Stochastic Variance-Reduced Cubic Regularization for Nonconvex Optimization

Zhe Wang

Yi Zhou

Ohio State University wang.10982@osu.edu

Duke University yi.zhou610@duke.edu Yingbin Liang

Guanghui Lan

Ohio State University liang.889@osu.edu

Georgia Institute of Technology george.lan@isye.gatech.edu

Abstract

Cubic regularization (CR) is an optimization method with emerging popularity due to its capability to escape saddle points and converge to second-order stationary solutions for nonconvex optimization. However, CR encounters a high sample complexity issue for finite-sum problems with a large data size. In this paper, we propose a stochastic variancereduced cubic-regularization (SVRC) method under random sampling, and study its convergence guarantee as well as sample complexity. We show that the iteration complexity of SVRC for achieving a second-order stationary solution within ϵ accuracy is $\mathcal{O}(\epsilon^{-3/2})$, which matches the state-of-art result on CR types of methods. Moreover, our proposed variance reduction scheme significantly reduces the periteration sample complexity. The resulting total Hessian sample complexity of our SVRC is $\tilde{\mathcal{O}}(N^{2/3}\epsilon^{-3/2})$, which outperforms the stateof-art result by a factor of $\tilde{\mathcal{O}}(N^{2/15})$. We also study our SVRC under random sampling without replacement scheme, which yields a lower per-iteration sample complexity, and hence justifies its practical applicability.

1 Introduction

Many machine learning problems are formulated as finite-sum nonconvex optimization problems that take the form

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \triangleq \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}), \tag{1}$$

Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS) 2019, Naha, Okinawa, Japan. PMLR: Volume 89. Copyright 2019 by the author(s).

where each component function f_i corresponds to the loss on the i-th data sample. While finding global optimal solutions of generic nonconvex optimization problems are challenging, various nonconvex problems in the form of eq. (1) have been shown to possess good landscape properties that facilitate convergence. For example, the square loss of a shallow linear neural network is shown to have only strict saddle points other than local minimum (Baldi and Hornik, 1989; Zhou and Liang, 2018). The same property also holds for some other nonconvex problems such as phase retrieval (Sun et al., 2017) and matrix factorization (Ge et al., 2016; Bhojanapalli et al., 2016). Such a remarkable property has motivated a growing research interest in designing algorithms that can escape strict saddle points and have guaranteed convergence to local minimum, and even to global minimum for problems without spurious local minimum.

Various algorithms have been designed to have the capability to escape strict saddle points in nonconvex optimization. Such a desired property requires that the obtained solution \mathbf{x}^* satisfies the second-order stationary conditions within an ϵ accuracy, i.e.,

$$\|\nabla F(\mathbf{x}^*)\| \leqslant \epsilon, \qquad \nabla^2 F(\mathbf{x}^*) \succcurlyeq -\sqrt{\epsilon} \mathbf{I}.$$
 (2)

Therefore, upon convergence, the gradient is guaranteed to be close to zero and the Hessian is guaranteed to be almost positive semidefinite, which thresh-out the possibility to converge to strict saddle points. Among these algorithms (which are reviewed in related work), the cubic-regularized Newton's method (also called cubic regularization or CR) (Nesterov and Polyak, 2006) is a popular method that provides the second-order stationary guarantee for the obtained solution. At each iteration k, CR solves a sub-problem that approximates the objective function in eq. (1) with a cubic-regularized second-order Taylor's expansion at the current iterate \mathbf{x}_k . In specific, the update rule of CR can be written

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

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_{k+1}.\tag{3}$$

It has been shown that CR converges to a point satisfying the second-order stationary condition (eq. (2)) within $\mathcal{O}(\epsilon^{-3/2})$ number of iterations. However, fully solving the exact cubic sub-problem in eq. (3) requires a high computation complexity, especially due to the computation of the Hessian matrices for loss functions on all the data samples. To evaluate the complexity of CR type algorithms, we define the stochastic Hessian oracle (SHO) as follows. Given a point x and the component number i, the oracle returns the corresponding Hessian $\nabla^2 f_i(\mathbf{x})$. Moreover, we define the subproblem oracle (SO) as a subroutine, which for a given a point \mathbf{x} , returns the minimizer of eq. (3). In Cartis et al. (2011), the authors proposed an inexact cubic-regularized (inexact-CR) Newton's method, which formulates the cubic sub-problem in eq. (3) with an inexact Hessian \mathbf{H}_k that satisfies

$$\|(\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k))\mathbf{s}_{k+1}\| \leqslant C \|\mathbf{s}_{k+1}\|^2, \qquad (4)$$

where $C\geqslant 0$ is a certain numerical constant. In particular, Cartis et al. (2011) showed that such an inexact method achieves the same order of theoretical guarantee as the original CR. This inexact condition has been explored in various situations (Kohler and Lucchi, 2017; Cartis et al., 2012a,b; Zhou et al., 2018). Especially, in order to satisfy the inexact Hessian condition in eq. (4), Kohler and Lucchi (2017) proposed a practical sub-sampling scheme (referred to SCR) to implement the inexact-CR. Specifically, at each iteration k, SCR collects two index sets $\xi_g(k), \xi_H(k)$ whose elements are sampled uniformly from $\{1,\ldots,N\}$ at random, and then evaluates respectively the gradients and Hessians of the corresponding component functions, i.e., $\mathbf{g}_k \triangleq \frac{1}{|\xi_g(k)|} \sum_{i \in \xi_g(k)} \nabla f_i(\mathbf{x}_k)$ and $\mathbf{H}_k \triangleq \frac{1}{|\xi_H(k)|} \sum_{i \in \xi_H(k)} \nabla^2 f_i(\mathbf{x}_k)$. Then, SCR solves the following cubic sub-problem at the k-th iteration.

$$\mathbf{s}_{k+1} = \operatorname*{argmin}_{\mathbf{s} \in \mathbb{R}^d} \mathbf{g}_k^{\top} \mathbf{s} + \tfrac{1}{2} \mathbf{s}^{\top} \mathbf{H}_k \mathbf{s} + \tfrac{M}{6} \| \mathbf{s} \|^3.$$

Kohler and Lucchi (2017) showed that if the mini-batch sizes to satisfy

$$|\xi_g(k)| \geqslant \mathcal{O}\left(\frac{1}{\|\mathbf{s}_{k+1}\|^4}\right), |\xi_H(k)| \geqslant \mathcal{O}\left(\frac{1}{\|\mathbf{s}_{k+1}\|^2}\right),$$
(5)

then the sub-sampled mini-batch of Hessians \mathbf{H}_k satisfies eq. (4) and the sub-sampled mini-batch of gradients \mathbf{g}_k satisfies

$$\|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\| \leqslant C_1 \|\mathbf{s}_{k+1}\|^2,$$
 (6)

where $C_1 \ge 0$ is a certain numerical constant, which further guarantee the same convergence rate for SCR as that the original exact CR.

Three important issues here motivate our design of a new sub-sampling CR algorithm.

- It can be seen from eq. (5) that as the algorithm converges, i.e., $\mathbf{s}_{k+1} \to \mathbf{0}$, the required sample size of SCR in Kohler and Lucchi (2017) grows polynomially fast, resulting significant increase in computational complexity. Thus, an important open issue here is to design an improved sub-sampling CR algorithm that reduces the sample complexity (and correspondingly computational complexity) particularly when the algorithm approaches to convergence.
- Another reason for the above pessimistic bound is because that Kohler and Lucchi (2017) analyzed the sample complexity for sampling with replacement, whereas in practice sampling without replacement can potentially have much lower sample complexity. As a clear evidence, the sample complexity for sampling with replacement to achieve a certain accuracy can be unbounded, whereas this for sampling without replacement can only be as large as the total sample size. Thus, the second open issue is to develop bounds for sampling without replacement in order to provide more precise guidance for sub-sampled CR methods.
- We also observe that eqs. (4) and (6) involve $\|\mathbf{s}_{k+1}\|$ (and hence \mathbf{x}_{k+1}), which is not available at iteration k. Kohler and Lucchi (2017) used s_k to replace s_{k+1} in experiments but not theory. A more recent study Wang et al. (2019) theoretically justified such a replacement with the convergence analysis, but not for stochastic sub-sampling scheme, for which the convergence analysis requires considerable efforts.

In this paper, we address the aforementioned open issues, and our contributions are summarized as follows.

Our Contributions

We propose a stochastic variance reduced cubic-regularized (SVRC) Newton's algorithm, which combines the variance reduced technique with concentration inequality under sub-sampling scheme. We show that the computation of the full Hessian and gradient can facilitate many steps of efficient inner-loop iteration as well as accurate approximation of Hessian and gradient under high probability perspective. SVRC can be associated with two sampling schemes, respectively with and without replacement.

We establish the convergence guarantee of SVRC with high probability under the implementable inexact condition similar with $\|\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)\| \leq C \|\mathbf{s}_k\|$. We show that the convergence of SVRC is at the same rate

 $^{^1{\}rm We}$ note that SVRC(ZSG) does not need the objective function and its gradient to be Lipschitz but we adopt such assumptions.

Algorithms		Total SHO	Total SO
CR	(Nesterov and Polyak, 2006)	$\mathcal{O}(N\epsilon^{-3/2})$	$\mathcal{O}(\epsilon^{-3/2})$
SCR	(Kohler and Lucchi, 2017)	$\mathcal{O}(\epsilon^{-5/2})$	$\mathcal{O}(\epsilon^{-3/2})$
Inexact CR	(Xu et al., 2017)	$\mathcal{O}(\epsilon^{-5/2})$	$\mathcal{O}(\epsilon^{-3/2})$
SVRC(ZXG)	(Zhou et al., 2018)	$\mathcal{O}(N^{4/5}\epsilon^{-3/2})$	$\mathcal{O}(\epsilon^{-3/2})$
SVRC	(This Work)	$\tilde{\mathcal{O}}(N^{2/3}\epsilon^{-3/2})$	$\mathcal{O}(\epsilon^{-3/2})$

Table 1: Comparison of total Hessian sample complexity

 $(O(\epsilon^{-3/2}))$ as the original CR (Nesterov and Polyak, 2006) or the other type of inexact-CR in Cartis et al. (2011, 011b); Kohler and Lucchi (2017).

We then develop the bounds on the total Hessian sample complexity of SVRC. We show that SVRC achieves $\tilde{\mathcal{O}}(N^{2/3}\epsilon^{-3/2})$ Hessian sample complexity (where we use $\tilde{\mathcal{O}}$ to hide the dependence on log factors), which outperforms CR (Nesterov and Polyak, 2006) by an order of $\tilde{\mathcal{O}}(N^{1/3})$ and outperform SCR (Kohler and Lucchi, 2017) in the regime of high accuracy requirement. Furthermore, our proposed SVRC order-wise outperforms the algorithm SVRC(ZSG) (Zhou et al., 2018) by an order of $\tilde{\mathcal{O}}(N^{2/15})$, which is also a variance reduced cubic regularized method concurrently proposed. A detailed comparison among these algorithms are summarized in Table 1.

We further provide an analysis for the case under sampling without replacement by developing a new concentration bound for sampling without replacement for random *matrices* by generalizing that for *scalar* random variables in Bardenet and Maillard (2015). Our result shows that sample replacement has lower sample complexity than that of with replacement in each iteration.

Related Works

Escaping saddle points: Various algorithms have been developed to escape strict saddle points and converge to local minimum for nonconvex optimization. The first-order such algorithms include the gradient descent algorithm with random initialization (Lee et al., 2016) and with injection of random noise (Rong et al., 2015; Chi et al., 2017). Various second-order algorithms were also proposed. In particular, Xu et al. (2017); Liu and Yang (2017); Carmon et al. (2016) proposed algorithms that exploit the negative curvature of Hessian to escape saddle points. The CR method as we describe below is another type of second-order algorithm that

has been shown to escape strict saddle points.

CR type of algorithms: The CR method was shown in Nesterov and Polyak (2006) that converges to a point that satisfies the first- and second-order optimality condition for nonconvex optimization. Its accelerated version was proposed in Nesterov (2008) and the convergence rate was characterized for convex optimization. Several methods have been proposed to solve the cubic sub-problem in CR more efficiently. Cartis et al. (2011) proposed to approximately solve the cubic sub-problem in Krylov space. Agarwal et al. (2017) proposed an alternative fast way to solve the sub-problem. Carmon and Duchi (2016) proposed a method based on gradient descent. Zhou et al. (2018) studied asymptotic convergence rate of CR under the nonconvex KŁ condition, and Wang et al. (2018) established convergence guarantee for CR with momentum in nonconvex optimization.

Inexact CR algorithms: Various inexact approaches were proposed to approximate Hessian and gradient in order to reduce the computational complexity for CR type of algorithms. In particular, Ghadimi et al. (2017) studied the inexact CR and accelerated CR for convex optimization, where the inexactness is fixed throughout the iterations. Tripuraneni et al. (2017) studied a similar inexact CR for nonconvex optimization. Alternatively, Cartis et al. (2011, 011b) studied the inexact CR for nonconvex optimization, where the inexact condition is adaptive during the iterations. Wang et al. (2019) established the convergence result of CR under a more reasonable inexact condition. Jiang et al. (2017) studied the adaptive inexact accelerated CR for convex optimization. In practice, sub-sampling is a very common approach to implement inexact algorithms. Kohler and Lucchi (2017) proposed a sub-sampling scheme that adaptively changes the sample complexity to guarantee the inexactness condition in Cartis et al. (2011, 011b). Xu et al. (2017) proposed uniform and nonuniform sub-sampling algorithms with fixed

inexactness condition for nonconvex optimization.

Stochastic variance reduced algorithms: Stochastic variance reduced algorithms have been applied to various first-order algorithms (known as SVRG algorithms), and the convergence rate has been studied for convex functions in, e.g., Johnson and Zhang (2013); Xiao and Zhang (2014) and for nonconvex functions in, e.g., Reddi et al. (2016); Li et al. (2017); Fang et al. (2018); Wang et al. (2018). Zhou et al. (2018) proposed a variance reduction version of CR. In this paper, we proposed another type of stochastic variance reduction to the second-order CR method to improve the state-of-art sample complexity result of approximating Hessian and gradient in probability perspective, and analyzed it in with and without replacement schemes.

