# Cubic Regularization with Momentum for Nonconvex Optimization

**Zhe Wang**
EECS Dept.
Ohio State University
wang.10982@osu.edu

**Yi Zhou**
EECS Dept.
Duke University
yi.zhou610@duke.edu

**Yingbin Liang**
EECS Dept.
Ohio State University
liang.889@osu.edu

**Guanghui Lan**
ISyE Dept.
Georgia Institute of Technology
george.lan@isye.gatech.edu

## Abstract

Momentum is a popular technique to accelerate the convergence in practical training, and its impact on convergence guarantee has been well-studied for first-order algorithms. However, such a successful acceleration technique has not yet been proposed for second-order algorithms in nonconvex optimization. In this paper, we apply the momentum scheme to cubic regularized (CR) Newton's method and explore the potential for acceleration. Our numerical experiments on various nonconvex optimization problems demonstrate that the momentum scheme can substantially facilitate the convergence of cubic regularization, and perform even better than the Nesterov's acceleration scheme for CR. Theoretically, we prove that CR under momentum achieves the best possible convergence rate to a second-order stationary point for nonconvex optimization. Moreover, we study the proposed algorithm for solving problems satisfying an error bound condition and establish a local quadratic convergence rate. Then, particularly for finite-sum problems, we show that the proposed algorithm can allow computational inexactness that reduces the overall sample complexity without degrading the convergence rate.

## 1 INTRODUCTION

In the era of machine learning, deep models such as neural networks have achieved great success in solving a variety of challenging tasks. However, training deep models is in general a difficult task and traditional first-order algorithms can easily get stuck at sub-optimal points such as saddle points, which have been shown to bottleneck the performance of practical training (Dauphin et al., 2014). Motivated by this, there is a rising interest in designing algorithms that can escape saddle points in general nonconvex optimization, and the cubic regularization (CR) Newton's method is such a type of popular optimization algorithm.

More specifically, consider the following generic nonconvex optimization problem.

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}), \tag{1}$$

where $f : \mathbb{R}^d \to \mathbb{R}$ is a twice-differentiable and nonconvex function. The CR algorithm (Nesterov and Polyak, 2006) takes an initialization $\mathbf{x}_0 \in \mathbb{R}^d$, a proper parameter $M > 0$, and generates a sequence $\{\mathbf{x}_k\}_k$ for solving eq. (1) via the following update rule.

$$\mathbf{s}_{k+1} = \operatorname*{argmin}_{\mathbf{s} \in \mathbb{R}^d} \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2}\mathbf{s}^\top \nabla^2 f(\mathbf{x}_k)\mathbf{s} + \frac{M}{6}\|\mathbf{s}\|^3,$$

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

Intuitively, the main step of CR solves a cubic minimization subproblem that is formulated by the second-order Taylor expansion at the current iterate with a cubic regularizer. Such a cubic subproblem can be efficiently solved by many dedicated solvers (Cartis et al., 2011a; Carmon and Duchi, 2016; Agarwal et al., 2017) that induce a low overall computation complexity (see Section 4.3 for further elaboration). By exploiting second order information (i.e., gradient and Hessian) of the objective function, the CR algorithm has been shown to produce a solution $\mathbf{x}$ that satisfies the $\epsilon$-second-order stationary condition, i.e.,

$$\|\nabla f(\mathbf{x})\| \leqslant \epsilon \quad \text{and} \quad \lambda_{\min}\big(\nabla^2 f(\mathbf{x})\big) \geqslant -\sqrt{\epsilon}, \tag{2}$$

where $\lambda_{\min}\big(\nabla^2 f(\mathbf{x})\big)$ denotes the minimum eigenvalue of the Hessian $\nabla^2 f(\mathbf{x})$. Unlike the first-order stationary condition (i.e., $\|\nabla f(\mathbf{x})\| \leqslant \epsilon$) which does not rule out the possibility of converging to a saddle point, the second-order stationary condition requires the corresponding Hessian to be almost positive semidefinite and hence can

avoid convergence to strict saddle points (i.e., at which Hessian has negative eigenvalue). In particular, a variety of nonconvex machine learning problems such as phase retrieval (Sun et al., 2017), dictionary learning (Sun et al., 2015) and tensor decomposition (Ge et al., 2015) have been shown to have only strict saddle points. Therefore, CR is guaranteed to escape all the saddle points and converge to a local minimum in solving these problems.

While most existing studies on the CR algorithm focus on reducing the computation complexity by various sampling schemes, e.g., mini-batch sampling (Xu et al., 2017), sub-sampling (Kohler and Lucchi, 2017), variance-reduced sampling (Wang et al., 2019b; Zhou et al., 2018), less attention has been paid to the design of new schemes for accelerating CR. The only exception is Nesterov (2008), where an acceleration scheme was proposed for CR, but has been shown to achieve a faster convergence rate than CR only for *convex* problems. Such an accelerated scheme consists of hyperparameters that are fine-tuned in the context of convex optimization, and hence may not guarantee to produce a second-order stationary solution in *nonconvex* optimization. There does not exist any accelerated CR algorithm that has provable convergence for nonconvex optimization. Therefore, the aim of this paper is to design a momentum-based scheme for CR with provable second-order stationary convergence guarantee for nonconvex optimization as well as yielding faster convergence in practical scenarios.

## 1.1 OUR CONTRIBUTIONS

Our major contribution lies in proposing the first CR algorithm that incorporates momentum technique, which has provable convergence guarantee to a second-order stationary point in nonconvex optimization. We also performed a comprehensive study of this algorithm from various aspects both in theory and experiments to demonstrate the appealing attributes of the proposed algorithm. Our specific contribution are listed as follows.

- We propose a CR type algorithm with momentum acceleration (referred to as CRm), which includes a cubic regularization step, a momentum step for acceleration and a monotone step. The momentum step introduces negligible computation complexity compared to that of the cubic regularization step in original CR, but can provide substantial advantage of acceleration.

- We establish the global convergence of CRm to a second-order stationary point in nonconvex optimization. The corresponding convergence rate is as fast as that of CR in the order-level, which is the best one can expect for nonconvex optimization. Our

experiments demonstrate that CRm substantially outperforms CR as well as Nesterov's accelerated CR (which does not have guaranteed performance for nonconvex optimization).

- We also show that CRm enjoys the local quadratic convergence property under a local error bound condition, which establishes the advantage of the second-order algorithms than the first-order algorithms in nonconvex optimization.

- We further show that the inexact variant of CRm significantly improves the computational complexity without losing the convergence rate. We also study the finite-sum problem, where we implement the inexact CRm via a subsampling approach, and established the total Hessian sample complexity to guarantee the convergence with high probability.

On the core of our proof technique, we rely on the delicate design of the adaptive momentum parameter in eq. (4), and the monotone step in the algorithm, which makes it possible to establish the convergence result under nonconvex optimization but with momentum acceleration. To the best of our knowledge, there is no result on accelerated CR type algorithms that have such good convergence property, or even the convergence property under nonconvex optimization.

## 1.2 RELATED WORKS

**Escaping saddle points:** A number of algorithms have been proposed to escape saddle points in order to find local minima. In general, There are three lines of research. It has been shown that with random perturbation, gradient descent algorithm (Jin et al., 2017), the stochastic gradient descent (Ge et al., 2015), the zero-th order method (Jin et al., 2018), and the accelerated gradient descent (Jin et al., 2017) can escape saddle points. The gradient descent has also been incorporated with the negative curvature descent in Carmon et al. (2016); Liu and Yang (2017); Xu et al. (2017) in order to converge to the second-order stationary points. Furthermore, the cubic regularized (CR) algorithm, which first appeared in Griewank (1981), has been shown by Nesterov and Polyak (2006) to converge to the second-order stationary points. Cartis et al. (2011a,b) then proposed an adaptive CR method with an approximate sub-problem solver. Agarwal et al. (2017) established an efficient sub-problem solver for CR by using the Hessian-vector product technique, and Carmon and Duchi (2016) showed that gradient descent can efficiently solve the sub-problem in CR. This paper further accelerates the CR algorithm with momentum and establishes its convergence rate to a second-order stationary point.

**Algorithms with momentum for nonconvex optimization:** Ghadimi and Lan (2016); Li and Lin (2015) proposed accelerated gradient descent type of algorithms for nonconvex optimization, which are guaranteed to converge as fast as gradient descent for nonconvex problems. Yao et al. (2017) proposed an efficient accelerated proximal gradient descent algorithm for nonconvex problems, which requires only one proximal step in each iteration as compared to the requirement of two proximal steps in each iteration in the algorithm proposed in Li and Lin (2015). Then Li et al. (2017) analyzed the algorithm in Yao et al. (2017) under the KL condition. While the existing studies analyzed only convergence to first-order stationary points, this paper proposes the CR algorithms with momentum that converge to a second-order stationary point.

**Inexact CR algorithms:** To reduce the computational complexity for the CR type of algorithms, various inexact Hessian and gradient approaches were proposed. In particular, Ghadimi et al. (2017) studied the inexact Hessian CR and accelerated CR for convex optimization, where the inexact level is fixed during iterations. Tripuraneni et al. (2017) studied an inexact CR for nonconvex optimization, which allows both the gradient and Hessian to be inexact. Alternatively, Cartis et al. (2011a,b) studied the inexact Hessian CR for nonconvex optimization, where the inexact condition is adaptive during iterations. Jiang et al. (2017) studied a unified scheme of inexact accelerated adaptive CR and gradient descent for convex optimization. Furthermore, Kohler and Lucchi (2017) proposed a subsampling CR (SCR) that adaptively changes the sample batch size to guarantee the inexactness condition in Cartis et al. (2011a,b), Wang et al. (2019a) relax the inexact condition in Kohler and Lucchi (2017); Cartis et al. (2011a,b), and Xu et al. (2017) proposed uniform and non-uniform sampling algorithms with fixed inexactness for nonconvex optimization. Wang et al. (2019b); Zhou et al. (2018) proposed stochastic variance reduced subsampling CR algorithms. This paper establishes the convergence rate for the inexact scenarios of the proposed CR algorithm with momentum.

**Local quadratic convergence:** The Newton's method and cubic regularized algorithm have been shown to converge quadratically to the global minimum under the strongly convex condition in Nesterov and Polyak (2006); Nesterov (2008), respectively. Furthermore, various Newton-type algorithms, i.e., the Levenberg-Marquardt method (Yamashita and Fukushima, 2001; Fan and Yuan, 2005), the regularized Newton method (Li et al., 2004), the regularized proximal Newton's method (Yue et al., 2016), and the CR algorithm (Yue et al., 2018), have been shown to have the local quadratic convergence under the more relaxed local error bound condition. This paper fur-

ther establishes such a property for the proposed CR with momentum algorithm.

## 2 CRm: CUBIC REGULARIZATION WITH MOMENTUM

In this section, we propose a CR-type algorithm that adopts a momentum scheme (referred to as CRm). The algorithm steps of CRm are summarized in Algorithm 1.

At each iteration, the proposed CRm conducts a cubic step (eq. (3)), a momentum step (eqs. (4) and (5)), and a monotone step (eq. (6)). In particular, the cubic step solves a subproblem of the second-order Taylor expansion with a cubic regularizer at the current iterate $\mathbf{x}_k$. The cubic step can be implemented efficiently by adopting the solver based on the Hessian-vector product approach (see Section 4.3 for details). The momentum step is an extrapolation step that aims to accelerate the algorithm. We note that the momentum step requires very little additional computation compared to the cubic step, but offers substantial advantage for accelerating the algorithm. The monotone step chooses the next iteration point between the cubic step and the momentum step to achieve the minimum function value. This guarantees that the algorithm outputs a desirable monotonically decreasing function value sequence, and helps to establish the convergence guarantee under nonconvex optimization.

---
**Algorithm 1** CRm
---
1: **Input:** Initialization $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbb{R}^d, \rho < 1, M > L_2$

2: **for** $k = 0, 1, \ldots$ **do**
3:     **Cubic step:**

$$\mathbf{s}_{k+1} = \operatorname*{argmin}_{\mathbf{s}} \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2}\mathbf{s}^\top \nabla^2 f(\mathbf{x}_k)\mathbf{s}$$
$$+ \frac{M}{6}\|\mathbf{s}\|^3$$
$$\mathbf{y}_{k+1} = \mathbf{x}_k + \mathbf{s}_{k+1} \qquad (3)$$

4:     **Momentum step:**

$$\beta_{k+1} = \min\{\rho, \|\nabla f(\mathbf{y}_{k+1})\|, \|\mathbf{y}_{k+1} - \mathbf{x}_k\|\} \qquad (4)$$
$$\mathbf{v}_{k+1} = \mathbf{y}_{k+1} + \beta_{k+1}(\mathbf{y}_{k+1} - \mathbf{y}_k) \qquad (5)$$

5:     **Monotone Step:**
$$\mathbf{x}_{k+1} = \operatorname*{argmin}_{\mathbf{x} \in \{\mathbf{y}_{k+1}, \mathbf{v}_{k+1}\}} f(\mathbf{x}) \qquad (6)$$

6: **end for**
---

We further highlight the ideas in the design of CRm. First,

we choose the momentum in the direction of $\mathbf{y}_{k+1} - \mathbf{y}_k$, which has been used for the first-order methods with momentum for nonconvex problems (Li et al., 2017; Yao et al., 2017). Second, the momentum parameter $\beta_{k+1}$ in eq. (4) is set to be adaptive (in fact proportional) to the norm of the progress made in the cubic regularization step and the norm of gradient, i.e., $\|\mathbf{y}_{k+1} - \mathbf{x}_k\|$ and $\|\nabla f(\mathbf{y}_{k+1})\|$. In this way, if the iterate is far away from a second-order stationary point, $\|\mathbf{y}_{k+1} - \mathbf{x}_k\|$ and $\|\nabla f(\mathbf{y}_{k+1})\|$ are large so that the momentum takes a large stepsize to make good progress. On the other hand, as the iterate is close to the stationary point, $\|\mathbf{y}_{k+1} - \mathbf{x}_k\|$ and $\|\nabla f(\mathbf{y}_{k+1})\|$ are small so that the momentum takes a small momentum stepsize in order not to miss the stationary point. It turns out that such a choice of the momentum parameter is critical to guarantee the convergence of CRm (as can be seen in the proof) as well as achieving acceleration. Our experiments (see Section 5) show that such a momentum scheme can substantially accelerate the convergence of CR in various nonconvex optimization problems. Therefore, the requirement of the adaptive step size $\beta_k$ in eq. (4) is not only intuitively reasonable but also theoretically sound.

