 \quad (5.12) \quad U(\mathbf{x}_{t+th}) - U(\tilde{\mathbf{x}}_t) \leq -f_{\text{th}} + \frac{16L_2 c^3 \mathcal{D}^3}{L_1^3}$$

$$325$$

326 with $t_{\text{th}} = \frac{\chi}{c^2} \frac{L_1}{\sqrt{L_2 \epsilon}}$. Next lemma states that addition of perturbation indeed results in escaping the
 327 saddle points.

328 **Lemma 5.4.** *Under the conditions similar to Lemma 5.3, let \mathbf{g}_1 is the minimum eigen vector of
 329 $\nabla^2 U(\tilde{\mathbf{x}})$. If we consider two sequences \mathbf{u}_t and \mathbf{w}_t such that*

$$330 \quad (5.13) \quad \|\mathbf{u}_0 - \tilde{\mathbf{x}}\| \leq r, \quad \mathbf{w}_0 = \mathbf{u}_0 + \mu r \mathbf{g}_1; \quad \mu \in [\delta/(2\sqrt{d}), 1],$$

332 then it holds that

$$333 \quad (5.14) \quad \min\{U(\mathbf{u}_T) - U(\mathbf{u}_0), U(\mathbf{w}_T) - U(\mathbf{w}_0)\} \leq -2.5F + 16L_2 \eta^3 \mathcal{D}^3$$

335 for any $\eta \leq \frac{c_{\max}}{L_1}$ and any $T \geq \frac{\mathcal{T}}{c_{\max}}$.

336 Lemma 5.4 states that for two points \mathbf{u}_0 and \mathbf{w}_0 , where \mathbf{w}_0 lies in the direction of the minimum
 337 eigenvector, one of both sequences \mathbf{u}_t and \mathbf{w}_t will result in a further decrement of the objective value

and hence escapes the saddle point. To prove the statement of Lemma 5.4, we need another two results stated next. In Lemma 5.5, we show that if the function value does not decrease after adding perturbation, then all the iterates belongs to a small ball around \mathbf{u}_0 . The next Lemma 5.6 shows that if the iterates starting from \mathbf{u}_0 get stuck, then the sequence \mathbf{w}_t will result in decreasing the function value and hence escapes the saddle point. Here, \mathbf{w}_t is the sequence obtained after applying SCA algorithm steps to \mathbf{w}_0 which by taking a step in the direction of minimum eigenvalue of the Hessian denoted by \mathbf{z}_1 and starting from \mathbf{u}_0 .

Lemma 5.5. *For any constant $\hat{c} \geq \frac{6}{\left(12 + \frac{2D}{\gamma L} \log \frac{d\kappa}{\delta}\right)^2}$, there exists a constant c_{\max} for any $\delta \in (0, \frac{d\kappa}{e}]$, $f(\cdot)$ satisfying the assumptions, and $\|f(\tilde{\mathbf{x}})\| \leq G$, for an initial point $\mathbf{u}_0 = \tilde{\mathbf{x}} + \zeta$ with $\|\mathbf{u}_0 - \tilde{\mathbf{x}}\| \leq 2r$, let*

$$(5.15) \quad T := \min\{\inf_t \{t | \tilde{f}_{\mathbf{u}_0}(\mathbf{u}_t) - f(\mathbf{u}_0) \leq -3F\}, \hat{c}F\}.$$

Then for any $\eta \leq c_{\max}/L_1$, it holds that for all $t < T$ we have $\|\mathbf{u}_t - \tilde{\mathbf{x}}\| \leq Z(L\hat{c})$, where $Z := 48 + \frac{8D}{\gamma L} \log\left(\frac{d\kappa}{\delta}\right)$ with $D = L_0(1 + \frac{1}{C})$.

Lemma 5.6. *There exists \hat{c} , c_{\max} for any $\delta \in (0, \frac{d\kappa}{e}]$, $f(\cdot)$ satisfying the assumptions, and $\|f(\tilde{\mathbf{x}})\| \leq G$, for an initial point $\mathbf{u}_0 = \tilde{\mathbf{x}} + \zeta$ and sequence $\{\mathbf{u}_t\}, \{\mathbf{w}_t\}$ where $\mathbf{w}_0 = \mathbf{u}_0 + \mu r \mathbf{z}_1$ with $\mu \in [\frac{\delta}{2\sqrt{2}}, 1]$, let us define*

$$(5.16) \quad T = \min\{\inf_t \{t | \tilde{f}_{\mathbf{w}_0}(\mathbf{w}_t) - f(\mathbf{w}_0) \leq -3F\}, \hat{c}F\}$$

then for any $\eta \leq c_{\max}/L_1$, if $\|\mathbf{u}_t - \tilde{\mathbf{x}}\| \leq Z(L\hat{c})$ for all $t < T$, it holds that $T < \hat{c}F$.

The proof of Lemmas 5.3 - Lemma 5.6 are provided in the Appendix B-E. These intermediate results are then used to prove Theorem 5.1; see Appendix F for the detailed proof.

6. Numerical Results. As a demonstration of the speedup afforded from SCA algorithms, we consider the multi-dimensional scaling (MDS) problem. Here, the goal is to embed relational data onto a low-dimensional euclidean space [18, 36, 4]. Given pairwise dissimilarities δ_{mn} where $(m, n) \in \mathcal{E}$ for all $\mathcal{E} \subset \{(m, n) | 1 \leq m < n \leq N\}$, between N objects, the embedding vectors are given by:

$$(6.1) \quad \{x_n\}_{n=1}^N = \arg \min_{\{\mathbf{x}_n\}_{n=1}^N} \sum_{1 \leq m < n \leq N} w_{mn}(\delta_{mn} - \|\mathbf{x}_m - \mathbf{x}_n\|)^2$$

where w_{mn} is the weight associated with the corresponding dissimilarity measurement, usually set to zero if it is not available. It can be seen that (6.1) is non-convex and even the global optimum is not unique, allowing rotational, translational, and reflectional ambiguity. Nevertheless, the objective in (6.1) can be expanded into a difference of convex terms, allowing us to apply the proposed SCA approach, where the surrogate is constructed by linearizing the concave term. The same approach has been widely used within the context of MDS and is referred to as the SMACOF algorithm. We consider a simple example with $N = 200$ and dissimilarities generated from calculating pairwise distances between N euclidean vectors in $[0, 1]^2$, keeping only 20% of the dissimilarities, and adding Gaussian noise to the distances with zero mean and variance 0.01. The non-zero weights are selected as inverse of the corresponding distances, which is the so-called Sammon mapping. The proposed

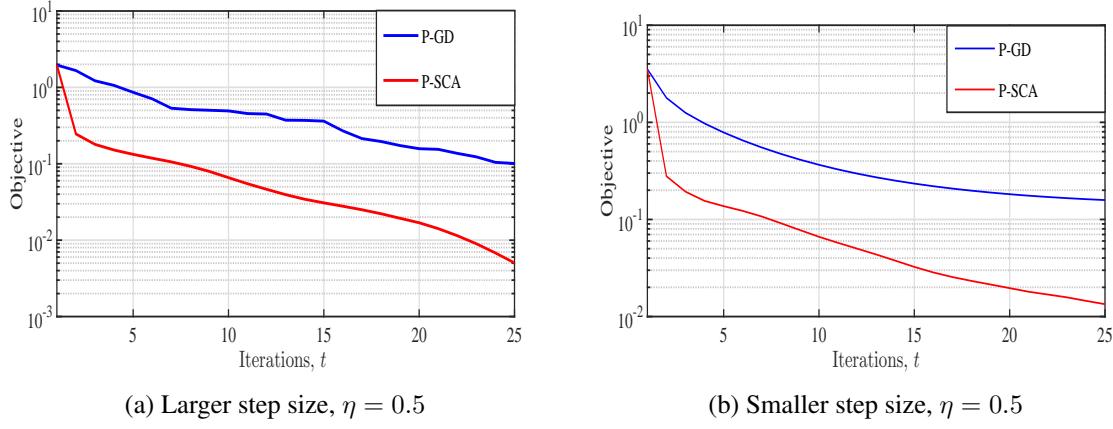


Figure 1: Performance comparisons for multidimensional scaling.

