

# A Momentum Accelerated Adaptive Cubic Regularization Method for Nonconvex Optimization

Yihang Gao\*

Michael K. Ng\*

## Abstract

The cubic regularization method (CR) and its adaptive version (ARC) are popular Newton-type methods in solving unconstrained non-convex optimization problems, due to its global convergence to local minima under mild conditions. The main aim of this paper is to develop a momentum accelerated adaptive cubic regularization method (ARCm) to improve the convergent performance. With the proper choice of momentum step size, we show the global convergence of ARCm and the local convergence can also be guaranteed under the KL property. Such global and local convergence can also be established when inexact solvers with low computational costs are employed in the iteration procedure. Numerical results for non-convex logistic regression and robust linear regression models are reported to demonstrate that the proposed ARCm significantly outperforms state-of-the-art cubic regularization methods (e.g., CR, momentum-based CR, ARC) and the trust region method. In particular, the number of iterations required by ARCm is less than 10% to 50% required by the most competitive method (ARC) in the experiments.

**Keywords.** Adaptive cubic regularization, momentum, KL property, inexact solvers, global and local convergence

## 1 Introduction

Most machine learning tasks involve solving challenging (non-convex) optimization problems

$$(1.1) \quad \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}).$$

Second-order critical points are usually the preferred solutions for (1.1) since many machine learning problems have been shown to have only strict saddle points and global minima without spurious local minima [8, 16]. Second-order methods that exploit Hessian information have been proposed to escape saddle points [15, 5] and enjoy faster local convergence [14, 19, 21]. Cubic regularization method (CR), first proposed by Griewank [9], and later independently by Nesterov and Polyak [15], and Weiser et al. [18], one of the second-order methods,

is widely applied in solving inverse problems [4, 12], regression models [19, 20] as well as minimax problems [11].

The vanilla CR method is formulated as

$$\begin{aligned} \mathbf{s}_k &\in \arg \min_{\mathbf{s}} \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s} + \frac{M_k}{6} \|\mathbf{s}\|_2^3, \\ \mathbf{x}_{k+1} &= \mathbf{x}_k + \mathbf{s}_k, \end{aligned}$$

where  $M_k := M$  is a fixed pre-defined constant for all iterations. In recent decades, various techniques are adopted to improve the performance of CR. Wang et al. [17] accelerated the vanilla CR by a momentum term (CRm), where the method works well mainly by enlarging the step size of  $\mathbf{s}_k$  when the cubic penalty parameter  $M$  is over-estimated.

Cartis et al. [4] proposed adaptive cubic regularization method (ARC), which adaptively assigns  $M_k$  based on the quality of the step  $\mathbf{s}_k$ , as an analogy to the trust region method (TR) [5]. Stochastic (subsampling) ARC (e.g., Kohler et al. [13] and Zhou et al. [20] etc.) were developed to reduce the overall computation, where the gradient  $\nabla f(\mathbf{x}_k)$  and the Hessian  $\nabla^2 f(\mathbf{x}_k)$  are inexactly evaluated. To the best of our knowledge, ARC behaves better in most of the applications compared with CR methods that fix  $M$  (e.g., vanilla CR and CRm). A natural question is how to further accelerate ARC with little extra cost in each iteration.

In this paper, we investigate a momentum-accelerated adaptive cubic regularization method (ARCm). Here are our contributions and the outline of the paper.

- We develop a momentum-accelerated adaptive cubic regularization method (ARCm). We adopt a general scheme for momentum, which is more suitable for ARC than CR (see in Theorem 2.1). The extra computation is cheap and ignorant as both the momentum and its step size are solely based on  $\{\mathbf{s}_k\}$  and are free of gradient or Hessian evaluations.
- With the proper setting of step size for momentum (see in the Algorithm 1), the global convergence of ARCm to second-order critical points is satisfied

\*Department of Mathematics, The University of Hong Kong, Pokfulam, Hong Kong SAR.

---

**Algorithm 1** Adaptive cubic regularization with momentum (ARCM)

---

**Input:**  $\mathbf{x}_0, \mathbf{v}_{-1} = \mathbf{0}, \gamma_1 > 1 \geq \gamma_2 > \gamma_3 > 0, 1 > \eta_2 > \eta_1 > 0, \sigma_0 > \sigma_{\min} > 0, 1 > \tau, \beta > 0$  and  $\alpha_1, \alpha_2 > 0$  in

**Output:**  $\{\mathbf{x}_k\}_{k=0}^T$  out

1: **for**  $k = 0$  to  $T - 1$  **do**

2: Solve the cubic subproblem:

$$\mathbf{s}_k \in \arg \min_{\mathbf{s}} m_k(\mathbf{s}) := \arg \min_{\mathbf{s}} f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s} + \frac{\sigma_k}{6} \|\mathbf{s}\|_2^3.$$

3: Compute  $f(\mathbf{x}_k + \mathbf{s}_k)$  and  $\rho_k = \frac{f(\mathbf{x}_k) - f(\mathbf{x}_k + \mathbf{s}_k)}{f(\mathbf{x}_k) - m_k(\mathbf{s}_k)}$ .

4: **if**  $\rho_k > \eta_1$  (successful update) **then**

5:  $\mathbf{y}_{k+1} = \mathbf{x}_k + \mathbf{s}_k$

6: **Momentum step:**

Select  $\beta_k \in [0, \min(\tau, \alpha_1 \|\mathbf{s}_k\|_2, \alpha_2 \|\mathbf{s}_k\|_2^2)]$  such that  $f(\mathbf{z}_{k+1}) \leq f(\mathbf{y}_{k+1})$  with  $\mathbf{v}_k = \beta_k \cdot \mathbf{v}_{k-1} + \mathbf{s}_k$  and  $\mathbf{z}_{k+1} = \mathbf{x}_k + \mathbf{v}_k$ .

7:  $\mathbf{x}_{k+1} = \mathbf{z}_{k+1}$

8: **if**  $\rho_k > \eta_2$  (very successful update) **then**

9:  $\sigma_{k+1} = \max(\sigma_{\min}, \gamma_3 \cdot \sigma_k)$

10: **else**

11:  $\sigma_{k+1} = \gamma_2 \cdot \sigma_k$

12: **end if**

13: **else**

14: (unsuccessful update)

$\mathbf{x}_{k+1} = \mathbf{x}_k$  and  $\mathbf{v}_k = \mathbf{v}_{k-1}$

15:  $\sigma_{k+1} = \gamma_1 \cdot \sigma_k$

16: **end if**

17: **end for**

---

(see in Theorem 2.2). We further show that the proposed ARCM enjoys local convergence under the KL property (see in Theorem 2.3), which is one of the advantages of second-order methods over first-order methods.

- We also study the global convergence of ARCM with the inexact cubic regularized subproblems (CRS) solutions as we usually approximately solve CRS in practice (see in Theorem 3.1). The local convergence is preserved if the error of CRS is decreasing with  $\|\mathbf{s}_k\|_2^3$  (see in Corollary 3.1) but may be destructed otherwise.
- We conduct experiments in solving high-dimensional and large-scale non-convex logistic regression and robust linear regression problems. Experimental results show that ARCM significantly outperforms CR, CRm, ARC and TR, where it is 10%-50% faster than ARC (the most competitive method among CR, CRm, ARC and TR) in terms of iterations for convergence.

We really appreciate CRm, which first accelerates CR by momentum. We would like to mention our main

difference with CRm [17]. Firstly, the scheme of the momentum for ARCM is different from that in CRm due to the adaptive selection strategy for  $M_k$ . Secondly, the step size for momentum in ARCM is free of gradient evaluation but CRm requires. Thirdly, besides the global convergence, we also study the local convergence of ARCM under the KL property, which is more general than the local error-bound condition studied for CRm. Furthermore, we analyzed the ARCM with inexact CRS solutions, which is more applicable in real applications.

## 2 The Proposed ARCM Algorithm

The detailed pseudocode for ARCM is shown in Algorithm 1. Here, we first assume that the cubic subproblem in Step 2 of Algorithm 1 is exactly computed. In Section 3, we will analyze the convergence property of ARCM with inexact solutions in Step 2.

Firstly, the momentum  $\mathbf{v}_{k-1}$  is an aggregation of previously accepted descent steps. It is always uniformly bounded if  $\beta_k \leq \tau < 1$  and  $\mathbf{s}_k$  is bounded (we will prove it in Lemma 2.3) for all  $k \leq T$ . At the beginning of the algorithm (i.e.,  $k \ll T$ ),  $\|\mathbf{s}_k\|_2$  is usually relative large. Therefore,  $\alpha_1 \|\mathbf{s}_k\|_2$  and

$\alpha_2\|\mathbf{s}_k\|_2^2$  dominate the term  $\min(\tau, \alpha_1\|\mathbf{s}_k\|_2, \alpha_2\|\mathbf{s}_k\|_2^2)$ , i.e.,  $\tau = \min(\tau, \alpha_1\|\mathbf{s}_k\|_2, \alpha_2\|\mathbf{s}_k\|_2^2)$  is a popular choice for step size of momentum. When  $\mathbf{x}_k$  approaches the local minima (then  $\|\mathbf{s}_k\|_2 \approx 0$ ), we may not expect momentum with large step size to work for second-order methods. Therefore, we adopt  $\alpha_1\|\mathbf{s}_k\|_2$  and  $\alpha_2\|\mathbf{s}_k\|_2^2$  in selecting  $\beta_k$  (where  $\beta_k \approx 0$ ) in order to preserve the local convergence property of the second-order method. Secondly, we require that  $f(\mathbf{z}_{k+1}) \leq f(\mathbf{y}_{k+1})$  since we do not hope to violate the sufficient descent of the objective in CR and ARC. Here, such  $\mathbf{z}_{k+1}$  must exist (e.g.,  $\beta_k = 0$ ) and we may use the bisection method to search appropriate  $\beta_k$ . The following Theorem 2.1 shows that  $\mathbf{z}_{k+1}$  with nonzero  $\beta_k$  may exist under two cases, where the momentum term helps the convergence of ARcm (i.e.,  $f(\mathbf{z}_{k+1}) < f(\mathbf{y}_{k+1})$  holds).

**THEOREM 2.1.** *Assume that the Hessian  $\nabla^2 f(\mathbf{x})$  of  $f(\mathbf{x})$  on the line segment  $[\mathbf{x}_k, \mathbf{x}_k + \mathbf{s}_k]$  is  $L_k$ -Lipschitz (i.e., for all  $\mathbf{x}, \mathbf{y} \in [\mathbf{x}_k, \mathbf{x}_k + \mathbf{s}_k]$  we have  $\|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{y})\|_2 \leq L_k \|\mathbf{x} - \mathbf{y}\|_2$ ), then in the following two cases the momentum may help the convergence, i.e., there exist small enough  $\beta_k$  such that  $f(\mathbf{z}_{k+1}) < f(\mathbf{y}_{k+1})$ :*

- (i)  $L_k < \sigma_k$  and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} > 0$ ;
- (ii)  $L_k > \sigma_k$  and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} < 0$ .