In the monotone step, the algorithm compares the function values between the cubic regularization step and the momentum step, and choose the better one to perform the next step. In this way, the proposed accelerated CR algorithm is guaranteed to be monotone, i.e., the generated function value sequences are monotonically decreasing. This monotone step is not required in convex optimization, but it seems crucial in nonconvex optimization due to the landscape of nonconvex function does not have strong structure as convex function. We further note that although the momentum step may not play a role in every iteration due to the monotone step, our experiments show that the momentum step does participate for most iterations during the course of convergence, validating its importance to accelerate the algorithm.

## 3 CONVERGENCE ANALYSIS OF CRm

In this section, we establish both the global and the local convergence rates of CRm to a second-order stationary point.

### 3.1 GLOBAL CONVERGENCE OF CRm

First recall that our goal is to minimize a twice-differentiable nonconvex function $f(\mathbf{x})$ (c.f. eq. (1)). We adopt the following standard assumptions on the objective function.

**Assumption 1.** *The objective function in eq. (1) satisfies:*

1. *$f$ is twice-continuously differentiable and bounded below, i.e., $f^\star \triangleq \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) > -\infty$;*
2. *For all $\alpha \in \mathbb{R}$, the sublevel set $\{\mathbf{x} : f(\mathbf{x}) \leqslant \alpha\}$ of $f$ is bounded;*
3. *The gradient $\nabla f(\cdot)$ and Hessian $\nabla^2 f(\cdot)$ are $L_1$ and $L_2$-Lipschitz continuous, respectively.*

Assumption 1 imposes standard conditions on the nonconvex objective function $f$. In particular, the bounded sublevel set condition in item 2 is satisfied whenever $f$ is coercive, i.e., $f(\mathbf{x}) \to +\infty$ as $\|\mathbf{x}\| \to +\infty$. This is true for many non-negative loss functions under mild conditions.

Based on Assumption 1, we characterize the global convergence rate of CRm to a second-order stationary point in the following result. We refer the readers to the supplementary materials for the proof.

**Theorem 1** (Global convergence rate). *Let Assumption 1 hold and fix any $\epsilon \leqslant 1$. Then, the sequence $\{\mathbf{x}_k\}_{k \geqslant 0}$ generated by CRm contains an $\epsilon$-second-order stationary point provided that the total number of iterations $k$ satisfies that*

$$k \geqslant \frac{C}{\epsilon^{3/2}}, \tag{7}$$

*where $C$ is a universal positive constant and is specified in the proof.*

Theorem 1 establishes the global convergence rate to an $\epsilon$-second-order stationary point for CRm. Although the obtained convergence rate of CRm achieves the same order as that of the original CR algorithm in Nesterov and Polyak (2006), which is in fact the best that one can expect for general nonconvex optimization, the technical proof critically exploits the design of the momentum scheme, and requires substantial machinery to handle the momentum step. Further in Section 5, we demonstrate via various experiments that CRm do enjoy the momentum acceleration and converge much faster than the original CR algorithm.

### 3.2 LOCAL CONVERGENCE OF CRm

It is well known that Newton-type second-order algorithms enjoy a local quadratic convergence rate for minimizing strongly convex functions. While strong convexity is a restrictive condition in nonconvex optimization, many nonconvex problems such as phase retrieval and low-rank matrix recovery have been shown to satisfy the following more relaxed local error bound condition (Yue et al., 2018).

**Assumption 2** (Local error bound). *Denote $\mathcal{X}$ as the set of second-order stationary points of $f$. There exists*

$\kappa, r > 0$ *such that for all* $\mathbf{x} \in \{\mathbf{x} : \text{dist}(\mathbf{x}, \mathcal{X}) \leqslant r\}$*, it holds that*

$$\text{dist}(\mathbf{x}, \mathcal{X}) \leqslant \kappa \|\nabla f(\mathbf{x})\|, \tag{8}$$

*where* $\text{dist}(\mathbf{x}, \mathcal{X})$ *denotes the point-to-set distance between* $\mathbf{x}$ *and* $\mathcal{X}$*.*

One can easily check that all strongly convex functions satisfy the above local error bound condition. Therefore, the local error bound condition is a more general geometry than strong convexity.

Next, we explore the local convergence property for CRm under the local error bound condition. Typically, such a property is due to the usage of the Hessian information in the algorithm. In CRm, the momentum step does not directly exploit the Hessian information. Hence, it is not clear *a priori* by including the momentum step whether CRm still enjoys the local quadratic convergence property. The following theorem provides an affirmative answer.

**Theorem 2.** *Let Assumptions 1 and 2 hold. Then, the sequence* $\{\mathbf{x}_k\}_{k \geqslant 0}$ *generated by CRm with* $M > L_2$ *converges quadratically to a point* $\mathbf{x}^\star \in \mathcal{X}$*, where* $\mathcal{X}$ *is the set of second-order stationary points of* $f$*. That is, there exists an integer* $k_1$ *such that for all* $k \geqslant k_1$*,*

$$\|\mathbf{x}_{k+1} - \mathbf{x}^\star\| \leqslant C\|\mathbf{x}_k - \mathbf{x}^\star\|^2, \tag{9}$$

*where* $C$ *is a universal positive constant and is specified in the proof.*

Under the local error bound condition, Theorem 2 shows that CRm enjoys a quadratic convergence rate as shown in eq. (9). To elaborate, note that Theorem 1 guarantees the convergence of CRm to a second-order stationary point, i.e., $\|\mathbf{x}_k - \mathbf{x}^\star\| \to 0$ as $k \to \infty$. Thus, the recursion in eq. (9) implies that $C\|\mathbf{x}_k - \mathbf{x}^\star\| \leqslant (C\|\mathbf{x}_{k_1} - \mathbf{x}^\star\|)^{2^{k-k_1}}$, which is at a quadratic converge rate. In particular, the region of quadratic convergence is defined by $\|\mathbf{x}_k - \mathbf{x}^\star\| \leqslant 1/C$. Such quadratic convergence achieves an $\epsilon$-accuracy second-order stationary point within $k = O(\log \log(1/\epsilon))$ number of iterations, which is much faster than the linear converge rate of fisrt-order methods in local region.

Local quadratic convergence has also been established for the original CR algorithm under the local error bound condition Yue et al. (2018). As a comparison, our proof of Theorem 2 for CRm exploits the proposed momentum scheme, which results in additional terms that requires extra effort to handle.

## 4 INEXACT VARIANTS OF CRm

The major computational load of CRm lies in the cubic step, which requires to solve a computationally costly optimization problem. In this section, we explore three imple-

mentation schemes that can efficiently perform the cubic step without sacrificing the acceleration performance.

### 4.1 CUBIC STEP WITH INEXACT HESSIAN

The cubic step requires the full Hessian information, which can be too costly in practice. Instead, we consider performing the following cubic step with an inexact approximation of the Hessian.

$$\hat{\mathbf{x}}_{k+1} = \underset{\mathbf{s} \triangleq \mathbf{x} - \mathbf{x}_k}{\arg\min} \nabla f(\mathbf{x}_k)^\top \mathbf{s} + \frac{1}{2}\mathbf{s}^\top \mathbf{H}_k \mathbf{s} + \frac{M}{6}\|\mathbf{s}\|^3, \tag{10}$$

where $\mathbf{H}_k$ denotes the inexact estimation of the full Hessian $\nabla^2 f(\mathbf{x}_k)$, and their difference is assumed to satisfy the following criterion. Section 4.2 proposes a subsampling scheme to achieve Assumption 3 for the finite-sum problem.

**Assumption 3.** *The inexact Hessian* $\mathbf{H}_k$ *in eq. (10) satisfies, for all* $k \geqslant 0$*,*

$$\|\mathbf{H}_k - \nabla^2 f(\mathbf{x}_k)\| \leqslant \epsilon_1.$$

Assumption 3 assumes that the inexact Hessian is close to the exact one in terms of a small operator norm gap. Such inexact criterion has been considered in Tripuraneni et al. (2017); Xu et al. (2017) to study the convergence property of the inexact CR algorithm.

Next, we study the convergence of the inexact variant of CRm by replacing the cubic step in eq. (3) with the inexact cubic step in eq. (10). Our main result is summarized as follows, and the proof is provided in the supplemental materials.

**Theorem 3.** *Let Assumptions 1 and 3 hold and fix any* $\epsilon \leqslant 1$*. Then, the sequence* $\{\mathbf{x}_k\}_{k \geqslant 0}$ *generated by the inexact CRm with* $M > 2L_2/3 + 2$ *and* $\epsilon_1 = \theta\sqrt{\epsilon}$ *contains an* $\epsilon$*-second-order stationary point provided that the total number of iterations* $k$ *satisfies that*

$$k \geqslant \frac{C}{\epsilon^{3/2}}, \tag{11}$$

*where* $C, \theta$ *are universal constants, and are specified in the proof.*

Theorem 3 shows that, under a proper inexact criterion, the iteration complexity of inexact CRm is on the same order as that of exact CRm for achieving an $\epsilon$-second-order stationary point. Since the inexact Hessian saves the computation in each iteration comparing to the full Hessian, it is clear that the overall computation complexity of the inexact CRm is less than that of the exact cases. In Appendix A.2, we verify through experiments that

the inexact algorithm do perform much better than the corresponding exact version.

We note that the proof of Theorem 3 suggests that the condition that $\|\mathbf{y}_{k+1} - \mathbf{x}_k\| \leqslant \epsilon_1$ implies the point $\mathbf{x}_{k+1}$ is an $\epsilon$-second-order stationary point, where $\epsilon_1 = \theta\sqrt{\epsilon}$. Thus, the implementation of the inexact CRm can terminate by checking the satisfaction of the condition $\|\mathbf{y}_{k+1} - \mathbf{x}_k\| \leqslant \epsilon_1$.

## 4.2 INEXACT CRm VIA SUBSAMPLING

In this subsection, we consider a general finite-sum optimization problem, where inexact CRm can be implemented via subsampling. More specifically, consider to solve the following optimization problem:

$$f(x) \triangleq \sum_{i=1}^{n} f_i(x), \tag{12}$$

where $f_i(\cdot)$ is possibly nonconvex. Furthermore, we assume that Assumption 1 holds for each $f_i(\cdot)$. For finite-sum problems, the full Hessian can be approximated by the Hessian of a mini-batch of data samples each uniformly randomly drawn from the dataset, i.e.,

$$\mathbf{H}_k = \frac{1}{|S_1|} \sum_{i \in S_1} \nabla^2 f_i(\mathbf{x}_k). \tag{13}$$

We use the subsampling technique introduced in Kohler and Lucchi (2017) to satisfy the inexact condition in Assumption 3. The following theorem provides our characterization of the overall Hessian sample complexity in order to guarantee the convergence of the subsampling algorithm with high probability over the entire iteration process.

**Theorem 4** (Total Hessian sample complexity). *Assuming that Assumption 1 holds for each $f_i(\cdot)$, and let the sub-sampled mini-batch of Hessians $\mathbf{H}_k, k = 0, 1, \ldots$ satisfies*

$$|S_1| = \left( \frac{8L_1^2}{\theta^2 \epsilon} + \frac{4L_1}{3\theta\sqrt{\epsilon}} \right) \log\left( \frac{4d}{\epsilon^{3/2}\delta} \right),$$

*then the sequence $\{\mathbf{x}_k\}_{k \geqslant 0}$ generated by the inexact CRm with $M > L_2 + 2$ outputs an $\epsilon$-second-order stationary point with probability at least $1 - \delta$ by taking at most the following number of Hessian samples in total:*

$$S \leqslant C \left( \frac{8L_1^2}{\theta^2 \epsilon^{5/2}} + \frac{4L_1}{3\theta\epsilon^2} \right) \log\left( \frac{4d}{\epsilon\delta} \right).$$

Theorem 4 characterizes the total Hessian sample complexity to guarantee the convergence of CRm with high

probability. This is the first such a type result for stochastic CR algorithms. Note that previous studies Kohler and Lucchi (2017); Xu et al. (2017) on subsampling CR provide only the Hessian sample complexity per iteration to guarantee inexactness condition with high probability. Our result indicates that even over the entire iteration process, the convergence is still guaranteed with high probability. In fact, if we let $N$ denote the total sample complexity, Theorem 4 implies that the failure probability $\delta$ decays exponentially fast as the total sample complexity $N$ becomes asymptotically large. Such a result by nature is stronger than those that characterize the convergence only in expectation, not in (high) probability, in existing literature.

## 4.3 EFFICIENT SOLVERS FOR CUBIC STEP

Since the cubic step does not have a closed form solution, an inexact solver is typically used for solving the cubic step. Various solvers have been proposed to approximately solve such a subproblem. The first type of solver is based on the Lanczos method (Cartis et al., 2011a,b), which solves the cubic subproblem in a Krylov subspace $\mathcal{K} = \text{span}\{\nabla f(\mathbf{x}_k), \nabla^2 f(\mathbf{x}_k)\nabla f(\mathbf{x}_k), \cdots\}$ instead of in the entire space. Each step of the solver can be implemented efficiently with a computation cost of $O(d)$ (Kohler and Lucchi, 2017). Moreover, building the subspace requires a Hessian-vector product, which introduces a cost of $O(nd)$ per additional subspace dimension for finite-sum problem with $n$ data samples. The second type of solver is proposed by Agarwal et al. (2017), which is based on the techniques of Hessian-vector product and binary search. The proposed solver can find an approximate solution of the cubic subproblem with a total cost of $O(nd/\epsilon^{1/4})$ for finite-sum problems, where $\epsilon$ is the desired accuracy. Carmon and Duchi (2016) proposed another solver based on the gradient descent method. The solver finds an approximate solution of the cubic subproblem within $O(\epsilon^{-1}\log(1/\epsilon))$ iterations for large $\epsilon$ and $O(\log(1/\epsilon))$ iterations for small $\epsilon$.