376 P-SCA as well as the P-GD algorithm have been applied to the MDS problem and the evolution of the
 377 objective function in 6.1, normalized by N^2 is shown in Fig. 1a and Fig. 1b for larger and smaller step
 378 sizes, respectively. It can be seen that the P-SCA is at least several times faster than P-GD algorithm,
 379 while retaining its saddle-point escaping properties. The performance of GD becomes worse for the
 380 smaller step sizes while SCA algorithm still performs better. It is remarked that for MDS problem, the
 381 per-iteration computational complexity of the two algorithms is exactly same.

382 **7. Conclusions and Future Directions.** This paper considers an algorithm for solving non-
 383 convex optimization formed by perturbation of successive convex approximation based method. The
 384 proposed perturbed successive convex approximation (P-SCA) algorithm is shown to converge to sec-
 385 ond order stationary point with a rate similar to the vanilla gradient descent (convergence to first order
 386 stationary point) up to a constant factor. Hence, the proposed algorithm can escape the saddle point
 387 efficiently utilizing the perturbed variant of the successive convex optimization algorithm. To perform
 388 the analysis, the proposed algorithm is proved to be a special case of the standard gradient descent
 389 algorithm with the gradient in error, which is called inexact gradient descent in literature. The idea of
 390 adding perturbation to the algorithms updates is utilized to escape the saddle points. As a future work,
 391 we are interested in considering the constrained version of the proposed algorithm which also escapes
 392 the saddle points.

393 **Appendix A. Proof of Lemma 5.2.** Note that $\hat{\mathbf{x}}(\mathbf{x}_t)$ is the solution of strongly convex problem
 394 in (4.4), hence from the first order optimality condition, it holds that

$$395 \quad (\mathbf{A}.1) \quad (\mathbf{y} - \hat{\mathbf{x}}(\mathbf{x}_t))^T \nabla \tilde{U}(\hat{\mathbf{x}}(\mathbf{x}_t); \mathbf{x}_t) \geq 0.$$

397 By selecting $\mathbf{y} = \mathbf{x}_t$ and add subtract $\nabla \tilde{U}(\mathbf{x}_t; \mathbf{x}_t)$ to the gradient term, we get

$$398 \quad (\mathbf{A}.2) \quad (\mathbf{x}_t - \hat{\mathbf{x}}(\mathbf{x}_t))^T \left(\nabla \tilde{U}(\hat{\mathbf{x}}(\mathbf{x}_t); \mathbf{x}_t) - \nabla \tilde{U}(\mathbf{x}_t; \mathbf{x}_t) + \nabla \tilde{U}(\mathbf{x}_t; \mathbf{x}_t) \right) \geq 0.$$

400 Following the inequality in Assumption 2, it holds that $\nabla U(\mathbf{x}_t) = \nabla \tilde{U}(\mathbf{x}_t; \mathbf{x}_t)$. This implies

$$401 \quad (\mathbf{A}.3) \quad (\mathbf{x}_t - \hat{\mathbf{x}}(\mathbf{x}_t))^T \nabla U(\mathbf{x}_t) \geq C \|\hat{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t\|^2.$$

403 From the statement of Assumption 1, it holds that the objective function $U(\mathbf{x})$ is Lipschitz continuous
 404 gradient with parameter L_1 , which implies that

$$405 \quad U(\mathbf{x}_{t+1}) \leq U(\mathbf{x}_t) + \nabla U(\mathbf{x}_t)^T (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L_1}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2.$$

$$406 \quad (A.4) \quad U(\mathbf{x}_{t+1}) \leq U(\mathbf{x}_t) + \eta \nabla U(\mathbf{x}_t)^T (\tilde{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t) + \frac{\eta^2 L_1}{2} \|\tilde{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t\|^2,$$

408 where (A.4) holds by utilizing the update for \mathbf{x}_{t+1} from Algorithm 4.2. Applying the upper bound in
 409 (A.3), we get

$$410 \quad (A.5) \quad U(\mathbf{x}_{t+1}) \leq U(\mathbf{x}_t) - \eta C \|\hat{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t\|^2 + \frac{L_1}{2} \|\eta(\tilde{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t)\|^2$$

$$411 \quad (A.6) \quad = U(\mathbf{x}_t) + \left(-\eta C + \frac{\eta^2 L_1}{2} \right) \|\tilde{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t\|^2 = U(\mathbf{x}_t) - \eta' \|\tilde{\mathbf{x}}(\mathbf{x}_t) - \mathbf{x}_t\|^2.$$

413 where $\eta' := \eta C - \frac{\eta^2 L_1}{2}$. Hence, one run of the Algorithm 4.2 results in a function value decrement
 414 with $\eta < \frac{2C}{L_1}$.

415 **Appendix B. Proof of Lemma 5.3.** From Lipschitz continuous gradient property

$$416 \quad (B.1) \quad U(\mathbf{x}_0) \leq U(\tilde{\mathbf{x}}) + \nabla U(\tilde{\mathbf{x}})^T (\mathbf{x}_0 - \tilde{\mathbf{x}}) + \frac{L_1}{2} \|\mathbf{x}_0 - \tilde{\mathbf{x}}\|^2.$$

418 From the statement of Lemma 5.3 and Cauchy-Schwarz inequality, it holds that

$$419 \quad (B.2) \quad U(\mathbf{x}_0) - U(\tilde{\mathbf{x}}) \leq \|\nabla U(\tilde{\mathbf{x}})\| \|\boldsymbol{\xi}\| + \frac{L_1}{2} \|\boldsymbol{\xi}\|^2$$

$$420 \quad \leq Gr + \frac{L_1}{2} r^2 \leq \frac{GL}{\kappa \log(\frac{d\kappa}{\delta})} + \frac{L_1}{2} \left(\frac{L}{\kappa \log(\frac{d\kappa}{\delta})} \right)^2 \leq \frac{3}{2} F.$$

422 The above mentioned bound represents the maximum value by which the function value can increase
 423 in the worst case after adding the perturbation. Next, following exactly the similar steps as performed
 424 in the proof of Lemma 14 in [28], we obtain

$$425 \quad U(\mathbf{x}_T) - U(\tilde{x}) = U(\mathbf{x}_T) - U(\mathbf{x}_0) + U(\mathbf{x}_0) - U(\tilde{x})$$

$$426 \quad (B.3) \quad \leq -2.5F + 16L_2\eta^3\mathcal{D}^3 + 1.5F \leq -F + 16L_2\eta^3\mathcal{D}^3.$$

428 **Appendix C. Proof of Lemma 5.4.** We consider $\tilde{\mathbf{x}} = 0$ without loss of generality, $T^* = \hat{c}\mathcal{T}$,
 429 and T' defined as $T' = \inf_t \{t | \tilde{U}_{\mathbf{u}_0}(\mathbf{u}_t) - U(\mathbf{u}_0) \leq -3F\}$. In order to prove the lemma statement, let
 430 us consider the two cases.

