

ZERO-ORDER RANDOMIZED SUBSPACE NEWTON METHODS

Erik Berglund, Sarit Khirirat, Xiaoyu Wang

KTH Royal Institute of Technology

ABSTRACT

Zeroth-order methods have become important tools for solving problems where we have access only to function evaluations. However, the zeroth-order methods only using gradient approximations are n times slower than classical first-order methods for solving n -dimensional problems. To accelerate the convergence rate, this paper proposes the zeroth order randomized subspace Newton (ZO-RSN) method, which estimates projections of the gradient and Hessian by random sketching and finite differences. This allows us to compute the Newton step in a lower dimensional subspace, with small computational costs. We prove that ZO-RSN can attain lower iteration complexity than existing zeroth order methods for strongly convex problems. Our numerical experiments show that ZO-RSN can perform black-box attacks under a more restrictive limit on the number of function queries than the state-of-the-art Hessian-aware zeroth-order method.

Index Terms— Zeroth-order optimization, sketching techniques, Newton-type method, adversarial black-box attacks, convolutional neural network.

1. INTRODUCTION

Several applications in machine learning, signal processing and communication networks can often be cast into optimization problems, where gradients are difficult or even infeasible to compute. Popular application examples include optimal hyper-parameter tuning for learning models [1, 2], black-box adversarial attacks on neural network models [3, 4, 5, 6] and sensor selection problems in smart grids or wireless networks [7, 8, 9]. This motivates the study of the zeroth-order methods. A prominent type of zeroth order methods uses function value differences to estimate the gradients [10, Section 3.4]. However, these methods are much slower than classical gradient descent [11], and also suffers from poor performance

particularly for ill-conditioned problems. An alternative way to improve their performance is to incorporate the second order information into zeroth-order methods. However, computing the full Hessian matrix can heavily increase the number of function evaluations and make the Newton step hard to compute, especially for high-dimensional problems. This necessitates us to approximate the Hessian matrix in a lower-dimensional subspace.

Ye *et al.*[12] developed the Hessian-aware zeroth order (ZOHA) methods, which integrate Hessian information into zeroth-order methods. The power-iteration based method ZOHA-PW has a lower query complexity than the gradient-estimating method by [11] when the eigenvalues of the Hessian decay sufficiently quickly. However, the power iteration method requires $O(n)$ function queries per iteration for n -dimensional problems, which is expensive when n is large. To decrease the query cost, they proposed the heuristic methods ZOHA-Gauss-DC and ZOHA-Diag-DC, which estimate the Hessian based on a limited number of random directions. However, no complexity bounds are provided for them.

Another approach to reduce the times of computing Hessian information for high-dimensional problems is to use randomized sketching techniques [13, 14, 15]. These sketching techniques construct lower dimensional sub-problems, which can be solved within small computation times, and enable classical optimization algorithms to have better scalability. For instance, a randomized subspace newton (RSN) method [14] exploits the sketching techniques on the Newton method to solve the problems with very large dimension and to achieve accelerated convergence rate.

In this paper, we propose Hessian-based zeroth-order algorithms using sketching techniques for huge-dimensional problems, called zeroth-order RSN (ZO-RSN). The methods exploit finite differences and sketching to approximate projections of the gradient and Hessian. We provide complexity bounds and prove that under certain conditions ZO-RSN attains lower query complexity than existing zeroth-order algorithms for strongly convex problems. Finally, our experiments with black-box attack problems on a convolutional neural network show that ZO-RSN has an overall competitive performance and higher success rate, compared to the ZOHA-Gauss-DC method in [12].

This work was supported partially by the Swedish Foundation for Strategic Research (SSF) project SoPhy and the Wallenberg Artificial Intelligence, Autonomous Systems and Software Program (WASP) funded by Knut and Alice Wallenberg Foundation. Emails: erberg1@kth.se, sarit@kth.se, wang10@kth.se. © 2022 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

1.1. Notation

For $x \in \mathbb{R}^n$ and $M \succ 0$, $\|x\|_2$ and $\|x\|_\infty$ are the ℓ_2 and ℓ_∞ norm, respectively, and $\|x\|_M^2 = x^T M x$. Given the sketching matrix $S \in \mathbb{R}^{n \times m}$, $s_1, s_2, \dots, s_m \in \mathbb{R}^n$ are its columns. For $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g(x) = \nabla f(x)$ and $H(x) = \nabla^2 f(x)$ are its gradient and Hessian. The function $f(x)$ is L -Lipschitz continuous if there exists a positive constant L such that $\|f(y) - f(x)\|_2 \leq L\|y - x\|_2$ for all $x, y \in \mathbb{R}^n$, and μ -strongly convex if there exists a positive constant μ such that $f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + (\mu/2)\|y - x\|_2^2$ for all $x, y \in \mathbb{R}^n$. We also state that the differentiable function $f(x)$ is L_s -smooth if its gradient $g(x)$ is L_s -Lipschitz continuous. Finally, for any $y \in \mathbb{R}^n$, $\Delta_y f(x) = f(x + y) - f(x)$.

2. PROBLEM FORMULATION

We consider the unconstrained optimization problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f(x), \quad (1)$$

where the dimension n could be very large. Here, $f(x)$ is a three times differentiable and μ -strongly convex function, which is bounded from below and has its minimum value f^* at the point x^* . $g(x)$ and $H(x)$ are also L_1 - and L_2 -Lipschitz continuous. To facilitate the analysis, we further make the following standard assumption on $f(x)$.