All of these solvers can be applied to solve the cubic step in CRm. Note that the momentum and monotone steps in CRm introduce order-level less computation complexity compared to these solvers for solving the cubic step. Thus, CRm have the same per-iteration computational complexity as CR when implementing the same solver, and have at least the same overall computational complexity as CR (in fact, much less overall computational complexity in practice as demonstrated by our experiments). As the solvers solve the cubic subproblem up to certain accuracy in practice, the total computation complexity of CRm to achieve a second-order stationary point can still be established.
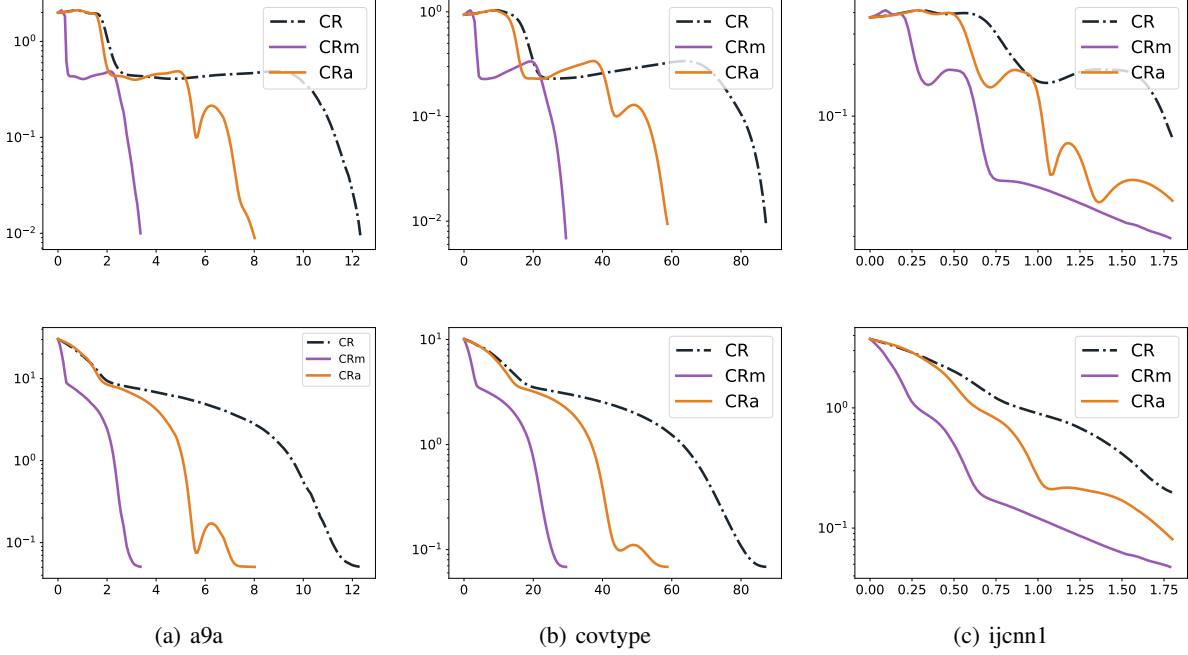
Figure 1: Nonconvex logistic regression. Top: gradient norm v.s. time. Bottom: function value gap v.s. time.

## 5 EXPERIMENTS

### 5.1 SETUP

We compare the performance among the following six algorithms: cubic regularization algorithm (**CR**) in [Nesterov and Polyak (2006)](#), accelerated cubic regularization algorithm (**CRa**) in [Nesterov (2008)](#) (whose convergence guarantee has not been established), cubic regularization algorithm (**CRm**) with momentum, cubic regularization algorithm with inexact Hessian (**CR_I**), accelerated cubic regularization with inexact Hessian (**CRa_I**), CRm with inexact Hessian (**CRm_I**). In this section, we present the comparison among the three exact algorithms. The comparisons among the inexact variants are presented in Appendix A due to space limitation. The details of the experiment settings can also be found in Appendix A.

We conduct two experiments. The first experiment solves the following logistic regression problem with a nonconvex regularizer

$$\min_{\mathbf{w}\in\mathbb{R}^d} -\left(\frac{1}{n}\sum_{i=1}^{n} y_i \log\left(\frac{1}{1+e^{-\mathbf{w}^T\mathbf{x}}}\right)\right.$$
$$\left. +(1-y_i)\log\left(\frac{e^{\mathbf{w}^T\mathbf{x}}}{1+e^{-\mathbf{w}^T\mathbf{x}}}\right)\right) + \alpha\sum_{i=1}^{d}\frac{w_i^2}{1+w_i^2},$$

where we set $\alpha = 0.1$ in our experiment. The second experiment solves the following nonconvex robust linear

regression problem

$$\min_{\mathbf{w}\in\mathbb{R}^d} \frac{1}{n}\sum_{i=1}^{n} \eta(y_i - \mathbf{w}^T\mathbf{x}_i), \qquad (14)$$

where $\eta(x) = \log(\frac{x^2}{2}+1)$. Each experiment is performed over three datasets, i.e., a9a, covtype, and ijcnn ([Chang and Lin, 2011](#)).

### 5.2 RESULTS

Figures 1 and 2 show the results of the two experiments for comparing the three exact algorithms, respectively. From both figures, it can be seen that CRm outperforms the vanilla CR, which demonstrates that the momentum step in CRm significantly accelerates the CR algorithm for nonconvex problems. Also, CRm outperforms CRa in the experiments with datasets a9a and covtype, while its performance is comparable to CRa in the experiments with dataset ijcnn1. Thus, our proposed momentum step achieves a faster convergence than the Nesterov's acceleration scheme for CR.

We note that similar comparisons are observed in the comparison of the corresponding three inexact variants of the algorithms (see Figures 3 and 4 in Appendix A), i.e., our momentum scheme with inexact Hessian outperforms other inexact CR algorithms. We also note that all inexact variants of the algorithms outperform their exact counterparts (see Figures 5 and 6 in Appendix A). Hence,

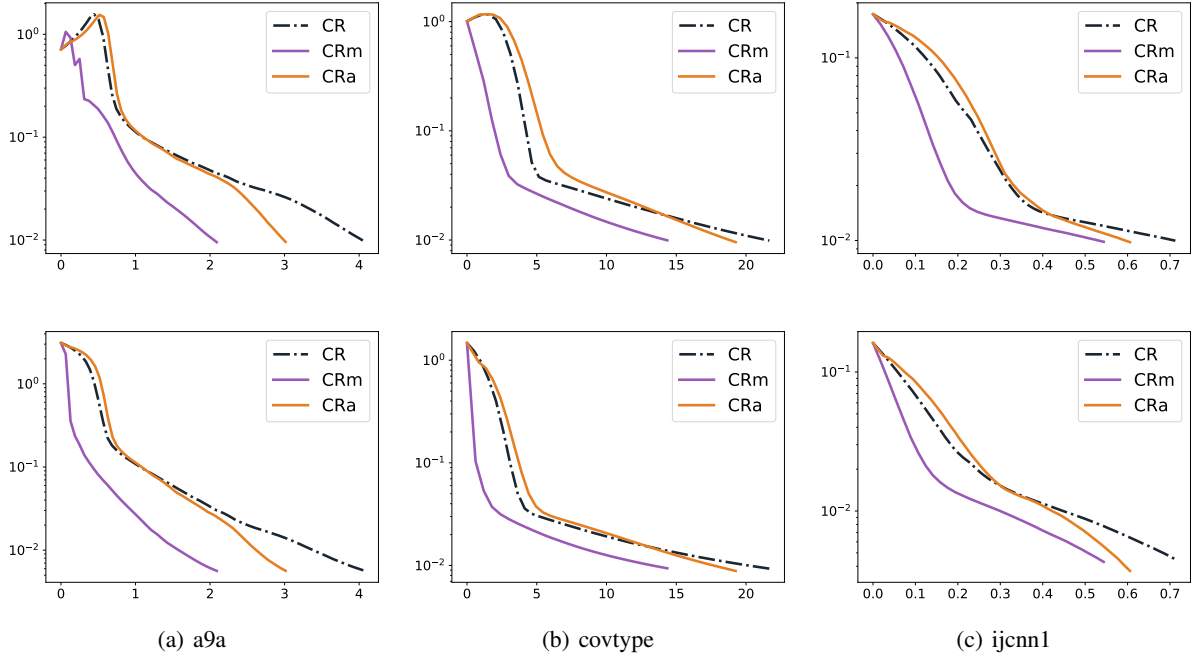(a) a9a      (b) covtype      (c) ijcnn1

Figure 2: Robust linear regression. Top: gradient norm v.s. time. Bottom: function value gap v.s. time.

the inexact implementation plays an important role in reducing the computation complexity of these CR type of algorithms in practice.

# 6 CONCLUSION

In this paper, we proposed a momentum scheme to accelerate the cubic regularization algorithm. We showed that the order of the global convergence rate of the proposed algorithm CRm is at least as fast as its vanilla version. We also established the local quadratic convergence property for the proposed algorithm, and extended our analysis for the proposed algorithm to the inexact Hessian case and established the total Hessian sample complexity to guarantee the convergence with high probability. We further conducted various experiments to demonstrate the advantage of applying momentum for accelerating the cubic regularized algorithm.

## Acknowledgement

# References

Agarwal, N., Allen-Zhu, Z., Bullins, B., Hazan, E., and Ma, T. (2017). Finding approximate local minima faster than gradient descent. In *Proc. 49th Annual ACM SIGACT Symposium on Theory of Computing (STOC)*, pages 1195–1199.

Carmon, Y. and Duchi, J. C. (2016). Gradient descent efficiently finds the cubic-regularized non-convex Newton step. *arXiv:1803.09357*.

Carmon, Y., Duchi, J. C., Hinder, O., and Sidford, A. (2016). Accelerated methods for non-convex optimization. *arXiv:1611.00756*.

Cartis, C., Gould, N. I. M., and Toint, P. L. (2011a). Adaptive cubic regularization methods for unconstrained optimization. Part I : Motivation, convergence and numerical results. *Mathematical Programming*, 127(2):245–295.

Cartis, C., Gould, N. I. M., and Toint, P. L. (2011b). Adaptive cubic regularization methods for unconstrained optimization. Part II worst-case function- and derivative-evaluation complexity. *Mathematical Programming*, 130(2):295–319.

Chang, C. and Lin, C. (2011). LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(3):1–27.

Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., and Bengio, Y. (2014). Identifying and attacking the saddle point problem in high-dimensional non-

convex optimization. In *Proc. International Conference on Neural Information Processing Systems (NIPS)*, pages 2933–2941.

Fan, J. and Yuan, Y. (2005). On the quadratic convergence of the levenberg-marquardt method without nonsingularity assumption. *Computing*, 74(1):23–39.

Ge, R., Huang, F., Jin, C., and Yuan, Y. (2015). Escaping from saddle points — online stochastic gradient for tensor decomposition. In *Proc. Conference on Learning Theory (COLT)*, volume 40, pages 797–842.

Ghadimi, S. and Lan, G. (2016). Gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 156(1):59–99.

Ghadimi, S., Liu, H., and Zhang, T. (2017). Second-order methods with cubic regularization under inexact information. *arXiv: 1710.05782*.

Griewank, A. (1981). The modification of newtons method for unconstrained optimization by bounding cubic terms. *Technical report, Department of Applied Mathematics and Theoretical Physics, University of Cambridge*.

Jiang, B., Lin, T., and Zhang, S. (2017). A unified scheme to accelerate adaptive cubic regularization and gradient methods for convex optimization. *arXiv:1710.04788*.

Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. (2017). How to escape saddle points efficiently. In *Proc. 34th International Conference on Machine Learning (ICML)*, volume 70, pages 1724–1732.

Jin, C., Liu, L., Ge, R., and Jordan, M. I. (2018). Minimizing nonconvex population risk from rough empirical risk. *arXiv:1803.09357*.

Jin, C., Netrapalli, P., and Jordan, M. I. (2017). Accelerated gradient descent escapes saddle points faster than gradient descent. *arXiv:1711.10456*.

Kohler, J. M. and Lucchi, A. (2017). Sub-sampled cubic regularization for non-convex optimization. In *Proc. 34th International Conference on Machine Learning (ICML)*, volume 70, pages 1895–1904.

Li, D.-H., Fukushima, M., Qi, L., and Yamashita, N. (2004). Regularized newton methods for convex minimization problems with singular solutions. *Computational Optimization and Applications*, 28(2):131–147.

Li, H. and Lin, Z. (2015). Accelerated proximal gradient methods for nonconvex programming. In *Proc. 28th Advances in Neural Information Processing Systems (NIPS)*, pages 379–387.

Li, Q., Zhou, Y., Liang, Y., and Varshney, P. K. (2017). Convergence analysis of proximal gradient with momentum for nonconvex optimization. In *Proc. 34th International Conference on Machine Learning (ICML)*, volume 70, pages 2111–2119.

Liu, M. and Yang, T. (2017). On noisy negative curvature descent: Competing with gradient descent for faster non-convex Optimization. *arXiv:1709.08571*.

Nesterov, Y. (2008). Accelerating the cubic regularization of Newton's method on convex problems. *Mathematical Programming*, 112(1):159–181.

Nesterov, Y. and Polyak, B. T. (2006). Cubic regularization of newton method and its global performance. *Mathematical Programming*, 108(1):177–205.