431 **Case I ($T' \leq T^*$):** From the statement of Lemma 5.5, we have $\|\mathbf{u}_{T'-1}\| \leq \mathcal{O}(L)$ which implies
 432 that

$$433 \quad \|\mathbf{u}_{T'}\| = \|\mathbf{u}_{T'-1}\| + \eta \|\nabla U(\mathbf{u}_{T'-1})\| + \eta \|\mathbf{e}_{T'-1}\|$$

$$434 \quad (C.1) \quad \leq \|\mathbf{u}_{T'-1}\| + \eta \|\nabla U(\tilde{\mathbf{x}})\| + \eta L_1 \|\mathbf{u}_{T'-1}\| + \eta \|\mathbf{e}_{T'-1}\| \leq (1 + \eta L_1) \hat{c}ZL + \eta G + \eta \mathcal{D}.$$

436 From the equality $\frac{G \log(\frac{d\kappa}{\delta})}{\gamma} = L$, we can write $\|\mathbf{u}_{T'}\| \leq Z_1 L + \eta \mathcal{D}$, where

437 $Z_1 := \left((1 + \eta L_1) \hat{c}Z + \eta \frac{\gamma}{\log(\frac{d\kappa}{\delta})} \right)$. By choosing a small c_{\max} and $\eta L_1 \leq c_{\max}$, for the Hessian
438 Lipschitz function, we have

$$439 \quad U(\mathbf{u}_{T'}) - U(\mathbf{u}_0) \leq \nabla U(\mathbf{u}_0)^T (\mathbf{u}_{T'} - \mathbf{u}_0) + \frac{1}{2} (\mathbf{u}_{T'} - \mathbf{u}_0)^T \nabla^2 U(\mathbf{u}_0) (\mathbf{u}_{T'} - \mathbf{u}_0) + \frac{L_2}{6} \|\mathbf{u}_{T'} - \mathbf{u}_0\|^3.$$

441 Consider the following quadratic approximation

$$442 \quad \tilde{U}_{\mathbf{y}}(\mathbf{x}) := U(\mathbf{y}) + \nabla U(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) + \frac{1}{2} (\mathbf{x} - \mathbf{y})^T \nabla^2 U(\tilde{\mathbf{x}}) (\mathbf{x} - \mathbf{y}).$$

444 For a Hessian Lipschitz function, it holds that

$$445 \quad U(\mathbf{u}_{T'}) - U(\mathbf{u}_0) \leq \nabla U(\mathbf{u}_0)^T (\mathbf{u}_{T'} - \mathbf{u}_0) + \frac{1}{2} (\mathbf{u}_{T'} - \mathbf{u}_0)^T \nabla^2 U(\mathbf{u}_0) (\mathbf{u}_{T'} - \mathbf{u}_0) + \frac{L_2}{6} \|\mathbf{u}_{T'} - \mathbf{u}_0\|^3$$

$$446 \quad \leq \nabla U(\mathbf{u}_0)^T (\mathbf{u}_{T'} - \mathbf{u}_0) + \frac{1}{2} (\mathbf{u}_{T'} - \mathbf{u}_0)^T \nabla^2 U(\tilde{\mathbf{x}}) (\mathbf{u}_{T'} - \mathbf{u}_0)$$

$$447 \quad + \frac{1}{2} (\mathbf{u}_{T'} - \mathbf{u}_0)^T (\nabla^2 U(\mathbf{u}_0) - \nabla^2 U(\tilde{\mathbf{x}})) (\mathbf{u}_{T'} - \mathbf{u}_0) + \frac{L_2}{6} \|\mathbf{u}_{T'} - \mathbf{u}_0\|^3$$

$$448 \quad (\text{C.2}) \quad \leq \tilde{U}_{\mathbf{u}_0}(\mathbf{u}_{T'}) - U(\mathbf{u}_0) + \frac{L_2}{2} \|\mathbf{u}_{T'} - \tilde{\mathbf{x}}\| \|\mathbf{u}_{T'} - \mathbf{u}_0\|^2$$

$$449 \quad + \frac{1}{2} (\mathbf{u}_{T'} - \mathbf{u}_0)^T \nabla^2 U(\mathbf{u}_0) (\mathbf{u}_{T'} - \mathbf{u}_0) + \frac{L_2}{6} \|\mathbf{u}_{T'} - \mathbf{u}_0\|^3$$

$$450 \quad \leq -3F + \mathcal{O}(L_2 L^3) + 16L_2 \eta^3 \mathcal{D}^3$$

$$451 \quad (\text{C.3}) \quad = -3F + \mathcal{O}(\sqrt{\eta L_1} F) + 16L_2 \eta^3 \mathcal{D}^3 \leq -2.5F + 16L_2 \eta^3 \mathcal{D}^3$$

453 by selecting $c_{\max} \leq \min\{1, \frac{1}{\hat{c}}\}$ and $\eta < \frac{1}{L_1}$. From Lemma 1, it holds that for any $T \geq \frac{1}{c_{\max}} \mathcal{T} \geq$
454 $\hat{c}\mathcal{T} = T^* \geq T'$, we can write

$$455 \quad U(\mathbf{u}_T) - U(\mathbf{u}_0) \leq U(\mathbf{u}_{T^*}) - U(\mathbf{u}_0) \leq U(\mathbf{u}_{T'}) - U(\mathbf{u}_0) \leq -2.5F + 16L_2 \eta^3 \mathcal{D}^3.$$

457 **Case II ($T' > T^*$):** For this case also, we know that $\|\mathbf{u}_t\| \leq \mathcal{O}(L) + \eta \mathcal{D}$ for all $t \leq T^*$. Let us
458 define T'' as $T'' = \inf_t \{t \mid \tilde{U}_{\mathbf{w}_0}(\mathbf{w}_t) - U(\mathbf{w}_0) \leq -3F\}$. From Lemma 5.6, we have $T'' \leq T^*$. Using
459 the similar arguments as in case 1, we conclude that for all $T \geq \frac{1}{c_{\max}} \mathcal{T}$, it holds that

$$460 \quad (\text{C.4}) \quad U(\mathbf{w}_T) - U(\mathbf{w}_0) \leq U(\mathbf{w}_{T^*}) - U(\mathbf{w}_0) \leq -2.5F + 16L_2 \eta^3 \mathcal{D}^3.$$

462 This concludes the proof.

463 **Appendix D. Proof of Lemma 5.5.** The idea here is to show that if function value does not
464 decrease after performing SCA iterations then it must be bounded in a small ball. This is performed
465 by analyzing the update dynamics of the proposed algorithm via decomposing d dimensional update
466 into two components: one along the subspace \mathcal{S} which is span of significantly negative eigenvalues
467 of the Hessian and second to the subspace compliment to it. Consider the second order Taylor series
468 approximation of the objective function $U(\mathbf{x})$ evaluated at $\tilde{\mathbf{x}}$ given by

$$469 \quad (\text{D.1}) \quad \tilde{U}_{\tilde{\mathbf{x}}}(\mathbf{x}) := U(\tilde{\mathbf{x}}) + \nabla U(\tilde{\mathbf{x}})^T (\mathbf{x} - \tilde{\mathbf{x}}) + \frac{1}{2} (\mathbf{x} - \tilde{\mathbf{x}})^T \mathcal{H} (\mathbf{x} - \tilde{\mathbf{x}})$$