Sampling without replacement: The sampling without replacement scheme for first-order methods has been studied by various papers. Recht and Re (2012) and Shamir (2016) studied stochastic gradient descent under sampling without replacement for least square problems. Gürbüzbalaban et al. (2015) provided convergence rate of the random reshuffling method. As for the sampling without replacement bounds. Hoeffding (1963) showed that the bound for sampling with replacement also holds for sampling without replacement. Friedlander and Schmidt (2012) provided deterministic bounds for without replacement sampling schemes for gradient approximations under certain assumptions. Bardenet and Maillard (2015) provided tight concentration bounds for sampling without replacement for scalar random variables, while bounds for random matrices remain unclear. We fill this gap, and provide a tight bound for random matrices under sampling without replacement in this paper.

2 Stochastic Variance Reduction Scheme for Cubic Regularization

In this paper, we are interested in solving the finite-sum problem given in eq. (1), which is rewritten below.

$$\min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) \triangleq \frac{1}{N} \sum_{i=1}^{N} f_i(\mathbf{x}), \tag{7}$$

where the component functions f_i , i = 1, ..., N correspond to the loss of the *i*-th data samples, respectively, and is nonconvex. More specifically, we adopt the following standard assumptions on the objective function in eq. (7) throughout the paper

Assumption 1. The objective function in eq. (7) satisfies

1. Function F is bounded below, i.e., $\inf_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}) > -\infty$;

2. For all component functions f_i , i = 1, ..., N, the function value f_i , the gradient ∇f_i , and the Hessian $\nabla^2 f_i$ are L_0 , L_1 and L_2 -Lipschitz, respectively.

Classical first-order stochastic optimization methods such as stochastic gradient descent has a low sample complexity per-iteration (Nemirovski et al., 2009). However, due to the variance of the stochastic gradients, the convergence rate is slow even with the incorporation of momentum (Lan, 2012; Ghadimi and Lan, 2016). A popular approach to maintain the sample complexity yet achieve a faster convergence rate that is comparable to that of the full batch first-order methods is the stochastic variance reduction scheme (Johnson and Zhang, 2013; Xiao and Zhang, 2014).

Motivated by the success of the variance reduction scheme in improving the sample complexity of first-order methods, we propose a stochastic variance reduced cubic-regularized Newton's method, and refer to it as SVRC. The detailed steps of SVRC are presented in Algorithm 1. To briefly elaborate the notation in Algorithm 1, we sequentially index the iterate variable \mathbf{x} across all inner loops by k for $k = 0, 1, \ldots$, so that for each \mathbf{x}_k , the initial variable of its inner loop is indexed as $\mathbf{x}_{\lfloor k/m \rfloor \cdot m}$ (where m is the number of iterations in each inner loop). For notational simplicity, we denote such an initial variable of each inner loop as $\tilde{\mathbf{x}}$ and denote its corresponding full gradient and Hessian as $\tilde{\mathbf{g}}$ and $\tilde{\mathbf{H}}$, whenever there is no confusion.

Algorithm 1 SVRC

```
Input: \mathbf{x}_0 \in \mathbb{R}^d, and \epsilon_1, m, M \in \mathbb{R}^+.
while k do
     if k \mod m = 0 then
           Set \mathbf{g}_k = \nabla F(\mathbf{x}_k), \mathbf{H}_k = \nabla^2 F(\mathbf{x}_k), \widetilde{\mathbf{g}} = \mathbf{g}_k, \widetilde{\mathbf{x}} =
          \mathbf{x}_k and \mathbf{H} = \mathbf{H}_k.
           Sample index sets \xi_q(k) and \xi_H(k) from
           \{1, ..., n\} uniformly at random.
           Compute
            \mathbf{g}_{k} = \frac{1}{|\xi_{g}(k)|} \left[ \sum_{i \in \xi_{g}(k)} \left( \nabla f_{i}(\mathbf{x}_{k}) - \nabla f_{i}(\tilde{\mathbf{x}}) \right) \right] + \tilde{\mathbf{g}},
          \mathbf{H}_k \!=\! \textstyle\frac{1}{|\xi_H(k)|} \! \left[ \sum_{i \in \xi_H(k)} \! \left( \nabla^2 f_i(\mathbf{x}_k) \!-\! \nabla^2 f_i(\tilde{\mathbf{x}}) \right) \! \right] \! + \! \widetilde{\mathbf{H}}.
     end if
     \mathbf{s}_{k+1} = \operatorname{argmin}_{\mathbf{s} \in \mathbb{R}^d} \mathbf{g}_k^{\top} \mathbf{s} + \frac{1}{2} \mathbf{s}^{\top} \mathbf{H}_k \mathbf{s} + \frac{M}{6} \|\mathbf{s}\|^3.
     \mathbf{x}_{k+1} = \mathbf{x}_k + \mathbf{s}_{k+1}.
     if \max\{\|s_{k+1}\|, \|s_k\|\} \le \epsilon_1 then
           return x_{k+1}
     end if
end while
```

To elaborate the algorithm, SVRC calculates a full gradient $\tilde{\mathbf{g}}$ and a full Hessian $\tilde{\mathbf{H}}$ in every outer loop

(i.e., for every m iterations), which are further used to construct the stochastic variance reduced gradients \mathbf{g}_k and Hessians \mathbf{H}_k in the inner loops. Note that the index sets $\xi_g(k), \xi_H(k)$ for the sampled gradients and Hessians are generated by a random sampling scheme. More specifically, we consider the following two types of sampling schemes in this paper.

Sampling with replacement: For k = 0, 1, ..., each element of the index sets $\xi_g(k)$ and $\xi_H(k)$ is sampled uniformly at random from $\{1, ..., N\}$.

Sampling without replacement: For k = 0, 1, ..., the index sets $\xi_g(k)$ and $\xi_H(k)$ are sampled uniformly at random from all subsets of $\{1, ..., N\}$ with cardinality $|\xi_g(k)|$ and $|\xi_H(k)|$, respectively.

To elaborate, the sampling with replacement scheme may sample the same index multiple times within each mini-batch, whereas the sampling without replacement scheme samples each index at most once within each mini-batch. Therefore, the sampling without replacement scheme has a smaller variance compared to that of the sampling with replacement scheme. Consequently, these sampling schemes lead to inexact gradients and inexact Hessians with different guarantees to meet the inexactness criterion.

3 Sample Complexity of SVRC

In this section, we study the sample complexity of SVRC for achieving a second-order stationary point via three technical steps, each corresponding to one subsection below.

3.1 Iteration Complexity under Modified Inexact Condition

In order to analyze the sample complexity of SVRC for achieving a second-order stationary point, it turns out that the inexact condition (Wang et al., 2019) on the estimated gradients and Hessians is not sufficient. Thus, we propose a modified inexact condition below, and then analyze the convergence to a second-order stationary point if SVRC satisfies such a condition.

Assumption 2. The approximate Hessian \mathbf{H}_k and approximate gradient \mathbf{g}_k satisfy, for all $k = 0, \dots$,

$$\|\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)\| \le \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$$
 (8)

$$\|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\| \le \beta \max \left\{ \|\mathbf{s}_k\|^2, \epsilon_1^2 \right\}$$
 (9)

where ϵ_1 , α and β are universal positive constants.

The inexact conditions in eqs. (8) and (9) introduce a slack variable ϵ_1 to avoid full batch sampling when $||s_k||$ is very close to zero upon convergence. It turns

out introduction of such a variable is essential for characterizing the total sample complexity of our proposed variance reduction scheme in Algorithm 1. Furthermore, since eqs. (8) and (9) are different from that in (Wang et al., 2019), and hence require the convergence analysis if SVRC satisfies such conditions. The following theorem presents the iteration complexity analysis under the modified conditions. The technical proof in fact requires considerable extra effort than that in Wang et al. (2019).

Theorem 1. Suppose Assumption 1 holds, and SVRC satisfies Assumption 2. Let

$$\tau \triangleq \min \left\{ \left(\frac{L+M}{2} + 2\beta + 2\alpha \right)^{-\frac{1}{2}}, \left(\frac{M+2L}{2} + 2\alpha \right)^{-1} \right\},$$

set

$$\epsilon_1 = \tau \sqrt{\epsilon},$$
(10)

and properly choose M, α and $\beta \in \mathbb{R}$ such that

$$\gamma \triangleq \left(\frac{3M - 2L_2}{24} - \frac{5}{2}\beta - \frac{5}{4}\alpha\right) > 0. \tag{11}$$

Then, the SVRC algorithm outputs an ϵ -approximate second-order stationary point, i.e.,

$$\|\nabla f(\mathbf{x}_{k+1})\| \leqslant \epsilon \quad and \quad \nabla^2 f(\mathbf{x}_{k+1}) \succcurlyeq -\sqrt{\epsilon} \mathbf{I} \quad (12)$$

within at most $k = O(\epsilon^{-3/2})$ number of iterations. Moreover, the following inequality holds

$$\sum_{i=1}^{k+1} \|\mathbf{s}_i\|^3 \leqslant C, \tag{13}$$

where
$$C \triangleq (f(\mathbf{x}_0) - f^* + (2\beta + \alpha + 2\gamma)\epsilon_1^3)/\gamma$$
.

As stated in Theorem 1, SVRC outputs an ϵ -approximate second-order stationary point with $k = O(\epsilon^{-3/2})$. Such an iteration complexity matches the state-of-art result and is the best result that one can expect on nonconvex optimization.

3.2 Per-iteration Sample Complexity

In this subsection, we bound the per-iteration sample complexity in order for SVRC (under sampling with replacement) to satisfy the inexact conditions in eqs. (8) and (9). We apply Bernstein's inequality and obtain the following theorem.

Theorem 2. Let Assumption 1 hold. Consider SVRC under the sampling with replacement scheme. Then, the sub-sampled mini-batch of gradients \mathbf{g}_k , k = 0, 1, ...

satisfies Assumption 2 with probability at least $1-\zeta$ provided that

$$|\xi_{g}(k)| \geqslant \left(\frac{8L_{1}^{2}}{\beta^{2} \max\{\|\mathbf{s}_{k}\|^{4}, \epsilon_{1}^{4}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|^{2} + \frac{4L_{1}}{3\beta \max\{\|\mathbf{s}_{k}\|^{2}, \epsilon_{1}^{2}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|\right) \log\left(\frac{2(d+1)}{\zeta}\right), \tag{14}$$

Furthermore, the sub-sampled mini-batch of Hessians $\mathbf{H}_k, k = 0, 1, \ldots$ of SVRC satisfies Assumption 2 with probability at least $1 - \zeta$ provided that

$$|\xi_{H}(k)| \geqslant \left(\frac{8L_{2}^{2}}{\alpha^{2} \max\{\|\mathbf{s}_{k}\|^{2}, \epsilon_{1}^{2}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|^{2} + \frac{4L_{2}}{3\alpha \max\{\|\mathbf{s}_{k}\|, \epsilon_{1}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|\right) \log\left(\frac{4d}{\zeta}\right).$$

$$(15)$$

We next compare the per-iteration Hessian sample complexity of SVRC under the sampling with replacement scheme (eq. (15)) with that of SCR under the same sampling scheme developed in Kohler and Lucchi (2017), which is rewritten below

$$|\xi_H(k)| \geqslant \mathcal{O}\left(\frac{1}{\|\mathbf{s}_{k+1}\|^2}\right). \tag{16}$$

To compare, our Theorem 2 requires a Hessian sample complexity of roughly the order

$$|\xi_H(k)| \geqslant \mathcal{O}\left(\frac{\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2}{\|\mathbf{s}_k\|^2}\right).$$
 (17)

It can be seen that the sample complexity bounds for SVRC in eq. (17) have an additional term $\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2$ in the numerators comparing to their corresponding bound for SCR in eq. (16). Intuitively, $\|\mathbf{x}_k - \tilde{\mathbf{x}}\| \to 0$ as the algorithm converges, and thus our variance reduction scheme requires a lower sample complexity than the stochastic sampling in SCR.

3.3 Total Sample Complexity of SVRC

Theorem 2 provides the sample complexity per iteration (each iteration in SVRC inner loop). We next provide our result on the sample complexity over the running process of SVRC, which is a key factor that impacts the computational complexity of SVRC.

Theorem 3. Let Assumptions 1 hold. For a given ϵ and δ , set $m = N^{1/3}$, then SVRC under the sampling with replacement scheme outputs an point \mathbf{x}_{k+1} such that satisfies $\|\nabla F(\mathbf{x}_{k+1})\| \leq \epsilon$ and $\nabla^2 F(\mathbf{x}_{k+1}) \geq -\epsilon \mathbf{I}$ with probability at least $1 - \delta$, and the total Hessian sample complexity of SVRC is bounded by

$$\sum_{i=1}^{K} |\xi_H(i)| \leqslant \frac{CN^{2/3}}{\epsilon^{3/2}} \log \left(\frac{8d}{\epsilon \delta} \right).$$

We next compare the total Hessian sample complexity of SVRC with that of other CR-type algorithms, which are given below.

SVRC:
$$\sum_{i=1}^{K} |\xi_H(i)| = \tilde{\mathcal{O}}\left(\frac{N^{2/3}}{\epsilon^{3/2}}\right), \quad (18)$$

SVRC (ZXG):
$$\sum_{i=1}^{K} |\xi_H(i)| = \mathcal{O}\left(\frac{N^{4/5}}{\epsilon^{3/2}}\right), \quad (19)$$

CR:
$$\sum_{i=1}^{K} |\xi_H(i)| \le \mathcal{O}\left(\frac{N}{\epsilon^{3/2}}\right), \qquad (20)$$

SCR:
$$\sum_{i=1}^{K} |\xi_H(i)| \leq \mathcal{O}\left(\frac{1}{\epsilon^{5/2}}\right). \tag{21}$$

Comparing eqs. (18) to (20). Clearly, our SVRC has lower total sample complexity than CR and SVRC(ZXG) by an order of $\mathcal{O}(N^{1/3})$ and $\mathcal{O}(N^{2/15})$, respectively. Therefore, our stochastic variance reduction scheme is sample efficient when applied to CR type of methods. Also, comparing the sample complexity of the two subsampled algorithms in eqs. (18) and (21), we observe that SVRC enjoys a lower-order complexity bound than SCR if $\epsilon = o(N^{-2/3})$, and hence performs better in the high accuracy regime.

4 SVRC under Sampling without Replacement Scheme

In this section, we explore the sample complexity of SVRC under the sampling without replacement scheme, which is commonly used in practice.