Sun, J., Qu, Q., and Wright, J. (2015). Complete dictionary recovery using nonconvex optimization. In *Proc. 32nd International Conference on Machine Learning (ICML)*, volume 37, pages 2351–2360.

Sun, J., Qu, Q., and Wright, J. (2017). A geometric analysis of phase retrieval. *Foundations of Computational Mathematics*, pages 2379–2383.

Tripuraneni, N., Stern, M., Jin, C., Regier, J., and Jordan, M. I. (2017). Stochastic cubic regularization for fast nonconvex optimization. *arXiv: 711.02838*.

Tropp, J. A. (2012). User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434.

Wang, Z., Zhou, Y., Liang, Y., and Lan, G. (2019a). A note on inexact gradient and Hessian conditions for cubic regularized Newtons method. *Operations Research Letters*.

Wang, Z., Zhou, Y., Liang, Y., and Lan, G. (2019b). Sample complexity of stochastic variance-reduced cubic regularization for nonconvex optimization. *In Proc. 22nd International Conference on Artificial Intelligence and Statistics (AISTATS)*.

Xu, P., Roosta-Khorasani, F., and Mahoney, M. W. (2017). Newton-type methods for non-convex optimization under inexact hessian information. *arXiv: 1708.07164*.

Xu, Y., Jin, R., and Yang, T. (2017). NEON+: Accelerated gradient methods for extracting negative curvature for non-convex optimization. *arXiv: 1712.01033*.

Yamashita, N. and Fukushima, M. (2001). On the rate of convergence of the levenberg-marquardt method. *Topics in Numerical Analysis*, pages 239–249.

Yao, Q., Kwok, J. T., Gao, F., Chen, W., and Liu, T.-Y. (2017). Efficient inexact proximal gradient algorithm for nonconvex problems. In *Proc. 26th International Joint Conference on Artificial Intelligence (IJCAI-17)*, pages 3308–3314.

Yue, M., Zhou, Z., and So, A. (2018). On the quadratic convergence of the cubic regularization method under a local error bound condition. *arXiv:1801.09387*.

Yue, M.-C., Zhou, Z., and Man-Cho So, A. (2016). A Family of inexact SQA methods for non-smooth convex minimization with provable convergence guarantees based on the Luo-Tseng error bound property. *arXiv:1605.07522*.

Zhou, D., Xu, P., and Gu, Q. (2018). Stochastic variance-reduced cubic regularized Newton method. *In Proc. 35th International Conference on Machine Learning (ICML)*.

# Supplementary Materials

We first note that in order for a clear illustration, we use a slightly difference notation such that $\mathbf{v}_{k+1}$ is replaced by $\tilde{\mathbf{x}}_{k+1}$ and $\mathbf{y}_{k+1}$ is replaced by $\hat{\mathbf{x}}_{k+1}$.

## A  Experiment Setting and Additional Result

**Parameters Setting:** The experiment specifications are as follows. For all algorithms, the subproblem in each iteration of the cubic step is solved by the Lanczos-type method as suggested by Cartis et al. (2011a). We set the parameter $M = 10$ through all the experiments. The momentum parameter $\beta_{k+1}$ in CRm is set to be $8 \times \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|$. Although in theory of CRm, we require $\beta_{k+1} < \min\{\rho, \|\nabla f(\hat{\mathbf{x}}_{k+1})\|, \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|\}$, in practice, a more relaxed value works well in our experiments. The initial point is set to be an all-two-vector for all datasets for the logistic regression problem and an all $0.5$ vector for all datasets for the robust linear regression problem.

### A.1  Comparison Among Inexact Algorithms

In this subsection, we present the comparison among the three algorithms with inexact Hessian. Namely, vanilla cubic regularized algorithm with inexact Hessian (**CR_I**), Nesterov accelerated cubic regularization algorithm with inexact Hessian (**CRa_I**), and the proposed CRm with inexact Hessian (**CRm_I**). For the implementation, since we solve a finite-sum problem $f(\mathbf{x}) = \sum_{i=1}^n f_i(\mathbf{x})$, we draw a mini-batch of data samples as the estimated Hessian, given by $\mathbf{H_k} = \sum_{i \in S_k} \nabla^2 f_i(\mathbf{x}_k)/|S_k|$. We take $n/20$ as the batch size in the logistic regression problem and $n/5$ in the robust linear regression problem. The results are shown in Figures 3 and 4. The performance comparison among the four algorithms with inexact Hessian is similar to that of their exact cases. However, we should note that the time used by inexact algorithms is much less than that of their corresponding exact cases.



(a) a9a ($n = 32651, d = 123$)  (b) covtype ($n = 581012, d = 54$)  (c) ijcnn1($n = 35000, d = 22$)
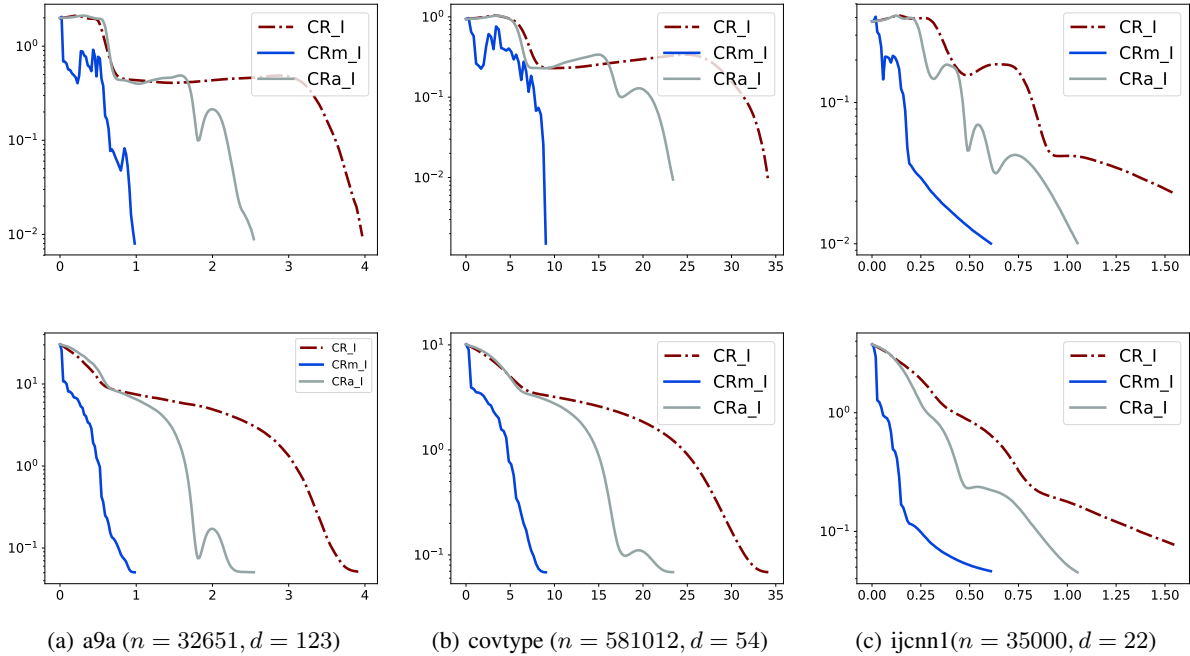
Figure 3: Logistic regression loss with a nonconvex regularizer: The top row presents the gradient norm versus runtime. The bottom row presents the function value gap versus runtime.

### A.2  Comparison between Exact and Inexact Algorithms

In this subsection, we present the comparison between the algorithms with exact Hessian and the algorithms with inexact Hessian. The results are shown in Figures 5 and 6. It is clear that all inexact algorithms significantly outperform

Figure 4: Robust linear regression loss: The top row presents the gradient norm versus runtime. The bottom row presents the function value gap versus runtime.

their corresponding exact algorithms. This demonstrates the efficiency of the inexact Hessian technique for second-order algorithms.



Figure 5: Nonconvex logistic regression. Top: gradient norm v.s. time. Bottom: function value gap v.s. time.

(a) a9a($n = 32651, d = 123$)  (b) covtype($n = 581012, d = 54$)  (c) ijcnn1($n = 35000, d = 22$)
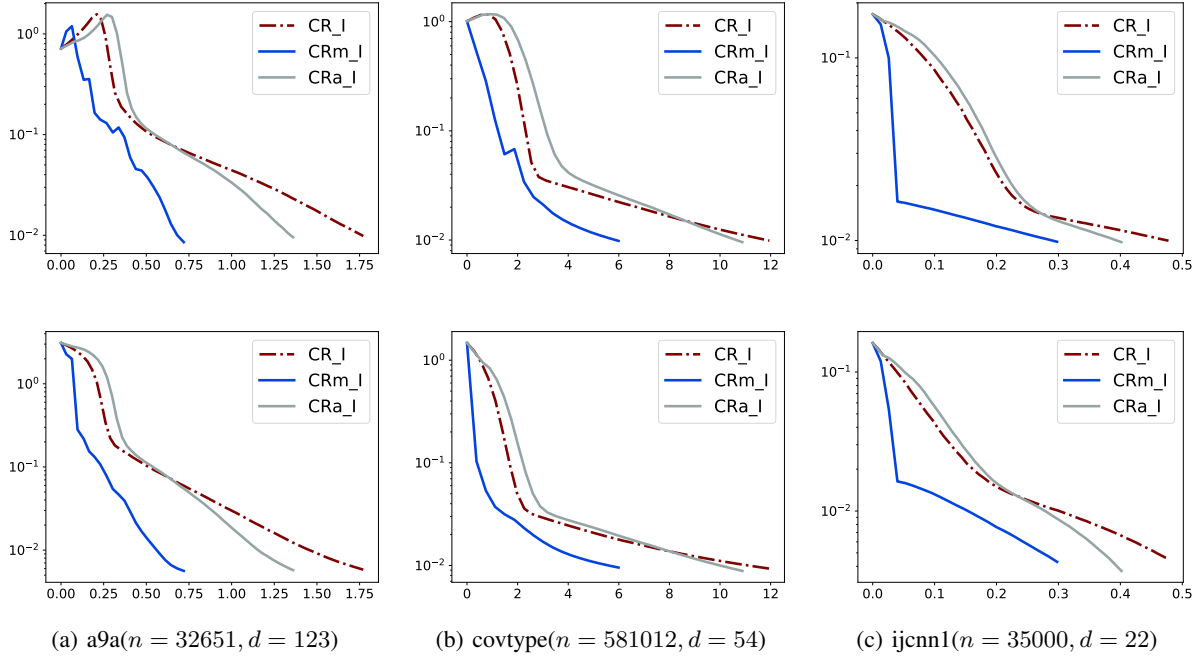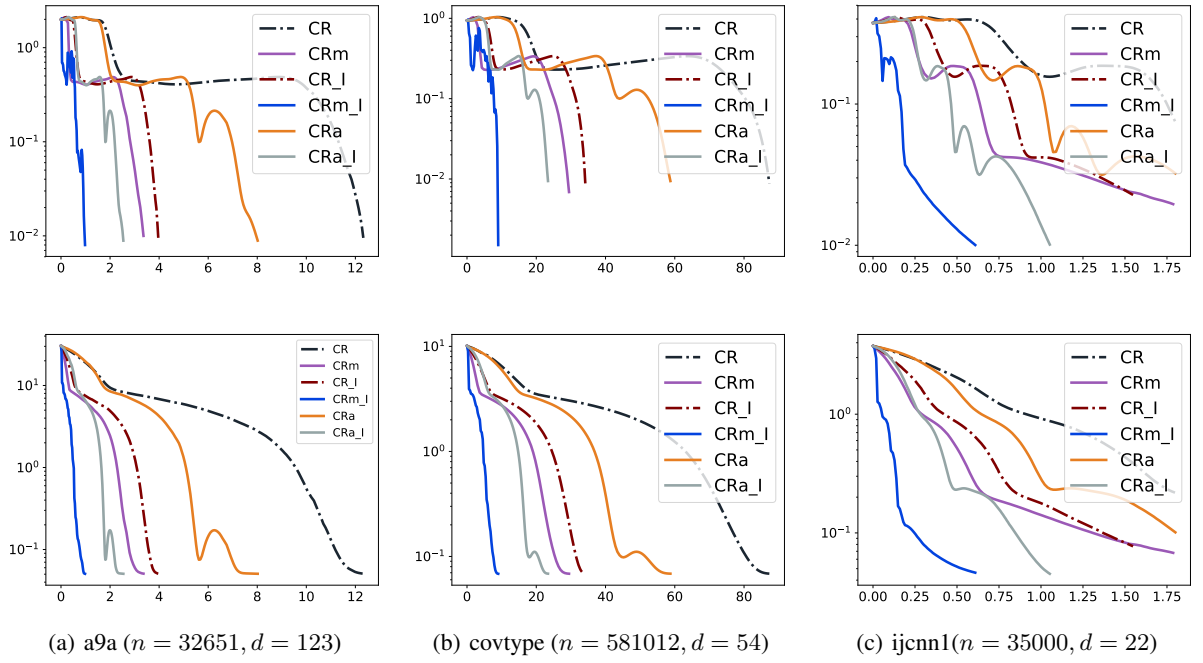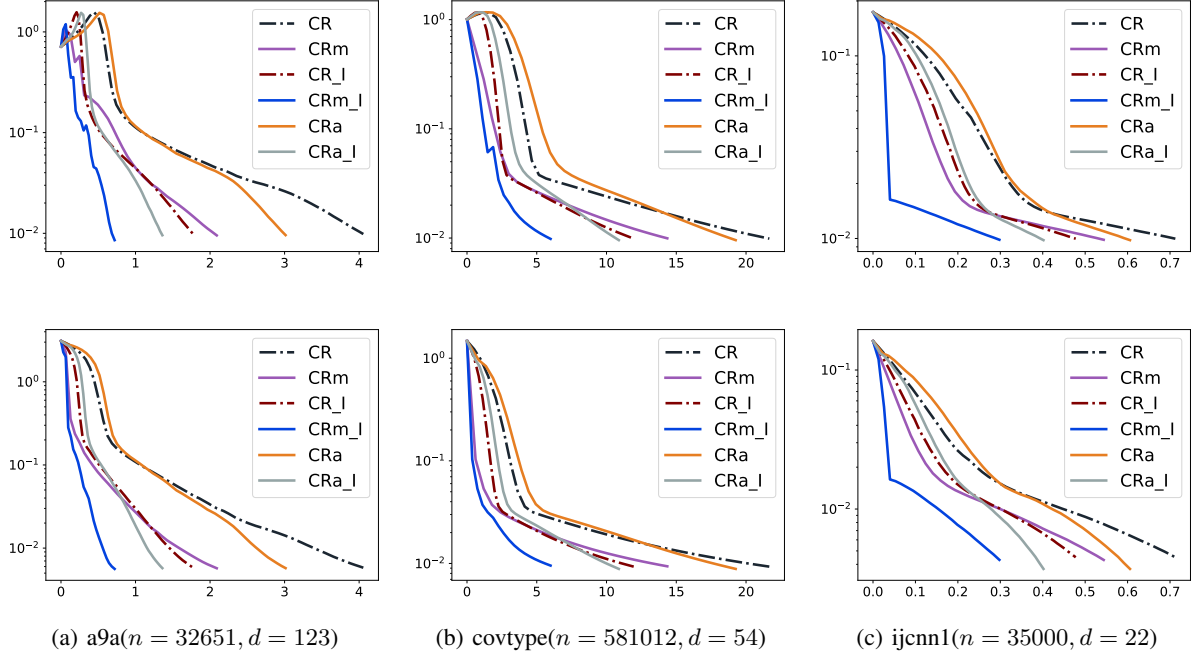
Figure 6: Robust linear regression loss: The top row presents the gradient norm versus runtime. The bottom row presents the function value gap versus runtime.