*Proof.* Suppose that the Hessian of  $f(\mathbf{x})$  is  $L$ -Lipschitz on a ball centered at  $\mathbf{y}_{k+1}$  with radius  $\tau\|\mathbf{v}_{k-1}\|_2$ , then we have

$$f(\mathbf{z}_{k+1}) - f(\mathbf{y}_{k+1}) \leq \beta \cdot \nabla f(\mathbf{y}_{k+1})^\top \mathbf{v}_{k-1} + \frac{1}{2}\beta^2 \cdot \mathbf{v}_{k-1}^\top \nabla^2 f(\mathbf{x}_{k+1}) \mathbf{v}_{k-1} + \frac{L}{6}\beta^3 \cdot \|\mathbf{v}_{k-1}\|_2^3,$$

by [15, Lemma 1] (we also provide the useful result in supplementary material (A.2)). If  $\nabla f(\mathbf{y}_{k+1})^\top \mathbf{v}_{k-1} < 0$ , then there exists a small enough  $\beta > 0$  such that  $f(\mathbf{z}_{k+1}) < f(\mathbf{y}_{k+1})$ . In the remaining part, we show that in the above two cases,  $\nabla f(\mathbf{y}_{k+1})^\top \mathbf{v}_{k-1} < 0$  may hold. Using the properties of the cubic regularization method [15, Lemma 1 & (2.5)], we have

$$\begin{aligned} & \nabla f(\mathbf{y}_{k+1})^\top \mathbf{v}_{k-1} \\ &= (\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k))^\top \mathbf{v}_{k-1} + \nabla f(\mathbf{x}_k)^\top \mathbf{v}_{k-1} \\ &= (\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k))^\top \mathbf{v}_{k-1} \\ & \quad - \left( \nabla^2 f(\mathbf{x}_k) \mathbf{s}_k + \frac{1}{2} \sigma_k \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k \right)^\top \mathbf{v}_{k-1} \\ &= -\frac{1}{2} \sigma_k \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1} \\ & \quad + (\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \mathbf{s}_k)^\top \mathbf{v}_{k-1} \\ &\leq -\frac{1}{2} \sigma_k \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1} + \frac{1}{2} L_k \|\mathbf{s}_k\|_2 \cdot |\mathbf{s}_k^\top \mathbf{v}_{k-1}|. \end{aligned} \tag{2.2}$$

If  $L_k < \sigma_k$  (i.e.,  $L_k \leq \sigma_k - \epsilon$  for some  $\epsilon > 0$ ) and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} > 0$ , then

$$\begin{aligned} \nabla f(\mathbf{y}_{k+1})^\top \mathbf{v}_{k-1} &\leq -\frac{1}{2} \sigma_k \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1} + \frac{1}{2} L_k \|\mathbf{s}_k\|_2 \cdot |\mathbf{s}_k^\top \mathbf{v}_{k-1}| \\ &\leq -\frac{1}{2} \epsilon \|\mathbf{s}_k\|_2 \cdot |\mathbf{s}_k^\top \mathbf{v}_{k-1}| < 0. \end{aligned}$$

In the last inequality of (2.2), we use the inequality that

$$\begin{aligned} & \left| (\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \mathbf{s}_k)^\top \mathbf{v}_{k-1} \right| \\ & \leq \frac{1}{2} L_k \|\mathbf{s}_k\|_2 \cdot |\mathbf{s}_k^\top \mathbf{v}_{k-1}|. \end{aligned}$$

If  $L_k > \sigma_k$ ,  $\mathbf{s}_k^\top \mathbf{v}_{k-1} < 0$  and

$$(\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \mathbf{s}_k)^\top \mathbf{v}_{k-1} = \frac{1}{2} \tilde{L}_k \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1}$$

for  $\sigma_k \leq \tilde{L}_k - \epsilon < L_k$  and  $\epsilon > 0$ , then

$$\begin{aligned} & \nabla f(\mathbf{y}_{k+1})^\top \mathbf{v}_{k-1} \\ &= -\frac{1}{2} \sigma_k \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1} \\ & \quad + (\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \mathbf{s}_k)^\top \mathbf{v}_{k-1} \\ &= -\frac{1}{2} \sigma_k \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1} + \frac{1}{2} \tilde{L}_k \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1} \\ &\leq \frac{1}{2} \epsilon \|\mathbf{s}_k\|_2 \cdot \mathbf{s}_k^\top \mathbf{v}_{k-1} < 0. \quad \square \end{aligned}$$

**REMARK 2.1.** *The two potential cases may happen since ARcm adaptively select  $M_k := \sigma_k$  by the criteria  $\rho_k$  where  $\sigma_k$  may be overestimated or underestimated. When  $\sigma_k > L_k$ , the step  $\mathbf{s}_k$  is too conservative and the new step  $\mathbf{x}_{k+1} - \mathbf{x}_k$  contributes to lower objective value if  $\mathbf{v}_{k-1}$  is also a descent direction (i.e.,  $\mathbf{s}_k^\top \mathbf{v}_{k-1} > 0$ ). Conversely, if  $\sigma_k < L_k$ , the step  $\mathbf{s}_k$  may be too aggressive, then the momentum with opposite direction (i.e.,  $\mathbf{s}_k^\top \mathbf{v}_{k-1} < 0$ ) may correct the imperfect step  $\mathbf{s}_k$ . In CRm, the momentum  $\mathbf{v}_{k-1}$  is of highly correlated to  $\mathbf{s}_k$  that  $\mathbf{s}_k^\top \mathbf{v}_{k-1} > 0$  and  $M > L_k$ , then only the first case will happen. Therefore, the momentum scheme in ARcm is more appropriate for ARC since ARC is more ambitious than vanilla CR.*

**2.1 Global Convergence** To prove global convergence, the following mild assumptions are essential.

**ASSUMPTION 2.1.** *We have the following assumptions for the objective  $f(\mathbf{x})$ :*

1.  $f(\mathbf{x})$  is globally second-order differentiable with respect to  $\mathbf{x}$ .
2.  $f(\mathbf{x})$  is bounded below, i.e.,  $f^* = \inf_{\mathbf{x}} f(\mathbf{x}) > -\infty$ .

3. For the given initial guess  $\mathbf{x}_0$ , there exists a closed convex set  $\mathcal{F}$  such that the level set  $\mathcal{L}(\mathbf{x}_0) := \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{x}_0)\} \subseteq \mathcal{F}$ . For all  $\mathbf{x}, \tilde{\mathbf{x}} \in \mathcal{F}$ , we have

$$(2.3) \quad \begin{aligned} \|\nabla^2 f(\mathbf{x})\|_2 &\leq \kappa_H, \\ \|\nabla f(\mathbf{x}) - \nabla f(\tilde{\mathbf{x}})\|_2 &\leq L_G \|\mathbf{x} - \tilde{\mathbf{x}}\|_2, \\ \|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\tilde{\mathbf{x}})\|_2 &\leq L_H \|\mathbf{x} - \tilde{\mathbf{x}}\|_2. \end{aligned}$$

The second and the third assumption hold if  $\lim_{\mathbf{x} \rightarrow \infty} f(\mathbf{x}) = +\infty$  (i.e.,  $f(\mathbf{x})$  is level bounded that  $\mathcal{L}(\mathbf{x}_0)$  is bounded) and  $f(\mathbf{x})$  is smooth enough, which is common in machine learning, e.g., non-negative (smooth) loss with ( $\ell_2$ ) regularization terms. For the objective that is not bounded below, global convergence is usually hard to be achieved theoretically. We then present some useful propositions and lemmas in preparation for deriving global convergence. Some proofs are placed in the supplement.

PROPOSITION 2.1. ([15, LEMMA 4], [4, LEMMA 3.3]) If  $\mathbf{s}_k \in \arg \min_{\mathbf{s}} m_k(\mathbf{s})$ , where  $m_k(\mathbf{s}) := f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2} \mathbf{s}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{s} + \frac{\sigma_k}{6} \|\mathbf{s}\|_2^3$ , then

$$(2.4) \quad f(\mathbf{x}_k) - m_k(\mathbf{s}_k) \geq \frac{1}{12} \sigma_k \|\mathbf{s}_k\|_2^3.$$

Furthermore, let  $\mathcal{S} := \{i : \rho_i > \eta_1\}$  denote the the set of index that successful update occurs, then for any  $k \in \mathcal{S}$  we have

$$(2.5) \quad f(\mathbf{x}_k) - f(\mathbf{y}_{k+1}) \geq \eta_1 \cdot (f(\mathbf{x}_k) - m_k(\mathbf{s}_k)) \geq \frac{\eta_1}{12} \sigma_k \|\mathbf{s}_k\|_2^3.$$

LEMMA 2.1. Under Assumption 2.1, the adaptive penalty parameter  $\sigma_k$  cannot be arbitrarily large, i.e.,

$$(2.6) \quad \sigma_k \leq \max\{L_H \gamma_1, \sigma_{\min}\} := \sigma_{\max}.$$

LEMMA 2.2. Denote  $\mathcal{S}_j := \{i < j : \rho_i > \eta_1\}$ , and  $\mathcal{U}_j := \{i < j : \rho_i \leq \eta_1\}$ . Then

$$(2.7) \quad |\mathcal{U}_j| \leq \left\lceil \log \frac{\max\{L_H \gamma_1, \sigma_{\min}\}}{\sigma_{\min}} \right\rceil \cdot |\mathcal{S}_j|,$$

and

$$(2.8) \quad |\mathcal{S}_T| \geq \frac{T}{1 + \left\lceil \log \frac{\max\{L_H \gamma_1, \sigma_{\min}\}}{\sigma_{\min}} \right\rceil}.$$

*Proof.* If the current update is successful, then we need at most  $\left\lceil \log \frac{\sigma_{\max}}{\sigma_{\min}} \right\rceil$  steps of unsuccessful updates to achieve the next successful update, according to Lemma 2.1. Combining with (2.7) and the fact that  $|\mathcal{S}_T| + |\mathcal{U}_T| = T$ , we have the relation (2.8).  $\square$

LEMMA 2.3. The followings hold for  $\mathbf{s}_k$  and  $\mathcal{S}_T$ :

$$\max_{i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2 \leq \left( \frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3},$$

and

$$\min_{i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2 \leq \left( \frac{12(f(\mathbf{x}_0) - f^*)}{|\mathcal{S}_T| \eta_1 \sigma_{\min}} \right)^{1/3}.$$

Therefore, we have

$$\|\mathbf{v}_k\|_2 \leq \frac{1}{1 - \tau} \left( \frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3},$$

for all  $k \in \mathcal{S}_T$ .