To this end, we first develop some technical concentration inequalities in the next subsection.

4.1 Concentration Inequality under Sampling without Replacement

The statistics of sampling without replacement is very different and more stable than that of sampling with replacement. However, theoretical analysis of sampling without replacement turns out to be very difficult. A common approach is to apply the concentration bound for sampling with replacement, which also holds for sampling without replacement (Tropp, 2012). However, such analysis can be too loose to capture the essence of the scheme of sampling without replacement. For example, the sample complexity for sampling with replacement to achieve a certain accuracy can be unbounded, whereas sampling without replacement can at most sample the total sample size.

Thus, in order to develop a tight sample complexity bound for SVRC under sampling without replacement, we first leverage a recently developed Hoeffding-type of concentration inequality for sampling without replacement (Bardenet and Maillard, 2015). There, the result is applicable only for scalar random variables, whereas our analysis here needs to deal with sub-sampled gradients and Hessians, which are vectors and matrices. This motivates us to first establish the matrix version of the Hoeffding-Serfling inequality. Such a concentration bound can be of independent interest in various other domains. The proof turns out to be very involved and is provided in the supplementary materials.

Theorem 4. Let $\mathcal{X} := \{\mathbf{A}_1, \cdots, \mathbf{A}_N\}$ be a collection of real-valued matrices in $\mathbb{R}^{d_1 \times d_2}$ with bounded spectral norm, i.e., $\|\mathbf{A}_i\| \le \sigma$ for all $i = 1, \dots, N$ and some $\sigma > 0$. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be n < N samples from \mathcal{X} under the sampling without replacement. Denote $\mu := \frac{1}{N} \sum_{i=1}^{N} \mathbf{A}_i$. Then, for any $\epsilon > 0$, the following bound holds

$$P\left(\left\|\frac{1}{n}\sum_{i=1}^{n}\mathbf{X}_{i}-\mu\right\| \geqslant \epsilon\right)$$

$$\leqslant 2(d_{1}+d_{2})\exp\left(-\frac{n\epsilon^{2}}{8\sigma^{2}(1+1/n)(1-n/N)}\right).$$

To further understand the above theorem, consider symmetric random matrix $\mathbf{X}_i \in \mathbb{R}^{d \times d}$. Suppose we want $\left\|\frac{1}{n}\sum_{i=1}^n \mathbf{X}_i - \mu\right\| \leq \epsilon$ to hold with probability $1 - \zeta$. Then the above theorem requires the sample size to satisfy

$$n_w \geqslant \left(\frac{1}{N} + \frac{\epsilon^2}{16\sigma^2 \log(4d)/\zeta}\right)^{-1}.$$
 (22)

We consider two regimes to understand the bound in eq. (22). (a) Low accuracy regime: Suppose ϵ is large enough so that the second term in eq. (22) dominates. In this case, we roughly have $n_w \geqslant \frac{16\sigma^2 \log(4d/\zeta)}{\epsilon^2}$, which has the same order as the suggested sample size by the matrix version of the Hoeffding inequality for sampling with replacement given below

$$n_b \geqslant \frac{8\sigma^2 \log(2d/\zeta)}{\epsilon^2}.$$
 (23)

Thus, the sample size is approximately the same for sampling with and without replacement to achieve a low accuracy concentration. (b) High accurary regime: Suppose ϵ is small enough so that the first term in eq. (22) dominates. Hence, eq. (22) roughly reduces to $n_w \geq N$, whereas the matrix version of the Hoeffding bound in eq. (23) for sampling with replacement requires infinite samples as $\epsilon \to 0$. Thus, the sample size is highly different for sampling with and without replacement to achieve a high accuracy concentration.

4.2 Per-iteration Sample Complexity

We apply Theorem 4 to analyze the sample complexity of SVRC under sampling without replacement. Our next theorem characterizes the sample size needed for SVRC in order to satisfy the inexact condition in Assumption 2.

Theorem 5. Let Assumption 1 hold. Consider SVRC under sampling without replacement. The sub-sampled mini-batches of gradients \mathbf{g}_k , $k = 0, 1, \ldots$ satisfy eq. (6) with probability at least $1 - \zeta$ provided that

$$|\xi_g(k)| \ge \left(\frac{1}{N} + \frac{\beta^2 \max\{\|\mathbf{s}_k\|^4, \epsilon_1^4\}}{64L_1^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2 \log(2(d+1)/\zeta)}\right)^{-1},$$
(24)

Furthermore, the sub-sampled mini-batches of Hessians $\mathbf{H}_k, k = 0, 1, \ldots$ satisfy eq. (4) with probability at least $1 - \zeta$ provided that

$$|\xi_H(k)| \ge \left(\frac{1}{N} + \frac{\alpha^2 \max\{\|\mathbf{s}_k\|^2, \epsilon_1^2\}}{64L_2^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2 \log(4d/\zeta)}\right)^{-1}.$$
(25)

In order to further understand the sample complexity in Theorem 5 and what improvement that SVRC makes in terms of sample complexity compared to the SCR algorithm in Kohler and Lucchi (2017), we next characterize the corresponding sample complexity for SCR under sampling without replacement below. (We note that the sample complexity for SCR under sampling with replacement was provided in Kohler and Lucchi (2017).)

Proposition 6. Let Assumptions 1 hold. Consider the SCR algorithm in Kohler and Lucchi (2017) under sampling without replacement. The sub-sampled minibatch of gradients \mathbf{g}_k , $k = 0, 1, \ldots$ satisfies eq. (6) with probability at least $1 - \zeta$ provided that for all k

$$|\xi_g(k)| \geqslant \left(\frac{1}{N} + \frac{C_1^2 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^4}{64L_0^2 \log(2(d+1)/\zeta)}\right)^{-1}.$$
 (26)

Furthermore, the sub-sampled mini-batch of Hessians $\mathbf{H}_k, k = 0, 1, \dots$ satisfies eq. (4) with probability at least $1 - \zeta$ provided that for all k

$$|\xi_H(k)| \geqslant \left(\frac{1}{N} + \frac{C_2^2 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2}{64L_1^2 \log(4d/\zeta)}\right)^{-1}.$$
 (27)

To compare the sample complexity for SVRC in Theorem 5 and SCR in Proposition 6, we take the sample complexity for mini-batch of gradients as an example. Comparing eq. (24) and eq. (26), the second term in the denominator in eq. (24) is additionally divided by $\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2$, which converges to zero as the algorithms

converge. Thus, $\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$ in eq. (26) converges to zero much faster than $\frac{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|^4}{\|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2}$ in eq. (24), so that the term 1/N dominates the denominator and results in the sample size close to the number of total samples much earlier in the iteration of SCR than SVRC.

We also note that Proposition 6 shows that as SCR approaches the convergence, the sample size goes to the total number of samples with technical rigor, whereas such a fact was only intuitively discussed in Kohler and Lucchi (2017).

4.3 Total Sample Complexity

We next characterize the total Hessian sample complexity of SVRC under sampling without replacement.

Theorem 7. Let Assumptions 1 hold. For a given ϵ and δ , set $m = N^{1/3}$, then SVRC under sampling without replacement outputs an point \mathbf{x}_{k+1} such that satisfies $\|\nabla F(\mathbf{x}_{k+1})\| \le \epsilon$ and $\nabla^2 F(\mathbf{x}_{k+1}) \ge -\epsilon \mathbf{I}$ with probability at least $1 - \delta$. Then the total sample complexity for Hessian used in SVRC is bounded by

$$\sum_{i=0}^{k} |\xi_H(k)| \leqslant \frac{CN^{2/3}}{\epsilon^{3/2}} \log \left(\frac{8d}{\epsilon \delta}\right). \tag{28}$$

In this theorem, we show that total sample complexity of SVRC under sampling without replacement is at least as good as SVRC under sampling with replacement. And the comparison of this bound with other bound follows similarly as we discuss in Section 3.3.

5 Discussion

Storage Issue: The proposed algorithm involves the storage of a Hessian, which requires $\mathcal{O}(d^2)$ space for storage. In this perspective, the proposed algorithm can be directly applied for solving small or medium scale machine learning problems. As for large scale problems, using PCA to store the main component of Hessian can be a possible solution.

With and Without replacement: We show that the total sample complexity of SVRC under sampling without replacement is at least as good as SVRC under sampling with replacement. Actually, if we compare the per iteration complexity of the two, i.e., we compare Theorem 5 with Theorem 2, the without replacement scheme has a better complexity than that with replacement in each iteration since there is a 1/N term in the denominator on the bound for the scheme without replacement. This does suggest the same total sample complexity for the two schemes is likely due to the technicality issue.

6 Conclusion

In this paper, we proposed a stochastic variance-reduced cubic regularization method. We characterized the per iteration sample complexity for Hessian and gradient that guarantees convergence of SVRC to a second-order optimality condition, under both sampling with and without replacement. We also developed the total sample size for Hessian. Our theoretic results imply that SVRC outperforms the state-of-art result by an factor of $O(N^{2/15})$. Moreover, Our study demonstrates that variance reduction can bring substantial advantage in sample size as well as computational complexity for second-order algorithms, along which direction we plan to explore further in the future.

Acknowledgement

The work of Z. Wang, Y. Zhou and Y. Liang was supported in part by the U.S. National Science Foundation under the grant CCF-1761506, and the work of G. Lan was supported in part by Army Research Office under the grant W911NF-18-1-0223 and U.S. National Science Foundation under the grant CMMI-1254446.

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.

Baldi, P. and Hornik, K. (1989). Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2(1):53 – 58.

Bardenet, R. and Maillard, O. (2015). Concentration inequalities for sampling without replacement. *Bernoulli*, 21(3):1361–1385.

Bhatia, R. (2007). *Positive Definite Matrices*. Princeton University Press.

Bhojanapalli, S., Neyshabur, B., and Srebro, N. (2016). Global optimality of local search for low rank matrix recovery. In *Proc. Advances in Neural Information Processing Systems (NIPS)*.

Carmon, Y. and Duchi, J. C. (2016). Gradient descent efficiently finds the cubic-regularized non-convex Newton step. ArXiv: 1612.00547.

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. (2011b). Adaptive cubic regularization methods for unconstrained optimization. part ii: worst-case function-

- and derivative-evaluation complexity. *Mathematical Programming*, 130(2):295–319.
- Cartis, C., Gould, N. I. M., and Toint, P. L. (2011). Adaptive cubic regularization methods for unconstrained optimization. part i: motivation, convergence and numerical results. *Mathematical Programming*.
- Cartis, C., Gould, N. I. M., and Toint, P. L. (2012a). An adaptive cubic regularization algorithm for non-convex optimization with convex constraints and its function-evaluation complexity. *IMA Journal of Numerical Analysis*, 32(4):1662 – 1695.
- Cartis, C., Gould, N. I. M., and Toint, P. L. (2012b). Complexity bounds for second-order optimality in unconstrained optimization. *Journal of Complexity*, 28(1):93 108.
- Chi, J., Rong, G., Praneeth, N., K., S. M., and Michael,
 I. J. (2017). How to escape saddle points efficiently.
 In Proc. 34th International Conference on Machine Learning (ICML), volume 70, pages 1724–1732.
- Fang, C., Li, C., Lin, Z., and Zhang, T. (2018). SPI-DER: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In Proc. Advances in Neural Information Processing Systems (NIPS), pages 687–697.
- Friedlander, M. and Schmidt, M. (2012). Hybrid deterministic-stochastic methods for data fitting. SIAM Journal on Scientific Computing, 34(3):A1380–A1405.
- Ge, R., Lee, J., and Ma, T. (2016). Matrix completion has no spurious local minimum. In *Proc. Advances* in Neural Information Processing Systems (NIPS), pages 2973–2981.
- Ghadimi, S. and Lan, G. (2016). Accelerated gradient methods for nonconvex nonlinear and stochastic programming. *Mathematical Programming*, 156:59–99.
- Ghadimi, S., Liu, H., and Zhang, T. (2017). Second-order methods with cubic regularization under inexact information. *ArXiv:* 1710.05782.
- Gürbüzbalaban, M., Ozdaglar, A., and Parrilo, P. (2015). Why Random Reshuffling Beats Stochastic Gradient Descent. *ArXiv:1510.08560*.
- Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *Journal of the Ameri*can Statistical Association, 58(301):13–30.
- 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.
- Johnson, R. and Zhang, T. (2013). Accelerating stochastic gradient descent using predictive variance reduc-

- tion. In Proc. 26th Advances in Neural Information Processing Systems (NIPS), pages 315–323.
- 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.
- Lan, G. (2012). An optimal method for stochastic composite optimization. *Mathematical Programming*, 133(1):365–397.
- Lee, J. D., Simchowitz, M., Jordan, M. I., and Recht, B. (2016). Gradient descent only converges to minimizers. In *Proc. 29th Annual Conference on Learning Theory (COLT)*, volume 49, pages 1246–1257.
- 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.
- Nemirovski, A. S., Juditsky, A., Lan, G., and Shapiro, A. (2009). Robust stochastic approximation approach to stochastic programming. SIAM Journal on Optimization, 19:1574–1609.
- 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. (2006). Cubic regularization of newton's method and its global performance. *Mathematical Programming*.
- Recht, B. and Re, C. (2012). Toward a noncommutative arithmetic-geometric mean inequality: conjectures, case-studies, and consequences. *Conference on Learning Theory*.
- Reddi, S., Hefny, A., Sra, S., B., P., and A., S. (2016). Stochastic variance reduction for nonconvex optimization. In *Proc. 33rd International Conference on Machine Learning (ICML)*, volume 48, pages 314–323.
- Rong, G.and Furong, H., Chi, J., and Yang, Y. (2015).
 Escaping from saddle points online stochastic gradient for tensor decomposition. In *Proc. 28th Conference on Learning Theory (COLT)*, volume 40, pages 797–842.
- Shamir, O. (2016). Without-replacement sampling for stochastic gradient methods. In *Proc. 29th Advances in Neural Information Processing Systems (NIPS)*, pages 46–54.
- Sun, J., Qu, Q., and Wright, J. (2017). A geometrical analysis of phase retrieval. Foundations of Computational Mathematics, pages 1–68.