Assumption 1 ([14, 16]). *There exists $\hat{L} \geq \hat{\mu} > 0$ such that for any $x, y \in \mathbb{R}^n$:*

$$f(x) \leq f(y) + g(y)^T(x - y) + (\hat{L}/2)\|x - y\|_{H(y)}^2, \quad (2)$$

$$f(x) \geq f(y) + g(y)^T(x - y) + (\hat{\mu}/2)\|x - y\|_{H(y)}^2. \quad (3)$$

Assumption 1 states the smoothness and strong convexity of $f(x)$ under the norm weighted by its Hessian $\|\cdot\|_{H(x)}$. Also, the \hat{L} -relative smoothness and $\hat{\mu}$ -relative convexity exist as a result of the L_1 -smoothness and μ -strong convexity assumption on $f(x)$, as shown below:

Proposition 2.1 ([14, 16]). *A function $f(x)$ is c -stable on a domain D if $\forall y, z \in D$, $\|z - y\|_{H(y)}^2$ and there exists a constant $c \geq 1$ such that $c = \|z - y\|_{H(z)}^2 / \|z - y\|_{H(y)}^2$. If $f(x)$ is μ -strongly convex and L_1 -smooth, then f is (L_1/μ) -stable. Furthermore, if $f(x)$ is c -stable, then Assumption 1 holds with $\hat{L} \leq c$ and $\hat{\mu} \geq 1/c$.*

2.1. RSN Methods

The randomized subspace Newton (RSN) method [14] is a popular inexact Newton method for solving huge-dimensional problems. This method solves an exact Newton system restricted to a random subspace. Given a fixed step-size $\gamma > 0$ and an initial point $x_0 \in \mathbb{R}^d$, the iterate x_k of the RSN method is updated via:

$$x_{k+1} = x_k + \gamma S_k \lambda_k, \quad S_k^T H(x_k) S_k \lambda_k = -S_k^T g(x_k), \quad (4)$$

where $S_k \in \mathbb{R}^{n \times m}$ stores m vectors that span the randomly selected subspace of \mathbb{R}^n . The next lemma characterizes the decrease in the function value from the ZO-RSN method (4).

Lemma 1. *Consider the RSN method (4) for solving Problem (1). If $\gamma \leq 1/\hat{L}$, then*

$$f(x_{k+1}) \leq f(x_k) - (\gamma/2)\|g(x_k)\|_{S_k(S_k^T H(x_k) S_k)^\dagger S_k^T}^2. \quad (5)$$

This descent lemma for the RSN method can be used to prove its linear convergence toward the exact optimum [14]. Furthermore, to implement the RSN method $S_k^T H(x_k) S_k$ and $S_k^T g(x_k)$ are computed efficiently by various sketching techniques such as sub-Gaussian sketches, randomized orthonormal system sketches, random sampling sketches and the Iterative Hessian Sketch [17] as well as the fast Johnson-Lindenstrauss sketch for problems with the appropriate structure [18]. These sketching techniques allow for computing λ_k with very small linear equation systems. If $m \ll n$, then λ_k in Eq. (4) can be solved quickly by inverting $S_k^T H(x_k) S_k \in \mathbb{R}^{m \times m}$.

3. ZERO-ORDER RSN METHODS

In this section, we introduce the zeroth-order randomized subspace Newton (ZO-RSN) method, which builds on the RSN method. The iterate x_k of the ZO-RSN algorithm is updated according to:

$$x_{k+1} = x_k + \gamma S_k \tilde{\lambda}_k, \quad \text{and} \quad \tilde{H}_{S_k}(x_k) \tilde{\lambda}_k = -\tilde{g}_{S_k}(x_k). \quad (6)$$

Here $\tilde{g}_{S_k}(x_k)$ and $\tilde{H}_{S_k}(x_k)$ are approximations of the sketched gradient and Hessian respectively. For a positive scalar α , they can be computed via:

$$[\tilde{g}_{S_k}(x_k)]_i := \Delta_{\alpha s_{i,k}} f(x_k) / \alpha \approx s_{i,k}^T g(x_k),$$

and

$$[\tilde{H}_{S_k}(x_k)]_{i,j} := \Delta_{\alpha s_{i,k}} \Delta_{\alpha s_{j,k}} f(x_k) / \alpha^2 \approx s_{i,k}^T H(x_k) s_{j,k},$$

for all $i = 1, \dots, m$. Similarly to Lemma 1, the ZO-RSN method can be proved to achieve the following bound:

$$f(x_{k+1}) \leq f(x_k) - \frac{\gamma}{2}\|g(x_k)\|_{S_k(S_k^T H(x_k) S_k)^\dagger S_k^T}^2 + O(\alpha). \quad (7)$$

This ensures function value improvement in Eq. (7) if α is sufficiently small and $\tilde{H}_{S_k}(x_k)$ is positive definite. In fact, we can ensure that positive definiteness of $\tilde{H}_{S_k}(x_k)$ follows from α being small enough if we choose S_k appropriately.

Lemma 2. *If $S_k^T S_k = I$ and $\|\tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k\|_2 < \mu$, then $\tilde{H}_{S_k}(x_k) \succ 0$.*

Based on this lemma, we set $S_k^T S_k = I$ to ensure that $\tilde{H}_{S_k}(x_k) \succ 0$. We also require $\mathbb{E}[S_k S_k^T] \succ 0$ so that the approximate sketching does not leave out any directions throughout every iteration. This requirement can be easily satisfied if $s_{1,k}, \dots, s_{m,k}$ are sampled from unit coordinate directions without replacement.

4. THEORETICAL RESULTS

We now provide a complexity bound for ZO-RSN methods.