*Proof.* According to Proposition 2.1, we have

$$(2.9) \quad \begin{aligned} &\sum_{k \in \mathcal{S}_T} \frac{1}{12} \eta_1 \sigma_k \|\mathbf{s}_k\|_2^3 \\ &\leq \sum_{k \in \mathcal{S}_T} f(\mathbf{x}_k) - f(\mathbf{y}_{k+1}) \leq \sum_{k \in \mathcal{S}_T} f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \\ &= \sum_{k=0}^T f(\mathbf{x}_k) - f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_0) - f^*, \end{aligned}$$

where the second inequality holds since  $f(\mathbf{x}_{k+1}) \leq f(\mathbf{y}_{k+1})$ , and the equality is satisfied because  $f(\mathbf{x}_{k+1}) = f(\mathbf{x}_k)$  for all  $k \notin \mathcal{S}_T$ . Then, we have

$$\sum_{k \in \mathcal{S}_T} \|\mathbf{s}_k\|_2^3 \leq \frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}}$$

as  $\sigma_k \geq \sigma_{\min}$ . Therefore, the first two inequalities hold by using  $\max_{i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2^3 \leq \sum_{i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2^3$  and  $\min_{i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2^3 \leq \frac{1}{|\mathcal{S}_T|} \sum_{i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2^3$ . The last inequality in the lemma holds since for any  $k \in \mathcal{S}_T$ ,

$$\|\mathbf{v}_k\|_2 \leq \frac{1}{1 - \tau} \cdot \max_{i \leq k, i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2 \leq \frac{1}{1 - \tau} \cdot \max_{i \in \mathcal{S}_T} \|\mathbf{s}_i\|_2. \quad \square$$

LEMMA 2.4. If  $k \in \mathcal{S}_T$ , we have

$$(2.10) \quad \|\nabla f(\mathbf{x}_{k+1})\|_2 \leq c_1 \|\mathbf{s}_k\|_2^2,$$

and

$$(2.11) \quad \lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1})) \geq -c_2 \|\mathbf{s}_k\|_2,$$

where  $c_1 = \frac{1}{2} \max\{L_H \gamma_1, \sigma_{\min}\} + \frac{1}{2} L_H + \frac{\alpha_2 L_G}{1 - \tau} \left( \frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3}$  and  $c_2 = \frac{1}{2} \max\{L_H \gamma_1, \sigma_{\min}\} + L_H + \frac{\alpha_1 L_H}{1 - \tau} \left( \frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3}$ .

*Proof.* If  $k \in \mathcal{S}_T$ , we have

$$\|\nabla f(\mathbf{y}_{k+1})\|_2 \leq \frac{1}{2}(\sigma_k + L_H) \|\mathbf{s}_k\|_2^2,$$

and

$$\lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) \geq -\left(\frac{1}{2}\sigma_k + L_H\right) \|\mathbf{s}_k\|_2.$$

The proof can be found in [15, Lemma 3 & 5]. We then derive the error bounds for  $\|\nabla f(\mathbf{x}_{k+1})\|_2$  and  $\lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1}))$ .

$$\begin{aligned} & \|\nabla f(\mathbf{x}_{k+1})\|_2 \\ & \leq \|\nabla f(\mathbf{y}_{k+1})\|_2 + \|\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{y}_{k+1})\|_2 \\ & \leq \|\nabla f(\mathbf{y}_{k+1})\|_2 + L_g \beta_k \|\mathbf{v}_k\|_2 \\ & \leq \frac{1}{2}(\sigma_k + L_H) \|\mathbf{s}_k\|_2^2 + L_g \alpha_2 \|\mathbf{s}_k\|_2^2 \cdot \|\mathbf{v}_k\|_2 \leq c_1 \|\mathbf{s}_k\|_2^2, \end{aligned}$$

where  $c_1 = \frac{1}{2}\sigma_{\max} + \frac{1}{2}L_H + \frac{\alpha_2 L_g}{1-\tau} \left(\frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}}\right)^{1/3}$ . Furthermore, we have

$$\begin{aligned} & \lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1})) \\ & \geq \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) - \|\nabla^2 f(\mathbf{x}_{k+1}) - \nabla^2 f(\mathbf{y}_{k+1})\|_2 \\ & \geq \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) - L_H \|\mathbf{x}_{k+1} - \mathbf{y}_{k+1}\|_2 \\ & \geq \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) - L_H \beta_k \|\mathbf{v}_k\|_2 \\ & \geq -\left(\frac{1}{2}\sigma_k + L_H\right) \|\mathbf{s}_k\|_2 - L_H \alpha_1 \|\mathbf{s}_k\|_2 \cdot \|\mathbf{v}_k\|_2 \\ & \geq -c_2 \|\mathbf{s}_k\|_2, \end{aligned}$$

where  $c_2 = \frac{1}{2}\sigma_{\max} + L_H + \frac{\alpha_1 L_H}{1-\tau} \left(\frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}}\right)^{1/3}$ .  $\square$

**THEOREM 2.2.** *We introduce the following measure of the local optimality:*

$$(2.12) \quad \mu(\mathbf{x}) = \max \left\{ \sqrt{\frac{1}{c_1} \|\nabla f(\mathbf{x})\|}, -\frac{1}{c_2} \lambda_{\min}(\nabla^2 f(\mathbf{x})) \right\},$$

where  $c_1 > 0$  and  $c_2 > 0$  are two universal constant defined in Lemma 2.4. Under Assumption 2.1, let the sequence  $\{\mathbf{x}_k\}_{k=1}^T$  be generated by Algorithm 1, then

$$(2.13) \quad \begin{aligned} & \min_{1 \leq k \leq T} \mu(\mathbf{x}_k) \\ & \leq \left( \frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min} T} \left( 1 + \left\lceil \log \frac{\max\{L_H \gamma_1, \sigma_{\min}\}}{\sigma_{\min}} \right\rceil \right) \right)^{1/3} \\ & = \mathcal{O}(T^{-1/3}). \end{aligned}$$

*Proof.* Lemma 2.4 implies that  $\mu(\mathbf{x}_{k+1}) \leq \|\mathbf{s}_k\|_2$  for all  $k \in \mathcal{S}_T$ . Then we have  $\min_{1 \leq k \leq T} \mu(\mathbf{x}_k) \leq \min_{k \in \mathcal{S}_T} \|\mathbf{s}_k\|_2$ . Combining it with Lemma 2.2 and Lemma 2.3, we finish the proof.  $\square$

**REMARK 2.2.** *Note that the proposed ARCM also satisfies the global convergence rate  $\mathcal{O}(T^{-1/3})$  for the measure  $\mu(\cdot)$ , as the vanilla CR [15] and ARC [4]. We do not expect to improve the convergence rate since ARCM is still a second-order method.*

**2.2 Local Convergence** In this subsection, we let  $T = +\infty$ . If the accumulation point  $\bar{\mathbf{x}}$  of the sequence  $\{\mathbf{x}_k\}_{k=0}^{\infty}$  generated by CR satisfies  $\lambda_{\min}(\nabla^2 f(\bar{\mathbf{x}})) > 0$ , then it achieves the local quadratic convergence [15]. However, the Hessian  $\nabla^2 f(\mathbf{x})$  at the local optima is usually not necessary to be positive definite. Therefore, the local quadratic convergence may not work. Yue et al. [19] proved that CR achieves the quadratic convergence if the objective  $f(\mathbf{x})$  satisfies the local error bound condition, which is weaker than the local positive definiteness of Hessian. Later, Zhou et al. [21] generalized the results in [19] to the objective that satisfies the KL property. Here, we show that the proposed ARCM also achieves the local convergence under the KL property. In the following analysis, sets  $\mathcal{S} := \{i : \rho_i > \eta_1\}$  and  $\mathcal{U} := \{i : \rho_i \leq \eta_1\}$  are in the ascending order. Lemma 2.2 implies that  $|\mathcal{S}|$  must be infinite but  $|\mathcal{U}|$  may be finite. We use  $k_j \in \mathcal{S}$  to denote the  $j$ -th element of  $\mathcal{S}$ .

**DEFINITION 2.1.** *A differentiable function  $f(\cdot)$  is said to satisfy the KL property if for any compact set  $\bar{\mathcal{X}}$  where  $f(\cdot)$  takes a constant value  $\bar{f}$ , there exist  $\epsilon_1, \epsilon_2 > 0$  such that for all  $\bar{\mathbf{x}} \in \bar{\mathcal{X}}$  and  $\mathbf{x} \in \{\mathbf{z} : \text{dist}(\mathbf{z}, \bar{\mathcal{X}}) < \epsilon_1, \bar{f} < f(\mathbf{z}) < \bar{f} + \epsilon_2\}$ ,*

$$(2.14) \quad \phi'(f(\mathbf{x}) - \bar{f}) \|\nabla f(\mathbf{x})\| \geq 1$$

holds, where  $\phi(t) = \frac{c}{\theta} t^\theta$  for some  $c > 0$  and  $\theta \in (0, 1)$ . Then, (2.14) is equivalent to

$$(2.15) \quad f(\mathbf{x}) - \bar{f} \leq c_0 \|\nabla f(\mathbf{x})\|^{\frac{1}{1-\theta}}$$

with  $c_0 = c^{1/(1-\theta)}$ .

Besides Assumption 2.1, we need further but mild assumptions for the local convergence. As is discussed following Assumption 2.1, the level boundedness (Assumption 2.2) and smoothness (e.g.,  $f(\mathbf{x})$  is third-order differentiable with  $\mathbf{x}$ ) imply the second and the third assumptions in Assumption 2.1.

**ASSUMPTION 2.2.** *We further assume that  $f(\mathbf{x})$  is level bounded, i.e., the level set  $\mathcal{L}(\tilde{\mathbf{x}})$  is bounded,  $\forall \tilde{\mathbf{x}} \in \mathcal{F}$ .*

Before deriving the local convergence for ARCM under the KL property, we first show that the sequence generated by  $\{\mathbf{x}_k\}_k$  is Cauchy and convergent to a second-order critical point. The proofs for Lemma 2.5

and Theorem 2.3 are extended from [21] and we put the tedious details in the supplementary material due to the page limit.