- 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., Ji, K., Zhou, Y., Liang, Y., and Tarokh, V. (2018). SpiderBoost: A class of faster variancereduced algorithms for nonconvex optimization. ArXiv:1810.10690.
- Wang, Z., Zhou, Y., Liang, Y., and Lan, G. (2018). Cubic Regularization with Momentum for Nonconvex Optimization. arXiv:1810.03763.
- Wang, Z., Zhou, Y., Liang, Y., and Lan, G. (2019).
 A note on inexact gradient and Hessian conditions for cubic regularized Newtons method. *Operations Research Letters*.
- Xiao, L. and Zhang, T. (2014). A proximal stochastic gradient method with progressive variance reduction. SIAM Journal on Optimization, 24(4):2057–2075.
- 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.
- Zhou, D., Xu, P., and Gu, Q. (2018). Stochastic Variance-Reduced Cubic Regularized Newton Method. In *Proc. 35th International Conference on Machine Learning (ICML)*.
- Zhou, Y. and Liang, Y. (2018). Critical points of linear neural networks: Analytical forms and landscape properties. In Proc. International Conference on Learning Representations (ICLR).
- Zhou, Y., Wang, Z., and Liang, Y. (2018). Convergence of cubic regularization for nonconvex optimization under KL property. In *Proc. 32nd Advances in Neural Information Processing Systems (NIPS)*.

Supplementary Materials

A Proof of Convergence

A.1 Lemmas

In this subsection, we introduce two useful lemmas, which will be used in the proof of convergence.

Lemma 8 (Nesterov and Polyak (2006), Lemma 1). Let the Hessian $\nabla^2 f(\cdot)$ of the function $f(\cdot)$ be L-Lipschitz continuous with L > 0. Then, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) - \nabla^2 f(\mathbf{x})(\mathbf{y} - \mathbf{x})\| \leqslant \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2,$$
 (29)

$$\left| f(\mathbf{y}) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) - \frac{1}{2} (\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x}) (\mathbf{y} - \mathbf{x}) \right| \leqslant \frac{L}{6} \|\mathbf{y} - \mathbf{x}\|^3.$$
 (30)

Lemma 9 (Wang et al. (2019), Lemma 3). Let $M \in \mathbb{R}$, $\mathbf{g} \in \mathbb{R}^d$, $\mathbf{H} \in \mathbb{S}^{d \times d}$, and

$$\mathbf{s} = \underset{\mathbf{u} \in \mathbb{R}^d}{\operatorname{argmin}} \mathbf{g}^{\top} \mathbf{u} + \frac{1}{2} \mathbf{u}^{\top} \mathbf{H} \mathbf{u} + \frac{M}{6} \| \mathbf{u} \|^{3}.$$
 (31)

Then, the following statements hold:

$$\mathbf{g} + \mathbf{H}\mathbf{s} + \frac{M}{2} \|\mathbf{s}\| \,\mathbf{s} = \mathbf{0},\tag{32}$$

$$\mathbf{H} + \frac{M}{2} \|\mathbf{s}\| \mathbf{I} \geq \mathbf{0},\tag{33}$$

$$\mathbf{g}^{\top}\mathbf{s} + \frac{1}{2}\mathbf{s}^{\top}\mathbf{H}\mathbf{s} + \frac{M}{6} \|\mathbf{s}\|^{3} \leqslant -\frac{M}{12} \|\mathbf{s}\|^{3}.$$
 (34)

A.2 Proof of Theorem 1

Proof. Since $\nabla^2 f(\mathbf{x})$ is L_2 -Lipschitz, thus we have

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_{k}) \stackrel{(i)}{\leqslant} \nabla f(\mathbf{x}_{k})^{\top} \mathbf{s}_{k+1} + \frac{1}{2} \mathbf{s}_{k+1}^{\top} \nabla f(\mathbf{x}_{k}) \mathbf{s}_{k+1} + \frac{L_{2}}{6} \| \mathbf{s}_{k+1} \|^{3}$$

$$\leqslant \mathbf{g}_{k}^{\top} \mathbf{s}_{k+1} + \frac{1}{2} \mathbf{s}_{k+1}^{\top} \mathbf{H}_{k} \mathbf{s}_{k+1} + \frac{M}{6} \| \mathbf{s}_{k+1} \|^{3} + (\nabla f(\mathbf{x}_{k}) - \mathbf{g}_{k})^{\top} \mathbf{s}_{k+1}$$

$$+ \frac{L_{2} - M}{6} \| \mathbf{s}_{k+1} \|^{3} + \frac{1}{2} \mathbf{s}_{k+1}^{\top} (\nabla^{2} f(\mathbf{x}_{k}) - \mathbf{H}_{k}) \mathbf{s}_{k+1}$$

$$\stackrel{(ii)}{\leqslant} -\frac{3M - 2L_{2}}{12} \| \mathbf{s}_{k+1} \|^{3} + (\nabla f(\mathbf{x}_{k}) - \mathbf{g}_{k})^{\top} \mathbf{s}_{k+1} + \frac{1}{2} \mathbf{s}_{k+1}^{\top} (\nabla f(\mathbf{x}_{k}) - \mathbf{H}_{k}) \mathbf{s}_{k+1}$$

$$(35)$$

where (i) follows from Lemma 8 with $\mathbf{y} = \mathbf{x}_{k+1}, \mathbf{x} = \mathbf{x}_k$ and $\mathbf{s}_{k+1} = \mathbf{x}_{k+1} - \mathbf{x}_k$, (ii) follows from eq. (34) in Lemma 9 with $\mathbf{g} = \mathbf{g}_k, \mathbf{H} = \mathbf{H}_k$ and $\mathbf{s} = \mathbf{s}_{k+1}$.

Next, we bound the terms $(\nabla f(\mathbf{x}_k) - \mathbf{g}_k)^{\top} \mathbf{s}_{k+1}$ and $\mathbf{s}_{k+1}^{\top} (\nabla f(\mathbf{x}_k) - \mathbf{H}_k) \mathbf{s}_{k+1}$. For the first term, we have that

$$(\nabla f(\mathbf{x}_{k}) - \mathbf{g}_{k})^{\top} \mathbf{s}_{k+1} \leq \|\nabla f(\mathbf{x}_{k}) - \mathbf{g}_{k}\| \|\mathbf{s}_{k+1}\| \stackrel{\text{(i)}}{\leq} \beta \left(\|\mathbf{s}_{k}\|^{2} + \epsilon_{1}^{2} \right) \|\mathbf{s}_{k+1}\| = \beta \left(\|\mathbf{s}_{k}\|^{2} \|\mathbf{s}_{k+1}\| + \epsilon_{1}^{2} \|\mathbf{s}_{k+1}\| \right)$$

$$\stackrel{\text{(ii)}}{\leq} \beta \left(\|\mathbf{s}_{k}\|^{3} + \|\mathbf{s}_{k+1}\|^{3} + \epsilon_{1}^{3} + \|\mathbf{s}_{k+1}\|^{3} \right) = \beta \left(\|\mathbf{s}_{k}\|^{3} + 2\|\mathbf{s}_{k+1}\|^{3} + \epsilon_{1}^{3} \right),$$
(36)

where (i) follows from Assumption 2, which gives that $\|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\| \le \beta \max \{\|\mathbf{s}_k\|^2, \epsilon_1^2\}$, and (ii) follows from the inequality that for $a, b \in \mathbb{R}^+$, $a^2b \le a^3 + b^3$, which can be verified by checking the cases with a < b and $a \ge b$, respectively. Similarly, we obtain that

$$\mathbf{s}_{k+1}^{\top}(\nabla f(\mathbf{x}_k) - \mathbf{H}_k)\mathbf{s}_{k+1} \leqslant \left\|\nabla^2 f(\mathbf{x}_k) - \mathbf{H}_k\right\| \left\|\mathbf{s}_{k+1}\right\|^2 \stackrel{\text{(i)}}{\leqslant} \alpha \left(\left\|\mathbf{s}_k\right\| + \epsilon_1\right) \left\|\mathbf{s}_{k+1}\right\|^2 = \alpha \left(\left\|\mathbf{s}_k\right\| \left\|\mathbf{s}_{k+1}\right\|^2 + \epsilon_1 \left\|\mathbf{s}_{k+1}\right\|^2\right)$$

$$\stackrel{\text{(ii)}}{\leq} \alpha \left(\|\mathbf{s}_{k}\|^{3} + \|\mathbf{s}_{k+1}\|^{3} + \epsilon_{1}^{3} + \|\mathbf{s}_{k+1}\|^{3} \right) = \alpha \left(\|\mathbf{s}_{k}\|^{3} + 2\|\mathbf{s}_{k+1}\|^{3} + \epsilon_{1}^{3} \right), \tag{37}$$

where (i) follows from Assumption 2, which gives that $\|\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)\| \le \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$, and (ii) follows from the inequality that for $a, b \in \mathbb{R}^+$, $a^2b \le a^3 + b^3$.

Plugging eqs. (36) and (37) into eq. (35) yields

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k) \leqslant -\frac{3M - 2L_2}{12} \|\mathbf{s}_{k+1}\|^3 + \beta \left(\|\mathbf{s}_k\|^3 + 2\|\mathbf{s}_{k+1}\|^3 + \epsilon_1^3\right) + \frac{\alpha}{2} \left(\|\mathbf{s}_k\|^3 + 2\|\mathbf{s}_{k+1}\|^3 + \epsilon_1^3\right)$$

$$= -\left(\frac{3M - 2L_2}{12} - 2\beta - \alpha\right) \|\mathbf{s}_{k+1}\|^3 + \left(\beta + \frac{\alpha}{2}\right) \|\mathbf{s}_k\|^3 + \left(\beta + \frac{\alpha}{2}\right) \epsilon_1^3$$
(38)

Summing Equation (38) for 0 to k, we obtain

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_{0}) \leqslant -\left(\frac{3M - 2L_{2}}{12} - 2\beta - \alpha\right) \sum_{i=1}^{k+1} \|\mathbf{s}_{i}\|^{3} + \left(\beta + \frac{\alpha}{2}\right) \sum_{i=0}^{k} \|\mathbf{s}_{i}\|^{3} + \left(\beta + \frac{\alpha}{2}\right) \sum_{i=0}^{k} \epsilon_{1}^{3}$$

$$\leqslant -\left(\frac{3M - 2L_{2}}{12} - 2\beta - \alpha\right) \sum_{i=1}^{k+1} \|\mathbf{s}_{i}\|^{3} + \left(\beta + \frac{\alpha}{2}\right) \sum_{i=0}^{k+1} \|\mathbf{s}_{i}\|^{3} + \left(\beta + \frac{\alpha}{2}\right) \sum_{i=0}^{k} \epsilon_{1}^{3}$$

$$\leqslant -\left(\frac{3M - 2L_{2}}{12} - 3\beta - \frac{3}{2}\alpha\right) \sum_{i=1}^{k+1} \|\mathbf{s}_{i}\|^{3} + \left(\beta + \frac{\alpha}{2}\right) \|\mathbf{s}_{0}\|^{3} + \left(\beta + \frac{\alpha}{2}\right) \sum_{i=0}^{k} \epsilon_{1}^{3}, \tag{39}$$

We next note that

$$\sum_{i=1}^{k+1} \|\mathbf{s}_i\|^3 = \frac{1}{2} \left(\sum_{i=1}^{k+1} \|\mathbf{s}_i\|^3 + \sum_{i=1}^{k+1} \|\mathbf{s}_i\|^3 \right) = \frac{1}{2} \left(\sum_{i=1}^{k+1} \|\mathbf{s}_i\|^3 + \sum_{i=0}^{k} \|\mathbf{s}_{i+1}\|^3 \right) \geqslant \frac{1}{2} \sum_{i=1}^{k} \left(\|\mathbf{s}_i\|^3 + \|\mathbf{s}_{i+1}\|^3 \right). \tag{40}$$

Plugging eq. (40) into eq. (39) yields that

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}_0) \leqslant -\sum_{i=1}^k \left(\frac{3M - 2L_2}{24} - \frac{3}{2}\beta - \frac{3}{4}\alpha \right) \left(\|\mathbf{s}_i\|^3 + \|\mathbf{s}_{i+1}\|^3 \right) + \left(\beta + \frac{\alpha}{2}\right) \|\mathbf{s}_0\|^3 + \left(\beta + \frac{\alpha}{2}\right) \sum_{i=0}^k \epsilon_1^3$$

$$\stackrel{\text{(i)}}{\leqslant} -\sum_{i=1}^k \left(\frac{3M - 2L_2}{24} - \frac{5}{2}\beta - \frac{5}{4}\alpha \right) \left(\|\mathbf{s}_i\|^3 + \|\mathbf{s}_{i+1}\|^3 \right) + \left(\beta + \frac{\alpha}{2}\right) \|\mathbf{s}_0\|^3 + \left(\beta + \frac{\alpha}{2}\right) \epsilon_1^3,$$

where (i) follows from the fact that before the algorithm terminates we always have that $||s_i|| \ge \epsilon_1$ or $||s_{i+1}|| \ge \epsilon_1$, which gives that $||s_i||^3 + ||s_{i+1}||^3 \ge \epsilon_1^3$. Therefore, we have

$$\sum_{i=1}^{k} \left(\frac{3M - 2L_2}{24} - \frac{5}{2}\beta - \frac{5}{4}\alpha \right) \left(\|\mathbf{s}_i\|^3 + \|\mathbf{s}_{i+1}\|^3 \right) \leqslant f(\mathbf{x}_0) - f^* + \left(\beta + \frac{\alpha}{2} \right) \|\mathbf{s}_0\|^3 + \left(\beta + \frac{\alpha}{2} \right) \epsilon_1^3$$

$$\stackrel{(i)}{=} f(\mathbf{x}_0) - f^* + (2\beta + \alpha) \epsilon_1^3$$

$$(41)$$

where (i) follows from the fact that $\|\mathbf{s}_0\| = \epsilon_1$. Thus, if the algorithm never terminates, then we always have that $\|s_i\| \ge \epsilon_1$ or $\|s_{i+1}\| \ge \epsilon_1$, which gives $\|s_i\|^3 + \|s_{i+1}\|^3 \ge \epsilon_1^3$. Following from Equation (41), we obtain that

$$k \times \gamma \epsilon_1^3 \leqslant f(\mathbf{x}_0) - f^* + (2\beta + \alpha) \epsilon_1^3, \tag{42}$$

where $\gamma \triangleq \left(\frac{3M-2L_2}{24} - \frac{5}{2}\beta - \frac{5}{4}\alpha\right)$. Therefore, we obtain

$$k \leqslant \frac{f(\mathbf{x}_0) - f^* + (2\beta + \alpha)\epsilon_1^3}{\gamma \epsilon_1^3},\tag{43}$$

which shows that the algorithm must terminates if the total number of iterations exceeds $O(\epsilon_1^{-3})$. With the choice of ϵ_1 in Theorem 1, we obtain that the algorithm terminates at most with total iteration $k = O(\epsilon^{-3/2})$.