Theorem 1. *Let the sketching matrix $S_k \in \mathbb{R}^{n \times m}$ satisfy $S_k^T S_k = I$ and $\mathbb{E}_{S_k \sim D}[S_k S_k^T] \succ 0$, and define $G(x) = \mathbb{E}_{S_k \sim D}[S_k (S_k^T H(x) S_k)^{-1} S_k^T]$,*

$$\rho(x) = \min_{v \in \mathbb{R}^n} \frac{v^T H(x)^{\frac{1}{2}} G(x) H(x)^{\frac{1}{2}} v}{\|v\|_2^2} \quad \text{and} \quad \rho = \min_{x \in \mathbb{R}^n} \rho(x).$$

Given $\varepsilon > 0$ and $\delta \in (0, 1)$, consider the ZO-RSN method (6) for Problem (1). If $\gamma \leq 1/\hat{L}$ and $\alpha \leq 0.3\mu/(mL_2)$ is small enough that

$$\frac{\alpha(C_1 + C_2\alpha)}{\rho\hat{\mu}\gamma - \alpha C_1 - \alpha^2 C_3} \leq \delta\varepsilon \quad \text{and} \quad \alpha C_1 + \alpha^2 C_3 < \rho\hat{\mu}\gamma,$$

then we can achieve $\mathbb{E}[f(x_k) - f^] \leq \varepsilon$ after*

$$k \geq \left\lceil \log \left(\frac{f(x_0) - f^*}{(1 - \delta)\varepsilon} \right) / \log \left(\frac{1}{1 - \rho\hat{\mu}\gamma + \alpha C_1 + \alpha^2 C_3} \right) \right\rceil$$

iterations where $C_1 = \gamma(\sqrt{m}L + B)/(2\mu)$, $C_2 = \gamma L_1^2[m + \sqrt{m}(1 + B)]/(2\mu^2)$, $C_3 = \gamma L_1[\sqrt{m}L_1(1 + B) + B(2 + B)]/(2\mu^2)$ and $B = 10mL_2/(3\mu)$.

Theorem 1 establishes a global, linear convergence for the ZO-RSN method toward an ε -accurate solution. The worst-case iteration complexity can be upper bounded as

$$k \geq \lceil \beta_1 \log([f(x_0) - f^*]/[(1 - \delta)\varepsilon]) \rceil. \quad (8)$$

where $\beta_1 = 1/(\rho\hat{\mu}\gamma - \alpha C_1 - \alpha^2 C_3)$. We can recover the convergence complexity for the RSN method [14] if α and δ approach zero. Furthermore, by choosing S_k properly, the iteration complexity for the ZO-RSN method in Eq. (8) can be lower than the complexities for existing zeroth-order methods. We show this with the following corollary:

Corollary 4.1. *Suppose all the conditions of Theorem 1 hold. If the columns of S_k are chosen randomly without replacement from a basis of orthonormal eigenvectors of $H(x_k)$, step-size $\gamma = 1/\hat{L}$, and $\alpha = (\sqrt{C_1^2/4 + (1 - \sigma)\rho\hat{\mu}\gamma} - C_1/2)/C_2$ for some $\sigma \in (0, 1)$, then $\rho = m/n$ and hence to achieve $\mathbb{E}[f(x_k) - f^*] \leq \varepsilon$, we need*

$$k \geq \left\lceil (n\hat{L}/[\sigma m\hat{\mu}]) \log([f(x_0) - f^*]/[(1 - \delta)\varepsilon]) \right\rceil. \quad (9)$$

Corollary 4.1 shows that the iteration complexity of the ZO-RSN methods depends on the subspace dimension m , the

problem dimension n and other parameters $\hat{\mu}, \hat{L}$. Since the ZO-RSN methods need $m(m + 1)/2$ function queries per iteration, we can obtain the total query complexity by multiplying Eq. (9) with this factor.

Now, we compare the complexity bounds for the ZO-RSN methods against the Hessian-aware zeroth-order method using the power iteration (ZOHA-PW) [12], which previously has been compared favourably to the zeroth-order method in [11]. Since the ZOHA-PW method also generates multiple random directions, here m refers to the number of the generated directions. For μ -strongly convex problems, the iteration complexity of ZOHA-PW is

$$k \geq \left\lceil \beta_2 \log \left([f(x_0) - f^*]/[(1 - \delta)\varepsilon] \right) \right\rceil, \quad (10)$$

where $\beta_2 = 64(n + 2)(\mu + 10\lambda_{s+1})/(\mu m)$, λ_{s+1} is an upper bound on the Hessian's $(s + 1)^{\text{th}}$ largest eigenvalue and $\hat{\delta}$ is a free parameter which is similar to δ in Eq. (9). Disregarding the function evaluations required to implement the power method, the total query complexity for ZOHA-PW is $2m$ times its iteration complexity. Consider the problem of minimizing a quadratic function. Then, $\hat{L} = \hat{\mu} = 1$. If $\delta, \hat{\delta}$ and m all are set to be equal for both methods, and also $\sigma = 0.5$, then the speedup in iteration complexity from using ZO-RSN instead of ZOHA-PW is

$$32(1 + 2/n)(1 + 10\lambda_{s+1}/\mu).$$

ZO-RSN is thus faster than ZOHA-PW by more than two orders of magnitude in iteration complexity, even for well-conditioned problems (when λ_{s+1}/μ is close to one). If function queries can be performed efficiently in parallel, then ZO-RSN has significantly lower run-time than ZOHA-PW. We can also prove that the speedup in query complexity for ZO-RSN compared to ZOHA-PW is

$$[128(1 + 2/n)(1 + 10\lambda_{s+1}/\mu)]/(m + 1).$$

Thus, as long as $m < 128(1 + 10\lambda_{s+1}/\mu) - 1$, the query complexity will be lower for ZO-RSN.

5. NUMERICAL EXPERIMENTS

We compare the performance of ZO-RSN against the existing Hessian-aware zeroth methods called ZOHA-Gauss-DC [12] that uses a descent-checking procedure to increase an attack success rate, and approximates Hessian according to

$$\tilde{H} = (2\alpha^2 b)^{-1} \sum_{i=1}^b |\Delta_{\alpha u_i} f(x) - \Delta_{\alpha u_i} f(x - \alpha u_i)| u_i u_i^T + \lambda I_d,$$

where λ is a positive constant and u_1, \dots, u_b are the vectors generated from the Gaussian distribution with zero mean and unit variance. In particular, we evaluate both methods on training un-targeted black box adversarial attacks over the MNIST data set [19, 12]. These attacks are carried out against

the trained convolutional neural network (CNN) model described in [12][Section 5.2]. For each example x_i^{nat} in the test set, the optimizer aims to generate an adversarial example x_i which differs from x_i^{nat} by at most ϵ in ℓ_∞ norm, while being classified differently with sufficient confidence. This is done by minimizing the following function [19]:

$$f(x) = \max \left\{ \max_{i \neq l} \log[Z(x)]_i - \log[Z(x)]_l, -\omega \right\}. \quad (11)$$

Here, $[Z(x)]_i$ represents the probability of an input x belonging to class i according to the trained neural network.