## B   Technical Lemmas

In this section, we introduce several technical lemmas that are useful for proving our main results.

**Lemma 5** (Nesterov and Polyak (2006), Lemma 1). *Let the Hessian $\nabla^2 f$ of the function $f$ be $L_2$-Lipschitz continuous with $L_2 > 0$. Then, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have*

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})\| \leqslant \frac{L_2}{2}\|\mathbf{y} - \mathbf{x}\|^2.$$

$$|f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) - \frac{1}{2}(\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})| \leqslant \frac{L_2}{6}\|\mathbf{y} - \mathbf{x}\|^3.$$

The following lemma provides bounds on $\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|$ as shown in (Yue et al., 2018, Lemma 1).

**Lemma 6** (Yue et al. (2018), Lemma 1). *Let Assumption 1 hold. Then, the sequences $\{\mathbf{x}_k\}_{k \geqslant 0}$ and $\{\hat{\mathbf{x}}_k\}_{k \geqslant 1}$ generated by Algorithm 1 satisfies, for all $k \geqslant 0$,*

$$\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| \leqslant c_1 \cdot dist(\mathbf{x}_k, \mathcal{X}), \tag{15}$$

*where $c_1 \triangleq \left(1 + \frac{L_2}{M} + \sqrt{\left(1 + \frac{L_2}{M}\right)^2 + \frac{L_2}{M}}\right)$.*

Next, we develop a number of useful bounds regarding the sequences $\{\mathbf{x}_k\}_{k \geqslant 0}$ and $\{\hat{\mathbf{x}}_k\}_{k \geqslant 1}$ that are generated by Algorithm 1. We refer to Appendix G for the details of the proofs.

**Lemma 7.** *Let Assumption 1 hold. Then, the sequences $\{\mathbf{x}_k\}_{k \geqslant 0}$ and $\{\hat{\mathbf{x}}_k\}_{k \geqslant 0}$ generated by Algorithms 1 with $M > \frac{2L_2}{3}$ satisfy, for all $k \geqslant 0$,*

$$f(\hat{\mathbf{x}}_{k+1}) - f(\mathbf{x}_k) \leqslant -\gamma\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^3, \tag{16}$$

$$\sum_{i=0}^{k} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^3 \leqslant \frac{f(\mathbf{x}_0) - f^\star}{\gamma}, \tag{17}$$

$$\|\nabla f(\hat{\mathbf{x}}_{k+1})\| \leqslant \frac{L_2 + M}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^2, \tag{18}$$

$$\lambda_{\min}(\nabla^2 f(\hat{\mathbf{x}}_{k+1})) \geqslant -\frac{M + 2L_2}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|, \tag{19}$$

*where* $\gamma = \frac{3M - 2L_2}{12}$. *Furthermore, there exists a* $k_0 \in \{0, \cdots k\}$ *such that*

$$\|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \leqslant \frac{1}{k^{1/3}}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}. \tag{20}$$

The following lemma establishes the corresponding bounds regarding $\{\mathbf{x}_k\}_{k \geqslant 0}$ and $\{\hat{\mathbf{x}}_k\}_{k \geqslant 1}$ that are generated by the *inexact* variants of Algorithms 1. We refer to Appendix G for the details of the proof.

**Lemma 8.** *Let Assumption 1 and Assumption 3 hold. Then, the sequences* $\{\mathbf{x}_k\}_{k \geqslant 0}$ *and* $\{\hat{\mathbf{x}}_k\}_{k \geqslant 0}$ *generated by the inexact variants of Algorithms 1 with* $M > \frac{2L_2}{3} + 2$ *satisfy, for all* $k \geqslant 0$ *and any* $\epsilon_1 > 0$,

$$f(\hat{\mathbf{x}}_{k+1}) - f(\mathbf{x}_k) \leqslant -\frac{3M - 2L_2}{12}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^3 + \frac{1}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^2 \epsilon_1, \tag{21}$$

$$\|\nabla f(\hat{\mathbf{x}}_{k+1})\| \leqslant \frac{L_2 + M}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^2 + \epsilon_1\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|, \tag{22}$$

$$\lambda_{\min}(\nabla^2 f(\hat{\mathbf{x}}_{k+1})) \geqslant -\frac{M + 2L_2}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| - \epsilon_1. \tag{23}$$

*Furthermore, if the total number of iterations* $k > \left(\frac{12}{3M - 2L_2 - 6}\right)\frac{f(\mathbf{x}_0) - f^\star}{\epsilon_1^3}$ *for any* $\epsilon_1 > 0$*, then there exists a* $k_0 \in \{0, \cdots k\}$ *such that*

$$\|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \leqslant \epsilon_1. \tag{24}$$

Note that in CRm, $\hat{\mathbf{x}}_{k+1}$ is generated by $\mathbf{x}_k$ through a cubic regularization step, and the output sequence $\{f(\mathbf{x}_k)\}_{k \geqslant 0}$ is monotone through a monotone step. Therefore, Lemmas 6 to 8 hold for CRm.

## C   Global Convergence: Proof of Theorem 1

We first prove a useful inequality. Note that

$$\|\mathbf{x}_{k+1} - \hat{\mathbf{x}}_{k+1}\| \overset{(i)}{\leqslant} \max\left(\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_{k+1}\|, \|\tilde{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_{k+1}\|\right)$$

$$= \|\tilde{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_{k+1}\| \overset{(ii)}{\leqslant} \beta_{k+1}\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\|, \tag{25}$$

where (i) follows from the definition of $\mathbf{x}_{k+1}$ (see eq. (6)) and (ii) follows from eq. (5).

We then present the following lemma that bounds $\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\|$, where $\{\hat{\mathbf{x}}_k\}_{k \geqslant 0}$ is generated by CRm.

**Lemma 9.** *Let Assumption 1 hold. Then, the sequence* $\{\hat{\mathbf{x}}_k\}$ *generated by CRm satisfies*

$$\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| \leqslant c_5, \tag{26}$$

*where* $c_5 \triangleq \frac{1}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}$.

*Proof.* See Appendix G.4. □

Next, we prove the main theorem for CRm. Note that

$$\|\nabla f(\mathbf{x}_{k+1})\| \leqslant \|\nabla f(\hat{\mathbf{x}}_{k+1})\| + \|\nabla f(\hat{\mathbf{x}}_{k+1}) - \nabla f(\mathbf{x}_{k+1})\|$$

$$\overset{(i)}{\leqslant} \|\nabla f(\hat{\mathbf{x}}_{k+1})\| + L_1\|\mathbf{x}_{k+1} - \hat{\mathbf{x}}_{k+1}\|$$

$$\overset{(ii)}{\leqslant} \|\nabla f(\hat{\mathbf{x}}_{k+1})\| + L_1\beta_{k+1}\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\|$$

$$\overset{(iii)}{\leqslant} \|\nabla f(\hat{\mathbf{x}}_{k+1})\| \left(1 + L_1\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\|\right) \tag{27}$$

$$\overset{(iv)}{\leqslant} \frac{L_2 + M}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^2 \left(1 + \frac{L_1}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right), \tag{28}$$

where (i) follows from the Lipschitz gradient assumption, (ii) follows from eq. (25), (iii) follows from eq. (4), which implies that $\beta_{k+1} \leqslant \|\nabla f(\hat{\mathbf{x}}_{k+1})\|$ and (iv) follows from Lemma 9 and eq. (18).

Next, we bound the minimum eigenvalue of the Hessian. Observe that

$$\lambda_{\min}\left(\nabla^2 f(\mathbf{x}_{k+1})\right) \overset{(i)}{\geqslant} \lambda_{\min}\left(\nabla^2 f(\hat{\mathbf{x}}_{k+1})\right) - \|\nabla^2 f(\mathbf{x}_{k+1}) - \nabla^2 f(\hat{\mathbf{x}}_{k+1})\|$$

$$\overset{(ii)}{\geqslant} \lambda_{\min}\left(\nabla^2 f(\hat{\mathbf{x}}_{k+1})\right) - L_2\|\mathbf{x}_{k+1} - \hat{\mathbf{x}}_{k+1}\|$$

$$\overset{(iii)}{\geqslant} \lambda_{\min}\left(\nabla^2 f(\hat{\mathbf{x}}_{k+1})\right) - L_2\beta_{k+1}\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\|$$

$$\overset{(iv)}{\geqslant} \lambda_{\min}\left(\nabla^2 f(\hat{\mathbf{x}}_{k+1})\right) - L_2\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| \tag{29}$$

$$\overset{(v)}{\geqslant} -\frac{M + 2L_2}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| - \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|\frac{L_2}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}$$

$$= -\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|\left(\frac{M + 2L_2}{2} + \frac{L_2}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right), \tag{30}$$

where (i) follows form Weyl's inequality, (ii) follows from the fact that $\nabla^2 f(\cdot)$ is $L_2$ Lipschitz, (iii) follows from eq. (25), (iv) follows from eq. (4), which implies that $\beta_{k+1} \leqslant \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|$ and (v) follows from Lemma 9 and eq. (19).

Then, by eq. (20), there exists a point $k_0 \in \{0, \cdots, k-1\}$ such that

$$\|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \leqslant \frac{1}{k^{1/3}}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}. \tag{31}$$

Plugging eq. (31) into eqs. (28) and (30) with $k = k_0$, we further obtain that

$$\|\nabla f(\mathbf{x}_{k_0+1})\| \leqslant \frac{1}{k^{2/3}}\frac{L_2 + M}{2}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{2/3}\left(1 + \frac{L_1}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right) \tag{32}$$

$$\lambda_{\min}\left(\nabla^2 f(\mathbf{x}_{k_0+1})\right) \geqslant -\frac{1}{k^{1/3}}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\left(\frac{M + 2L_2}{2} + \frac{L_2}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right). \tag{33}$$

Thus, in order to guarantee $\|\nabla f(\mathbf{x}_{k_0+1})\| \leqslant \epsilon$ in eq. (32) and $\lambda_{\min}\left(\nabla^2 f(\mathbf{x}_{k_0+1})\right) \geqslant -\sqrt{\epsilon}$ in eq. (33), we require

$$k \geqslant \frac{1}{\epsilon^{3/2}}\left(\frac{L_2 + M}{2}\right)^{3/2}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)\left(1 + \frac{L_1}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right)^{3/2}. \tag{34}$$

$$k \geqslant \frac{1}{\epsilon^{3/2}}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)\left(\frac{M + 2L_2}{2} + \frac{L_2}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right)^3. \tag{35}$$

Combining eqs. (34) and (35), we obtain that CRm must pass an $\epsilon$-approximate second-order stationary point if

$$k \geqslant \frac{1}{\epsilon^{3/2}}\max\{c_3, c_4\},$$

where

$$c_3 \triangleq \left(\frac{L_2 + M}{2}\right)^{3/2} \left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right) \left(1 + \frac{L_1}{1 - \rho} \left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right)^{3/2}$$

$$c_4 \triangleq \left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right) \left(\frac{M + 2L_2}{2} + \frac{L_2}{1 - \rho} \left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right)^3.$$

## D   Local Quadratic Convergence: Proof of Theorem 2

We first present the following lemma that characterizes the properties of the sequence $\{\mathbf{x}_k\}_{k \geqslant 0}$ generated by CRm.