Suppose that the algorithm terminates at iteration k, then according to the analysis in eq. (41), we have that

$$\sum_{i=1}^{k-1} \gamma \left(\|\mathbf{s}_i\|^3 + \|\mathbf{s}_{i+1}\|^3 \right) \leqslant f(\mathbf{x}_0) - f^* + (2\beta + \alpha) \,\epsilon_1^3. \tag{44}$$

On the other hand, according to eq. (44) and the terminal condition that $||s_i|| \le \epsilon_1$ and $||s_{i+1}|| \le \epsilon_1$, we obtain

$$\sum_{i=1}^{k} \gamma \left(\|\mathbf{s}_{i}\|^{3} + \|\mathbf{s}_{i+1}\|^{3} \right) \leqslant f(\mathbf{x}_{0}) - f^{*} + (2\beta + \alpha + 2\gamma) \epsilon_{1}^{3},$$

which gives that

$$\sum_{i=1}^{k+1} \|\mathbf{s}_i\|^3 \leqslant \frac{f(\mathbf{x}_0) - f^* + (2\beta + \alpha + 2\gamma)\epsilon_1^3}{\gamma}.$$
 (45)

We next consider the convergence of $\|\nabla f(x_k)\|$ and $\|\nabla^2 f(x_k)\|$. Next, we prove the convergence rate of $\nabla f(\cdot)$ and $\nabla^2 f(\cdot)$. We first derive

$$\|\nabla f(\mathbf{x}_{k+1})\| \stackrel{\text{(i)}}{=} \|\nabla f(\mathbf{x}_{k+1}) - \left(\mathbf{g}_{k} + \mathbf{H}_{k}\mathbf{s}_{k+1} + \frac{M}{2} \|\mathbf{s}_{k+1}\| \mathbf{s}_{k+1}\right) \|$$

$$\leq \|\nabla f(\mathbf{x}_{k+1}) - (\mathbf{g}_{k} + \mathbf{H}_{k}\mathbf{s}_{k+1})\| + \frac{M}{2} \|\mathbf{s}_{k+1}\|^{2}$$

$$\leq \|\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_{k}) - \nabla^{2} f(\mathbf{x}_{k}) \mathbf{s}_{k+1} \| + \|\nabla f(\mathbf{x}_{k}) - \mathbf{g}_{k}\| + \|(\nabla^{2} f(\mathbf{x}_{k}) - \mathbf{H}_{k}) \mathbf{s}_{k+1}\| + \frac{M}{2} \|\mathbf{s}_{k+1}\|^{2}$$

$$\stackrel{\text{(ii)}}{\leq} \frac{L_{2}}{2} \|\mathbf{s}_{k+1}\|^{2} + \beta(\|\mathbf{s}_{k}\|^{2} + \epsilon_{1}^{2}) + \alpha(\|\mathbf{s}_{k}\| + \epsilon_{1}) \|\mathbf{s}_{k+1}\| + \frac{M}{2} \|\mathbf{s}_{k+1}\|^{2}$$

$$\stackrel{\text{(iii)}}{\leq} \left(\frac{L+M}{2} + 2\beta + 2\alpha\right) \epsilon_{1}^{2} \stackrel{\text{(iv)}}{\leq} \epsilon,$$

where (i) follows from eq. (32) with $\mathbf{g} = \mathbf{g}_k$, $\mathbf{H} = \mathbf{H}_k$ and $\mathbf{s} = \mathbf{s}_{k+1}$, (ii) follows from eq. (29) in Lemma 8 and Assumption 2, (iii) follows from the terminal condition of the algorithm, and (iv) follows from eq. (10).

Similarly, we have

$$\nabla^{2} f(\mathbf{x}_{k+1}) \stackrel{(i)}{\succcurlyeq} \mathbf{H}_{k} - \left\| \mathbf{H}_{k} - \nabla^{2} f(\mathbf{x}_{k+1}) \right\| \mathbf{I}$$

$$\stackrel{(ii)}{\succcurlyeq} - \frac{M}{2} \left\| \mathbf{s}_{k+1} \right\| \mathbf{I} - \left\| \mathbf{H}_{k} - \nabla^{2} f(\mathbf{x}_{k+1}) \right\| \mathbf{I}$$

$$\succcurlyeq - \frac{M}{2} \left\| \mathbf{s}_{k+1} \right\| \mathbf{I} - \left\| \mathbf{H}_{k} - \nabla^{2} f(\mathbf{x}_{k}) \right\| \mathbf{I} - \left\| \nabla^{2} f(\mathbf{x}_{k}) - \nabla^{2} f(\mathbf{x}_{m+1}) \right\| \mathbf{I}$$

$$\stackrel{(iii)}{\succcurlyeq} - \frac{M}{2} \left\| \mathbf{s}_{k+1} \right\| \mathbf{I} - \alpha (\left\| \mathbf{s}_{k} \right\| + \epsilon_{1}) \mathbf{I} - L_{2} \left\| \mathbf{s}_{k+1} \right\| \mathbf{I}$$

$$\stackrel{(iv)}{\succcurlyeq} - \left(\frac{M + 2L_{2}}{2} + 2\alpha \right) \epsilon_{1} \mathbf{I} \stackrel{(v)}{\succcurlyeq} \epsilon \mathbf{I},$$

where (i) follows from Weyl's inequality, (ii) follows from eq. (33) with $\mathbf{H} = \mathbf{H}_m$ and $\mathbf{s} = \mathbf{s}_{m+1}$, (iii) follows from Assumption 2 and the fact that $\nabla^2 f(\cdot)$ is L_2 -Lipschitz, (iv) follows from the terminal condition of the algorithm, and (v) follows from eq. (10).

B Proofs for SVRC under Sampling with Replacement

B.1 Proof of Theorem 2

The idea of the proof is to apply the following matrix Bernstein inequality Tropp (2012) for sampling with replacement to characterize the sample complexity in order to satisfy the inexactness condition $\|\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)\| \le \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$ with the probability at least $1 - \zeta$.

Lemma 10 (Matrix Bernstein Inequality). Consider a finite sequence $\{X_k\}$ of independent, random matrices with dimensions $d_1 \times d_2$. Assume that each random matrix satisfies

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

Define

$$\sigma^{2} \triangleq \max\left(\left\|\sum_{k} \mathbb{E}(\mathbf{X}_{k} \mathbf{X}_{k}^{*})\right\|, \left\|\sum_{k} \mathbb{E}(\mathbf{X}_{k}^{*} \mathbf{X}_{k})\right\|\right). \tag{46}$$

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

$$P\left(\left\|\sum_{k} \mathbf{X}_{k}\right\| \geqslant \epsilon\right) \leqslant 2(d_{1} + d_{2}) \exp\left(-\frac{\epsilon^{2}/2}{\sigma^{2} + R\epsilon/3}\right).$$

Let $\xi_H(k)$ be the collection of index that uniformly picked from $1, \dots, N$ with replacement, and \mathbf{X}_i be

$$\mathbf{X}_i = \frac{1}{|\xi_H(k)|} \left(\nabla^2 f_i(\mathbf{x}_k) - \nabla^2 f_i(\tilde{\mathbf{x}}) + \nabla^2 F(\tilde{\mathbf{x}}) - \nabla^2 F(\mathbf{x}_k) \right),$$

then we have

$$\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k) = \sum_{i \in \xi_H(k)} \mathbf{X}_i. \tag{47}$$

Moreover, we have $\mathbb{E}\mathbf{X}_i = \mathbf{0}$, and

$$R \triangleq \|\mathbf{X}_i\| = \frac{1}{|\xi_H(k)|} \|\nabla^2 f_{\xi_i}(\mathbf{x}_k) - \nabla^2 f_{\xi_i}(\tilde{\mathbf{x}}) + \nabla^2 F(\tilde{\mathbf{x}}) - \nabla^2 F(\mathbf{x}_k)\|$$

$$\stackrel{\text{(i)}}{\leq} \frac{2L_2}{|\xi_H(k)|} \|\mathbf{x}_k - \tilde{\mathbf{x}}\|, \tag{48}$$

where (i) follows because $\nabla^2 f_i(\cdot)$ is L_2 Lipschitz, for $1 \leq i \leq N$.

The variance also can be bounded by

$$\sigma^{2} \triangleq \max\left(\left\|\sum_{k \in \xi_{H}(k)} \mathbb{E}(\mathbf{X}_{k} \mathbf{X}_{k}^{*})\right\|, \left\|\sum_{k \in \xi_{H}(k)} \mathbb{E}(\mathbf{X}_{k}^{*} \mathbf{X}_{k})\right\|\right)$$

$$\stackrel{(i)}{\leq} \left\|\sum_{k \in \xi_{H}(k)} \mathbb{E}(\mathbf{X}_{k}^{2})\right\| \stackrel{(ii)}{\leq} \sum_{k \in \xi_{H}(k)} \mathbb{E}\left\|\mathbf{X}_{k}^{2}\right\| \leq \sum_{k \in \xi_{H}(k)} \mathbb{E}\left\|\mathbf{X}_{k}^{2}\right\|$$

$$\stackrel{(ii)}{\leq} \frac{4L_{2}^{2}}{|\xi_{H}(k)|} \left\|\mathbf{x}_{k} - \tilde{\mathbf{x}}\right\|^{2}$$

$$(49)$$

where (i) follows from the fact that \mathbf{X}_k is real and symmetric, (ii) follows from Jensen's inequality, and (iii) follows from eq. (48).

Therefore, in order to satisfy $\|\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)\| \le \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$ with probability at least $1 - \zeta$, by eq. (47), it is equivalent to require $\left\|\sum_{i \in \xi_H(k)} \mathbf{X}_i\right\| \le \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$ with probability at least $1 - \zeta$. We now apply Lemma 10 for \mathbf{X}_i , and it is sufficient to have:

$$2(d_1+d_2)\exp\left(\frac{-\epsilon^2/2}{\sigma^2+R\epsilon/3}\right) \leqslant \zeta$$

which is equivalent to have

$$\frac{1}{\sigma^2 + R\epsilon/3} \geqslant \frac{2}{\epsilon^2} \log \left(\frac{2(d_1 + d_2)}{\zeta} \right). \tag{50}$$

Plugging eqs. (48) and (49) into eq. (50) yields

$$\frac{1}{\frac{4L_{2}^{2}}{\left|\xi_{H}(\tilde{k})\right|}\left\|\mathbf{x}_{k}-\tilde{\mathbf{x}}\right\|^{2}+\frac{2L_{2}}{\left|\xi_{H}(\tilde{k})\right|}\left\|\mathbf{x}_{k}-\tilde{\mathbf{x}}\right\|\epsilon/3}\geqslant\frac{2}{\epsilon^{2}}\log\left(\frac{4d}{\zeta}\right),$$

which gives

$$|\xi_H(k)| \geqslant \left(\frac{8L_2^2}{\epsilon^2} \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2 + \frac{4L_2}{3\epsilon} \|\mathbf{x}_k - \tilde{\mathbf{x}}\|\right) \log\left(\frac{4d}{\zeta}\right). \tag{51}$$

Substituting $\epsilon = \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$, we obtain the required sample size to be bounded by

$$|\xi_{H}(k)| \geqslant \left(\frac{8L_{2}^{2}}{\alpha^{2} \max\{\|\mathbf{s}_{k}\|^{2}, \epsilon_{1}^{2}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|^{2} + \frac{4L_{2}}{3\alpha \max\{\|\mathbf{s}_{k}\|, \epsilon_{1}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|\right) \log\left(\frac{4d}{\zeta}\right). \tag{52}$$

We next bound the sample size $|\xi_g(k)|$ for the gradient in the similar procedure. We first define $\mathbf{X}_i \in \mathbb{R}^{d \times 1}$ as

$$\mathbf{X}_{i} = \frac{1}{|\xi_{g}(k)|} \left(\nabla f_{\xi_{i}}(\mathbf{x}_{k}) - \nabla f_{\xi_{i}}(\tilde{\mathbf{x}}) + \nabla F(\tilde{\mathbf{x}}) - \nabla F(\mathbf{x}_{k}) \right), \tag{53}$$

then we have

$$\mathbf{g}_k - \nabla f(\mathbf{x}_k) = \sum_{i \in \xi_g(k)} \mathbf{X}_i \tag{54}$$

Furthermore,

$$R = \|\mathbf{X}_i\| = \frac{1}{|\xi_g(k)|} \|\nabla f_{\xi_i}(\mathbf{x}_k) - \nabla f_{\xi_i}(\tilde{\mathbf{x}}) + \nabla F(\tilde{\mathbf{x}}) - \nabla F(\mathbf{x}_k)\| \stackrel{(i)}{\leqslant} \frac{2L_1}{|S_{g,k}|} \|\mathbf{x}_k - \tilde{\mathbf{x}}\|,$$
 (55)

where (i) follows because $\nabla f_i(\cdot)$ is L_1 Lipschitz, for $i = 1, \ldots, N$, and

$$\sigma^{2} \triangleq \max\left(\left\|\sum_{k \in \xi_{g}(k)} \mathbb{E}(\mathbf{X}_{k} \mathbf{X}_{k}^{*})\right\|, \left\|\sum_{k \in \xi_{g}(k)} \mathbb{E}(\mathbf{X}_{k}^{*} \mathbf{X}_{k})\right\|\right) \leqslant \sum_{k \in \xi_{H}(k)} \mathbb{E}\left\|\mathbf{X}_{k}\right\|^{2}$$

$$\stackrel{\text{(ii)}}{\leqslant} \frac{4L_{1}^{2}}{|\xi_{g}(k)|} \left\|\mathbf{x}_{k} - \tilde{\mathbf{x}}\right\|^{2}$$