LEMMA 2.5. *Under Assumption 2.1 and 2.2, the followings hold for the sequence  $\{\mathbf{x}_k\}_{k=0}^{+\infty}$  generated by ARCM (Algorithm 1):*

1.  $\bar{f} := \lim_{k \rightarrow \infty} f(\mathbf{x}_k)$  exists.
2. The sequence  $\{\mathbf{x}_k\}_{k=0}^{+\infty}$  is bounded and  $\lim_{k \rightarrow \infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 = 0$ . Moreover, the set  $\bar{\mathcal{X}}$  of accumulation points of the sequence is non-empty, satisfying

$$f(\bar{\mathbf{x}}) = \bar{f}, \quad \nabla f(\bar{\mathbf{x}}) = \mathbf{0}, \quad \lambda_{\min}(\nabla^2 f(\bar{\mathbf{x}})) \geq \mathbf{0},$$

for all  $\bar{\mathbf{x}} \in \bar{\mathcal{X}}$ .

3. If  $f(\mathbf{x})$  satisfies the KL property, then  $\bar{\mathcal{X}} = \{\bar{\mathbf{x}}\}$  is a singleton.

THEOREM 2.3. *Let the objective  $f(\mathbf{x})$  satisfies Assumption 2.1, 2.2 and the KL property, then there exists a large enough  $j_0 \in \mathbb{N}$  such that the sequence  $\{\mathbf{x}_k\}_{k=0}^{+\infty}$  (or  $\{\mathbf{x}_{k_j+1}\}_{j=0}^{+\infty}$ ,  $k_j \in \mathcal{S}$ ) generated by ARCM (Algorithm 1) satisfies*

1. If  $\theta \in (\frac{1}{3}, 1)$ , then the local convergence is super-linear with

$$\|\mathbf{x}_{k_j+1} - \bar{\mathbf{x}}\|_2 \leq \mathcal{O} \left( \exp \left( - \left( \frac{2\theta}{1-\theta} \right)^{j-j_0} \right) \right)$$

and

$$\|\mathbf{x}_{k+1} - \bar{\mathbf{x}}\|_2 \leq \mathcal{O} \left( \exp \left( - \left( \frac{2\theta}{1-\theta} \right)^{\left\lceil \frac{k-k_{j_0}}{c_3} \right\rceil} \right) \right),$$

where  $c_3 = 1 + \left\lceil \log \frac{\max\{L_H \gamma_1, \sigma_{\min}\}}{\sigma_{\min}} \right\rceil$ .

2. If  $\theta = \frac{1}{3}$ , then the local convergence is linear with

$$\|\mathbf{x}_{k_j+1} - \bar{\mathbf{x}}\|_2 \leq \mathcal{O}(\exp(-c_4(j-j_0)))$$

and

$$\|\mathbf{x}_{k+1} - \bar{\mathbf{x}}\|_2 \leq \mathcal{O} \left( \exp \left( -c_4 \left\lceil \frac{k-k_{j_0}}{c_3} \right\rceil \right) \right),$$

for some constants  $c_4 > 0$ .

3. If  $\theta \in (0, \frac{1}{3})$ , then the local convergence is sub-linear with

$$\|\mathbf{x}_{k_j+1} - \bar{\mathbf{x}}\|_2 \leq \mathcal{O} \left( (j-j_0)^{-\frac{2\theta}{1-3\theta}} \right)$$

and

$$\|\mathbf{x}_{k+1} - \bar{\mathbf{x}}\|_2 \leq \mathcal{O} \left( \left\lceil \frac{k-k_{j_0}}{c_3} \right\rceil^{-\frac{2\theta}{1-3\theta}} \right).$$

### 3 ARCM with Inexact Solutions

A popular approach to the exact solution of CRS in Step 2 of Algorithm 1, is solving the corresponding secular equation [15], where full eigendecomposition for the Hessian matrix is required (with computational complexity  $\mathcal{O}(d^3)$ ). Cartis et al. [4] designed a Newton-Cholesky iteration method for solving CRS, however, its computational cost is still of  $\mathcal{O}(d^3)$ . Exact solutions of CRS are usually very computationally expensive for high-dimensional problems. In practice, inexact solvers are more popular. Cartis et al. [4] projected the CRS to a Krylov subspace in order to lower the dimension. The convergence of the Krylov subspace method for inexactly solving CRS was analyzed by Carmon and Duchi [3]. Later, they showed that the simple gradient descent with proper learning rates finds the global solution of CRS [2]. Jiang et al. [12] proposed an accelerated first-order method that reformulates the CRS into a constrained convex problem. Recently, Gao et al. [7] suggested solving an approximate secular equation rather than the exact secular equation, where partial eigendecomposition is required. Suppose that we solve the CRS by inexact solvers that satisfy Condition 3.1:

$$\tilde{\mathbf{s}}_k \approx \arg \min_{\mathbf{s}} m_k(\mathbf{s}),$$

$$\mathbf{y}_{k+1} = \mathbf{x}_{k+1} + \tilde{\mathbf{s}}_k,$$

$$\tilde{\mathbf{v}}_k = \tilde{\beta}_k \tilde{\mathbf{v}}_{k-1} + \tilde{\mathbf{s}}_k \text{ and } \mathbf{z}_{k+1} = \mathbf{x}_k + \tilde{\mathbf{v}}_k$$

$$\text{with } \tilde{\beta}_k \in [0, \min\{\tau, \alpha_1 \|\tilde{\mathbf{s}}_k\|_2, \alpha_2 \|\tilde{\mathbf{s}}_k\|_2^2\}]$$

$$\text{and } f(\mathbf{z}_{k+1}) \leq f(\mathbf{y}_{k+1}),$$

$$\mathbf{x}_{k+1} = \mathbf{z}_{k+1}.$$

CONDITION 3.1. *Suppose that  $\tilde{\mathbf{s}}_k$  is the inexact solution of the CRS in the  $k$ -th iteration, satisfying the following  $\delta_k$ -conditions:*

$$1. m_k(\tilde{\mathbf{s}}_k) - f(\mathbf{x}_k) \leq -\frac{\sigma_k}{12} \|\tilde{\mathbf{s}}_k\|_2^3 + \delta_k < 0;$$

$$2. \nabla m_k(\tilde{\mathbf{s}}_k) \leq \delta_k^{2/3};$$

$$3. \|\tilde{\mathbf{s}}_k\| - \|\mathbf{s}_k\| \leq \delta_k^{1/3}.$$

Note that the first item  $m_k(\tilde{\mathbf{s}}_k) - f(\mathbf{x}_k) < 0$  is easy to be satisfied since the Cauchy point method guaranteed the sufficient decrease except at a stable point. Before analyzing the global convergence of ARCM with inexact CRS solutions, we first provide some useful lemmas where some proofs are in the supplement.

LEMMA 3.1. *Under Assumption 2.1 and Condition 3.1, the adaptive penalty parameter  $\sigma_k$  cannot be arbitrary large, i.e.,*

$$\sigma_k \leq \max\{L_H \gamma_1, \sigma_{\min}\} = \sigma_{\max}.$$

Therefore, (2.7) and (2.8) still hold for  $|\mathcal{U}_T|$  and  $|\mathcal{S}_T|$ .

LEMMA 3.2. *Without the loss of generality, we assume that  $\delta_k < f(\mathbf{x}_0) - f^*$ . Suppose that Assumption 2.1 and Condition 3.1 hold, then we have*

$$\max_{k \in \mathcal{S}_T} \|\tilde{\mathbf{s}}_k\|_2 \leq \left( \frac{24(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3},$$

and

$$\min_{k \in \mathcal{S}_T} \|\tilde{\mathbf{s}}_k\|_2^3 - \frac{12\delta_k}{\sigma_k} \leq \frac{12(f(\mathbf{x}_0) - f^*)}{|\mathcal{S}_T| \eta_1 \sigma_{\min}}.$$

Therefore, we have

$$\|\tilde{\mathbf{v}}_k\|_2 \leq \frac{1}{1-\tau} \left( \frac{24(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3},$$

for all  $k \in \mathcal{S}_T$

LEMMA 3.3. *Under Assumption 2.1 and Condition 3.1, if  $k \in \mathcal{S}_T$ , we have*

$$(3.16) \quad \|\nabla f(\mathbf{x}_{k+1})\|_2 \leq c_5 \|\tilde{\mathbf{s}}_k\|_2^2 + \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2,$$

and

$$(3.17) \quad \lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1})) \geq -c_6 \|\tilde{\mathbf{s}}_k\|_2 - \frac{\sigma_{\max}}{2} \|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2,$$

where  $c_5 = \frac{1}{2} \max\{L_H \gamma_1, \sigma_{\min}\} + \frac{1}{2} L_H + \frac{\alpha_2 L_g}{1-\tau} \left( \frac{24(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3}$  and  $c_6 = \frac{1}{2} \sigma_{\max} + L_H + \frac{\alpha_1 L_H}{1-\tau} \left( \frac{24(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3}$

*Proof.* We first derive the error bound for  $\|\nabla f(\mathbf{y}_{k+1})\|$  and  $\lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1}))$ :

$$(3.18) \quad \begin{aligned} & \|\nabla f(\mathbf{y}_{k+1})\|_2 \\ \leq & \|\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \tilde{\mathbf{s}}_k\|_2 \\ & + \left\| \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) \tilde{\mathbf{s}}_k + \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2 \tilde{\mathbf{s}}_k \right\|_2 + \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2^2 \\ \leq & \frac{L_H}{2} \|\tilde{\mathbf{s}}_k\|_2^2 + \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2 + \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2^2 \\ \leq & \left( \frac{\sigma_{\max}}{2} + \frac{L_H}{2} \right) \|\tilde{\mathbf{s}}_k\|_2^2 + \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2, \end{aligned}$$

where the second inequality is due to (A.1); and

$$\begin{aligned} & \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) \\ \geq & \lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) - \|\nabla^2 f(\mathbf{y}_{k+1}) - \nabla^2 f(\mathbf{x}_k)\|_2 \\ \geq & -\frac{\sigma_k}{2} \|\mathbf{s}_k\|_2 - L_H \|\tilde{\mathbf{s}}_k\|_2 \\ \geq & -\frac{\sigma_k}{2} \|\|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2\| - \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2 - L_H \|\tilde{\mathbf{s}}_k\|_2 \\ = & -\left( \frac{\sigma_{\max}}{2} - L_H \right) \|\tilde{\mathbf{s}}_k\|_2 - \frac{\sigma_{\max}}{2} \|\|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2\|. \end{aligned}$$

where the second inequality comes from a well-known result of cubic regularization [15, Proposition 1]. Using similar arguments as in Lemma 2.4, we complete the proof. We put the remaining details in the supplement.  $\square$

THEOREM 3.1. *We introduce the following measure of the local optimality:*

(3.19)

$$\tilde{\mu}(\mathbf{x}) = \max \left\{ \sqrt{\frac{1}{c_5} \|\nabla f(\mathbf{x})\|}, -\frac{1}{c_6} \lambda_{\min}(\nabla^2 f(\mathbf{x})) \right\}.$$