**Lemma 10.** *Let Assumption 1 hold and assume that $\mathcal{L}(f(\mathbf{x}_k))$ is bounded for some $k \geqslant 0$. Then, the sequence $\{\mathbf{x}_k\}_{k \geqslant 0}$ generated by CRm and its set of accumulation points $\bar{\mathcal{X}}$ satisfy*

*(i) $v \triangleq \lim_{k \to \infty} f(\mathbf{x}_k)$ exists.*

*(ii) $\lim_{k \to \infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\| = 0$.*

*(iii) The sequence $\{\mathbf{x}_k\}_{k \geqslant 0}$ is bounded.*

*(iv) The set $\bar{\mathcal{X}}$ is bounded and non-empty. Moreover, every $\bar{\mathbf{x}} \in \bar{\mathcal{X}}$, satisfies*

$$f(\bar{\mathbf{x}}) = v, \quad \nabla f(\bar{\mathbf{x}}) = 0, \quad \nabla^2 f(\bar{\mathbf{x}}) \succcurlyeq 0.$$

*Proof.* See Appendix G.5. $\qquad\square$

We next prove the main theorem for CRm. Denote $\bar{\mathbf{x}}_k \in \mathrm{argmin}_{\mathbf{z} \in \mathcal{X}} \|\mathbf{x}_k - \mathbf{z}\|^2$ as the projection of $\mathbf{x}_k$ onto $\mathcal{X}$. Since $\{\mathbf{x}_k\}_{k \geqslant 0}$ is bounded (Lemma 10, (iii)) and $\bar{\mathcal{X}}$ is non-empty and bounded (Lemma 10, (iv)), we conclude that $\lim_{k \to \infty} \mathrm{dist}(\mathbf{x}_k, \bar{\mathcal{X}}) = 0$. By (iv) of Lemma 10 and the definition of $\mathcal{X}$, we have $\bar{\mathcal{X}} \subseteq \mathcal{X}$. Thus, we obtain that

$$\lim_{k \to \infty} \|\mathbf{x}_k - \bar{\mathbf{x}}_k\| = \lim_{k \to \infty} \mathrm{dist}(\mathbf{x}_k, \mathcal{X}) \leqslant \lim_{k \to \infty} \mathrm{dist}(\mathbf{x}_k, \bar{\mathcal{X}}) = 0, \tag{36}$$

which implies that

$$\lim_{k \to \infty} \mathrm{dist}(\mathbf{x}_k, \mathcal{X}) = 0. \tag{37}$$

Therefore, for any $r > 0$, there exists a $k_1 \geqslant 0$ such that $\mathrm{dist}(\mathbf{x}_k, \mathcal{X}) \leqslant r$ for all $k \geqslant k_1$. Combining this with Assumption 2, we obtain that

$$\mathrm{dist}(\mathbf{x}_k, \mathcal{X}) \leqslant \kappa \|\nabla f(\mathbf{x}_k)\|, \quad \forall k \geqslant k_1. \tag{38}$$

Hence, for all $k \geqslant k_1$, we obtain that

$$\begin{aligned}
\mathrm{dist}(\mathbf{x}_{k+1}, \mathcal{X}) &\leqslant \kappa \|\nabla f(\mathbf{x}_{k+1})\| \\
&\leqslant \kappa \|\nabla f(\mathbf{x}_{k+1}) - \nabla f(\hat{\mathbf{x}}_{k+1})\| + \kappa \|\nabla f(\hat{\mathbf{x}}_{k+1})\| \\
&\overset{(i)}{\leqslant} \kappa L_1 \|\mathbf{x}_{k+1} - \hat{\mathbf{x}}_{k+1}\| + \kappa \|\nabla f(\hat{\mathbf{x}}_{k+1})\| \\
&\overset{(ii)}{\leqslant} \kappa L_1 \beta_{k+1} \|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| + \kappa \|\nabla f(\hat{\mathbf{x}}_{k+1})\| \\
&\overset{(iii)}{\leqslant} \kappa \|\nabla f(\hat{\mathbf{x}}_{k+1})\| (L_1 \|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| + 1) \\
&\overset{(iv)}{\leqslant} \kappa \|\nabla f(\hat{\mathbf{x}}_{k+1})\| (L_1 c_5 + 1) \\
&\overset{(v)}{\leqslant} \kappa \left(\frac{L_2 + M}{2}\right) (L_1 c_5 + 1) \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^2
\end{aligned} \tag{39}$$

where (i) follows from the Lipschitz gradient property, (ii) follows from eq. (25), (iii) follows from eq. (5), which implies that $\beta_{k+1} \leqslant \|\nabla f(\hat{\mathbf{x}}_{k+1})\|$, (iv) follows from Lemma 9 and (v) follows from eq. (18). Combining eq. (39) with Lemma 6, we obtain that, for all $k \geqslant k_1$,

$$
\begin{aligned}
\text{dist}(\mathbf{x}_{k+1}, \mathcal{X}) &\leqslant \kappa \left( \frac{L_2 + M}{2} \right) (L_1 c_5 + 1) \, c_1^2 \cdot \text{dist}(\mathbf{x}_k, \mathcal{X})^2 \\
&= c_6 \cdot \text{dist}(\mathbf{x}_k, \mathcal{X})^2,
\end{aligned}
\tag{40}
$$

where $c_6 \triangleq \kappa \left( \frac{L_2 + M}{2} \right) (L_1 c_5 + 1) \, c_1^2$.

Next, we prove that $\{\mathbf{x}_k\}_{k \geqslant 0}$ is Cauchy, and hence is a convergent sequence. For any $\epsilon > 0$, by eq. (37), there exists $k_2 \geqslant 0$ such that

$$
\text{dist}(\mathbf{x}_k, \mathcal{X}) \leqslant \min \left( \frac{1}{2c_6}, \frac{\epsilon}{2c_1(c_5 + 1)} \right), \quad \forall k \geqslant k_2.
\tag{41}
$$

Therefore, for any $k \geqslant \max(k_1, k_2)$ and any $j \geqslant 0$, we have

$$
\begin{aligned}
\|\mathbf{x}_{k+j} - \mathbf{x}_k\| &\leqslant \sum_{i=k}^{k+j-1} \|\mathbf{x}_{i+1} - \mathbf{x}_i\| \leqslant \sum_{i=k}^{k+j-1} \left( \|\mathbf{x}_{i+1} - \hat{\mathbf{x}}_{i+1}\| + \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| \right) \\
&\overset{(i)}{\leqslant} \sum_{i=k}^{k+j-1} \left( \beta_{i+1} \|\hat{\mathbf{x}}_{i+1} - \hat{\mathbf{x}}_i\| + \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| \right) \\
&\overset{(ii)}{\leqslant} \sum_{i=k}^{k+j-1} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| \left( \|\hat{\mathbf{x}}_{i+1} - \hat{\mathbf{x}}_i\| + 1 \right) \\
&\overset{(iii)}{\leqslant} \sum_{i=k}^{k+j-1} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| (c_5 + 1) \overset{(iv)}{\leqslant} \sum_{i=k}^{k+j-1} c_1 \cdot \text{dist}(\mathbf{x}_i, \mathcal{X})(c_5 + 1) \\
&= c_1(c_5 + 1) \sum_{i=k}^{k+j-1} \text{dist}(\mathbf{x}_i, \mathcal{X}) \\
&\overset{(v)}{\leqslant} c_1(c_5 + 1)\text{dist}(\mathbf{x}_k, \mathcal{X}) \sum_{i=0}^{\infty} \frac{1}{2^i} \\
&= 2c_1 (c_5 + 1) \cdot \text{dist}(\mathbf{x}_k, \mathcal{X}) \\
&\overset{(iv)}{\leqslant} \epsilon,
\end{aligned}
\tag{42}
$$

where (i) follows from eq. (25), (ii) follows from eq. (5), which implies that $\beta_{i+1} \leqslant \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|$, (iii) follows from Lemma 9, (iv) follows from Lemma 6, (v) follows from eq. (40) and eq. (41), which implies that $\text{dist}(\mathbf{x}_{k+1}, \mathcal{X}) \leqslant \text{dist}(\mathbf{x}_k, \mathcal{X})/2$ and (iv) follows from eq. (41). Then, we conclude that $\{\mathbf{x}_k\}_{k \geqslant 0}$ is a Cauchy sequence, and thus converges.

Next, we study the convergence rate of $\{\mathbf{x}_k\}_{k \geqslant 0}$. Let $\mathbf{x}^\star \triangleq \lim_{k \to \infty} \mathbf{x}_k$. By (iv) of Lemma 10, we have $\mathbf{x}^\star \in \mathcal{X}$. Then, for all $k \geqslant \max\{k_1, k_2\}$, we obtain that

$$
\begin{aligned}
\|\mathbf{x}^\star - \mathbf{x}_{k+1}\| &= \lim_{j \to \infty} \|\mathbf{x}_{k+1+j} - \mathbf{x}_{k+1}\| \overset{(i)}{\leqslant} 2c_1 (c_5 + 1) \cdot \text{dist}(\mathbf{x}_{k+1}, \mathcal{X}) \\
&\overset{(ii)}{\leqslant} 2c_1 (c_5 + 1) c_6 \cdot \text{dist}(\mathbf{x}_k, \mathcal{X})^2 \overset{(iii)}{\leqslant} 2c_1 (c_5 + 1) c_6 \|\mathbf{x}^\star - \mathbf{x}_k\|^2
\end{aligned}
\tag{43}
$$

where (i) follows from eq. (42), (ii) follows from eq. (40), and (iii) follows from the fact that $\text{dist}(\mathbf{x}_k, \mathcal{X}) \leqslant \|\mathbf{x}^\star - \mathbf{x}_k\|$.

Note that eq. (43) implies that

$$
\frac{\|\mathbf{x}^\star - \mathbf{x}_{k+1}\|}{\|\mathbf{x}^\star - \mathbf{x}_k\|^2} \leqslant 2c_1 (c_5 + 1) c_6, \quad \forall k \geqslant \max\{k_1, k_2\}.
\tag{44}
$$

Hence, $\{\mathbf{x}_k\}_{k\geqslant 0}$ converges at least Q-quadratically to $\mathbf{x}^\star$. In particular, the Q-quadratic convergence region of CRm is given by

$$\|\mathbf{x}_k - \mathbf{x}^\star\| \leqslant \frac{1}{2c_1 c_6 \left(c_5 + 1\right)}. \tag{45}$$

# E  Inexact CRm Convergence: Proof of Theorem 3

We first present the following lemma, which bounds the term $\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\|$, where sequence $\{\mathbf{x}_k\}_{k\geqslant 0}$ is generated by the inexact variant of CRm.

**Lemma 11.** *Let Assumption 1 and Assumption 3 hold. Set $M > 2L_2/3 + 2, \beta_k \leqslant \rho, \epsilon_1 \leqslant 1$. Then the sequence $\{\hat{\mathbf{x}}_k\}_{k\geqslant 0}$ generated by the inexact variant of CRm satisfies*

$$\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| \leqslant c_8 \tag{46}$$

*where $c_8 \triangleq \frac{1}{1-\rho} \left( \left( \frac{3M-2L_2-6}{12} \right)^{-1/3} (f(\mathbf{x}_0) - f^\star)^{1/3} + 1 \right).$*

*Proof.* See Appendix G.6. $\qquad\square$

We next prove the main theorem for CRm. The proof of this theorem is similar to that of the exact case, and hence we only highlight the main difference for simplicity. Define $\epsilon_1 = \theta\sqrt{\epsilon}$ where

$$\theta \triangleq \min \left\{ \sqrt{\frac{2}{\left(1 + L_1 c_8\right)\left(L_2 + M + 2\right)}}, \frac{2}{M + 2L_2 + 2 + L_2 c_8} \right\}.$$

Since $\frac{2}{M+2L_2+2+L_2 c_8} \leqslant 1$ and $\epsilon \leqslant 1$, we obtain that $\epsilon_1 \leqslant 1$.

Note that Lemma 8 implies that if the total number of iterations $k > \left( \frac{12}{3M-2L_2-6} \right) \frac{f(\mathbf{x}_0)-f^\star}{\theta^3 \epsilon^{3/2}}$, then there exists a $k_0 \in \{0, \cdots, k\}$ such that

$$\|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \leqslant \epsilon_1. \tag{47}$$

Following the reasoning similar to that for proving eq. (27), we obtain that

$$\begin{aligned}
\|\nabla f(\mathbf{x}_{k_0+1})\| &\leqslant \|\nabla f(\hat{\mathbf{x}}_{k_0+1})\| \left(1 + L_1 \|\hat{\mathbf{x}}_{k_0+1} - \hat{\mathbf{x}}_{k_0}\|\right) \\
&\overset{(i)}{\leqslant} \|\nabla f(\hat{\mathbf{x}}_{k_0+1})\| \left(1 + L_1 c_8\right) \\
&\overset{(ii)}{\leqslant} \left(1 + L_1 c_8\right) \left( \frac{L_2 + M}{2} \|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\|^2 + \epsilon_1 \|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \right) \\
&\overset{(iii)}{\leqslant} \left(1 + L_1 c_8\right) \left( \frac{L_2 + M + 2}{2} \right) \epsilon_1^2 \overset{(iv)}{\leqslant} \epsilon,
\end{aligned} \tag{48}$$

where (i) follows from Lemma 11, (ii) follows from eq. (22), (iii) follows from eq. (47), and (iv) follows from the definition of $\epsilon_1$.

Then, following the reasoning similar to that for proving eq. (29), we further obtain that

$$\begin{aligned}
\lambda_{\min}\left(\nabla^2 f(\mathbf{x}_{k_0+1})\right) &\geqslant \lambda_{\min}\left(\nabla^2 f(\hat{\mathbf{x}}_{k_0+1})\right) - L_2 \|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \|\hat{\mathbf{x}}_{k_0+1} - \hat{\mathbf{x}}_{k_0}\| \\
&\overset{(i)}{\geqslant} \lambda_{\min}\left(\nabla^2 f(\hat{\mathbf{x}}_{k_0+1})\right) - L_2 c_8 \|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \\
&\overset{(ii)}{\geqslant} -\frac{M + 2L_2}{2} \|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| - \epsilon_1 - L_2 c_8 \|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \\
&\overset{(iii)}{\geqslant} -\frac{M + 2L_2 + 2 + L_2 c_8}{2} \epsilon_1 \overset{(iv)}{\geqslant} -\sqrt{\epsilon},
\end{aligned} \tag{49}$$

where (i) follows from Lemma 11, (ii) follows from eq. (19), (iii) follows from eq. (47), and (iv) follows from the definition of $\epsilon_1$. Combining eqs. (48) and (49), we obtain the main statement of Theorem 3.

# F  Subsampling Technique

To prove Theorem 4, we first establish a useful Proposition (Proposition 1) to characterize the per iteration complexity of Hessian in Appendix F.1, and then establish the overall convergence guarantee in Appendix F.2.