In order to satisfy $\|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\| \le \beta \max \left\{ \|\mathbf{s}_k\|^2, \epsilon_1^2 \right\}$ with the probability at least $1 - \zeta$, by eq. (54), it is equivalent to require $\left\| \sum_{i \in \xi_g(k)} |\mathbf{X}_i| \right\| \le \beta \max \left\{ \|\mathbf{s}_k\|^2, \epsilon_1^2 \right\}$ with the probability at least $1 - \zeta$. We then apply Lemma 10 for \mathbf{X}_i in the way similar to that for bounding the sample size for Hessian, with $R = \frac{2L_1}{|S_{g,k}|} \|\mathbf{x}_k - \tilde{\mathbf{x}}\|$, $\epsilon = \beta \max \left\{ \|\mathbf{s}_k\|^2, \epsilon_1^2 \right\}$, and $\sigma^2 = \frac{4L_1^2}{|\xi_g(k)|} \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2$, and obtain the required sample size to satisfy

$$|\xi_{g}(k)| \geqslant \left(\frac{8L_{1}^{2}}{\beta^{2} \max\{\|\mathbf{s}_{k}\|^{4}, \epsilon_{1}^{4}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|^{2} + \frac{4L_{1}}{3\beta \max\{\|\mathbf{s}_{k}\|^{2}, \epsilon_{1}^{2}\}} \|\mathbf{x}_{k} - \tilde{\mathbf{x}}\|\right) \log\left(\frac{2(d+1)}{\zeta}\right). \tag{56}$$

B.2 Proof of Theorem 3

First, by eq. (13), we have

$$\sum_{i=1}^{k+1} \|\mathbf{x}_i - \mathbf{x}_{i-1}\|^3 \leqslant C. \tag{57}$$

We then derive

$$\sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \|\mathbf{x}_{i \cdot m+j} - \mathbf{x}_{i \cdot m}\|^{2}$$

$$\leq \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\|\mathbf{x}_{i \cdot m+j} - \mathbf{x}_{i \cdot m+j-1}\| + \dots + \|\mathbf{x}_{i \cdot m+1} - \mathbf{x}_{i \cdot m}\| \right)^{2}
\leq \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\|\mathbf{x}_{i \cdot m+m-1} - \mathbf{x}_{i \cdot m+m-2}\| + \dots + \|\mathbf{x}_{i \cdot m+1} - \mathbf{x}_{i \cdot m}\| \right)^{2}
= \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\sum_{l=1}^{m-1} \|\mathbf{x}_{i \cdot m+l} - \mathbf{x}_{i \cdot m+l-1}\| \right)^{2} \leq \sum_{i=0}^{i} \sum_{j=1}^{k/m-1} \sum_{l=1}^{m-1} \|\mathbf{x}_{i \cdot m+l} - \mathbf{x}_{i \cdot m+l-1}\|^{2}
\leq m^{2} \sum_{i=0}^{k/m-1} \sum_{l=1}^{m-1} \|\mathbf{x}_{i \cdot m+l} - \mathbf{x}_{i \cdot m+l-1}\|^{2} \leq m^{2} \sum_{i=1}^{k} \|\mathbf{x}_{i} - \mathbf{x}_{i-1}\|^{2}
\leq m^{2} k^{1/3} \left(\sum_{i=1}^{k} \|\mathbf{x}_{i} - \mathbf{x}_{i-1}\|^{3} \right)^{2/3} \leq m^{2} k^{1/3} C^{2/3}, \tag{58}$$

where (i) follows from the Cauthy-Schwaz inequality (ii) follows because j is not a variable in the inner summation, (iii) follows from Holder's inequality, and (iv) follows from eq. (57).

Similarly, we have that

$$\sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \|\mathbf{x}_{i \cdot m+j} - \mathbf{x}_{i \cdot m}\| \leqslant \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\|\mathbf{x}_{i \cdot m+j} - \mathbf{x}_{i \cdot m+j-1}\| + \dots + \|\mathbf{x}_{i \cdot m+1} - \mathbf{x}_{i \cdot m}\| \right)
\leqslant \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\|\mathbf{x}_{i \cdot m+m-1} - \mathbf{x}_{i \cdot m+m-2}\| + \dots + \|\mathbf{x}_{i \cdot m+1} - \mathbf{x}_{i \cdot m}\| \right)
= \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\sum_{l=1}^{m-1} \|\mathbf{x}_{i \cdot m+l} - \mathbf{x}_{i \cdot m+l-1}\| \right) \leqslant m \sum_{i=0}^{k/m-1} \sum_{l=1}^{m-1} \|\mathbf{x}_{i \cdot m+l} - \mathbf{x}_{i \cdot m+l-1}\|
\leqslant m \sum_{i=1}^{k} \|\mathbf{x}_{i} - \mathbf{x}_{i-1}\| \leqslant m k^{2/3} \left(\sum_{i=1}^{k} \|\mathbf{x}_{i} - \mathbf{x}_{i-1}\|^{3} \right)^{1/3} \leqslant m k^{2/3} C^{1/3}, \tag{59}$$

where (i) follows because j is not a variable in the inner summation, (ii) follows from Holder's inequality, and (iii) follows from eq. (57).

Thus, the total sample size for Hessian is given by

$$\begin{split} m + \frac{kN}{m} + \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left| \xi_H(k) \right| \\ \leqslant \frac{CkN}{m} + \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\frac{8L_2^2}{\alpha^2 \max\{\left\|\mathbf{s}_k\right\|^2, \epsilon_1^2\}} \left\| \mathbf{x}_{i \cdot m + j} - \mathbf{x}_{i \cdot m} \right\|^2 + \frac{4L_2}{3\alpha \max\{\left\|\mathbf{s}_k\right\|, \epsilon_1\}} \left\| \mathbf{x}_{i \cdot m + j} - \mathbf{x}_{i \cdot m} \right\| \right) \log \left(\frac{4d}{\zeta} \right) \\ \leqslant \frac{CkN}{m} + \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \left(\frac{8L_2^2}{\alpha^2 \epsilon_1^2} \left\| \mathbf{x}_{i \cdot m + j} - \mathbf{x}_{i \cdot m} \right\|^2 + \frac{4L_2}{3\alpha \epsilon_1} \left\| \mathbf{x}_{i \cdot m + j} - \mathbf{x}_{i \cdot m} \right\| \right) \log \left(\frac{4d}{\zeta} \right) \\ \leqslant \frac{CkN}{m} + \left(\frac{8L_2^2}{\alpha^2 \epsilon_1^2} m^2 k^{1/3} C^{2/3} + \frac{4L_2}{3\alpha \epsilon_1} m k^{2/3} C^{1/3} \right) \log \left(\frac{4d}{\zeta} \right) \\ \leqslant \log \left(\frac{4d}{\zeta} \right) \left(\frac{N}{m \epsilon^{3/2}} + \frac{C}{\epsilon^{3/2}} m^2 + \frac{C}{\epsilon^{3/2}} m \right) = \log \left(\frac{4d}{\zeta} \right) \frac{C}{\epsilon^{3/2}} \left(\frac{N}{m} + m^2 \right) \end{split}$$

where (i) follows form Theorem 2, and (ii) follows form eqs. (58) and (59), (iii) follows from the fact that $\zeta \leq 1$ and $d \geq 1$ which gives $\log\left(\frac{4d}{\zeta}\right) > 1$, and $\epsilon_1 = O(\epsilon^{1/2})$ such that $k = O(\epsilon^{-3/2})$ according to Theorem 1

We minimize the above bound over m, substitute the minimizer $m^* = N^{1/3}$, and obtain

$$\sum_{i=0}^{k} |\xi_H(k)| \leqslant \frac{CN^{2/3}}{\epsilon^{3/2}} \log \left(\frac{4d}{\zeta}\right).$$

Next, according to Theorem 2, Assumption 2 is satisfies with probability at least $1-\zeta$ for gradient and $1-\zeta$ for Hessian . Thus, according to the union bound, the probability of a failure satisfaction per iteration is at most 2ζ . Then, for k iteration, the probability of failure satisfaction of Assumption 2 is at most $2k\zeta$ according to the union bound. To obtain Assumption 2 holds for the total k iteration with probability least $1-\delta$, we require

$$1 - 2k\zeta \geqslant 1 - \delta$$
,

which yields

$$\zeta \leqslant \frac{\delta}{2k}$$
.

Thus, with probability $1 - \delta$, the algorithms successfully outputs an ϵ approximated second-order stationary point, with the total Hessian sample complexity is bounded by

$$\sum_{i=0}^{k} |\xi_H(k)| \leqslant \frac{CN^{2/3}}{\epsilon^{3/2}} \log \left(\frac{8d}{\epsilon^{3/2} \delta} \right) \leqslant \frac{CN^{2/3}}{\epsilon^{3/2}} \log \left(\frac{8d}{\epsilon \delta} \right). \tag{60}$$

which gives

$$\sum_{i=0}^{k} |\xi_H(k)| = \tilde{O}\left(\frac{N^{2/3}}{\epsilon^{3/2}}\right). \tag{61}$$

C Proof of Concentration Inequality for Sampling without replacement

The proof generalizes the Hoeffding-Serfling inequality for scalar random variables in Bardenet and Maillard (2015) to that for random matrices. We also apply various properties for handling random matrices in Tropp (2012).

C.1 Definitions and Useful Lemmas

We first introduce the definition of the matrix function following Tropp (2012), and then introduce a number of Lemmas that are useful in the proof.

Given a symmetric matrix **A**, suppose its eigenvalue decomposition is given by $\mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \in \mathbb{R}^{d\times d}$, where $\mathbf{\Lambda} = diag(\lambda_1, \dots, \lambda_d)$. Then a function $f: \mathbb{R} \to \mathbb{R}$ of **A** is defined as:

$$f(\mathbf{A}) \triangleq \mathbf{U}f(\mathbf{\Lambda})\mathbf{U}^T,\tag{62}$$

where $f(\mathbf{\Lambda}) = diag(f(\lambda_1), \dots, f(\lambda_d))$, i.e., $f(\mathbf{\Lambda})$ applies the function $f(\cdot)$ to each diagonal entry of the matrix $\mathbf{\Lambda}$.

The trace exponential function tr $\exp: \mathbf{A} \to \operatorname{tre}^{\mathbf{A}}$, i.e., tr $\exp(\mathbf{A})$, is defined to first apply the exponential matrix function $\exp(\mathbf{A})$, and then take the trace of $\exp(\mathbf{A})$. Such a function is monotone with respect to the semidefinite order:

$$\mathbf{A} \preceq \mathbf{H} \implies \operatorname{tr} \exp(\mathbf{A}) \preceq \operatorname{tr} \exp(\mathbf{H}),$$
 (63)

which follows because for two symmetric matrices **A** and **H**, if $\mathbf{A} \leq \mathbf{H}$, then $\lambda_i(\mathbf{A}) \leq \lambda_i(\mathbf{H})$ for every i, where $\lambda_i(\mathbf{A})$ is the i-th largest eigenvalue of **A**. Furthermore, the matrix function $\log(\cdot)$ is monotone with respect to the semidefinite order (see the exercise 4.2.5 in Bhatia (2007)):

$$\mathbf{0} \prec \mathbf{A} \preccurlyeq \mathbf{H} \implies \log(\mathbf{A}) \preccurlyeq \log(\mathbf{H}).$$
 (64)

The next three lemmas follow directly from Bardenet and Maillard (2015) because the proofs are applicable for matrices.

Lemma 11. [Bardenet and Maillard (2015)] Let $\mathbf{Z}_k \triangleq \frac{1}{k} \sum_{i=1}^k \mathbf{X}_i$. The following reverse martingale structure holds for $\{\mathbf{Z}_k\}_{k \leq N}$:

$$\mathbb{E}[\mathbf{Z}_k|\mathbf{Z}_{k+1},\cdots\mathbf{Z}_{N-1}] = \mathbf{Z}_{k+1}.\tag{65}$$

Lemma 12. [Bardenet and Maillard (2015)] Let $\mathbf{Y}_k \triangleq \mathbf{Z}_{N-k}$ for $1 \leqslant k \leqslant N-1$. For any $\lambda > 0$, the following equality holds for $2 \leqslant k \leqslant n$,

$$\lambda \mathbf{Y}_k = \lambda \mathbf{Y}_{k-1} - \lambda \frac{\mathbf{X}_{N-k+1} - \mu - \mathbf{Y}_{k-1}}{N-k}.$$
 (66)

Lemma 13. [Bardenet and Maillard (2015)] Let $\mathbf{Y}_k \triangleq \mathbf{Z}_{N-k}$ for $1 \leqslant k \leqslant N-1$. For $2 \leqslant k \leqslant N$, the following equality holds

$$\mathbb{E}[\mathbf{X}_{N-k+1} - \mu - \mathbf{Y}_{k-1} | Y_1, \cdots, \mathbf{Y}_{k-1}] = 0, \tag{67}$$

where $\mu = \frac{1}{N} \sum_{t=1}^{N} \mathbf{X}_t$.

The following lemma is an extension of Hoeffding's inequality for scalars to matrices. We include a brief proof for completeness.

Lemma 14 (Hoeffding's Inequality for Matrix). For a random symmetric matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$, suppose

$$\mathbb{E}[\mathbf{X}] = 0 \quad and \quad a\mathbf{I} \preccurlyeq \mathbf{X} \preccurlyeq b\mathbf{I}.$$

where a and b are real constants. Then for any $\lambda > 0$, the following inequality holds

$$\mathbb{E}[e^{\lambda \mathbf{X}}] \preceq \exp\left(\frac{1}{8}\lambda^2 (b-a)^2 \mathbf{I}\right). \tag{68}$$

Proof. The proof follows from the standard reasoning for scalar version. We emphasize only the difference in handling matrices. Suppose the eigenvalue decomposition of the symmetric random matrix **X** can be written as $\mathbf{X} = \mathbf{U}\Lambda\mathbf{U}^T$, where $\mathbf{U} = [\mathbf{u}_1, \cdots, \mathbf{u}_d]$ and $\mathbf{\Lambda} = diag(\lambda_1, \cdots, \lambda_d)$. Therefore, we obtain $e^{\lambda \mathbf{X}} = \sum_{i=1}^d e^{\lambda \lambda_i} \mathbf{u}_i \mathbf{u}_i^T$.