*Under Assumption 2.1 and Condition 3.1, let the sequence  $\{\mathbf{x}_k\}_{k=1}^T$  be generated by Algorithm 1, we have*

1. *If  $0 < \delta_k \leq \delta$  for some  $\delta > 0$ , then*

$$(3.20) \quad \min_{1 \leq k \leq T} \tilde{\mu}(\mathbf{x}_k) \leq \mathcal{O}\left(T^{-1/3} + \delta^{1/3}\right).$$

2. *If there exists  $0 < \varepsilon < \frac{1}{12}$  such that  $0 < \delta_k \leq \varepsilon \sigma_k \|\tilde{\mathbf{s}}_k\|_2^3$ , then*

$$(3.21) \quad \min_{1 \leq k \leq T} \tilde{\mu}(\mathbf{x}_k) \leq \mathcal{O}\left(T^{-1/3}\right).$$

*Proof.* Lemma 3.3 implies that  $\tilde{\mu}(\mathbf{x}_{k+1}) \leq \|\tilde{\mathbf{s}}_k\|_2 + \frac{1}{\sqrt{c_5}} \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2^{1/2} + \frac{\sigma_{\max}}{2c_6} \|\|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2\|$  for all  $k \in \mathcal{S}_T$ . Then we have  $\min_{1 \leq k \leq T} \tilde{\mu}(\mathbf{x}_k) \leq \min_{k \in \mathcal{S}_T} \|\tilde{\mathbf{s}}_k\|_2 + \frac{1}{\sqrt{c_5}} \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2^{1/2} + \frac{\sigma_{\max}}{2c_6} \|\|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2\|$ . Combining it with Lemma 3.1, Lemma 3.2 and Condition 3.1, we finish the proof.  $\square$

COROLLARY 3.1. *Under Assumption 2.1 and Condition 3.1, let the sequence  $\{\mathbf{x}_k\}_{k=1}^{+\infty}$  be generated by Algorithm 1 and  $f(\mathbf{x})$  satisfies the KL property. If there exists  $0 < \varepsilon < \frac{1}{12}$  such that  $0 < \delta_k \leq \varepsilon \sigma_k \|\tilde{\mathbf{s}}_k\|_2^3$ , then ARCm with inexact CRS solutions still hold local convergence property as in Theorem 2.3.*

## 4 Numerical Experiments

In this section, we conduct experiments on the proposed ARCm and some state-of-the-art second-order methods (e.g., CR, CRm, ARC, and TR) in solving non-convex logistic regression and robust linear regression models.

**Settings.** All these algorithms involve solving CRS in each iteration (note that trust region subproblems are similar to cubic regularization subproblems in Step 2 of Algorithm 1). We adopt the Krylov subspace method [4, 5] with at most 50 subspaces to approximately solve CRS. All hyperparameters are tuned to achieve nearly optimal results in terms of the iterations for convergence. For ARCm, we set the starting cubic penalty

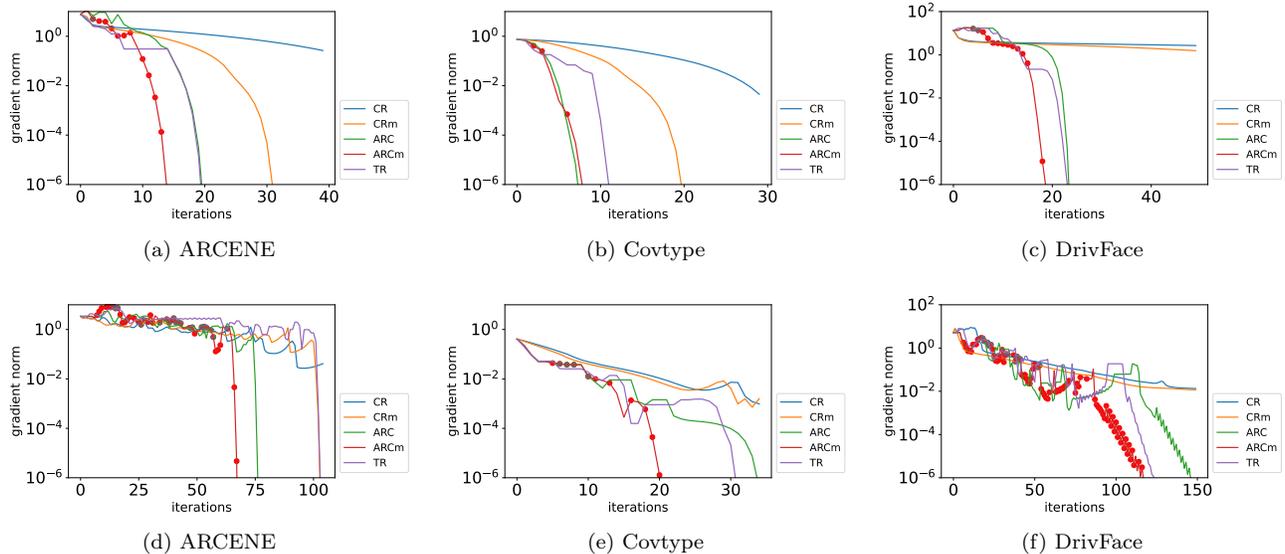


Figure 1: Gradient norm versus the iteration steps on logistic regression (the first row) and robust linear regression (the second row) models for three datasets. Red dots:  $\beta_k > 0$  and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} > 0$ ; gray dots:  $\beta_k > 0$  and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} < 0$ .

parameter  $\sigma_0 = 1.0$  and the momentum parameters  $(\tau, \alpha_1, \alpha_2) = (0.5, 0.1, 1.0)$ . Except for the momentum parameters, all other parameters of ARCm are the same as ARC (e.g.,  $(\eta_1, \eta_2) = (0.1, 0.9)$ ).

**Datasets.** For both two models, we test algorithms on ARCENE [10], Covtype [1] and DrivFace [6] datasets. Each training data in all three datasets is in the form of  $(\mathbf{a}_i, b_i)$ , where  $\mathbf{a}_i \in \mathbb{R}^d$  is a multi-variate attribute vector and  $b_i \in \{0, 1\}$  is the corresponding label. The detailed information for these two datasets is shown in Table 1.

Table 1: The overview of datasets.

Dataset	sample size $n$	dimension $d$
ARCENE	100	10000
Covtype	581012	54
DrivFace	606	6400

**Models.** For a given dataset  $\{(\mathbf{a}_i, b_i)\}_{i=1}^n$ , the empirical loss for the non-convex logistic regression is

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n b_i \log [\psi(\mathbf{w}^\top \mathbf{a}_i)] + (1 - b_i) \log [1 - \psi(\mathbf{w}^\top \mathbf{a}_i)] + \chi \sum_{j=1}^d \frac{w_j^2}{1 + w_j^2},$$

where  $\psi(x) = \frac{1}{1 + \exp(-x)}$ ,  $\mathbf{w} = (w_1 \cdots w_d)^\top$  and  $\chi = 0.1$ . The empirical loss for the non-convex robust

linear regression is formulated as

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \varrho(b_i - \mathbf{w}^\top \mathbf{a}_i),$$

where  $\varrho(x) = \log\left(\frac{x^2}{2} + 1\right)$ .

The trajectories of gradient norm during the iterations are displayed in Figure 1, where gradient norm is the most popular measurement for convergence [4, 17]. We would like to emphasize that the Hessian is (nearly) semi-positive definite in all experiments if algorithms converge. Due to limited spaces, we put trajectories of the empirical loss in the supplement. We merely show the gradient norm versus iteration steps rather than time (seconds), since the main computation in each step is from the CRS. The proposed ARCm achieves the best performances in these examples, being 2-10 times faster than CRm and 10%-50% faster than ARC in terms of iterations for convergence, except for the logistic regression on the Covtype dataset where ARCm and ARC have similar results. We do not expect the proposed ARCm significantly outperforms ARC if the problem is less challenging and ARC can easily find the optima (Figure 1 (b)). Moreover, we also highlight the two cases of momentum in Theorem 2.1. When the momentum term helps the convergence of ARCm (i.e.,  $f(\mathbf{z}_{k+1}) < f(\mathbf{y}_{k+1})$ ), we point it in red and gray if  $\mathbf{s}_k^\top \mathbf{v}_{k-1} > 0$  and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} < 0$  respectively. We observe that the momentum in opposite direction helps the convergence of ARCm and it usually occurs at the

beginning of the algorithm, which is within our expectation.

## 5 Conclusion and Future Works

In this paper, we propose the momentum accelerated ARC (ARCM) that improves the performance of ARC. The global convergence and the local convergence under the KL property are theoretically studied. Furthermore, we analyze a more practical case where inexact CRS solutions are obtained in ARCM. Experimental results show that the proposed ARCM significantly outperforms some state-of-the-art second-order methods in solving non-convex logistic regression and robust linear regression models. There are still many unsolved problems for cubic regularization methods. Firstly, better strategies for CR-based and ARC-based algorithms. For example, some works modify the criterion  $\rho_k$  in ARC. Secondly, fast solvers of CRS. Some methods enforce structured Hessian in CRS to reduce the computation.