## F.1  Per iteration Complexity

In order to satisfy the inexact criterion $\epsilon_1$ in Assumption 3, the mini-batch size should be large enough to guarantee statistical concentration with high probability (Xu et al., 2017; Tripuraneni et al., 2017; Kohler and Lucchi, 2017; Wang et al., 2019b; Zhou et al., 2018). Such an approach is referred to as the subsampling technique, and has been used in Kohler and Lucchi (2017) to implement the inexact CR. Here, we apply such an approach to CRm, and the following theorem characterizes the sample complexity to guarantee the inexact criterion for each iteration.

**Proposition 1** (Per iteration Hessian sample complexity)**.** *Assuming that Assumption 1 holds for each $f_i(\cdot)$, then sub-sampled mini-batch of Hessians $\mathbf{H}_k, k = 0, 1, \ldots$ satisfies Assumption 3 with probability at least $1 - \zeta$ provided that*

$$|S_1| \geqslant \left( \frac{8L_1^2}{\epsilon_1^2} + \frac{4L_1}{3\epsilon_1} \right) \log \left( \frac{4d}{\zeta} \right). \tag{50}$$

The idea of the proof is to apply the following matrix Bernstein inequality (Tropp, 2012) to characterize the sample complexity in order to satisfy the inexactness condition in Assumption 3 with the probability at least $1 - \zeta$.

**Lemma 12** (Tropp (2012), Theorem 1.6.2)**.** *Consider a finite sequence $\{\mathbf{X}_k\}$ of independent, random matrices with dimensions $d_1 \times d_2$. Assume that each random matrix satisfies*

$$\mathbb{E}\mathbf{X}_k = \mathbf{0} \quad and \quad \|\mathbf{X}_k\| \leqslant R \quad almost\ surely.$$

*Define*

$$\sigma^2 \triangleq \max \left( \left\| \sum_k \mathbb{E}(\mathbf{X}_k \mathbf{X}_k^*) \right\|, \left\| \sum_k \mathbb{E}(\mathbf{X}_k^* \mathbf{X}_k) \right\| \right). \tag{51}$$

*Then, for all $\epsilon \geqslant 0$,*

$$P \left( \left\| \sum_k \mathbf{X}_k \right\| \geqslant \epsilon \right) \leqslant 2(d_1 + d_2) \exp \left( - \frac{\epsilon^2/2}{\sigma^2 + R\epsilon/3} \right).$$

With this lemma in hand, we are ready to prove our main result.

*Proof of Proposition 1.* In order to apply Lemma 12, we first define

$$\mathbf{X}_i = \frac{1}{|S_1|} \left( \nabla^2 f_i(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \right).$$

Then, we obtain that

$$\mathbb{E}\mathbf{X}_i = 0 \tag{52}$$

and

$$\|\mathbf{X}_i\| \leqslant \frac{\|\nabla^2 f_i(\mathbf{x}_k)\| + \|\nabla^2 f(\mathbf{x}_k)\|}{|S_1|} = \frac{2L_1}{|S_1|} \triangleq R. \tag{53}$$

where (i) follows from item 3 of Assumption 1 that $\nabla f_i(\cdot)$ is $L1$-Lipschitz which implies that $\|\nabla^2 f_i(\cdot)\| \leqslant L_1$ and $\|\nabla^2 f(\cdot)\| \leqslant L_1$.

Moreover, we have that

$$\sigma^2 = \max \left( \left\| \sum_{i \in S_1} \mathbb{E}(\mathbf{X}_i \mathbf{X}_i^*) \right\|, \left\| \sum_{i \in S_1} \mathbb{E}(\mathbf{X}_i^* \mathbf{X}_i) \right\| \right)$$

$$\overset{\text{(i)}}{=} \left\| \sum_{i \in S_1} \mathbb{E}(\mathbf{X}_i^2) \right\| \leqslant \sum_{i \in S_1} \left\| \mathbb{E}(\mathbf{X}_i^2) \right\| \overset{\text{(ii)}}{\leqslant} \sum_{i \in S_1} \mathbb{E} \left\| \mathbf{X}_i^2 \right\| \leqslant \sum_{i \in S_1} \mathbb{E} \left\| \mathbf{X}_i \right\|^2 \overset{\text{(iii)}}{\leqslant} \frac{4L_1^2}{|S_1|}, \tag{54}$$

where (i) follows from the fact that $\mathbf{X}_i$ is real and symmetric, (ii) follows from Jasen's inequality, and (iii) follows from eq. (53).

Plugging eqs. (54), (52) and (53) into Lemma 12, we obtain

$$P\left( \left\| \sum_{i \in S_1} \mathbf{X}_k \right\| \geqslant \epsilon_1 \right) \leqslant 4d \exp\left( -\frac{\epsilon_1^2/2}{\frac{4L_1^2}{|S_1|} + \frac{2L_1\epsilon_1}{3|S_1|}} \right). \tag{55}$$

Thus, in order to satisfies $\left\| \sum_{i \in S_1} \mathbf{X}_k \right\| \geqslant \epsilon_1$ with probability at least $1 - \zeta$, it is sufficient to require

$$4d \exp\left( -\frac{\epsilon_1^2/2}{\frac{4L_1^2}{|S_1|} + \frac{2L_1\epsilon_1}{3|S_1|}} \right) \leqslant \zeta, \tag{56}$$

which gives that

$$|S_1| \geqslant \left( \frac{8L_1^2}{\epsilon_1^2} + \frac{4L_1}{3\epsilon_1} \right) \log\left( \frac{4d}{\zeta} \right). \tag{57}$$

$\square$

## F.2 Overall Complexity: Proof of Theorem 4

*Proof.* We first note that Theorem 3 shows that let $\epsilon_1 = \theta\sqrt{\epsilon}$, then the sequence $\{\mathbf{x}_k\}_{k \geqslant 0}$ generated by the inexact CRm contains an $\epsilon$-second-order stationary point if the total number $k$ of iterations satisfies

$$k = \frac{C}{\epsilon^{3/2}}. \tag{58}$$

Next, according to Proposition 1, Assumption 3 is satisfies with probability at least $1 - \zeta$ for Hessian . Thus, according to the union bound, for $k$ iterations, the probability of failure satisfaction of Assumption 3 is at most $k\zeta$. To obtain Assumption 3 holds for the total $k$ iteration with probability least $1 - \delta$, we require

$$1 - k\zeta \geqslant 1 - \delta,$$

which yields

$$\zeta \leqslant \frac{\delta}{k}.$$

Thus, with probability $1 - \delta$, the algorithms successfully outputs an $\epsilon$ approximated second-order stationary point if we set $\zeta = \delta/k$. Therefore, according to Proposition 1, Assumption 3 with $\epsilon_1 = \theta\sqrt{\epsilon}$ holds with probability at least $1 - \zeta$ given that

$$|S_1| = \left( \frac{8L_1^2}{\theta^2\epsilon} + \frac{4L_1}{3\theta\sqrt{\epsilon}} \right) \log\left( \frac{4dk}{\delta} \right), \tag{59}$$

and the total Hessian sample complexity is bounded by

$$S = k \times |S_1| = k \times \left( \frac{8L_1^2}{\theta^2\epsilon} + \frac{4L_1}{3\theta\sqrt{\epsilon}} \right) \log\left( \frac{4dk}{\delta} \right) \overset{\text{(i)}}{\leqslant} C\left( \frac{8L_1^2}{\theta^2\epsilon^{5/2}} + \frac{4L_1}{3\theta\epsilon^2} \right) \log\left( \frac{4d}{\epsilon\delta} \right). \tag{60}$$

where (i) follows from eq. (58).

$\square$

# G Proof of Technical Lemmas

In this section, we provide the proofs of the technical lemmas.

## G.1 Useful Inequality

**Lemma 13.** *For $x, \Lambda \in \mathbb{R}$, $0 < x \leqslant 1$, and $0 < \Lambda < 1$, the following inequality holds*

$$x(1 - \Lambda) \log \left( \frac{1}{1 - \Lambda} \right) \leqslant 1 - (1 - \Lambda)^x. \tag{61}$$

*Proof.* Let

$$f(x) = 1 - (1 - \Lambda)^x - x(1 - \Lambda) \log \left( \frac{1}{1 - \Lambda} \right),$$

it is sufficient to prove $f(x) \geqslant 0$ for $0 < x \leqslant 1$, and $0 < \Lambda < 1$. We first note that

$$\nabla f(x) = -(1 - \Lambda)^x \log(1 - \Lambda) - (1 - \Lambda) \log \left( \frac{1}{1 - \Lambda} \right), \tag{62}$$

which is decreasing with respect to $x$. Thus, we have, for $0 < x \leqslant 1$

$$\nabla f(x) \geqslant \nabla f(1) = 0,$$

which implies $f(x)$ increasing within $0 < x \leqslant 1$. Thus,

$$f(x) \geqslant f(0) = 0. \tag{63}$$

Therefore, we complete our proof. $\qquad \square$

## G.2 Proof of Lemma 7

*Proof.* We first prove eq. (16). Define $\mathbf{s}_k \triangleq \hat{\mathbf{x}}_{k+1} - \mathbf{x}_k$. Then, we obtain that

$$
\begin{aligned}
f(\hat{\mathbf{x}}_{k+1}) &\overset{(i)}{\leqslant} f(\mathbf{x}_k) + \nabla f(\mathbf{x})^T \mathbf{s}_k + \frac{1}{2} \mathbf{s}_k^T \nabla^2 f(\mathbf{x}) \mathbf{s}_k + \frac{L_2}{6} \|\mathbf{s}_k\|^3 \\
&\overset{(ii)}{\leqslant} f(\mathbf{x}_k) - \frac{M}{12} \|\mathbf{s}_k\|^3 + \frac{L_2 - M}{6} \|\mathbf{s}_k\|^3 \\
&= f(\mathbf{x}_k) - \frac{3M - 2L_2}{12} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^3,
\end{aligned}
$$

where (i) follows from lemma 5, and (ii) follows from Lemma 4 in Nesterov and Polyak (2006) and the definition of $\hat{\mathbf{x}}_{k+1}$. Then, we further obtain that

$$f(\hat{\mathbf{x}}_{k+1}) - f(\mathbf{x}_k) \leqslant -\gamma \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^3, \tag{64}$$

which gives eq. (16).

Next, we prove eq. (17). Note that eq. (16) implies that, for all $i \geqslant 0$,

$$\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \overset{(i)}{\leqslant} \frac{f(\mathbf{x}_i) - f(\hat{\mathbf{x}}_{i+1})}{\gamma} \overset{(ii)}{\leqslant} \frac{f(\mathbf{x}_i) - f(\mathbf{x}_{i+1})}{\gamma}, \tag{65}$$

where (i) follows from eq. (64), and (ii) follows from the definition of $\mathbf{x}_{i+1}$. Summing eq. (65) over $i$ from 0 to $k$, we obtain that

$$\sum_{i=0}^{k} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \leqslant \frac{f(\mathbf{x}_0) - f(\mathbf{x}_{k+1})}{\gamma} \leqslant \frac{f(\mathbf{x}_0) - f^\star}{\gamma},$$

which gives eq. (17).

To prove eqs. (18) and (19), note that

$$\hat{\mathbf{x}}_{k+1} = \operatorname*{argmin}_{\mathbf{s} \triangleq \mathbf{x} - \mathbf{x}_k} \nabla f(\mathbf{x}_k)^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \nabla^2 f(\mathbf{x}_k) \mathbf{s} + \frac{M}{6} \|\mathbf{s}\|^3.$$

Then, Lemma 5 in Nesterov and Polyak (2006) directly implies eqs. (18) and (19).

Next, we prove eq. (20). Note that

$$\min_{0 \leqslant i \leqslant k} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \leqslant \frac{1}{k} \sum_{i=0}^{k} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \overset{(i)}{\leqslant} \frac{1}{k} \frac{f(\mathbf{x}_0) - f^\star}{\gamma}, \tag{66}$$

where (i) follows from eq. (17). Then, eq. (20) follows by taking the cubic root on both sides of eq. (66). □

### G.3  Proof of Lemma 8

*Proof.* We first prove eq. (21). Note that

$$
\begin{aligned}
f(\hat{\mathbf{x}}_{k+1}) \\
&\overset{(i)}{\leqslant} f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{s}_k + \frac{1}{2} \mathbf{s}_k^T \nabla^2 f(\mathbf{x}_k) \mathbf{s}_k + \frac{L_2}{6} \|\mathbf{s}_k\|^3 \\
&= f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{s}_k + \frac{1}{2} \mathbf{s}_k^T \mathbf{H}_k \mathbf{s}_k + \frac{M}{6} \|\mathbf{s}_k\|^3 + \frac{L_2 - M}{6} \|\mathbf{s}_k\|^3 + \frac{1}{2} \mathbf{s}_k^T (\nabla^2 f(\mathbf{x}_k) - \mathbf{H}_k) \mathbf{s}_k \\
&\overset{(ii)}{\leqslant} f(\mathbf{x}_k) - \frac{M}{12} \|\mathbf{s}_k\|^3 + \frac{L_2 - M}{6} \|\mathbf{s}_k\|^3 + \frac{1}{2} \mathbf{s}_k^T (\nabla^2 f(\mathbf{x}_k) - \mathbf{H}_k) \mathbf{s}_k \\
&\overset{(iii)}{\leqslant} f(\mathbf{x}_k) - \frac{3M - 2L_2}{12} \|\mathbf{s}_k\|^3 + \frac{1}{2} \|\mathbf{s}_k\|^2 \epsilon_1,
\end{aligned}
$$

where (i) follows from Lemma 5, (ii) follows from Lemma 4 in Nesterov and Polyak (2006) and the fact that $\hat{\mathbf{x}}_{k+1} = \operatorname{argmin}_{\mathbf{s} \triangleq \mathbf{x} - \mathbf{x}_k} \nabla f(\mathbf{x}_k)^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H}_k \mathbf{s} + \frac{M}{6} \|\mathbf{s}\|^3$, and (iii) follows from Assumption 3.