Since the problem is constrained and does not have guarantees for μ -strong convexity or L_1 -smoothness, we need to modify the ZO-RSN algorithm. Firstly, we artificially ensure positive definiteness and boundedness of $\tilde{H}_{S_k}(x_k)$ by applying the operator $\Pi_{[\lambda_{\min}, \lambda_{\max}]}(\cdot)$ that projects its eigenvalues onto an interval $[\lambda_{\min}, \lambda_{\max}]$ to get a modified matrix $\hat{H}_{S_k}(x_k)$. Secondly, we consider ℓ_∞ -norm constraints by determining $\tilde{\lambda}_k$ that solves the following minimization problem

$$\begin{aligned} & \underset{\lambda \in \mathbb{R}^m}{\text{minimize}} && f(x_k) + \gamma \tilde{g}_{S_k}(x_k)^T \lambda + \frac{\gamma}{2} \|\lambda\|_{\tilde{H}_{S_k}(x_k)}^2 \\ & \text{subject to} && -\gamma S_k \lambda \leq x_k - x_i^{nat} - \mathbf{1}\epsilon \\ & && \gamma S_k \lambda \leq \mathbf{1}\epsilon + x_i^{nat} - x_k. \end{aligned} \quad (12)$$

This approach corresponds to using sequential quadratic programming (SQP) for nonlinear problems with linear constraints, but with the step to the next iterate being restricted to lie in a specific subspace. To solve the auxiliary problem (12) quickly with a standard `cvxopt` solver [20], we generate S_k by choosing its columns to be unit coordinate vectors. This enables us to formulate the problem with only m constraints. This adapted ZO-RSN algorithm is called ZO-RSN-SQP. Finally, we use the descent-checking technique corresponding to that for ZOHA-Gauss-DC. The full description of ZO-RSN-SQP is given in Algorithm 1.

We trained the network model until its accuracy reached 98.84%, and also set $\alpha = 0.1, \gamma = 1, m = 3$ and $m_{\max} = 20$ for ZO-RSN-SQP and the same parameters for ZOHA-Gauss-DC for the un-targeted black box attacks described in [12]. In the experiments, we either ended a test run if the algorithm managed to find a point with function value at $\omega = -1$, or if the algorithm called queried the neural network for a prediction 50000 times. We labelled the former result as a success and the latter result as a failure.

The results of our black box attack experiments were summarized in Table 1. Firstly, ZO-RSN-SQP has a more stable performance than ZOHA-Gauss-DC. Even though both algorithms implement the same decent checking technique, only ZO-RSN-SQP succeeds in the attacks for all cases. Secondly, the mean number of queries for ZO-RSN-SQP is lower than that for ZOHA-Gauss-DC. This results from a minority of the problems, where ZOHA-Gauss-DC requires a large number of queries to solve. In contrast, ZOHA-Gauss-DC has a lower median value than ZO-RSN-SQP. As ZO-RSN requires

Algorithm 1 ZO-RSN-SQP for black-box attack

```

Initialize  $x_0 \leftarrow x_i^{nat}, \alpha, \gamma, m, m_{\max}$ 
for  $k = 0, 1, \dots, k_{\max}$  do
    Generate  $S_k = [s_{1,k}, \dots, s_{m,k}]$ 
    Compute  $\tilde{g}_{S_k}$  and  $\tilde{H}_{S_k}$ 
     $\hat{H}_{S_k} \leftarrow \Pi_{[\lambda_{\min}, \lambda_{\max}]}(\tilde{H}_{S_k})$ 
     $\tilde{\lambda}_k \leftarrow$  Solution to (12) with  $\hat{H}_{S_k}(x_k)$  and  $\tilde{g}_{S_k}(x_k)$ 
     $x_{\text{trial}} \leftarrow x_k + S_k \tilde{\lambda}_k$ 
    while  $f(x_{\text{trial}}) \geq f(x_k)$  and  $\bar{m} < m_{\max}$  do
         $\bar{m} \leftarrow \bar{m} + 1$ 
        Generate  $s_{\bar{m},k}$  such that  $[S_k, s_{\bar{m},k}]^T [S_k, s_{\bar{m},k}] = I$ 
         $S_k \leftarrow [s_{1,k}, \dots, s_{\bar{m},k}]$ 
         $[\tilde{g}_{S_k}(x_k)]_{\bar{m}} \leftarrow \Delta_{\alpha s_{i,k}} f(x_k) / \alpha$ 
        for  $j = 1, 2, \dots, \bar{m}$  do
             $[\tilde{H}_{S_k}(x_k)]_{\bar{m},j} \leftarrow \Delta_{\alpha s_{i,k}} \Delta_{\alpha s_{j,k}} f(x_k) / \alpha^2$ 
             $[\hat{H}_{S_k}(x_k)]_{j,\bar{m}} \leftarrow [\tilde{H}_{S_k}(x_k)]_{\bar{m},j}$ 
        end for
         $\hat{H}_{S_k} \leftarrow \Pi_{[\lambda_{\min}, \lambda_{\max}]}(\hat{H}_{S_k})$ 
         $\tilde{\lambda}_k \leftarrow$  Solution to (12) with  $\hat{H}_{S_k}(x_k)$  and  $\tilde{g}_{S_k}(x_k)$ 
         $x_{\text{trial}} \leftarrow x_k + \gamma S_k \tilde{\lambda}_k$ 
    end while
    if  $f(x_{\text{trial}}) \leq f(x_k)$  then
         $x_{k+1} \leftarrow x_{\text{trial}}$ 
    else
         $x_{k+1} \leftarrow x_k$ 
    end if
end for