471 where $\mathcal{H} = \nabla^2 U(\tilde{\mathbf{x}})$. It is emphasized that the gradient of U can be approximated by gradient of $\tilde{U}_{\tilde{\mathbf{x}}}$
 472 provided \mathbf{x} and $\tilde{\mathbf{x}}$ are close. This result is stated as [21]

$$473 \quad (D.2) \quad \left\| \nabla U(\mathbf{x}) - \nabla \tilde{U}_{\tilde{\mathbf{x}}}(\mathbf{x}) \right\| \leq \frac{L_2}{2} \|\mathbf{x} - \tilde{\mathbf{x}}\|^2$$

475 for a L_2 -Hessian Lipschitz function U . Set $\mathbf{u}_0 = 0$, from the SCA steps in (4.7), we can write

$$476 \quad (D.3) \quad \mathbf{u}_{t+1} = \mathbf{u}_t - \eta \nabla U(\mathbf{u}_t) + \eta \mathbf{e}_t$$

$$477 \quad (D.4) \quad = \mathbf{u}_t - \eta \nabla U(0) - \eta \left[\int_0^1 \nabla^2 U(\theta \mathbf{u}_t) d\theta \right] \mathbf{u}_t + \eta \mathbf{e}_t$$

$$478 \quad (D.5) \quad = \mathbf{u}_t - \eta \nabla U(0) - \eta (\mathcal{H} + \Delta_t) \mathbf{u}_t + \eta \mathbf{e}_t$$

$$479 \quad (D.6) \quad = (\mathbf{I} - \eta \mathcal{H} - \eta \Delta_t) \mathbf{u}_t - \eta (\nabla U(0) - \mathbf{e}_t)$$

481 where $\Delta_t := \int_0^1 \nabla^2 U(\theta \mathbf{u}_t) d\theta - \mathcal{H}$. Next, from the Hessian Lipschitz in (5.5), we have $\|\Delta_t\| \leq$
 482 $L_2(\|\mathbf{u}_t\| + \|\tilde{\mathbf{x}}\|)$. From the smoothness of the gradient in Assumption 1, we have

$$483 \quad (D.7) \quad \|\nabla U(0) - \mathbf{e}_t\| \leq \|\nabla U(0) - \nabla U(\tilde{\mathbf{x}}) + \nabla U(\tilde{\mathbf{x}}) + \mathbf{e}_t\|$$

$$484 \quad (D.8) \quad \leq \|\nabla U(\tilde{\mathbf{x}})\| + L_1 \|\tilde{\mathbf{x}}\| + \mathcal{D} \leq G + L_1 \cdot 2 \frac{L}{\kappa \cdot \log \left(\frac{d\kappa}{\delta} \right)} + \mathcal{D} \leq 3G + \mathcal{D}.$$

486 Next, we will calculate projections of \mathbf{u}_t on to the subspaces S and S_c , where S is the subspace of eigen
 487 vectors whose corresponding eigen values are less than $-\frac{\gamma}{\hat{c} \log \frac{d\kappa}{\delta}}$. Further, S_c defines the subspace of
 488 the remaining eigen vectors. Next, from the definition of T , we know that $\tilde{U}_{\mathbf{u}_0}(\mathbf{u}_t) - U(\mathbf{u}_0) > -3F$.
 489 Let $\mathbf{u}_0 = 0$, we get

$$490 \quad (D.9) \quad -3F < \tilde{U}_0(\mathbf{u}_t) - U(0) = \nabla U(0)^T \mathbf{u}_t + \frac{1}{2} \mathbf{u}_t^T \mathcal{H} \mathbf{u}_t$$

492 Now decomposing the vector \mathbf{u}_t into α_t (projection on to S) and β_t (projection on to S_c), we get

$$493 \quad (D.10) \quad \alpha_{t+1} = (\mathbf{I} - \eta \mathcal{H}) \alpha_t - \eta \mathcal{P}_S \Delta_t \mathbf{u}_t - \eta \mathcal{P}_S (\nabla U(0) - \mathbf{e}_t)$$

$$494 \quad (D.11) \quad \beta_{t+1} = (\mathbf{I} - \eta \mathcal{H}) \beta_t - \eta \mathcal{P}_{S^c} \Delta_t \mathbf{u}_t - \eta \mathcal{P}_{S^c} (\nabla U(0) - \mathbf{e}_t).$$

496 Utilizing these definitions, we get $-3F \leq \nabla U(0)^T \mathbf{u}_t - \frac{\gamma}{2} \frac{\|\alpha_t\|^2}{\hat{c} \log \frac{d\kappa}{\delta}} + \frac{1}{2} \beta_t^T \mathcal{H} \beta_t$, where the negative sign
 497 comes for the second term comes from the upper bound on the eigen values in subspace S . Note that
 498 $\|\mathbf{u}_t\|^2 = \|\alpha_t\|^2 + \|\beta_t\|^2$, we get

$$499 \quad (D.12) \quad -3F \leq \nabla U(0)^T \mathbf{u}_t - \frac{\gamma}{2} \frac{\|\mathbf{u}_t\|^2 - \|\beta_t\|^2}{\hat{c} \log \frac{d\kappa}{\delta}} + \frac{1}{2} \beta_t^T \mathcal{H} \beta_t.$$

501 After rearranging the terms, we get

$$502 \quad \|\mathbf{u}_t\|^2 \leq \frac{2\hat{c} \log \frac{d\kappa}{\delta}}{\gamma} \left(3F + \nabla U(0)^T \mathbf{u}_t + \frac{1}{2} \beta_t^T \mathcal{H} \beta_t \right) + \|\beta_t\|^2$$

$$503 \quad \leq 4 \max \left\{ \frac{6F \hat{c} \log \frac{d\kappa}{\delta}}{\gamma}, \frac{6G \hat{c} \log \frac{d\kappa}{\delta}}{\gamma} \|\mathbf{u}_t\|, \frac{\beta_t^T \mathcal{H} \beta_t \hat{c} \log \frac{d\kappa}{\delta}}{\gamma}, \|\beta_t\|^2 \right\}.$$

505 Taking square root on both sides results in

$$506 \quad \| \mathbf{u}_t \| \leq 4 \max \left\{ \sqrt{\frac{6F\hat{c} \log \frac{d\kappa}{\delta}}{\gamma}}, \frac{6G\hat{c} \log \frac{d\kappa}{\delta}}{\gamma}, \sqrt{\frac{\beta_t^T \mathcal{H} \beta_t \hat{c} \log \frac{d\kappa}{\delta}}{\gamma}}, \|\beta_t\| \right\}. \\ 507$$

508 In order to prove the lemma statement

$$509 \quad (\text{D.13}) \quad \| \mathbf{u}_t \| \leq \hat{c}LZ$$

511 we proceed by induction. It holds trivially for $t = 0$ since $\mathbf{u}_0 = 0$. Assume that (D.13) holds for
512 $\tau \leq t$, and then it remains to show that (D.13) holds for $t+1 < T$. It is sufficient to bound the last two
513 terms of (D) which are $\beta_t^T \mathcal{H} \beta_t$ and $\|\beta_t\|$. Consider the update equation for β_t from (D.11), we have

$$514 \quad (\text{D.14}) \quad \beta_t = (\mathbf{I} - \eta \mathcal{H}) \beta_t + \eta \delta_t$$