Next, we prove eq. (22). Note that

$$\hat{\mathbf{x}}_{k+1} = \operatorname*{argmin}_{\mathbf{s} \triangleq \mathbf{x} - \mathbf{x}_k} \nabla f(\mathbf{x}_k)^T \mathbf{s} + \frac{1}{2} \mathbf{s}^T \mathbf{H}_k \mathbf{s} + \frac{M}{6} \|\mathbf{s}\|^3. \tag{67}$$

By the first-order optimality condition, we obtain that

$$\nabla f(\mathbf{x}_k) + \mathbf{H}_k \mathbf{s}_k + \frac{M}{2} \mathbf{s}_k \|\mathbf{s}_k\| = 0. \tag{68}$$

Then, we further obtain that

$$
\begin{aligned}
\|\nabla f(\hat{\mathbf{x}}_{k+1})\| &\overset{(i)}{=} \left\| \nabla f(\hat{\mathbf{x}}_{k+1}) - \nabla f(\mathbf{x}_k) - \mathbf{H}_k \mathbf{s}_k - \frac{M}{2} \mathbf{s}_k \|\mathbf{s}_k\| \right\| \\
&\leqslant \|\nabla f(\hat{\mathbf{x}}_{k+1}) - \nabla f(\mathbf{x}_k) - \mathbf{H}_k \mathbf{s}_k\| + \frac{M}{2} \|\mathbf{s}_k\|^2 \\
&\leqslant \|\nabla f(\hat{\mathbf{x}}_{k+1}) - \nabla f(\mathbf{x}_k) - \nabla^2 f(\mathbf{x}_k) \mathbf{s}_k\| + \|(\nabla^2 f(\mathbf{x}_k) - \mathbf{H}_k) \mathbf{s}_k\| + \frac{M}{2} \|\mathbf{s}_k\|^2 \\
&\overset{(ii)}{\leqslant} \frac{L_2 + M}{2} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^2 + \epsilon_1 \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|,
\end{aligned}
$$

where (i) follows from eq. (68), and (ii) follows from Lemma 5, Assumption 3, and the fact that $\mathbf{s}_k \triangleq \hat{\mathbf{x}}_{k+1} - \mathbf{x}_k$.

Next, we prove eq. (23). By eq. (67) and Proposition 1 in Nesterov and Polyak (2006), we obtain that

$$\mathbf{H}_k \succcurlyeq -\frac{M}{2} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| \mathbf{I}. \tag{69}$$

Then, we further obtain that

$$
\begin{aligned}
\lambda_{\min}(\nabla^2 f(\mathbf{x}_{k+1})) &\overset{(i)}{\geqslant} \lambda_{\min}(\mathbf{H}_k) - \|\nabla^2 f(\mathbf{x}_{k+1}) - \mathbf{H}_k\| \\
&\overset{(ii)}{\geqslant} -\frac{M}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| - \|\nabla^2 f(\mathbf{x}_{k+1}) - \nabla^2 f(\mathbf{x}_k)\| - \|\nabla^2 f(\mathbf{x}_k) - \mathbf{H}_k\| \\
&\overset{(iii)}{\geqslant} -\frac{M}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| - L_2\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| - \epsilon_1 \\
&= -\frac{M + 2L_2}{2}\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| - \epsilon_1,
\end{aligned}
$$

where (i) follows from Wely's inequality, (ii) follows from eq. (69) and (iii) follows from Assumption 3 and the fact that $\nabla^2 f$ is $L_2$-Lipschitz.

Next, we prove eq. (24) by contradiction. Suppose for every $i \in \{0, \cdots, k\}$ it holds that

$$
\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| > \epsilon_1. \tag{70}
$$

Then, eq. (21) further implies that

$$
\begin{aligned}
f(\hat{\mathbf{x}}_{i+1}) - f(\mathbf{x}_i) &\leqslant -\frac{3M - 2L_2}{12}\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 + \frac{1}{2}\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^2 \epsilon_1 \\
&\overset{(i)}{\leqslant} -\left(\frac{3M - 2L_2 - 6}{12}\right)\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3,
\end{aligned}
$$

where (i) follows from eq. (70). Therefore, we have

$$
\left(\frac{3M - 2L_2 - 6}{12}\right)\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \leqslant f(\mathbf{x}_i) - f(\hat{\mathbf{x}}_{i+1}) \overset{(i)}{\leqslant} f(\mathbf{x}_i) - f(\mathbf{x}_{i+1}), \tag{71}
$$

where (i) follows from the definition of $\mathbf{x}_{k+1}$. Summing up eq. (71) over $i$ from 0 to $k$, we obtain that

$$
\sum_{i=0}^{k} \left(\frac{3M - 2L_2 - 6}{12}\right)\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \leqslant f(\mathbf{x}_0) - f^\star. \tag{72}
$$

Combining eq. (72) with the fact that $\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| > \epsilon_1$ for $i \in \{0, \cdots, k\}$ and $M > 2L_2/3 + 2$, we have

$$
k\left(\frac{3M - 2L_2 - 6}{12}\right)\epsilon_1^3 \leqslant \sum_{i=0}^{k} \left(\frac{3M - 2L_2 - 6}{12}\right)\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \leqslant f(\mathbf{x}_0) - f^\star,
$$

which gives

$$
k \leqslant \left(\frac{12}{3M - 2L_2 - 6}\right)\frac{f(\mathbf{x}_0) - f^\star}{\epsilon_1^3}.
$$

This contradicts with our assumption that $k > \left(\frac{12}{3M - 2L_2 - 6}\right)\frac{f(\mathbf{x}_0) - f^\star}{\epsilon_1^3}$. Therefore, there must exist an integer $k_0 \in \{0, \cdots k\}$ such that

$$
\|\hat{\mathbf{x}}_{k_0+1} - \mathbf{x}_{k_0}\| \leqslant \epsilon_1, \tag{73}
$$

and the proof is complete. □

### G.4 Proof of Lemma 9

*Proof.* For $i \geqslant 1$, note that

$$
\|\hat{\mathbf{x}}_{i+1} - \hat{\mathbf{x}}_i\| \overset{(i)}{\leqslant} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| + \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|
$$

$$\overset{(ii)}{\leqslant} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| + \beta_i\|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_{i-1}\|$$

$$\overset{(iii)}{\leqslant} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| + \rho\|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_{i-1}\|, \tag{74}$$

where (i) follows from the triangle inequality, (ii) follows from eq. (25), and (iii) follows from the fact that $\beta_{k+1} \leqslant \rho$ for all $k \geqslant 0$.

Recursively applying eq. (74), we obtain that

$$\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| \leqslant \rho^k\|\hat{\mathbf{x}}_1 - \hat{\mathbf{x}}_0\| + \sum_{i=1}^{k} \rho^{k-i}\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|$$

$$\overset{(i)}{\leqslant} \sum_{i=0}^{k} \rho^{k-i}\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|, \tag{75}$$

where (i) follows because $\hat{\mathbf{x}}_0 = \mathbf{x}_0$ in Algorithm 1. Note that eq. (75) is also true for $k = 0$. Then, by eq. (75), we further obtain that

$$\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| \leqslant \sum_{i=0}^{k} \rho^{k-i}\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|$$

$$\leqslant \max_{i\in\{0,\cdots,k\}} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| \sum_{i=0}^{k} \rho^{k-i}$$

$$\overset{(i)}{\leqslant} \max_{i\in\{0,\cdots,k\}} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| \frac{1}{1-\rho} \tag{76}$$

$$\overset{(ii)}{\leqslant} \frac{1}{1-\rho} \left( \frac{f(\mathbf{x}_0) - f^\star}{\gamma} \right)^{1/3},$$

where (i) follows from the fact that $\rho < 1$ and (ii) follows from eq. (17). The proof of Lemma 9 is complete. $\qquad\square$

### G.5 Proof of Lemma 10

*Proof.* To prove item $(i)$, it suffices to show that $\{f(\mathbf{x}_k)\}_{k\geqslant 0}$ is a decreasing sequence with a lower bound. By Assumption 1, $f$ is bounded below. Thus, $\{f(\mathbf{x}_k)\}_{k\geqslant 0}$ is bounded below. Also, note that

$$f(\mathbf{x}_{k+1}) \overset{(i)}{\leqslant} f(\hat{\mathbf{x}}_{k+1}) \overset{(ii)}{\leqslant} f(\mathbf{x}_k) \tag{77}$$

where (i) follows from eq. (6), and (ii) follows from eq. (16). Thus, $\{f(\mathbf{x}_k)\}_{k\geqslant 0}$ is a decreasing sequence and is bounded below, which further imply that $\{f(\mathbf{x}_k)\}_{k\geqslant 0}$ converges. We denote the corresponding limit as $v$.

To prove item $(ii)$, note that

$$\lim_{k\to\infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leqslant \lim_{k\to\infty} \|\mathbf{x}_{k+1} - \hat{\mathbf{x}}_{k+1}\| + \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|$$

$$\overset{(i)}{\leqslant} \lim_{k\to\infty} \beta_{k+1}\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| + \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|$$

$$\overset{(ii)}{\leqslant} \lim_{k\to\infty} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| \left( \|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| + 1 \right)$$

$$\overset{(iii)}{\leqslant} \lim_{k\to\infty} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| \left( \frac{1}{1-\rho} \left( \frac{f(\mathbf{x}_0) - f^\star}{\gamma} \right)^{1/3} + 1 \right)$$

$$\overset{(iv)}{=} 0 \tag{78}$$

where (i) follows from eq. (25), (ii) follows from eq. (5), which implies that $\beta_{k+1} \leqslant \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|$, (iii) follows from Lemma 9 and (iv) follows from eq. (17), which implies that

$$\lim_{k\to\infty} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| = 0. \tag{79}$$

Then, we conclude that $\lim_{k\to\infty} \|\mathbf{x}_{k+1} - \mathbf{x}_k\| = 0$.

To prove item $(iii)$, note that eq. (77) implies that $\mathbf{x}_k \in \mathcal{L}(f(\mathbf{x}_0))$ for all $k$. By the assumption that $\mathcal{L}(f(\mathbf{x}_k))$ is bounded for some $k \geqslant 0$, we conclude that $\{\mathbf{x}_k\}_{k\geqslant 0}$ is also bounded.

To prove item $(iv)$, note that the Bolzano-Weierstarss theorem and item $(iii)$ of Lemma 10 imply that $\{\mathbf{x}_k\}_{k\geqslant 0}$ has a convergent subsequence. Also, the set of its accumulation points $\bar{\mathcal{X}}$ is bounded. Moreover, for every accumulation point $\bar{\mathbf{x}}$, by eqs. (28), (30) and (79), we obtain that

$$\|\nabla f(\bar{\mathbf{x}})\| \leqslant \limsup_{k\to\infty} \|\nabla f(\mathbf{x}_{k+1})\|$$

$$\leqslant \limsup_{k\to\infty} \frac{L_2}{2} \|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\|^2 \left(1 + \frac{L_1}{1-\rho} \left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right) = 0,$$

and

$$\lambda_{\min}\left(\nabla^2 f(\bar{\mathbf{x}})\right) \geqslant \liminf_{k\to\infty} \lambda_{\min}\left(\nabla^2 f(\mathbf{x}_{k+1})\right)$$

$$\geqslant \liminf_{k\to\infty} -\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| \left(\frac{M + 2L_2}{2} + \frac{L_2}{1-\rho}\left(\frac{f(\mathbf{x}_0) - f^\star}{\gamma}\right)^{1/3}\right) = 0.$$

Thus, we conclude that $\nabla f(\bar{\mathbf{x}}) = 0, \nabla^2 f(\bar{\mathbf{x}}) \succcurlyeq 0$. Furthermore, item $(i)$ of **??** implies that $f(\mathbf{x}_k)_{k\geqslant 0}$ converges to its limit $v$.

$\square$

## G.6 Proof of Lemma 11

*Proof.* Following the proof similar to that of Lemma 9, one can show that eq. (76) also holds for the inexact algorithm, i.e.,

$$\|\hat{\mathbf{x}}_{k+1} - \hat{\mathbf{x}}_k\| \leqslant \max_{i\in\{0,\cdots,k\}} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| \frac{1}{1-\rho}. \tag{80}$$

Then, it suffices to bound $\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|$. Suppose that the inexact variant of CRm2 terminates at iteration $k$. By the termination criterion, we have

$$\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| > \epsilon_1 \quad \text{for} \quad 0 \leqslant i \leqslant k-1, \tag{81}$$

and

$$\|\hat{\mathbf{x}}_{k+1} - \mathbf{x}_k\| \leqslant \epsilon_1 \leqslant 1. \tag{82}$$

For $0 \leqslant i \leqslant k-1$, eq. (21) implies that

$$f(\hat{\mathbf{x}}_{i+1}) - f(\mathbf{x}_i) \leqslant -\frac{3M - 2L_2}{12} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 + \frac{1}{2}\|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^2 \epsilon_1$$

$$\overset{(i)}{\leqslant} -\left(\frac{3M - 2L_2 - 6}{12}\right) \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \tag{83}$$

where (i) follows from eq. (81). Summing eq. (83) over $i$ from 0 to $k-1$, we obtain that

$$\sum_{i=0}^{k-1} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\|^3 \leqslant \left(\frac{3M - 2L_2 - 6}{12}\right)^{-1} (f(\mathbf{x}_0) - f^\star), \tag{84}$$

which further implies that

$$\max_{i\in\{0,\cdots,k-1\}} \|\hat{\mathbf{x}}_{i+1} - \mathbf{x}_i\| \leqslant \left(\frac{3M - 2L_2 - 6}{12}\right)^{-1/3} (f(\mathbf{x}_0) - f^\star)^{1/3}. \tag{85}$$

Combining eqs. (80), (82) and (85), we obtain the statement of Lemma 11. $\square$