```

Algorithm	ZO-RSN-SQP	ZOHA-Gauss-DC
Success rate (%)	100	95.33
Median queries	2336	815
Mean queries	2510	4164
Max queries	8239	50000
$f_{\text{est}2000} - f^*$	1.94	$3.70 \cdot 10^{-1}$
$f_{\text{est}4000} - f^*$	$1.89 \cdot 10^{-1}$	$1.86 \cdot 10^{-1}$
$f_{\text{est}6000} - f^*$	$2.42 \cdot 10^{-2}$	$1.47 \cdot 10^{-1}$

Table 1. Comparison of ℓ_∞ norm based black-box attacks on a CNN model trained on the MNIST data.

more function queries per iteration and subspace dimension than ZOHA-Gauss-DC, one can hypothesize this extra effort is worthwhile mainly for the harder-to-attack test examples.

To investigate the speed of convergence, we also ran a separate experiment where we made estimates of the average objective value after 2000, 4000 and 6000 queries, $f_{\text{est}2000}$, $f_{\text{est}4000}$, $f_{\text{est}6000}$, using the first 100 MNIST examples. The suboptimality based on these results are also shown in Table 1. As we can see, ZOHA-Gauss-DC is initially faster, but ZO-RSN-SQP becomes more accurate towards the end.

6. CONCLUSIONS

We have proposed the ZO-RSN method, a Hessian-based zeroth-order method that approximates sketched gradients and Hessians by finite differences. Our results display a lower iteration complexity of the ZO-RSN method than exist-

ing zeroth-order methods for strongly convex problems. The experiments with un-targeted adversarial attacks on a CNN model illustrate that the modified ZO-RSN method named ZO-RSN-SQP attains an overall competitive performance and a higher stability, compared to ZOHA-Gauss-DC.

7. REFERENCES

- [1] Jasper Snoek, Hugo Larochelle, and Ryan P Adams, “Practical bayesian optimization of machine learning algorithms,” *Advances in neural information processing systems*, vol. 25, 2012.
- [2] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl, “Algorithms for hyper-parameter optimization,” *Advances in neural information processing systems*, vol. 24, 2011.
- [3] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh, “Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models,” in *Proceedings of the 10th ACM workshop on artificial intelligence and security*, 2017, pp. 15–26.
- [4] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami, “Practical black-box attacks against machine learning,” in *Proceedings of the 2017 ACM on Asia conference on computer and communications security*, 2017, pp. 506–519.
- [5] Weiwei Hu and Ying Tan, “Generating adversarial malware examples for black-box attacks based on GAN,” *arXiv preprint arXiv:1702.05983*, 2017.
- [6] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin, “Black-box adversarial attacks with limited queries and information,” in *International Conference on Machine Learning*. PMLR, 2018, pp. 2137–2146.
- [7] Sijia Liu, Sundeep Prabhakar Chepuri, Makan Fardad, Engin Maşazade, Geert Leus, and Pramod K Varshney, “Sensor selection for estimation with correlated measurement noise,” *IEEE Transactions on Signal Processing*, vol. 64, no. 13, pp. 3509–3522, 2016.
- [8] Alfred O Hero and Douglas Cochran, “Sensor management: Past, present, and future,” *IEEE Sensors Journal*, vol. 11, no. 12, pp. 3064–3075, 2011.
- [9] Sijia Liu, Jie Chen, Pin-Yu Chen, and Alfred Hero, “Zeroth-order online alternating direction method of multipliers: Convergence analysis and applications,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2018, pp. 288–297.
- [10] Boris Polyak, *Introduction to Optimization*, 07 2020.
- [11] Yurii Nesterov and Vladimir Spokoiny, “Random gradient-free minimization of convex functions,” *Foundations of Computational Mathematics*, vol. 17, no. 2, pp. 527–566, Apr 2017.

- [12] Haishan Ye, Zhichao Huang, Cong Fang, Chris Junchi Li, and Tong Zhang, “Hessian-aware zeroth-order optimization for black-box adversarial attack,” *arXiv preprint arXiv:1812.11377*, 2018.
- [13] Junqi Tang, Mohammad Golbabaee, and Mike E Davies, “Gradient projection iterative sketch for large-scale constrained least-squares,” in *International Conference on Machine Learning*. PMLR, 2017, pp. 3377–3386.
- [14] Robert Gower, Dmitry Kovalev, Felix Lieder, and Peter Richtarik, “RSN: randomized subspace newton,” in *Advances in Neural Information Processing Systems*. 2019, vol. 32, Curran Associates, Inc.
- [15] Mert Pilanci and Martin J Wainwright, “Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence,” *SIAM Journal on Optimization*, vol. 27, no. 1, pp. 205–245, 2017.
- [16] Sai Praneeth Karimireddy, Sebastian U Stich, and Martin Jaggi, “Global linear convergence of Newton’s method without strong-convexity or lipschitz gradients,” *arXiv preprint arXiv:1806.00413*, 2018.
- [17] Mert Pilanci and Martin J. Wainwright, “Iterative Hessian sketch: Fast and accurate solution approximation for constrained least-squares,” *J. Mach. Learn. Res.*, vol. 17, no. 1, pp. 1842–1879, Jan. 2016.
- [18] Nir Ailon and Bernard Chazelle, “The fast johnson–lindenstrauss transform and approximate nearest neighbors,” *SIAM Journal on Computing*, vol. 39, no. 1, pp. 302–322, 2009.
- [19] Nicholas Carlini and David Wagner, “Towards evaluating the robustness of neural networks,” in *2017 IEEE Symposium on Security and Privacy (SP)*. 2017, pp. 39–57, IEEE.
- [20] M Andersen, J Dahl, and L Vandenberghe, “CVXOPT: python software for convex optimization, version 1.2.6,” URL <https://cvxopt.org>, 2021.
- [21] R.A. Horn and C.R. Johnson, *Matrix Analysis*, Matrix Analysis. Cambridge University Press, 2013.