516 where $\|\delta_t\|$ is bounded as

$$517 \quad (\text{D.15}) \quad \|\delta_t\| \leq \|\Delta_t\| \|\mathbf{u}_t\| + \|\nabla U(0) - \mathbf{e}_t\|$$

$$518 \quad (\text{D.16}) \quad \leq L_2 (\|\mathbf{u}_t\| + \|\tilde{\mathbf{x}}\|) \|\mathbf{u}_t\| + \|\nabla U(0) - \mathbf{e}_t\|$$

$$519 \quad \leq L_2 \hat{c}Z \left(\hat{c}Z + \frac{2}{\kappa \cdot \log(\frac{d\kappa}{\delta})} \right) L^2 + 3G + \mathcal{D}$$

$$520 \quad (\text{D.17}) \quad \leq [\hat{c}Z(\hat{c}Z + 2)\sqrt{\eta L_1} + 3] G + \mathcal{D} \leq 4G + \mathcal{D}$$

522 where the last step follows by choosing $c_{\max} \leq \frac{1}{Z\hat{c}(Z\hat{c}+2)}$ and step size $\eta < c_{\max}/L$.

523 **Bounding $\|\beta_{t+1}\|$:** We have

$$524 \quad (\text{D.18}) \quad \|\beta_{t+1}\| \leq \left(1 + \frac{\eta\gamma}{\hat{c} \log(\frac{d\kappa}{\delta})} \right) \|\beta_t\| + \eta(4G + \mathcal{D})$$

526 Applying recursively, we get

$$527 \quad (\text{D.19}) \quad \|\beta_{t+1}\| \leq \sum_{\tau=0}^t (4G + \mathcal{D}) \left(1 + \frac{\eta\gamma}{\hat{c} \log(\frac{d\kappa}{\delta})} \right)^{\tau} \eta \leq 3(4G + \mathcal{D})\eta T \leq (12L + \eta\mathcal{D}\mathcal{T})\hat{c}.$$

529 Now, from the definition we have $T \leq \hat{c}F$, so that $\left(1 + \frac{\eta\gamma}{\hat{c} \log(\frac{d\kappa}{\delta})} \right)^T \leq 3$.

530 **Bounding $\beta_{t+1}^T \mathcal{H} \beta_{t+1}$:** From (D.14), we can write $\beta_t = \sum_{\tau=0}^{t-1} (\mathbf{I} - \eta \mathcal{H})^\tau \delta_\tau$, which implies that

$$531 \quad (\text{D.20}) \quad \beta_{t+1}^T \mathcal{H} \beta_{t+1} = \eta^2 \sum_{\tau_1=0}^t \sum_{\tau_2=0}^t \delta_{\tau_1}^T (\mathbf{I} - \eta \mathcal{H})^{\tau_1} \mathcal{H} (\mathbf{I} - \eta \mathcal{H})^{\tau_2} \delta_{\tau_2}$$

$$532 \quad (\text{D.21}) \quad \leq \eta^2 \sum_{\tau_1=0}^t \sum_{\tau_2=0}^t \|\delta_{\tau_1}\| \|(\mathbf{I} - \eta \mathcal{H})^{\tau_1} \mathcal{H} (\mathbf{I} - \eta \mathcal{H})^{\tau_2}\| \|\delta_{\tau_2}\|$$

$$533 \quad (\text{D.22}) \quad \leq (4G + \mathcal{D})^2 \eta^2 \sum_{\tau_1=0}^t \sum_{\tau_2=0}^t \|(\mathbf{I} - \eta \mathcal{H})^{\tau_1} \mathcal{H} (\mathbf{I} - \eta \mathcal{H})^{\tau_2}\|.$$

534

If $\{\lambda_i\}$ denotes eigen values of \mathcal{H} , then $\lambda_i(1 - \eta\lambda_i^{\tau_1+\tau_2})$ denotes the corresponding eigen value of $(\mathbf{I} - \eta\mathcal{H})^{\tau_1}\mathcal{H}(\mathbf{I} - \eta\mathcal{H})^{\tau_2}$. The maxima is achieved for $\lambda_{i,t}^* = \frac{1}{(1+t)\eta}$. Note that the function $\lambda_i(1 - \eta\lambda_i^{\tau_1+\tau_2})$ is monotonically increasing between $(-\infty, \lambda_{i,t}^*]$. Hence,

$$\begin{aligned} 538 \quad & \|(\mathbf{I} - \eta\mathcal{H})^{\tau_1}\mathcal{H}(\mathbf{I} - \eta\mathcal{H})^{\tau_2}\| = \max_i \lambda_i(1 - \eta\lambda_i)^{\tau_1+\tau_2} \\ 539 \quad & (\text{D.23}) \quad \leq \hat{\lambda}(1 - \eta\hat{\lambda})^{\tau_1+\tau_2} \leq \frac{1}{(1 + \tau_1 + \tau_2)\eta} \\ 540 \end{aligned}$$

541 where $\hat{\lambda} = \min\{L_1, \lambda_{\tau_1+\tau_2}^*\}$. Hence, we get

$$\begin{aligned} 542 \quad & \beta_{t+1}^T \mathcal{H} \beta_{t+1} \leq 2(4G + \mathcal{D})^2 \eta T \\ 543 \quad & (\text{D.24}) \quad \leq 4(4G^2 \eta \hat{c} \mathcal{T} + \eta \hat{c} \mathcal{T} \mathcal{D}^2) \leq 2(4\hat{c}L^2 \gamma \log^{-1}\left(\frac{d\kappa}{\delta}\right) + \eta \hat{c} \mathcal{T} \mathcal{D}^2) \\ 544 \end{aligned}$$

545 which follows from

$$\begin{aligned} 546 \quad & \sum_{\tau_1=0}^t \sum_{\tau_2=0}^t \frac{1}{(1 + \tau_1 + \tau_2)\eta} = \sum_{\tau=0}^{2t} \min\{1 + \tau, 2t + 1 - \tau\} \cdot \frac{1}{1 + \tau} \\ 547 \quad & (\text{D.25}) \quad \leq 2t + 1 < 2T. \\ 548 \end{aligned}$$

549 Finally, substituting the upper bounds in (D.19) and (D.24) into the right hand side of (D), and further
550 simplifying the bounds, we get $\|\mathbf{u}_t\| \leq ZL\hat{c}$ where $Z := 48 + \frac{8\mathcal{D}}{\gamma L} \log\left(\frac{d\kappa}{\delta}\right)$ is a constant and \hat{c} is
551 selected such that $\hat{c} \geq \frac{6}{\left(12 + \frac{2\mathcal{D}}{\gamma L} \log\frac{d\kappa}{\delta}\right)^2}$. This concludes the proof.

Appendix E. Proof of Lemma 5.6.