## References

- [1] Jock A Blackard and Denis J Dean. Comparative accuracies of artificial neural networks and discriminant analysis in predicting forest cover types from cartographic variables. *Computers and Electronics in Agriculture*, 24(3):131–151, 1999.
- [2] Yair Carmon and John Duchi. Gradient descent finds the cubic-regularized nonconvex newton step. *SIAM Journal on Optimization*, 29(3):2146–2178, 2019.
- [3] Yair Carmon and John C Duchi. Analysis of krylov subspace solutions of regularized non-convex quadratic problems. *Advances in Neural Information Processing Systems*, 31, 2018.
- [4] Coralia Cartis, Nicholas IM Gould, and Philippe L Toint. Adaptive cubic regularisation methods for unconstrained optimization. part i: motivation, convergence and numerical results. *Mathematical Programming*, 127(2):245–295, 2011.
- [5] Andrew R Conn, Nicholas IM Gould, and Philippe L Toint. *Trust region methods*. SIAM, 2000.
- [6] Katerine Diaz-Chito, Aura Hernández-Sabaté, and Antonio M López. A reduced feature set for driver head pose estimation. *Applied Soft Computing*, 45:98–107, 2016.
- [7] Yihang Gao, Man-chung Yue, and Michael K. Ng. Approximate secular equations for the cubic regularization subproblem. *OpenReview, to appear in Advances in Neural Information Processing Systems*, 2022.
- [8] Rong Ge, Jason D Lee, and Tengyu Ma. Matrix completion has no spurious local minimum. *Advances in Neural Information Processing Systems*, 29, 2016.
- [9] Andreas Griewank. The modification of newton’s method for unconstrained optimization by bounding cubic terms. Technical report, Technical Report NA/12, 1981.
- [10] Isabelle Guyon, Steve Gunn, Asa Ben-Hur, and Gideon Dror. Result analysis of the nips 2003 feature selection challenge. *Advances in Neural Information Processing Systems*, 17, 2004.
- [11] Kevin Huang, Junyu Zhang, and Shuzhong Zhang. Cubic regularized newton method for the saddle point models: A global and local convergence analysis. *Journal of Scientific Computing*, 91(2):1–31, 2022.
- [12] Rujun Jiang, Man-Chung Yue, and Zhishuo Zhou. An accelerated first-order method with complexity analysis for solving cubic regularization subproblems. *Computational Optimization and Applications*, 79(2):471–506, 2021.
- [13] Jonas Moritz Kohler and Aurelien Lucchi. Subsampled cubic regularization for non-convex optimization. In *International Conference on Machine Learning*, pages 1895–1904. PMLR, 2017.
- [14] Dong-Hui Li, Masao Fukushima, Liqun Qi, and Nobuo Yamashita. Regularized newton methods for convex minimization problems with singular solutions. *Computational Optimization and Applications*, 28(2):131–147, 2004.
- [15] Yurii Nesterov and Boris T Polyak. Cubic regularization of newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- [16] Ju Sun, Qing Qu, and John Wright. A geometrical analysis of phase retrieval. In *International Symposium on Information Theory*, 2016.
- [17] Zhe Wang, Yi Zhou, Yingbin Liang, and Guanghui Lan. Cubic regularization with momentum for non-convex optimization. In *Uncertainty in Artificial Intelligence*, pages 313–322. PMLR, 2020.
- [18] Martin Weiser, Peter Deuffhard, and Bodo Erdmann. Affine conjugate adaptive Newton methods for nonlinear elastomechanics. *Optimisation Methods and Software*, 22(3):413–431, 2007.
- [19] Man-Chung Yue, Zirui Zhou, and Anthony Man-Cho So. On the quadratic convergence of the cubic regularization method under a local error bound condition. *SIAM Journal on Optimization*, 29(1):904–932, 2019.
- [20] Dongruo Zhou, Pan Xu, and Quanquan Gu. Stochastic variance-reduced cubic regularized newton methods. In *International Conference on Machine Learning*, pages 5990–5999. PMLR, 2018.
- [21] Yi Zhou, Zhe Wang, and Yingbin Liang. Convergence of cubic regularization for nonconvex optimization under kl property. *Advances in Neural Information Processing Systems*, 31, 2018.

## Supplementary Materials

### A Some Technical Proofs

**A.1 Proof for Lemma 2.1** For any  $\mathbf{x}, \mathbf{y} \in \mathcal{F}$ , we have

$$\begin{aligned}
 (A.1) \quad & \left\| \nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \right\|_2 \\
 &= \left\| \int_0^1 [\nabla^2 f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla^2 f(\mathbf{x})] (\mathbf{y} - \mathbf{x}) dt \right\|_2 \\
 &\leq L_H \|\mathbf{y} - \mathbf{x}\|_2^2 \int_0^1 t dt = \frac{L_H}{2} \|\mathbf{y} - \mathbf{x}\|_2^2,
 \end{aligned}$$

and

$$\begin{aligned}
 (A.2) \quad & \left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x}_k)^\top (\mathbf{y} - \mathbf{x}) \right. \\
 & \quad \left. - \frac{1}{2} (\mathbf{y} - \mathbf{x})^\top \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \right| \\
 &= \left| \int_0^1 (\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}) \right. \\
 & \quad \left. - t \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x}) \right)^\top (\mathbf{y} - \mathbf{x}) dt \right| \\
 &\leq \frac{L_H}{2} \|\mathbf{y} - \mathbf{x}\|_2^3 \cdot \int_0^1 t^2 dt = \frac{L_H}{6} \|\mathbf{y} - \mathbf{x}\|_2^3.
 \end{aligned}$$

Then

$$f(\mathbf{x}_k + \mathbf{s}_k) - m_k(\mathbf{s}_k) \leq \left( \frac{L_H}{6} - \frac{\sigma_k}{6} \right) \|\mathbf{s}_k\|_2^3.$$

Let  $r_k := f(\mathbf{x}_k + \mathbf{s}_k) - m_k(\mathbf{s}_k) + (1 - \eta_2)(m_k(\mathbf{s}_k) - f(\mathbf{x}_k))$ . The very successful update occurs (i.e.,  $\rho_k > \eta_2$ ) if and only if  $r_k < 0$ . Note that  $m_k(\mathbf{s}_k) - f(\mathbf{x}_k) < 0$  if  $\mathbf{x}_k$  is not a local minimizer. If  $\sigma_k \geq L_H$ , then  $f(\mathbf{x}_k + \mathbf{s}_k) - m_k(\mathbf{s}_k) \leq \left( \frac{L_H}{6} - \frac{\sigma_k}{6} \right) \|\mathbf{s}_k\|_2^3 \leq 0$  and  $r_k < 0$ . Suppose that  $k-1 \notin \mathcal{S}$  and  $\sigma_{k-1} \approx L_H$  (but  $\sigma_{k-1} < L_H$ ), then  $L_H < \sigma_k = \gamma_1 \sigma_{k-1} \leq \gamma_1 L_H$  and  $k \in \mathcal{S}$ .

**A.2 Proof for Lemma 2.5** Note that the sequence  $\{f(\mathbf{x}_k)\}$  is lower bounded by  $f^*$  and is monotonically decreasing, then it must be convergent. Assumption 2.2 and the non-increasing property of  $\{f(\mathbf{x}_k)\}$  imply the boundedness of  $\{\mathbf{x}_k\}_{k=0}^{+\infty}$ . Here the boundedness of the sequence  $\{\mathbf{x}_k\}_{k=0}^{+\infty}$  implies that the set  $\bar{\mathcal{X}}$  is non-empty.

Using Lemma 2.3, we have for all  $k \in \mathcal{S}$

$$\begin{aligned}
 & \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 \\
 &\leq \|\mathbf{x}_{k+1} - \mathbf{y}_{k+1}\|_2 + \|\mathbf{y}_{k+1} - \mathbf{x}_k\|_2 \\
 &\leq \beta_k \|\mathbf{v}_k\|_2 + \|\mathbf{s}_k\|_2 \\
 &\leq \alpha_1 \|\mathbf{s}_k\|_2 \cdot \|\mathbf{v}_k\|_2 + \|\mathbf{s}_k\|_2 \\
 &\leq \left( \alpha_1 \frac{1}{1 - \tau} \left( \frac{12(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3} + 1 \right) \|\mathbf{s}_k\|_2 \\
 &:= c_7 \|\mathbf{s}_k\|_2.
 \end{aligned}$$

Lemma 2.2 and 2.3 jointly show that  $\lim_{k \rightarrow \infty, k \in \mathcal{S}} \|\mathbf{s}_k\|_2 = 0$  and thus  $\lim_{k \rightarrow \infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 = 0$  (note that  $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 = 0$  if  $k \notin \mathcal{S}$ ). Using Lemma 2.4, we have

$$\begin{aligned}
 \|\nabla f(\bar{\mathbf{x}})\| &\leq \limsup_{k \rightarrow \infty, k \in \mathcal{S}} \|\nabla f(\mathbf{x}_{k+1})\| \\
 &\leq \limsup_{k \rightarrow \infty, k \in \mathcal{S}} c_1 \|\mathbf{s}_k\|_2^2 = 0
 \end{aligned}$$

and

$$\begin{aligned}
 \lambda_{\min}(\nabla^2 f(\bar{\mathbf{x}})) &\geq \limsup_{k \rightarrow \infty, k \in \mathcal{S}} \lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1})) \\
 &\geq \limsup_{k \rightarrow \infty, k \in \mathcal{S}} -c_2 \|\mathbf{s}_k\|_2 = 0.
 \end{aligned}$$

Define  $e_j = f(\mathbf{x}_{k_j+1}) - \bar{f}$ , where  $k_j \in \mathcal{S}$  is the  $j$ -th element in  $\mathcal{S}$ . If  $f(\cdot)$  satisfies the KL property, then there exist a large enough  $j_0$  such that  $\phi'(e_j) \geq \frac{1}{\|\nabla f(\mathbf{x}_{k_j+1})\|}$  and  $1 - \alpha_2 \|\mathbf{s}_{k_j}\|_2 \|\mathbf{v}_{k_j}\|_2 \geq \alpha_3$  for all  $j \geq j_0$  and some  $1 > \alpha_3 > 0$ , since  $\lim_{j \rightarrow \infty} \|\mathbf{s}_{k_j}\|_2 = 0$  and  $\|\mathbf{v}_{k_j}\|_2$  is bounded. Furthermore,

$$\begin{aligned}
 (A.3) \quad & \|\mathbf{x}_{k_j+1} - \mathbf{x}_{k_j}\|_2 \geq \|\mathbf{s}_{k_j}\|_2 - \beta_{k_j} \|\mathbf{v}_{k_j}\|_2 \\
 &\geq \|\mathbf{s}_{k_j}\|_2 - \alpha_2 \|\mathbf{s}_{k_j}\|_2^2 \|\mathbf{v}_{k_j}\|_2 \\
 &= (1 - \alpha_2 \|\mathbf{s}_{k_j}\|_2 \|\mathbf{v}_{k_j}\|_2) \|\mathbf{s}_{k_j}\|_2 \\
 &\geq \alpha_3 \|\mathbf{s}_{k_j}\|_2,
 \end{aligned}$$

$$\begin{aligned}
 \|\mathbf{x}_{k_j+1} - \mathbf{x}_{k_j}\|_2 &\leq \|\mathbf{s}_{k_j}\|_2 + \beta_{k_j} \|\mathbf{v}_{k_j}\|_2 \\
 &\leq \|\mathbf{s}_{k_j}\|_2 + \alpha_2 \|\mathbf{s}_{k_j}\|_2^2 \|\mathbf{v}_{k_j}\|_2 \\
 &= (1 + \alpha_2 \|\mathbf{s}_{k_j}\|_2 \|\mathbf{v}_{k_j}\|_2) \|\mathbf{s}_{k_j}\|_2 \\
 &\leq (2 - \alpha_3) \|\mathbf{s}_{k_j}\|_2,
 \end{aligned}$$

and

$$\phi'(e_j) \geq \frac{1}{\|\nabla f(\mathbf{x}_{k_j+1})\|} \geq \frac{1}{c_1 \|\mathbf{s}_{k_j}\|_2^2} \geq \frac{\alpha_3^2}{c_1 \|\mathbf{x}_{k_j+1} - \mathbf{x}_{k_j}\|_2^2}.$$