Since scalar function $e^{\lambda x}$ is convex for any $\lambda > 0$, for $1 \le i \le d$, we have

$$e^{\lambda \lambda_i} \leqslant \left(\frac{b - \lambda_i}{b - a}e^{\lambda a} + \frac{\lambda_i - a}{b - a}e^{\lambda b}\right),$$
 (69)

which implies that

$$e^{\lambda \lambda_i} \mathbf{u}_i \mathbf{u}_i^T \preceq \left(\frac{b - \lambda_i}{b - a} e^{\lambda a} + \frac{\lambda_i - a}{b - a} e^{\lambda b} \right) \mathbf{u}_i \mathbf{u}_i^T.$$
 (70)

Then,

$$\begin{split} \mathbb{E}[e^{\lambda \mathbf{X}}] &= \mathbb{E}\bigg[\sum_{i=1}^{d} e^{\lambda \lambda_{i}} \mathbf{u}_{i} \mathbf{u}_{i}^{T}\bigg] \overset{\text{(i)}}{\preccurlyeq} \mathbb{E}\bigg[\sum_{i=1}^{d} \left(\frac{b-\lambda_{i}}{b-a} e^{\lambda a} + \frac{\lambda_{i}-a}{b-a} e^{\lambda b}\right) \mathbf{u}_{i} \mathbf{u}_{i}^{T}\bigg] \\ &= \mathbb{E}\bigg[\sum_{i=1}^{d} \frac{b}{b-a} e^{\lambda a} \mathbf{u}_{i} \mathbf{u}_{i}^{T} - \sum_{i=1}^{d} \frac{\lambda_{i}}{b-a} e^{\lambda a} \mathbf{u}_{i} \mathbf{u}_{i}^{T} + \sum_{i=1}^{d} \frac{\lambda_{i}}{b-a} e^{\lambda b} \mathbf{u}_{i} \mathbf{u}_{i}^{T} - \sum_{i=1}^{d} \frac{a}{b-a} e^{\lambda b} \mathbf{u}_{i} \mathbf{u}_{i}^{T}\bigg] \\ &\stackrel{\text{(ii)}}{=} \mathbb{E}\bigg[\sum_{i=1}^{d} \frac{b}{b-a} e^{\lambda a} \mathbf{u}_{i} \mathbf{u}_{i}^{T} - \frac{e^{\lambda a}}{b-a} \mathbf{X} + \frac{e^{\lambda b}}{b-a} \mathbf{X} - \sum_{i=1}^{d} \frac{a}{b-a} e^{\lambda b} \mathbf{u}_{i} \mathbf{u}_{i}^{T}\bigg] \\ &\stackrel{\text{(iii)}}{=} \mathbb{E}\bigg[\sum_{i=1}^{d} \frac{b}{b-a} e^{\lambda a} \mathbf{u}_{i} \mathbf{u}_{i}^{T} - \sum_{i=1}^{d} \frac{a}{b-a} e^{\lambda b} \mathbf{u}_{i} \mathbf{u}_{i}^{T}\bigg] \\ &\stackrel{\text{(iv)}}{=} \mathbb{E}\bigg[\frac{b}{b-a} e^{\lambda a} \mathbf{I} - \frac{a}{b-a} e^{\lambda b} \mathbf{I}\bigg] = \bigg(\frac{b}{b-a} e^{\lambda a} - \frac{a}{b-a} e^{\lambda b}\bigg) \mathbf{I} \end{split}$$

$$\preccurlyeq \exp\left(\frac{1}{8}\lambda^2(b-a)^2\right)\mathbf{I} \stackrel{\text{(v)}}{=} \exp\left(\frac{1}{8}\lambda^2(b-a)^2\mathbf{I}\right),$$
(71)

where (i) follows from eq. (70) and the fact that the expectation of random matrix preserves the semi-definite order, (ii) follows from $\mathbf{X} = \sum_{i=1}^{d} \lambda_i \mathbf{u}_i \mathbf{u}_i^T$, (iii) follows because $\mathbb{E}[\mathbf{X}] = 0$, (iv) follows because $\mathbf{I} = \mathbf{U}\mathbf{U}^T = \sum_{i=1}^{d} \mathbf{u}_i \mathbf{u}_i^T$, and (v) follows from the standard steps in the proof of the scalar version of Hoeffding's inequality.

Lemma 15. Tropp (2012)[Corollary 3.3] Let \mathbf{H} be a fixed self-adjoint matrix, and let \mathbf{X} be a random self-adjoint matrix. The following inequality holds

$$\mathbb{E} \operatorname{tr} \exp(\mathbf{H} + \mathbf{X}) \leqslant \operatorname{tr} \exp(\mathbf{H} + \log(\mathbb{E}e^{\mathbf{X}})). \tag{72}$$

Lemma 16. Bardenet and Maillard (2015) For integer $n \leq N$, the following inequality holds

$$\sum_{t=1}^{n} \left(\frac{1}{N-t} \right)^{2} \leqslant \frac{n}{(N-n)^{2}} \left(1 - \frac{n-1}{N} \right)$$

C.2 Proof of Theorem 4

First, it suffices to show the theorem only for symmetric matrices, due to the technique of *dilations* in Tropp (2012) that transforms the asymmetric matrix to a symmetric matrix while keeping the spectral norm to be the same.

Second, it also suffices to show that for $1 \le i \le N$, \mathbf{X}_i are symmetric and bounded, i.e., $a\mathbf{I} \le \mathbf{X}_i \le b\mathbf{I}$, and $1 \le n \le N-1$, the following inequality holds

$$P\bigg(\lambda_{\max}\bigg(\frac{1}{n}\sum_{i=1}^{n}\mathbf{X}_{i}-\mu\bigg)\geqslant\epsilon\bigg)\leqslant d\exp\bigg(-\frac{n\epsilon^{2}}{2(b-a)^{2}(1+1/n)(1-n/N)}\bigg).$$

This is because the above result, with X_i being replaced with $-X_i$, implies

$$P\left(\lambda_{\min}\left(\frac{1}{n}\sum_{i=1}^{n}\mathbf{X}_{i}-\mu\right)\leqslant-\epsilon\right)\leqslant d\exp\left(-\frac{n\epsilon^{2}}{2(b-a)^{2}(1+1/n)(1-n/N)}\right). \tag{73}$$

Then the combination of the two results completes the desired theorem.

We start the proof by applying the matrix version of Chernoff inequality as follows. Let $\mathbf{Z}_k \triangleq \frac{1}{k} \sum_{i=1}^k \mathbf{X}_i$, for any $\lambda > 0$, we obtain

$$P\left(\lambda_{\max}(\mathbf{Z}_{n}) \geqslant \epsilon\right) = P\left(\exp(\lambda \lambda_{\max}(\mathbf{Z}_{n})) \geqslant \exp(\lambda \epsilon)\right)$$

$$\stackrel{(i)}{\leqslant} \exp(-\lambda \epsilon) \mathbb{E} \exp\left(\lambda \lambda_{\max}(\mathbf{Z}_{n})\right)$$

$$\stackrel{(iii)}{\leqslant} \exp(-\lambda \epsilon) \mathbb{E} \lambda_{\max}\left(\exp(\lambda \mathbf{Z}_{n})\right)$$

$$\stackrel{(iiii)}{\leqslant} \exp(-\lambda \epsilon) \mathbb{E} \operatorname{tr} \exp(\lambda \mathbf{Z}_{n})$$

$$\stackrel{(iv)}{\leqslant} \exp(-\lambda \epsilon) \operatorname{tr} \exp\left(\frac{\lambda^{2}}{2}(b-a)^{2}\frac{(n+1)}{n^{2}}\left(1-\frac{n}{N}\right)I\right)$$

$$\stackrel{(v)}{\leqslant} d \exp\left(\frac{\lambda^{2}}{2}(b-a)^{2}\frac{(n+1)}{n^{2}}\left(1-\frac{n}{N}\right)\right) \exp(-\lambda \epsilon)$$

$$= d \exp\left(\frac{\lambda^{2}}{2}(b-a)^{2}\frac{(n+1)}{n^{2}}\left(1-\frac{n}{N}\right) - \lambda \epsilon\right)$$

$$(74)$$

where (i) follows from the matrix version of Chernoff inequality, (ii) follows from the fact that $\exp(\cdot)$ is an increasing function, thus $\exp(\lambda \lambda_{\max}(\mathbf{Z}_n)) = \lambda_{\max}(\exp(\lambda \mathbf{Z}_n))$, and (iii) follows from the fact that $\lambda_{\max}(\mathbf{A}) \leq \operatorname{tr}(\mathbf{A})$, with $\mathbf{A} = \exp(\lambda \mathbf{Z}_n)$, we get the desire result.

We next bound \mathbb{E} tr $\exp(\lambda \mathbf{Z}_n)$. Let $Y_k \triangleq Z_{N-k}$ for $1 \leqslant k \leqslant N-1$, and $\mathbb{E}_k[\cdot] \triangleq \mathbb{E}[\cdot|\mathbf{Y}_1,\cdots,\mathbf{Y}_k]$. Thus,

$$\mathbb{E} \operatorname{tr} \exp(\lambda \mathbf{Y}_{n}) \stackrel{\text{(i)}}{=} \mathbb{E} \operatorname{tr} \exp\left(\lambda \mathbf{Y}_{n-1} - \lambda \frac{\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1}}{N-n}\right)$$

$$\stackrel{\text{(ii)}}{=} \mathbb{E} \mathbb{E}_{n-1} \operatorname{tr} \exp\left(\lambda \mathbf{Y}_{n-1} - \lambda \frac{\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1}}{N-n}\right)$$

$$\stackrel{\text{(iii)}}{\leq} \mathbb{E} \operatorname{tr} \exp\left(\lambda \mathbf{Y}_{n-1} + \log \mathbb{E}_{n-1} \exp\left(-\lambda \frac{\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1}}{N-n}\right)\right), \tag{75}$$

where (i) follows from Lemma 12, (ii) follows from the tower property of expectation, (iii) follows by applying Lemma 15, where $\lambda \mathbf{Y}_{n-1}$ is deterministic given $\mathbf{Y}_1, \dots, \mathbf{Y}_k$, and $-\lambda (\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1})/(N-n)$ is a random variable matrix.

In order to apply Lemma 14 to bound $\mathbb{E}_{n-1} \exp(-\lambda (\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1})/(N-n))$, we first bound $\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1}$ as follows:

$$\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1} \stackrel{\text{(i)}}{=} \mathbf{X}_{N-n+1} - \mu - \mathbf{Z}_{N-n+1}$$

$$\stackrel{\text{(ii)}}{=} \mathbf{X}_{N-n+1} - \mu - \frac{1}{N-n+1} \sum_{i=1}^{N-n+1} \left(\mathbf{X}_i - \mu \right)$$

$$= \mathbf{X}_{N-n+1} - \frac{1}{N-n+1} \sum_{i=1}^{N-n+1} \mathbf{X}_i, \tag{76}$$

where (i) follows from the definition of \mathbf{Y}_{n-1} and (ii) follows from the definition of \mathbf{Z}_{N-n+1} . Since $a\mathbf{I} \preceq \mathbf{X}_i \preceq b\mathbf{I}$, the above equality implies

$$-\frac{(b-a)}{N-n}\mathbf{I} \preccurlyeq \frac{\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1}}{N-n} \preccurlyeq \frac{(b-a)}{N-n}\mathbf{I}.$$
 (77)

By applying Lemma 14, and the fact $\mathbb{E}_{n-1}[\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1}] = 0$ due to Lemma 13, we obtain

$$\mathbb{E}_{n-1} \exp\left(\mathbf{X}_{N-n+1} - \mu - \mathbf{Y}_{n-1}\right) \preceq \exp\left(\frac{1}{8}\lambda^2 \left(\frac{2(b-a)}{N-n}\right)^2 \mathbf{I}\right) = \exp\left(\frac{1}{2}\lambda^2 \left(\frac{b-a}{N-n}\right)^2 \mathbf{I}\right),\tag{78}$$

Substituting eq. (78) into eq. (75), we obtain

$$\mathbb{E} \operatorname{tr} \exp(\lambda \mathbf{Y}_{n}) \overset{\text{(i)}}{\leqslant} \mathbb{E} \operatorname{tr} \exp\left(\lambda \mathbf{Y}_{n-1} + \log \exp\left(\frac{1}{2}\lambda^{2} \left(\frac{b-a}{N-n}\right)^{2} \mathbf{I}\right)\right)$$

$$= \mathbb{E} \operatorname{tr} \exp\left(\lambda \mathbf{Y}_{n-1} + \frac{\lambda^{2}}{2} \left(\frac{b-a}{N-n}\right)^{2} \mathbf{I}\right)$$

$$\dots$$

$$\overset{\text{(ii)}}{\leqslant} \operatorname{tr} \exp\left(\log \mathbb{E}[e^{\lambda \mathbf{Y}_{1}}] + \sum_{t=2}^{n} \frac{\lambda^{2}}{2} \left(\frac{b-a}{N-t}\right)^{2} \mathbf{I}\right). \tag{79}$$

where (i) follows from eqs. (63) and (64), and (ii) follows by applying the steps similar to obtain eq. (78) for n-2 times.

To bound $\mathbb{E}[e^{\lambda \mathbf{Y}_1}]$, we first note that

$$\mathbf{Y}_{1} = \mathbf{Z}_{N-1} = \frac{1}{N-1} \sum_{i=1}^{N-1} \left(\mathbf{X}_{i} - \mu \right) \stackrel{\text{(i)}}{=} \frac{1}{N-1} \left(N\mu - \mathbf{X}_{N} - (N-1)\mu \right) = \frac{1}{N-1} \left(\mu - \mathbf{X}_{N} \right),$$

where (i) follows because $N\mu = \sum_{i=1}^{N} \mathbf{X}_{i}$. Thus with $a\mathbf{I} \leq \mathbf{X}_{i} \leq b\mathbf{I}$ and $a\mathbf{I} \leq \mu \leq b\mathbf{I}$, we obtain

$$-\frac{(b-a)}{N-1}\mathbf{I} \preccurlyeq \mathbf{Y}_1 \preccurlyeq \frac{(b-a)}{N-1}\mathbf{I}.$$
(80)

Applying the matrix Hoeffding lemma with eq. (80) and $\mathbb{E}[Y_1] = \mathbb{E}[Z_{N-1}] = 0$, we obtain

$$\mathbb{E}[e^{\lambda \mathbf{Y}_1}] \preccurlyeq \exp\left(\frac{1}{2}\lambda^2 \left(\frac{b-1}{N-1}\right)^2 \mathbf{I}\right). \tag{81}$$

Substituting eq. (81) into eq. (79), we obtain

$$\mathbb{E} \operatorname{tr} \exp(\lambda \mathbf{Y}_{n}) \leqslant \operatorname{tr} \exp\left(\sum_{t=1}^{n} \frac{\lambda^{2}}{2} \left(\frac{b-a}{N-t}\right)^{2} \mathbf{I}\right)$$

$$= \operatorname{tr} \exp\left(\frac{\lambda^{2}}{2} (b-a)^{2} \sum_{t=1}^{n} \left(\frac{1}{N-t}\right)^{2} \mathbf{I}\right)$$

$$\stackrel{\text{(i)}}{\leqslant} \operatorname{tr} \exp\left(\frac{\lambda^{2}}{2} (b-a)^{2} \frac{n}{(N-n)^{2}} \left(1 - \frac{n-1}{N}\right) \mathbf{I}\right), \tag{82}$$

where (i) follows from lemma 16.