A. PROOF OF LEMMA 1

If $\gamma \leq 1/\hat{L}$, then from Eq. (2) with $x = x_{k+1}, y = x_k$

$$f(x_{k+1}) \leq f(x_k) + g(x_k)^T(x_{k+1} - x_k) + \frac{1}{2\gamma} \|x_{k+1} - x_k\|_{H(x_k)}^2. \quad (13)$$

Utilizing the updates from Eq. (4) that $x_{k+1} - x_k = \gamma S_k \lambda_k$ and $\lambda_k = -(S_k^T H_k S_k)^\dagger S_k^T g(x_k)$, we complete the proof.

B. PROOF OF LEMMA 2

If $S_k^T S_k = I$ and also $v \in \mathbb{R}^m$ has norm 1, then $\mu \leq v^T S_k^T H(x_k) S_k v$. This condition implies that $S_k^T H(x_k) S_k$ is positive definite and its lowest eigenvalue is bounded by μ . Then, $\tilde{H}_{S_k}(x_k) \succ 0$ iff $\tilde{H}_{S_k}(x_k)(S_k^T H(x_k) S_k)^{-1} \succ 0$. Since

$$\begin{aligned} \tilde{H}_{S_k}(x_k)(S_k^T H(x_k) S_k)^{-1} \\ = I + (\tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k)(S_k^T H(x_k) S_k)^{-1}, \end{aligned}$$

positive definiteness of $\tilde{H}_{S_k}(x_k)$ is ensured if

$$\|(\tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k)(S_k^T H(x_k) S_k)^{-1}\|_2 < 1. \quad (14)$$

Since

$$\begin{aligned} & \|(\tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k)(S_k^T H(x_k) S_k)^{-1}\|_2 \\ & \leq \|\tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k\|_2 \|(S_k^T H(x_k) S_k)^{-1}\|_2 \\ & \leq \|\tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k\|_2 / \mu, \end{aligned}$$

a sufficient condition ensuring that Eq. (14) holds is

$$\|(\tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k)\|_2 < \mu.$$

C. LEMMA 3

To facilitate the analysis, we establish key error bounds due to finite difference estimations.

Lemma 3. Consider the ZO-RSN method (6) for solving Problem (1). Let $e_k = \tilde{g}_{S_k}(x_k) - S_k^T g(x_k)$ and $E_k = \tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k$. Then,

$$\|e_k\|_2 \leq \sqrt{m}\alpha L_1/2, \quad \text{and} \quad \|E_k\|_2 \leq 5m\alpha L_2/3.$$

Proof. Define $e_k = \tilde{g}_{S_k}(x_k) - S_k^T g(x_k)$ and $E_k = \tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k$. To prove the upper-bound for $\|e_k\|_2$, consider the first-order Taylor expansion of $f(x_k + \alpha s_{i,k})$ with the error term on Lagrange form: For $\theta \in [0, 1]$

$$\begin{aligned} & f(x_k + \alpha s_{i,k}) \\ & = f(x_k) + \alpha g(x_k)^T s_{i,k} + \frac{\alpha^2}{2} s_{i,k}^T \nabla^2 H(x_k + \theta \alpha s_{i,k}) s_{i,k}. \end{aligned}$$

Therefore,

$$\begin{aligned} |\Delta_{\alpha s_{i,k}} f(x_k)/\alpha - g(x_k)^T s_{i,k}| &= \frac{\alpha}{2} |s_{i,k}^T H(x_k + \theta \alpha s_{i,k}) s_{i,k}| \\ &\leq \alpha \frac{L_1}{2}, \end{aligned}$$

which implies a bound for each component of e_k . We can conclude that $\|e_k\|_2 \leq \sqrt{m} \alpha L_1/2$.

Next, denote the third derivative tensor of $f(x)$ by $f'''(x)$, where

$$[f'''(x)]_{ijk} = \frac{\partial}{\partial[x]_i} \frac{\partial}{\partial[x]_j} \frac{\partial}{\partial[x]_k} f(x).$$

For $u, v, w \in \mathbb{R}^d$, let

$$f'''(x)[u][v][w] = \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n [f'''(x)]_{ijk} [u]_i [v]_j [w]_k.$$

The second-order Taylor expansion of $f(x_k + \alpha u)$ is

$$\begin{aligned} f(x_k + \alpha u) &= f(x_k) + \alpha g(x_k)^T u + \frac{\alpha^2}{2} u^T H(x_k) u \\ &\quad + \frac{\alpha^3}{6} f'''(x_k + \theta \alpha u)[u][u][u], \quad \theta \in [0, 1]. \end{aligned}$$

The L_2 -Lipschitz continuity assumption on $H(x)$ implies that $f'''(x)[u][v][w]$ is bounded by $L_2 \|u\|_2 \|v\|_2 \|w\|_2$. Thus, for $\theta_1, \theta_2, \theta_3 \in [0, 1]$

$$\begin{aligned} |\Delta_{\alpha s_{i,k}} \Delta_{\alpha s_{j,k}} f(x_k)/\alpha^2 - s_{i,k}^T H(x_k) s_{j,k}| &= \frac{\alpha}{6} |f'''(c_k^1)[s_{i,k} + s_{j,k}][s_{i,k} + s_{j,k}][s_{i,k} + s_{j,k}] \\ &\quad - f'''(c_k^2)[s_{i,k}][s_{i,k}][s_{i,k}] - f'''(c_k^3)[s_{j,k}][s_{j,k}][s_{j,k}]| \\ &\leq \frac{5}{3} \alpha L_2, \end{aligned}$$

where $c_k^1 = x_k + \theta_1 \alpha (s_{i,k} + s_{j,k})$, $c_k^2 = x_k + \theta_2 \alpha s_{i,k}$ and $c_k^3 = x_k + \theta_3 \alpha s_{j,k}$. This gives a bound on each component of E_k . To finish the analysis, we invoke the following theorem:

Theorem 2 (Geršgorin theorem, [21]). *Let $A = [a_{ij}] \in M_n$ and let $R'_i(A) = \sum_{j \neq i} |a_{ij}|$, $i \in \{1, \dots, n\}$ denote the deleted absolute row sums of A . Consider the n Geršgorin discs $\{z \in \mathbb{C} : |z - a_i| \leq R'_i(A)\}$. The eigenvalues of A are in the union of the Geršgorin discs.*

Finally, by applying Theorem 2 on $E_k^T E_k$, we have $\rho(E_k^T E_k) \leq (25/9) L_2^2 m^2 \alpha^2$. We can hence conclude that $\|E_k\|_2 = \sqrt{\rho(E_k^T E_k)} \leq (5/3) L_2 m \alpha$. \square

D. PROOF OF THEOREM 1

The RSN method chooses x_{k+1} by minimizing the right hand side of (13) with respect to x and subject to the condition that $x - x_k$ is a linear combination of the columns of S_k . This

constraint can directly be taken into account by the following change of variables to an m -dimensional variable vector λ :

$$x = x_k + \gamma S_k \lambda.$$

Denote $T_k(\lambda)$ as the upper bound of (13), i.e.

$$T_k(\lambda) = f(x_k) + \gamma g(x_k)^T S_k \lambda + \frac{\gamma}{2} \|\lambda\|_{S_k^T H(x_k) S_k}^2.$$

Here, $\lambda_k = -(S_k^T H_k S_k)^\dagger S_k^T g(x_k)$ from Eq. (4) is the λ minimizing $T_k(\lambda)$.

Since the ZO-RSN method only accesses approximations of the sketched gradient and Hessian $\tilde{g}_{S_k}(x_k)$ and $\tilde{H}_{S_k}(x_k)$, it tries to minimize $\tilde{T}_k(\lambda)$, where

$$\tilde{T}_k(\lambda) = f(x_k) + \gamma \tilde{g}_{S_k}(x_k)^T \lambda + \frac{\gamma}{2} \|\lambda\|_{\tilde{H}_{S_k}(x_k)}^2. \quad (15)$$

Let $\tilde{\lambda}_k$ be the minimizer of $\tilde{T}_k(\lambda)$. By setting $x_{k+1} = x_k + \gamma S_k^T \tilde{\lambda}_k$, we get

$$\begin{aligned} f(x_{k+1}) &\leq T(\tilde{\lambda}_k) = T(\lambda_k) + T(\tilde{\lambda}_k) - T(\lambda_k) \\ &= f(x_k) - \frac{\gamma}{2} \|g(x_k)\|_{S_k^T H(x_k) S_k}^2 \\ &\quad + \gamma g(x_k)^T S_k (\tilde{\lambda}_k - \lambda_k) \\ &\quad + \frac{\gamma}{2} (\|\tilde{\lambda}_k\|_{S_k^T H(x_k) S_k}^2 - \|\lambda_k\|_{S_k^T H(x_k) S_k}^2) \quad (16) \\ &= f(x_k) - \frac{\gamma}{2} \|g(x_k)\|_{S_k^T H(x_k) S_k}^2 \\ &\quad + \gamma g(x_k)^T S_k (\tilde{\lambda}_k - \lambda_k) \\ &\quad + \frac{\gamma}{2} (\tilde{\lambda}_k + \lambda_k)^T S_k^T H(x_k) S_k (\tilde{\lambda}_k - \lambda_k). \end{aligned}$$

To complete the proof, we need to determine upper-bounds for $\|\lambda_k\|_2$, $\|\tilde{\lambda}_k - \lambda_k\|_2$ and $\|\lambda_k + \tilde{\lambda}_k\|_2$. We first prove the upper-bound for $\|\lambda_k\|_2$. Since $f(x)$ is L_1 -smooth and μ -strongly convex, $\mu I \preceq S_k^T H(x_k) S_k \preceq L_1 I$. By the fact that $S_k^T H(x_k) S_k \lambda_k = -S_k^T g(x_k)$,

$$\|\lambda_k\|_2 = \|(S_k^T H(x_k) S_k)^{-1} S_k^T g(x_k)\|_2 \leq \|g(x_k)\|_2 / \mu. \quad (17)$$

We next find the upper-bound for $\|\tilde{\lambda}_k - \lambda_k\|_2$. Define $e_k = \tilde{g}_{S_k}(x_k) - S_k^T g(x_k)$ and $E_k = \tilde{H}_{S_k}(x_k) - S_k^T H(x_k) S_k$. If $\alpha \leq 3\mu/(10L_2m)$, then $\|E_k\|_2 \|(S_k^T H(x_k) S_k)^{-1}\|_2 \leq 1/2$. We can then use the following lemma:

Lemma 4 ([21]). *Let $A \in M_n$ be non-singular with condition number $\kappa(A)$, let $b, \Delta b \in \mathbb{R}^n$ and let $\Delta A \in M_n$ be such that $\|A^{-1}\|_2 \|\Delta A\|_2 < 1$. If $x = A^{-1}b$, there exists a Δx such that*

$$(A + \Delta A)(x + \Delta x) = b + \Delta b,$$

and

$$\|\Delta x\|_2 \leq \frac{\|A^{-1}\|_2}{1 - \kappa(A) \frac{\|\Delta A\|_2}{\|A\|_2}} (\|\Delta b\|_2 + \|\Delta A\|_2 \|x\|_2),$$