Note that  $\phi(\cdot)$  is concave when  $\theta < 1$  and  $e_{j+1} < e_j$ , we

have

$$\begin{aligned}
& \phi(e_j) - \phi(e_{j+1}) \\
& \geq \phi'(e_j)(e_j - e_{j+1}) \\
& \geq \frac{\alpha_3^2}{c_1 \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2^2} \cdot \frac{1}{12} \eta_1 \sigma_{\min} \|\mathbf{s}_{k_{j+1}}\|_2^3 \\
& \geq \frac{\alpha_3^2 \eta_1 \sigma_{\min}}{12 c_1 (2 - \alpha_3)} \frac{\|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2^3}{\|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2^2} \\
& := c_8 \frac{\|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2^3}{\|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2^2}.
\end{aligned} \tag{A.4}$$

Then by the Hölder's inequality, we have

$$\begin{aligned}
& \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \\
& \leq c_9 \cdot \sum_{j=j_0}^J (\phi(e_j) - \phi(e_{j+1}))^{1/3} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2^{2/3} \\
& \leq c_9 \left( \sum_{j=j_0}^J \phi(e_j) - \phi(e_{j+1}) \right)^{1/3} \left( \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 \right)^{2/3} \text{ and equivalently} \\
& \leq c_9 \phi(e_{j_0})^{1/3} \cdot \left( \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 \right)^{2/3},
\end{aligned}$$

where  $c_9 = c_8^{-1/3}$ . If  $\sum_{j=j_0}^{+\infty} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 = +\infty$ , then we may let  $J > j_0$  to be large enough such that  $\sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 \approx \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \gg 1$ . However the relation (A.2) is violated. Therefore  $\sum_{j=j_0}^{+\infty} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 < +\infty$  and thus  $\sum_{k=0}^{+\infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\|_2 < +\infty$ , which implies that  $\{\mathbf{x}_k\}_{k=0}^{+\infty}$  is a Cauchy sequence and  $\mathcal{X} = \{\bar{\mathbf{x}}\}$  is a singleton.

**A.3 Proof for Theorem 2.3** First note that if  $f(\mathbf{x})$  satisfies the KL property, then it also satisfies the local KL -error bound [21, Proposition 1], which is a generalization of the local error bound in [19], i.e.,

$$\|\mathbf{x}_{k_j+1} - \bar{\mathbf{x}}\|_2 \leq \kappa \|\nabla f(\mathbf{x}_{k_{j+1}})\|_2^{\frac{\theta}{1-\theta}} \tag{A.5}$$

holds for  $j > j_0$  with large enough  $j_0 > 0$  and some  $\kappa > 0$ . Combining (A.5) with (A.3) and Lemma 2.4, we have

$$\begin{aligned}
\|\mathbf{x}_{k_j+1} - \bar{\mathbf{x}}\|_2 & \leq \kappa \|\nabla f(\mathbf{x}_{k_{j+1}})\|_2^{\frac{\theta}{1-\theta}} \\
& \leq \kappa (c_1 \|\mathbf{s}_{k_j}\|_2^2)^{\frac{\theta}{1-\theta}}
\end{aligned} \tag{A.6}$$

Note that [19, Lemma 1] shows that there exists  $c_{10} > 0$  such that

$$\|\mathbf{s}_{k_j}\|_2 \leq c_{10} \|\mathbf{x}_{k_j} - \bar{\mathbf{x}}\|_2, \tag{A.7}$$

where  $c_{10} = \left( 1 + \frac{L_H}{\sigma_{\min}} + \sqrt{\left( 1 + \frac{L_H}{\sigma_{\min}} \right)^2 + \frac{L_H}{\sigma_{\min}}} \right)$ .

Equations (A.6) and (A.7) jointly imply that

$$\|\mathbf{x}_{k_{j+1}} - \bar{\mathbf{x}}\|_2 \leq c_{11} \|\mathbf{x}_{k_j} - \bar{\mathbf{x}}\|_2^{\frac{2\theta}{1-\theta}} = c_{11} \|\mathbf{x}_{k_{j-1}+1} - \bar{\mathbf{x}}\|_2^{\frac{2\theta}{1-\theta}},$$

where the convergence is super linear since  $\frac{2\theta}{1-\theta} > 1$  when  $\theta \in (\frac{1}{3}, 1)$ . Without the loss of generality we assume that  $c_{11} \leq 1$  and  $\|\mathbf{x}_{k_{j_0}+1} - \bar{\mathbf{x}}\|_2 \leq \exp(-1)$ , we conclude the results in item 1 for  $\theta \in (\frac{1}{3}, 1)$ .

We then discuss the cases when  $\theta = \frac{1}{3}$  and  $\theta \in (0, \frac{1}{3})$ . Here we adopt similar techniques in [21, Theorem 4]. Fix  $\nu \in (0, 1)$  and  $j > j_0$  for a large enough  $j_0 \in \mathcal{S}$ . If  $\|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \geq \nu \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2$ , then using (A.4) we have

$$\begin{aligned}
\phi(e_j) - \phi(e_{j+1}) & \geq c_8 \frac{\|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2^3}{\|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2^2} \\
& \geq c_8 \nu^2 \|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2
\end{aligned}$$

$$\|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \leq \frac{1}{c_8 \nu^2} (\phi(e_j) - \phi(e_{j+1})).$$

Otherwise,  $\|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \leq \nu \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2$ . Summing up two inequalities, we have

$$\begin{aligned}
& \|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \\
& \leq \frac{1}{c_8 \nu^2} (\phi(e_j) - \phi(e_{j+1})) + \nu \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2.
\end{aligned} \tag{A.8}$$

Let  $\Delta_j = \sum_{l=j}^{\infty} \|\mathbf{x}_{k_{l+1}} - \mathbf{x}_{k_l}\|_2$ , then  $\|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 = \Delta_j - \Delta_{j+1}$ . Summing up (A.8) from  $j_0$  to  $J$  yields that

$$\begin{aligned}
& \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \\
& \leq \nu \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 + \frac{1}{c_8 \nu^2} \phi(e_{j_0}).
\end{aligned}$$

and then

$$\begin{aligned}
& \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 \\
& \leq \nu \left( \sum_{j=j_0}^J \|\mathbf{x}_{k_{j+1}+1} - \mathbf{x}_{k_{j+1}}\|_2 + \|\mathbf{x}_{k_{j_0}+1} - \mathbf{x}_{k_{j_0}}\|_2 \right) \\
& \quad + \frac{1}{c_8 \nu^2} \phi(e_{j_0}).
\end{aligned}$$

Let  $J \rightarrow \infty$  and we can further simplify the above

inequality as

$$\begin{aligned}
& \Delta_{j+1} \\
& \leq \frac{\nu}{1-\nu} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 + \frac{1}{c_8 \nu^2 (1-\nu)} \phi(e_j) \\
& = \frac{\nu}{1-\nu} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 + \frac{c}{c_8 \nu^2 (1-\nu) \theta} e_j^\theta \\
& \leq \frac{\nu}{1-\nu} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 + \frac{c}{c_8 \nu^2 (1-\nu) \theta} \left( c_0 \|\nabla f(\mathbf{x}_{k_{j+1}})\|^{1-\frac{1}{\theta}} \right)^\theta \\
& \leq \frac{\nu}{1-\nu} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 + \frac{c}{c_8 \nu^2 (1-\nu) \theta} \left( c_0 \left( c_1 \|\mathbf{s}_{k_j}\|_2^2 \right)^{\frac{1}{1-\theta}} \right)^\theta \\
& \leq \frac{\nu}{1-\nu} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2 + c_{12} \|\mathbf{x}_{k_{j+1}} - \mathbf{x}_{k_j}\|_2^{\frac{2\theta}{1-\theta}} \\
& = \frac{\nu}{1-\nu} (\Delta_j - \Delta_{j+1}) + c_{12} (\Delta_j - \Delta_{j+1})^{\frac{2\theta}{1-\theta}}.
\end{aligned}$$

If  $\theta = \frac{1}{3}$ , then  $\Delta_{j+1} \leq c_{13} (\Delta_j - \Delta_{j+1})$  and thus  $\Delta_{j+1} \leq \frac{c_{13}}{1+c_{13}} \Delta_j$ . Therefore,  $\Delta_j \leq \left( \frac{c_{13}}{1+c_{13}} \right)^{j-j_0} \Delta_{j_0}$  and  $\|\mathbf{x}_{k_{j+1}} - \bar{\mathbf{x}}\|_2 \leq \Delta_{j+1} = \mathcal{O}(\exp(-c_{14}(j-j_0)))$ , where  $c_{14} = \log\left(\frac{1+c_{13}}{c_{13}}\right)$ . If  $\theta \in (0, \frac{1}{3})$ , then  $(\Delta_j - \Delta_{j+1})^{\frac{2\theta}{1-\theta}}$  is dominant term of the right hand side if  $j_0$  is large enough. Then there exists  $c_{15} > 0$  such that  $\Delta_{j+1} \leq c_{15} (\Delta_j - \Delta_{j+1})$ . Define  $h(t) = t^{-\frac{1-\theta}{2\theta}}$  and a fixed constant  $\omega > 1$ . If  $h(\Delta_{k+1}) \leq \omega h(\Delta_k)$ , then

$$\begin{aligned}
1 & \leq c_{15} (\Delta_j - \Delta_{j+1}) h(\Delta_{j+1}) \leq c_{15} \omega (\Delta_j - \Delta_{j+1}) h(\Delta_j) \\
& \leq c_{15} \omega \int_{\Delta_{j+1}}^{\Delta_j} h(t) dt = c_{15} \omega \frac{2\theta}{3\theta-1} \left( \Delta_j^{\frac{3\theta-1}{2\theta}} - \Delta_{j+1}^{\frac{3\theta-1}{2\theta}} \right)
\end{aligned}$$

and

$$(A.9) \quad \Delta_{j+1}^{\frac{3\theta-1}{2\theta}} - \Delta_j^{\frac{3\theta-1}{2\theta}} \geq \frac{1-3\theta}{2\omega\theta c_{15}}.$$

On the other hand, if  $h(\Delta_{k+1}) > \omega h(\Delta_k)$ , then  $\Delta_{j+1}^{\frac{3\theta-1}{2\theta}} > \omega^{\frac{1-3\theta}{1-\theta}} \Delta_j^{\frac{3\theta-1}{2\theta}}$  and  $\Delta_{j+1}^{\frac{3\theta-1}{2\theta}} - \Delta_j^{\frac{3\theta-1}{2\theta}} > \left( \omega^{\frac{1-3\theta}{1-\theta}} - 1 \right) \Delta_j^{\frac{3\theta-1}{2\theta}}$ .