Now let m = N - n, where $1 \leq m \leq N - 1$, and hence $\mathbf{Y}_n = \mathbf{Z}_{N-n}$. Thus, eq. (82) implies

$$\mathbb{E} \text{ tr } \exp(\lambda \mathbf{Z}_m) \leqslant \text{tr } \exp\left(\frac{\lambda^2}{2}(b-a)^2 \frac{(m+1)}{m^2} \left(1 - \frac{m}{N}\right) \mathbf{I}\right).$$

Substituting the above bound into eq. (74), we obtain

$$P\left(\lambda_{\max}(\mathbf{Z}_n) \geqslant \epsilon\right) \leqslant \exp(-\lambda \epsilon) \operatorname{tr} \exp\left(\frac{\lambda^2}{2} (b-a)^2 \frac{(n+1)}{n^2} \left(1 - \frac{n}{N}\right) \mathbf{I}\right)$$
$$= d \exp\left(\frac{\lambda^2}{2} (b-a)^2 \frac{(n+1)}{n^2} \left(1 - \frac{n}{N}\right) - \lambda \epsilon\right), \tag{83}$$

where the last step follows form the equation $\operatorname{tr}(a\mathbf{I}) = da$ for $\mathbf{I} \in \mathbb{R}^{d \times d}$. The proof is completed by minimizing the above bound with respect to $\lambda > 0$, and then substituting the minimizer $\lambda^* = \frac{n\epsilon}{(b-a)^2(1+\frac{1}{n})(1-\frac{n}{N})}$.

D Proofs for SVRC under Sampling without Replacement

D.1 Proof of Theorem 5

Proof. The idea of the proof is to apply the matrix concentration inequality for sampling without replacement that we developed in Theorem 4 to characterize the sample complexity in order to satisfy the inexactness condition $\|\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)\| \leq \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$ with the probability at least $1 - \zeta$.

We first note that

$$\begin{split} \mathbf{H}_k - \nabla^2 F(\mathbf{x}_k) &\stackrel{\text{(i)}}{=} \frac{1}{|\xi_H(k)|} \left[\sum_{i \in \xi_H(k)} (\nabla^2 f_i(\mathbf{x}_k) - \nabla^2 f_i(\tilde{\mathbf{x}})) \right] + \nabla^2 F(\tilde{\mathbf{x}}_k) - \nabla^2 F(\mathbf{x}_k) \\ &= \frac{1}{|\xi_H(k)|} \sum_{i \in \xi_H(k)} \left(\nabla^2 f_i(\mathbf{x}_k) - \nabla^2 f_i(\tilde{\mathbf{x}}) + \nabla^2 F(\tilde{\mathbf{x}}) - \nabla^2 F(\mathbf{x}_k) \right) \end{split}$$

where (i) follows from the definition of \mathbf{H}_k in Algorithm 1. In order to apply the concentration inequality (Theorem 4) to bound $\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)$, we define, for $1 \leq i \leq N$,

$$\mathbf{X}_i = \nabla^2 f_i(\mathbf{x}_k) - \nabla^2 f_i(\tilde{\mathbf{x}}) + \nabla^2 F(\tilde{\mathbf{x}}) - \nabla^2 F(\mathbf{x}_k),$$

which gives

$$\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k) = \frac{1}{|\xi_H(k)|} \sum_{i \in \xi_H(k)} \mathbf{X}_i.$$
(84)

Moreover, we have $\mu \triangleq \frac{1}{N} \sum_{i=1}^{N} \mathbf{X}_i = \mathbf{0}$, and

$$\sigma \triangleq \|\mathbf{A}_i\| = \|\nabla^2 f_i(\mathbf{x}_k) - \nabla^2 f_i(\tilde{\mathbf{x}}) + \nabla^2 F(\tilde{\mathbf{x}}) - \nabla^2 F(\mathbf{x}_k)\| \stackrel{\text{(i)}}{\leqslant} 2L_2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|,$$

where (i) follows because $\nabla^2 f_i(\cdot)$ is L_2 Lipschitz, for $1 \leq i \leq N$.

Thus, in order to satisfy $\|\mathbf{H}_k - \nabla^2 F(\mathbf{x}_k)\| \le \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$ with probability at least $1 - \zeta$, by eq. (84), it is equivalent to satisfy $\left\|\frac{1}{|\xi_H(k)|}\sum_{i\in\xi_H(k)}\mathbf{X}_i - \mu\right\| \le \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$ with probability at least $1 - \zeta$. We now apply Theorem 4 for \mathbf{X}_i , and it is sufficient to have:

$$2(d_1 + d_2) \exp\left(-\frac{n\epsilon^2}{8\sigma^2(1 + 1/n)(1 - n/N)}\right) \leqslant \zeta,$$

which implies

$$\frac{n\epsilon^2}{8\sigma^2(1+1/n)(1-n/N)} \geqslant \log(\frac{2(d_1+d_2)}{\zeta}).$$

Using $(1+1/n) \leq 2$, it is sufficient to have:

$$\frac{n\epsilon^2}{16\sigma^2(1-n/N)} \geqslant \log(\frac{2(d_1+d_2)}{\zeta}),$$

which implies

$$n \geqslant \frac{1}{\frac{1}{N} + \frac{\epsilon^2}{16\sigma^2 \log(2(d_1 + d_2)/\zeta)}}.$$
(85)

We then substitute $\sigma = 2L_2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|$, $\epsilon = \alpha \max\{\|\mathbf{s}_k\|, \epsilon_1\}$, and $n = |\xi_H(k)|$, and obtain the required sample size to satisfy

$$|\xi_H(k)| \geqslant \frac{1}{\frac{1}{N} + \frac{\alpha^2 \max\{\|\mathbf{s}_k\|^2, \epsilon_1^2\}}{64L_2^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2 \log(4d/\zeta)}}.$$
 (86)

We next bound the sample size $|\xi_g(k)|$ for the gradient, the proof follows the same procedure. We first define $\mathbf{X}_i \in \mathbb{R}^{d \times 1}$ as

$$\mathbf{X}_{i} = \nabla f_{i}(\mathbf{x}_{k}) - \nabla f_{i}(\tilde{\mathbf{x}}) + \nabla F(\tilde{\mathbf{x}}) - \nabla F(\mathbf{x}_{k}), \tag{87}$$

and hence

$$\mathbf{g}_k - \nabla F(\mathbf{x}_k) = \frac{1}{|\xi_g(k)|} \sum_{i \in \xi_g(k)} \mathbf{X}_i.$$
(88)

Moreover, we have $\mu = \frac{1}{N} \sum_{i \in \xi_q(k)} \mathbf{A}_i = \mathbf{0}$, and

$$\sigma \triangleq \|\mathbf{A}_i\| = \|\nabla f_i(\mathbf{x}_k) - \nabla f_i(\tilde{\mathbf{x}}) + \nabla F(\tilde{\mathbf{x}}) - \nabla F(\mathbf{x}_k)\| \stackrel{\text{(i)}}{\leqslant} 2L_1 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|,$$

where (i) follows because $\nabla f_i(\cdot)$ is L_1 Lipschitz, for $1 \leq i \leq N$.

In order to satisfy $\|\mathbf{g}_k - \nabla F(\mathbf{x}_k)\| \le \beta \max\{\|\mathbf{s}_k\|^2, \epsilon_1^2\}$ with probability at least $1 - \zeta$, by eq. (88), it is equivalent to satisfy $\left\|\frac{1}{|\xi_g(k)|}\sum_{i\in\xi_g(k)}\mathbf{X}_i - \mu\right\| \le \beta \max\{\|\mathbf{s}_k\|^2, \epsilon_1^2\}$ with probability at least $1 - \zeta$. We then apply Theorem 4 for \mathbf{X}_i in the way similar to that for bounding the sample size for Hessian, with $\sigma = 2L_1 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|$, $\mu = 0$, $\epsilon = \beta \max\{\|\mathbf{s}_k\|^2, \epsilon_1^2\}$, and $n = |\xi_g(k)|$, and obtain the required sample size to satisfy

$$|\xi_g(k)| \geqslant \frac{1}{\frac{1}{N} + \frac{\beta^2 \max\{\|\mathbf{s}_k\|^4, \epsilon_1^4\}}{64L^2 \|\mathbf{x}_k - \tilde{\mathbf{x}}\|^2 \log(2(d+1)/\zeta)}}.$$
(89)

D.2 Proof of Proposition 6

Proof. The proof of Proposition 6 is similar to the proof of Theorem 5. We first define $\mathbf{A}_i \in \mathbb{R}^{d \times d}$ as

$$\mathbf{A}_i = \nabla^2 f_i(\mathbf{x}_k) - \nabla^2 F(\mathbf{x}_k), \tag{90}$$

and hence $\mu = \frac{1}{N} \sum_{i \in \xi_H(k)} \mathbf{A}_i = \mathbf{0}$. Furthermore,

$$\sigma \triangleq \|\mathbf{A}_i\| = \|\nabla^2 f_i(\mathbf{x}_k) - \nabla^2 F(\mathbf{x}_k)\| \stackrel{\text{(i)}}{\leqslant} 2L_1,$$

where (i) follows from Assumption 1.

Let $\{\mathbf{X}_i\}_{i=1}^{|\xi_g(k)|} = \{\mathbf{A}_i : i \in \xi_H(k)\}, \text{ and we have }$

$$\frac{1}{|\xi_H(k)|} \sum_{i \in \xi_H(k)} \mathbf{X}_i - \mu \stackrel{\text{(i)}}{=} \frac{1}{|\xi_H(k)|} \sum_{i \in \xi_H(k)} \mathbf{A}_i \stackrel{\text{(ii)}}{=} \mathbf{H}_k - \nabla^2 F(\mathbf{x}_k), \tag{91}$$

where (i) follows from the fact that $\mu = 0$ and (ii) follows from the definition of \mathbf{H}_k in Algorithm 1.

We then apply Theorem 4 for \mathbf{X}_i with $\sigma = 2L_1$, $\mu = 0$, $\epsilon = C_2 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|$, and $n = |\xi_H(k)|$, and obtain the require sampled size to satisfy

$$|\xi_H(k)| \geqslant \frac{1}{\frac{1}{N} + \frac{C_2^2 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2}{64L_1^2 \log(4d/\zeta)}}.$$
 (92)

To bound the sample size of gradient, i.e., $|\xi_q(k)|$, we follow the similar proof by constructing

$$\mathbf{A}_i = \nabla f_i(\mathbf{x}_k) - \nabla F(\mathbf{x}_k),\tag{93}$$

and applying Theorem 4 with $\sigma = 2L_0$, $\mu = 0$, $\epsilon = C_1 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^2$, and $n = |\xi_g(k)|$, and obtain the required sample size to satisfy

$$|\xi_g(k)| \geqslant \frac{1}{\frac{1}{N} + \frac{C_1^2 \|\mathbf{x}_{k+1} - \mathbf{x}_k\|^4}{64L_2^2 \log(2(d+1)/\zeta)}}.$$
 (94)

D.3 Proof of Theorem 7

Proof. Assume the algorithm terminates at iteration k, then the total Hessian complexity is given by

$$m + \frac{kN}{m} + \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} |\xi_H(k)| \stackrel{\text{(i)}}{\leqslant} \frac{CkN}{m} + \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \frac{1}{\frac{1}{N} + \frac{\alpha^2 \max\{\|\mathbf{s}_k\|^2, \epsilon_1^2\}}{64L_2^2 \|\mathbf{x}_{i \cdot m + j} - \mathbf{x}_{i \cdot m}\|^2 \log(4d/\zeta)}}$$

$$\leqslant \frac{CkN}{m} + \sum_{i=0}^{k/m-1} \sum_{j=1}^{m-1} \frac{64L_2^2 \|\mathbf{x}_{i \cdot m + j} - \mathbf{x}_{i \cdot m}\|^2 \log(4d/\zeta)}{\alpha^2 \epsilon_1^2}$$

$$\stackrel{\text{(ii)}}{\leqslant} \frac{CkN}{m} + \frac{64L_2^2}{\alpha^2 \epsilon_1^2} \left(m^2 k^{1/3} C^{2/3}\right) \log\left(\frac{4d}{\zeta}\right)$$

$$\stackrel{\text{(iii)}}{\leqslant} C \log\left(\frac{4d}{\zeta}\right) \left(\frac{N}{m\epsilon^{3/2}} + \frac{m^2}{\epsilon^{3/2}}\right) = \frac{C}{\epsilon^{3/2}} \log\left(\frac{4d}{\zeta}\right) \left(\frac{N}{m} + m^2\right)$$

where (i) follows form Theorem 5, and (ii) follows form eq. (58), (iii) follows from the fact that $\zeta < 1$ and $d \ge 1$ which gives $\log\left(\frac{4d}{\zeta}\right) > 1$, and the fact that $\epsilon_1 = O(\epsilon^{1/2})$ such that $k = O(\epsilon^{-3/2})$ according to Theorem 1.

We minimize the above bound over m, substitute the minimizer $m^* = N^{1/3}$, and follows the similar procedure in the proof of eq. (13) to ensure a successful event overall iteration with at least $1 - \delta$, which gives that

$$\sum_{i=0}^{k} |\xi_H(k)| \leqslant \frac{CN^{2/3}}{\epsilon^{3/2}} \log \left(\frac{8d}{\epsilon \delta}\right). \tag{95}$$

Thus, we have

$$\sum_{i=0}^{k} |\xi_H(k)| = \tilde{O}\left(\frac{N^{3/2}}{\epsilon^{3/2}}\right). \tag{96}$$