552 Consider the update equation for \mathbf{w}_t , we have

$$\begin{aligned} 553 \quad & (\text{E.1}) \quad \mathbf{u}_{t+1} + \mathbf{v}_{t+1} = \mathbf{w}_{t+1} - \eta \nabla U(\mathbf{w}_t) + \eta \mathbf{e}_t = \mathbf{u}_t + \mathbf{v}_t - \eta \nabla U(\mathbf{u}_t + \mathbf{v}_t) \\ 554 \quad & = \mathbf{u}_t + \mathbf{v}_t - \eta \nabla U(\mathbf{u}_t) - \eta \left[\int_0^1 \nabla^2 U(\mathbf{u}_t + \theta \mathbf{v}_t) d\theta \right] \mathbf{v}_t + \eta \mathbf{e}_t \\ 555 \quad & = \mathbf{u}_t + \mathbf{v}_t - \eta \nabla U(\mathbf{u}_t) - \eta (\mathcal{H} + \Delta'_t) \mathbf{v}_t + \eta \mathbf{e}_t \\ 556 \quad & (\text{E.2}) \quad = \mathbf{u}_t + \mathbf{v}_t - \eta \nabla U(\mathbf{u}_t) - \eta (\mathcal{H} + \Delta'_t) \mathbf{v}_t + \eta \mathbf{e}_t \\ 557 \quad & (\text{E.3}) \quad = \mathbf{u}_t - \eta \nabla U(\mathbf{u}_t) + \eta \mathbf{e}_t + \eta (\mathbf{I} - \eta \mathcal{H} - \Delta'_t) \mathbf{v}_t. \\ 558 \end{aligned}$$

559 From now onwards, the proof is exactly similar to that of Lemma 17 in [28] but provided here for
560 completeness. As in the earlier proofs, , we have $\|\Delta'_t\| \leq L_2(\|\mathbf{u}_t\| + \|\mathbf{v}_t\| + \|\tilde{\mathbf{x}}\|)$. Now, we can write
561 the update for \mathbf{v}_{t+1} as

$$\begin{aligned} 562 \quad & (\text{E.4}) \quad \mathbf{v}_{t+1} = (\mathbf{I} - \eta \mathcal{H} - \Delta'_t) \mathbf{v}_t \\ 563 \end{aligned}$$

564 Note that we have $\|\mathbf{w}_0 - \tilde{\mathbf{x}}\| \leq \|\mathbf{u}_0 - \tilde{\mathbf{x}}\| + \|\mathbf{v}_0\| \leq \frac{L}{\kappa \cdot \log\left(\frac{d\kappa}{\delta}\right)}$, from Lemma (5.5), we have $\|\mathbf{w}_t\| \leq$
565 $ZL\hat{c}$ for all $t \leq T$. From the condition in Lemma 5.6, we note that $\|\mathbf{u}_t\| \leq ZL\hat{c}$ for all $t < T$. This
566 implies that

$$\begin{aligned} 567 \quad & (\text{E.5}) \quad \|\mathbf{v}_t\| \leq \|\mathbf{u}_t\| + \|\mathbf{w}_t\| \leq 2ZL\hat{c} \\ 568 \end{aligned}$$

569 for all $t < T$. From here, we can derive that

570 (E.6)
$$\|\triangle'_t\| \leq L_2(\|\mathbf{u}_t\| + \|\mathbf{v}_t\| + \|\tilde{\mathbf{x}}\|) \leq L_2 \left(2ZL\hat{c} + \frac{L}{\kappa \log \frac{d\kappa}{\delta}} \right) \leq L_2 L (2Z\hat{c} + 1).$$

571

572 Next, we develop a lower bound on $\|\mathbf{v}_t\|$. Let ψ_t be the norm of the projection of \mathbf{v}_t onto \mathbf{z}_1 direction
573 and ϕ_t be the norm of the projection of \mathbf{v}_t on to the remaining subspace. From (E.4), we can write

574 (E.7)
$$\psi_t \geq (1 + \gamma\eta)\psi_t - \mu\sqrt{\psi_t^2 + \phi_t^2}$$

575 (E.8)
$$\phi_t \leq (1 + \gamma\eta)\Psi_t + \mu\sqrt{\psi_t^2 + \phi_t^2}$$

577 where $\mu = \eta L_2 L (2Z\hat{c} + 1)$. By induction, next we prove that the following inequality holds for all
578 $t < T$

579 (E.9)
$$\phi_t \leq 4\mu t \psi_t.$$

581 For $t = 1$ the base case of the induction holds for $t = 0$ from the definition of \mathbf{v}_0 . Now, let us assume
582 that (E.9) holds for $t \leq T$, we need to show that it holds for $t + 1 \leq T$, we have

583 (E.10)
$$4\mu(t+1)\psi_t \geq 4\mu(t+1) \left((1 + \gamma\eta)\psi_t - \mu\sqrt{\psi_t^2 + \phi_t^2} \right)$$

584 (E.11)
$$\phi_{t+1} \leq 4\mu t(1 + \gamma\eta)\psi_t + \mu\sqrt{\psi_t^2 + \phi_t^2}.$$

586 Next, we only need to show that

587 (E.12)
$$(1 + 4\mu(t+1))\sqrt{\psi_t^2 + \phi_t^2} \leq 4(1 + \gamma\eta)\psi_t.$$

589 By selecting $\sqrt{c_{\max}} \leq \frac{1}{2Z\hat{c}+1} \min\{\frac{1}{2\sqrt{2}}, \frac{1}{4\hat{c}}\}$ and $\eta \leq \frac{c_{\max}}{L_1}$, we get

590 (E.13)
$$4\mu(t+1) \leq 4\mu T \leq 4\eta L_2 L (2Z\hat{c} + 1) \hat{c} T = 4\sqrt{\eta L_1} (2Z\hat{c} + 1) \hat{c} \leq 1.$$

592 This implies that

593 (E.14)
$$4(1 + \gamma\eta)\psi_t \geq 4\psi_t \leq 2\sqrt{2\phi_t^2} \geq (1 + 4\mu(t+1))\sqrt{\psi_t^2 + \phi_t^2}$$

595 which proves the induction. From (E.9), we have $\phi_t \leq 4\mu t \psi_t \leq \psi_t$, which gives

596 (E.15)
$$\psi_{t+1} \geq (1 + \gamma\eta)\psi_t - \sqrt{2}\mu\psi_t \geq \left(1 + \frac{\gamma\eta}{2}\right) \psi_t$$

598 where the last inequality follows from $\mu = \eta L_2 L (2Z\hat{c} + 1) \leq \sqrt{c_{\max}} (2Z\hat{c} + 1) \gamma\eta \log^{-1} \frac{d\kappa}{\delta} < \frac{\gamma\eta}{2\sqrt{2}}$.

599 From (E.5) and (E.15), we obtain for all $t < T$

600
$$2ZL\hat{c} \geq \|\mathbf{v}_t\| \geq \psi_t \geq \left(1 + \frac{\gamma\eta}{2}\right)^t \psi_0$$

601 (E.16)
$$= \left(1 + \frac{\gamma\eta}{2}\right)^t c_0 \frac{L}{\kappa} \log^{-1} \frac{d\kappa}{\delta} \geq \left(1 + \frac{\gamma\eta}{2}\right)^t \frac{\delta}{2\sqrt{d}} \frac{L}{\kappa} \log^{-1} \left(\frac{d\kappa}{\delta}\right).$$

602

603 The above result implies that

$$604 \quad (E.17) \quad T \leq \frac{1}{2} \frac{\log \left(4Z \frac{\kappa\sqrt{d}}{\delta} \cdot \hat{c} \log \left(\frac{d\kappa}{\delta} \right) \right)}{\log \left(1 + \frac{\gamma\eta}{2} \right)} \leq \frac{1}{2} \frac{\log \left(4Z \frac{\kappa\sqrt{d}}{\delta} \cdot \hat{c} \log \left(\frac{d\kappa}{\delta} \right) \right)}{\gamma\eta} \leq (2 + \log(4Z\hat{c}))\mathcal{T}. \\ 605$$

606 The inequality follows from the selection $\delta \in (0, \frac{d\kappa}{\delta}]$ and we have $\log \left(\frac{d\kappa}{\delta} \right) \geq 1$. We choose \hat{c} large
607 enough such that $2 + \log(4Z\hat{c}) \leq \hat{c}$, we get $T < \hat{c}\mathcal{T}$, which concludes the proof.