Note that  $\omega^{\frac{1-3\theta}{1-\theta}} - 1 > 0$  and  $\lim_{j \rightarrow \infty} \Delta_j^{\frac{3\theta-1}{2\theta}} = +\infty$ , then there must exist large enough  $j_0$  such that  $\left( \omega^{\frac{1-3\theta}{1-\theta}} - 1 \right) \Delta_j^{\frac{3\theta-1}{2\theta}} \geq \frac{1-3\theta}{2\omega\theta c_{15}}$ . Therefore, the inequality (A.9) holds for all  $j > j_0$ . Summing up (A.9) from  $j = j_0$  we have that

$$\Delta_{j+1}^{\frac{3\theta-1}{2\theta}} - \Delta_{j_0}^{\frac{3\theta-1}{2\theta}} \geq \frac{1-3\theta}{2\omega\theta c_{15}} (j+1-j_0)$$

and

$$\begin{aligned}
\Delta_{j+1} & \leq \left( \Delta_{j_0}^{\frac{3\theta-1}{2\theta}} + \frac{1-3\theta}{2\omega\theta c_{15}} (j+1-j_0) \right)^{-\frac{2\theta}{1-3\theta}} \\
& = \mathcal{O}\left( (c_{16} (j-j_0))^{-\frac{2\theta}{1-3\theta}} \right),
\end{aligned}$$

which completes the proof since  $\|\mathbf{x}_{k_{j+1}} - \bar{\mathbf{x}}\|_2 \leq \Delta_{j+1}$ .

**A.4 Proof for Lemma 3.1** A similar analysis is conducted here as Lemma 2.1. Let  $\tilde{r}_k = f(\mathbf{x}_k + \tilde{\mathbf{s}}_k) - m_k(\tilde{\mathbf{s}}_k) + (1-\eta_2)(m_k(\tilde{\mathbf{s}}_k) - f(\mathbf{x}_k))$ . According to (A.1), we have

$$f(\mathbf{x}_k + \tilde{\mathbf{s}}_k) - m_k(\tilde{\mathbf{s}}_k) \leq \left( \frac{L_H}{6} - \frac{\sigma_k}{6} \right) \|\tilde{\mathbf{s}}_k\|_2^3.$$

The proof can be similarly developed by using the argument in Lemma 2.1.

**A.5 Proof for Lemma 3.2** As an analogy to Lemma 2.3, we provide only the key steps here and omit some details. The inequality (2.9) can be similarly developed as

$$\sum_{k \in \mathcal{S}_T} \eta_1 \left( \frac{\sigma_k}{12} \|\tilde{\mathbf{s}}_k\|_2^3 - \delta_k \right) \leq f(\mathbf{x}_0) - f^*,$$

then

$$\max_{k \in \mathcal{S}_T} \frac{\sigma_k}{12} \|\tilde{\mathbf{s}}_k\|_2^3 \leq \frac{f(\mathbf{x}_0) - f^*}{\eta_1} + \delta_k \leq \frac{2(f(\mathbf{x}_0) - f^*)}{\eta_1}$$

and

$$\min_{k \in \mathcal{S}_T} \frac{\sigma_k}{12} \|\tilde{\mathbf{s}}_k\|_2^3 - \delta_k \leq \frac{f(\mathbf{x}_0) - f^*}{\eta_1 |\mathcal{S}_T|}.$$

**A.6 Proof for Lemma 3.3** We first derive the error bound for  $\|\nabla f(\mathbf{y}_{k+1})\|$  and  $\lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1}))$ :

(A.10)

$$\begin{aligned}
& \|\nabla f(\mathbf{y}_{k+1})\|_2 \\
& \leq \|\nabla f(\mathbf{y}_{k+1}) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \tilde{\mathbf{s}}_k\|_2 \\
& \quad + \left\| \nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) \tilde{\mathbf{s}}_k + \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2 \tilde{\mathbf{s}}_k \right\|_2 + \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2^2 \\
& \leq \frac{L_H}{2} \|\tilde{\mathbf{s}}_k\|_2^2 + \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2 + \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2^2 \\
& \leq \left( \frac{\sigma_{\max}}{2} + \frac{L_H}{2} \right) \|\tilde{\mathbf{s}}_k\|_2^2 + \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2,
\end{aligned}$$

where the second inequality is due to (A.1); and

$$\begin{aligned}
& \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) \\
& \geq \lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) - \|\nabla^2 f(\mathbf{y}_{k+1}) - \nabla^2 f(\mathbf{x}_k)\|_2 \\
& \geq -\frac{\sigma_k}{2} \|\mathbf{s}_k\|_2 - L_H \|\tilde{\mathbf{s}}_k\|_2 \\
& \geq -\frac{\sigma_k}{2} \|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2 - \frac{\sigma_k}{2} \|\tilde{\mathbf{s}}_k\|_2 - L_H \|\tilde{\mathbf{s}}_k\|_2 \\
& = -\left( \frac{\sigma_{\max}}{2} - L_H \right) \|\tilde{\mathbf{s}}_k\|_2 - \frac{\sigma_{\max}}{2} \|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2.
\end{aligned}$$

where the second inequality comes from a well-known result of cubic regularization [15, Proposition 1]. Now,

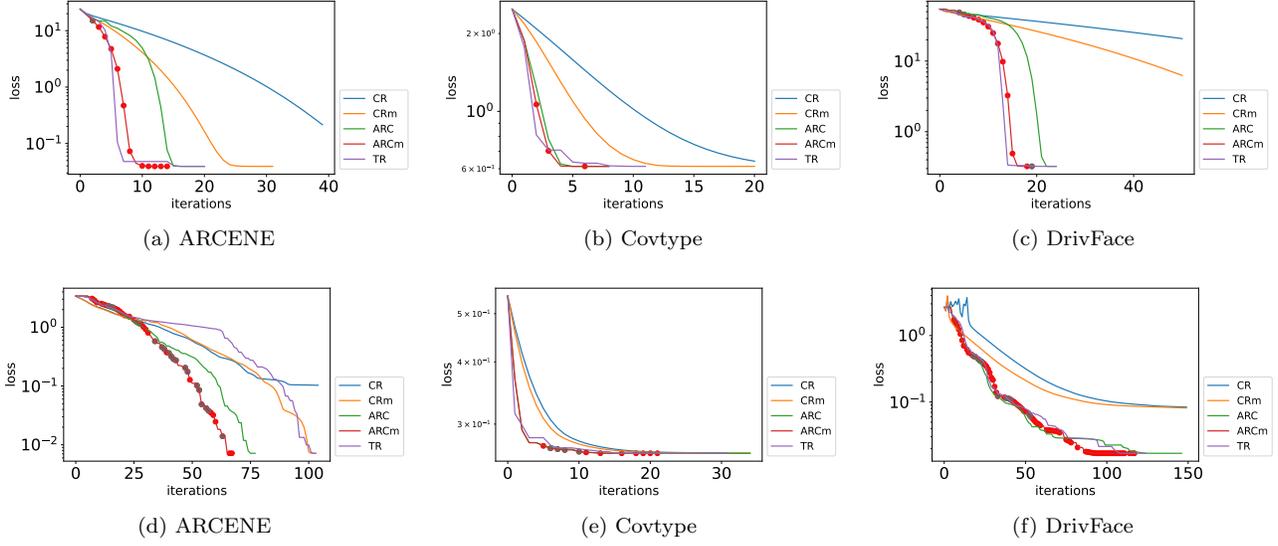


Figure 2: Loss versus the iteration steps on logistic regression (the first row) and robust linear regression (the second row) models for three datasets. Red dots:  $\beta_k > 0$  and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} > 0$ ; gray dots:  $\beta_k > 0$  and  $\mathbf{s}_k^\top \mathbf{v}_{k-1} < 0$ .

we are ready to develop the error bound  $\|\nabla f(\mathbf{x}_{k+1})\|$  and  $\lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1}))$ . Using similar arguments as in Lemma 2.4, we have

$$\begin{aligned}
& \|\nabla f(\mathbf{x}_{k+1})\|_2 \\
& \leq \|\nabla f(\mathbf{y}_{k+1})\|_2 + \|\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{y}_{k+1})\|_2 \\
& \leq \|\nabla f(\mathbf{y}_{k+1})\|_2 + L_g \tilde{\beta}_k \|\tilde{\mathbf{v}}_k\|_2 \\
& \leq \left( \frac{\sigma_{\max}}{2} + \frac{L_H}{2} \right) \|\tilde{\mathbf{s}}_k\|_2^2 + \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2 \\
& \quad + L_g \alpha_2 \|\tilde{\mathbf{s}}_k\|_2^2 \cdot \|\tilde{\mathbf{v}}_k\|_2 \\
& \leq c_5 \|\tilde{\mathbf{s}}_k\|_2^2 + \|\nabla m_k(\tilde{\mathbf{s}}_k)\|_2,
\end{aligned}$$

and

$$\begin{aligned}
& \lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1})) \\
& \geq \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) - \|\nabla^2 f(\mathbf{x}_{k+1}) - \nabla^2 f(\mathbf{y}_{k+1})\|_2 \\
& \geq \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) - L_H \|\mathbf{x}_{k+1} - \mathbf{y}_{k+1}\|_2 \\
& \geq \lambda_{\min}(\nabla^2 f(\mathbf{y}_{k+1})) - L_H \tilde{\beta}_k \|\tilde{\mathbf{v}}_k\|_2 \\
& \geq - \left( \frac{1}{2} \sigma_{\max} + L_H \right) \|\tilde{\mathbf{s}}_k\|_2 - \frac{\sigma_{\max}}{2} \left| \|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2 \right| \\
& \quad - L_H \alpha_1 \|\tilde{\mathbf{s}}_k\|_2 \cdot \|\tilde{\mathbf{v}}_k\|_2 \\
& \geq -c_6 \|\tilde{\mathbf{s}}_k\|_2 - \frac{\sigma_{\max}}{2} \left| \|\mathbf{s}_k\|_2 - \|\tilde{\mathbf{s}}_k\|_2 \right|,
\end{aligned}$$

where  $c_5 = \frac{1}{2} \max\{L_H \gamma_1, \sigma_{\min}\} + \frac{1}{2} L_H + \frac{\alpha_2 L_g}{1-\tau} \left( \frac{24(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3}$  and  $c_6 = \frac{1}{2} \sigma_{\max} + L_H + \frac{\alpha_1 L_H}{1-\tau} \left( \frac{24(f(\mathbf{x}_0) - f^*)}{\eta_1 \sigma_{\min}} \right)^{1/3}$ .

## B Additional Experimental Results

The trajectories of the empirical loss during the iterations are displayed in Figure 2, where we have found in Figure 1 that the momentum in opposite direction helps the convergence of ARCm at the beginning of the algorithm.