608 **Appendix F. Proof of Theorem 5.1.** Let $c < \min\{c_{\max}, 2C\}$, the goal of this proof is to
609 achieve a point \mathbf{x} for which it holds that

$$610 \quad (F.1) \quad \|\nabla U(\mathbf{x})\| \leq g_{\text{th}} = \frac{\epsilon\sqrt{c}}{\chi^2}, \quad \lambda_{\min}(\nabla^2 U(\mathbf{x})) \geq -\sqrt{\epsilon L_2}. \\ 611$$

612 Note that $c \leq 1$, $\chi \geq 1$ as defined in Algorithm 4.2, which implies that $\frac{\sqrt{c}}{\chi^2} \leq 1$ and hence any \mathbf{x}
613 satisfying (F.1) is an ϵ second order stationary point. To proceed with the proof, let us start from \mathbf{x}_0 , if
614 \mathbf{x}_0 does not satisfy (F.1), then there are two possible outcomes

- 615 1. $\|\nabla U(\mathbf{x}_0)\| > g_{\text{th}}$. This means that the first order stationary point is not reached yet and hence
616 Algorithm 2 will not perform the perturbation step. From Lemma 5.2, we have

$$617 \quad (F.2) \quad U(\mathbf{x}_1) - U(\mathbf{x}_0) \leq -\frac{\eta}{2}g_{\text{th}}^2 = -\frac{c^2\epsilon^2}{2\chi^4 L_1}. \\ 618$$

- 619 2. $\|\nabla U(\mathbf{x}_0)\| \leq g_{\text{th}}$. Under this condition, a perturbation step is performed by Algorithm 2, then
620 SCA steps are performed for the next t_{th} iterations, and then termination condition is checked.
621 If termination condition is not satisfied, it holds that

$$622 \quad (F.3) \quad U(\mathbf{x}_{\text{th}}) - U(\mathbf{x}_0) \leq -f_{\text{th}} + \frac{16L_2c^3\mathcal{D}^3}{L_1^3}. \\ 623$$

624 where we have substituted step size $\eta = \frac{c}{L_1}$. Now we select $c^3 = \frac{f_{\text{th}}}{16L_2\mathcal{D}^3}$ with $s \in (0, 1)$ and introduce
625 s as an adjustment parameter to make $c < \min\{c_{\max}, 2C\}$. Note that we cannot make s arbitrarily
626 small because that would result in a smaller step size η which is proportional to c . Next, we obtain

$$627 \quad (F.4) \quad U(\mathbf{x}_{\text{th}}) - U(\mathbf{x}_0) \leq -(1-s)f_{\text{th}}.$$

629 The above results implies that on an average with each step of the algorithm, the function value
630 reduces by

$$631 \quad (F.5) \quad \frac{U(\mathbf{x}_{\text{th}}) - U(\mathbf{x}_0)}{t_{\text{th}}} \leq -\frac{c^3\epsilon^2}{\chi^4 L_1}.$$

633 The function values decreases at least by $\frac{c^3\epsilon^2}{\chi^4 L_1}$ on an average after running for t_{th} iterations. Since
634 the maximum amount by which the function value may decrease is upper bound by $U(\mathbf{x}_0) - U^*$, the
635 algorithm must terminate in the following number of iterations

$$636 \quad (F.6) \quad \frac{U(\mathbf{x}_0) - U^*}{\frac{c^3\epsilon^2}{\chi^4 L_1}} = \frac{\chi^4}{c^3} \cdot \frac{L_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2} = \mathcal{O} \left(\frac{L_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2} \log \left(\frac{dL_1\Delta_U}{\epsilon^2\delta} \right) \right). \\ 637$$

638 From Algorithm 4.2, note that the perturbation is added at iteration t if the gradient is small or
 639 $\nabla U(\tilde{x}_t) \leq g_{\text{th}}$. From Lemma 5.3, the probability of happening this at each time is at least $1 - \delta$
 640 where $\delta = \frac{dL_1}{\sqrt{L_2}\epsilon} \exp(-\chi)$. In other words, the number of times perturbation is added for one run of
 641 the algorithm is given by

$$642 \quad (F.7) \quad \frac{1}{t_{\text{th}}} \cdot \frac{\chi^4}{c^3} \cdot \frac{L_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2} = \frac{\chi^3}{c} \cdot \frac{L_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2}.$$

644 Using the union bound, it holds that Lemma 5.3 holds for each of the iterations after adding perturba-
 645 tions. This result makes sure that the Algorithm 4.2 converges to the ϵ second order stationary point
 646 with probability

$$647 \quad (F.8) \quad 1 - \frac{dL_1}{\sqrt{L_2}\epsilon} \exp(-\chi) \cdot \frac{\chi^3}{c} \frac{L_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2} = 1 - \frac{\chi^3 \exp(-\chi)}{c} \cdot \frac{dL_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2}.$$

649 Note that we choose $\chi = 3 \max\{\log\left(\frac{dL_1 \Delta_U}{c\epsilon^2}\right), 4\}$, which implies that $\chi \geq 12$, and $\chi^3 \exp(-\chi) \leq$
 650 $\exp^{-\chi/3}$. Therefore, we have

$$651 \quad (F.9) \quad \frac{\chi^3 \exp(-\chi)}{c} \cdot \frac{dL_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2} \leq \frac{\exp^{-\chi/3}}{c} \cdot \frac{dL_1(U(\mathbf{x}_0) - U^*)}{\epsilon^2} \leq \delta.$$

653 which completes the proof.

654

REFERENCES

- 655 [1] A. BECK, A. BEN-TAL, AND L. TETRUASHVILI, *A sequential parametric convex approximation method with appli-*
 656 *cations to nonconvex truss topology design problems*, J. Global Opt., 47 (2010), pp. 29–51.
- 657 [2] A. S. BEDI, P. SARMA, AND K. RAJAWAT, *Tracking moving agents via inexact online gradient descent algorithm*,
 658 IEEE J. Sel. Topics Signal Process., 12 (2018), pp. 202–217.
- 659 [3] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Gradient convergence in gradient methods with errors*, SIAM J. Opt., 10
 660 (2000), pp. 627–642.
- 661 [4] I. BORG AND P. GROENEN, *Modern multidimensional scaling: Theory and applications*, Journal of Educational
 662 Measurement, 40 (2003), pp. 277–280.
- 663 [5] E. J. CANDÁS, X. LI, AND M. SOLTANOLKOTABI, *Phase retrieval via wirtinger flow: Theory and algorithms*, IEEE
 664 Trans. Inf. Theory, 61 (2015), pp. 1985–2007.
- 665 [6] F. E. CURTIS, D. P. ROBINSON, AND M. SAMADI, *A trust region algorithm with a worst-case iteration complexity*
 666 *of $\mathcal{O}(\epsilon^{-3/2})$ for nonconvex optimization*, Math. Prog., 162 (2017), pp. 1–32.
- 667 [7] H. DANESHMAND, J. KOHLER, A. LUCCHI, AND T. HOFMANN, *Escaping saddles with stochastic gradients*, in
 668 International Conference on Machine Learning, 2018, pp. 1163–1172.
- 669 [8] R. DIXIT, A. S. BEDI, R. TRIPATHI, AND K. RAJAWAT, *Online learning with inexact proximal online gradient*
 670 *descent algorithms*, IEEE Trans. Signal Process., 67 (2019), pp. 1338–1352.
- 671 [9] R. GE, F. HUANG, C. JIN, AND Y. YUAN, *Escaping from saddle points — online stochastic gradient for tensor*
 672 *decomposition*, in Proceedings of The 28th Conference on Learning Theory, 2015, pp. 797–842.
- 673 [10] S. B. GELFAND AND S. K. MITTER, *Recursive stochastic algorithms for global optimization in r^d* , SIAM J. Control
 674 Opt., 29 (1991), pp. 999–1018.
- 675 [11] S. S. HAYKIN ET AL., *Neural networks and learning machines/Simon Haykin.*, New York: Prentice Hall., 2009.
- 676 [12] P. JAIN, P. KAR, ET AL., *Non-convex optimization for machine learning*, Foundations and Trends® in Machine
 677 Learning, 10 (2017), pp. 142–336.
- 678 [13] C. JIN, R. GE, P. NETRAPALLI, S. M. KAKADE, AND M. I. JORDAN, *How to escape saddle points efficiently*, in
 679 Proc. ICML, 2017, pp. 1724–1732.

- 680 [14] K. Y. LEVY, *The power of normalization: Faster evasion of saddle points*, arXiv preprint arXiv:1611.04831, (2016).
- 681 [15] A. LIU, V. LAU, AND B. KANANIAN, *Stochastic successive convex approximation for non-convex constrained sto-*
682 *chastic optimization*, arXiv preprint arXiv:1801.08266, (2018).
- 683 [16] X.-Y. LIU, S. AERON, V. AGGARWAL, AND X. WANG, *Low-tubal-rank tensor completion using alternating mini-*
684 *mization*, arXiv preprint arXiv:1610.01690, (2016).
- 685 [17] W.-K. K. MA, *Semidefinite relaxation of quadratic optimization problems and applications*, IEEE Signal Processing
686 Magazine, 1053 (2010).
- 687 [18] L. V. D. MAATEN AND G. HINTON, *Visualizing data using t-sne*, Journal of machine learning research, 9 (2008),
688 pp. 2579–2605.
- 689 [19] J. MAIRAL, *Stochastic majorization-minimization algorithms for large-scale optimization*, in Adv. Neural Inf.
690 Process. Syst., 2013, pp. 2283–2291.
- 691 [20] Y. NESTEROV, *Introductory lectures on convex programming volume i: Basic course*, Lecture notes, (1998).
- 692 [21] Y. NESTEROV, *Introductory lectures on convex optimization: A basic course*, vol. 87, Springer Science & Business
693 Media, 2013.
- 694 [22] Y. NESTEROV AND B. POLYAK, *Cubic regularization of newton method and its global performance*, Mathematical
695 Programming, 108 (2006), pp. 177–205.
- 696 [23] P. NETRAPALLI, P. JAIN, AND S. SANGHAVI, *Phase retrieval using alternating minimization*, in Adv. Neural Inf.
697 Process. Sys., 2013, pp. 2796–2804.
- 698 [24] R. PEMANTLE ET AL., *Nonconvergence to unstable points in urn models and stochastic approximations*, Annals of
699 Prob., 18 (1990), pp. 698–712.
- 700 [25] M. RAGINSKY, A. RAKHLIN, AND M. TELGARSKY, *Non-convex learning via stochastic gradient langevin dynam-*
701 *ics: a nonasymptotic analysis*, in Proc. of COLT, S. Kale and O. Shamir, eds., vol. 65 of Proceedings of Machine
702 Learning Research, Amsterdam, Netherlands, 07–10 Jul 2017, PMLR, pp. 1674–1703.
- 703 [26] M. RAZAVIYAYN, M. HONG, AND Z.-Q. LUO, *A unified convergence analysis of block successive minimization*
704 *methods for nonsmooth optimization*, SIAM J. Opt., 23 (2013), pp. 1126–1153.
- 705 [27] M. SCHMIDT, N. L. ROUX, AND F. R. BACH, *Convergence rates of inexact proximal-gradient methods for convex*
706 *optimization*, in Advances in neural information processing systems, 2011, pp. 1458–1466.
- 707 [28] G. SCUTARI, F. FACCHINEI, AND L. LAMPARIELLO, *Parallel and distributed methods for constrained nonconvex*
708 *optimization?part i: Theory*, IEEE Trans. Signal Process., 65 (2017), pp. 1929–1944.
- 709 [29] G. SCUTARI, F. FACCHINEI, P. SONG, D. P. PALOMAR, AND J. PANG, *Decomposition by partial linearization:*
710 *Parallel optimization of multi-agent systems*, IEEE Trans. Signal Process., 62 (2014), pp. 641–656.
- 711 [30] G. SCUTARI, F. FACCHINEI, P. SONG, D. P. PALOMAR, AND J.-S. PANG, *Decomposition by partial linearization:*
712 *Parallel optimization of multi-agent systems*, IEEE Trans. Signal Process., 62 (2013), pp. 641–656.
- 713 [31] A. M.-C. SO AND Z. ZHOU, *Non-asymptotic convergence analysis of inexact gradient methods for machine learning*
714 *without strong convexity*, Opt. Methods Soft., 32 (2017), pp. 963–992.
- 715 [32] A. M.-C. SO AND Z. ZHOU, *Non-asymptotic convergence analysis of inexact gradient methods for machine learning*
716 *without strong convexity*, Opt. Methods Software, 32 (2017), pp. 963–992.
- 717 [33] J. SUN, Q. QU, AND J. WRIGHT, *A geometric analysis of phase retrieval*, Foundations of Computational Mathematics, 18 (2018), pp. 1131–1198.
- 718 [34] R. SUN AND Z. LUO, *Guaranteed matrix completion via non-convex factorization*, IEEE Trans. Inf. Theory, 62
719 (2016), pp. 6535–6579.
- 720 [35] Y. SUN, P. BABU, AND D. P. PALOMAR, *Majorization-minimization algorithms in signal processing, communica-*
721 *tions, and machine learning*, IEEE Trans. Signal Process., 65 (2017), pp. 794–816.
- 722 [36] J. TZENG, H. H.-S. LU, AND W.-H. LI, *Multidimensional scaling for large genomic data sets*, BMC bioinformatics,
723 9 (2008), p. 179.
- 724 [37] W. WANG, V. AGGARWAL, AND S. AERON, *Efficient low rank tensor ring completion*, in Proc. IEEE Int. Conf.
725 Computer Vision, 2017, pp. 5697–5705.
- 726 [38] W. WANG, V. AGGARWAL, AND S. AERON, *Tensor train neighborhood preserving embedding*, IEEE Trans. Signal
727 Process., 66 (2018), pp. 2724–2732.
- 728 [39] X. WANG, S. MA, D. GOLDFARB, AND W. LIU, *Stochastic quasi-newton methods for nonconvex stochastic opti-*
729 *mization*, SIAM J. Opt., 27 (2017), pp. 927–956.
- 730 [40] L. XIAO AND T. ZHANG, *A proximal stochastic gradient method with progressive variance reduction*, SIAM J. Opt.,
731 24 (2014), pp. 2057–2075.