By Lemma 4, and by the fact that

$$\kappa(S_k^T H(x_k) S_k) / \|S_k^T H(x_k) S_k\|_2 = \|(S_k^T H(x_k) S_k)^{-1}\|_2,$$

we have

$$\begin{aligned} & \|\tilde{\lambda}_k - \lambda_k\|_2 \\ & \leq \frac{\|(S_k^T H(x_k) S_k)^{-1}\|_2}{1 - \kappa(S_k^T H(x_k) S_k) \frac{\|E_k\|_2}{\|S_k^T H(x_k) S_k\|_2}} (\|e_k\|_2 + \|E_k\|_2 \|\lambda_k\|_2) \\ & \leq \frac{1}{\mu} (\sqrt{m} L_1 + \frac{10}{3} m L_2 \|\lambda_k\|_2) \alpha \\ & \leq \frac{1}{\mu} (\sqrt{m} L_1 + m \frac{10 L_2}{3 \mu} \|g(x_k)\|_2) \alpha. \end{aligned} \quad (18)$$

We finally can prove the upper-bound for $\|\lambda_k + \tilde{\lambda}_k\|_2$:

$$\begin{aligned} \|\lambda_k + \tilde{\lambda}_k\|_2 & \leq 2\|\lambda_k\|_2 + \|\tilde{\lambda}_k - \lambda_k\|_2 \\ & \leq \frac{1}{\mu} \left(\sqrt{m} L_1 + \left(2 + m \frac{10 L_2}{3 \mu} \right) \|g(x_k)\|_2 \right) \alpha. \end{aligned} \quad (19)$$

Next, plugging in inequalities (17), (18) and (19) into (16), and then using the fact that $\|g(x_k)\|_2 \leq (1 + \|g(x_k)\|_2^2)/2$

$$\begin{aligned} f(x_{k+1}) & \leq f(x_k) - \frac{\gamma}{2} \|g(x_k)\|_{S_k(S_k^T H(x_k) S_k)^{-1} S_k^T}^2 \\ & \quad + \alpha(C_1 + C_2 \alpha + \|g(x_k)\|_2^2 (C_1 + C_3, \alpha)) \end{aligned} \quad (20)$$

where $C_1 = \gamma(\sqrt{m}L + B)/(2\mu)$, $C_2 = \gamma L_1^2[m + \sqrt{m}(1 + B)]/(2\mu^2)$, $C_3 = \gamma L_1[\sqrt{m}L_1(1 + B) + B(2 + B)]/(2\mu^2)$ and $B = 10mL_2/(3\mu)$. Taking the expectation with respect to x_k on both sides of Inequality (20), we have

$$\begin{aligned} \mathbb{E}[f(x_{k+1})|x_k] & \leq f(x_k) - \frac{\gamma}{2} \|g(x_k)\|_{G(x_k)}^2 \\ & \quad + \alpha(C_1 + C_2 \alpha + \|g(x_k)\|_2^2 (C_1 + C_3 \alpha)), \end{aligned} \quad (21)$$

where $G(x) = \mathbb{E}_{S_k \sim \mathcal{D}}[S_k(S_k^T H(x) S_k)^{-1} S_k^T]$.

To prove the linear convergence of the ZO-RSN method from Eq. (21), we need to bound $\|g(x_k)\|_2^2$ and $\|g(x_k)\|_{G(x_k)}^2$. We first prove the upper bound for $\|g(x_k)\|_2$ by the L_1 -smoothness assumption of $f(x)$, i.e.

$$f(x) \leq f(y) + g(y)^T(x - y) + \frac{L_1}{2} \|x - y\|_2^2.$$

Setting $y = x_k$ and minimizing both sides with respect to x separately results in

$$\|g(x_k)\|_2^2 \leq 2L_1(f(x_k) - f^*). \quad (22)$$

We next show the lower bound for $\|g(x_k)\|_{G(x_k)}^2$. If $H(x_k)$ is non-singular, then

$$\begin{aligned} \|g(x_k)\|_{G(x_k)}^2 & = \|H(x_k)^{\frac{1}{2}} H(x_k)^{-\frac{1}{2}} g(x_k)\|_{G(x_k)}^2 \\ & \geq \rho \|g(x_k)\|_{H(x_k)^{-1}}^2. \end{aligned}$$

Setting $y = x_k$ in (3) and minimizing both sides of the equation with respect to x separately gives

$$f^* \geq f(x_k) - \frac{1}{2\hat{\mu}} \|g(x_k)\|_{H(x_k)^{-1}}^2.$$

Therefore,

$$2\rho\hat{\mu}(f(x_k) - f^*) \leq \rho \|g(x_k)\|_{H(x_k)^{-1}}^2 \leq \|g(x_k)\|_{G(x_k)}^2. \quad (23)$$

Next, by plugging (22) and (23) into (21), then by subtracting f^* from both sides of the inequality, and after that by taking the total expectation, we get

$$V_{k+1} \leq [1 - \rho\hat{\mu}\gamma + \alpha(C_1 + \alpha C_3)]V_k + \alpha(C_1 + C_2\alpha).$$

where $V_k = \mathbb{E}[f(x_k) - f^*]$.

If α satisfies $\alpha C_1 + \alpha^2 C_3 < \rho\hat{\mu}\gamma$, then by applying the inequality recursively and by using the fact $\sum_{l=0}^{k-1} \beta^l \leq \sum_{l=0}^{\infty} \beta^l = 1/(1 - \beta)$ for $\beta \in (0, 1)$

$$V_k \leq (1 - \rho\hat{\mu}\gamma + \alpha(C_1 + \alpha C_3))^k V_0 + \frac{\alpha(C_1 + C_2\alpha)}{\rho\hat{\mu}\gamma - \alpha C_1 - \alpha^2 C_3}. \quad (24)$$

If α also satisfies

$$\frac{\alpha(C_1 + C_2\alpha)}{\rho\hat{\mu}\gamma - \alpha C_1 - \alpha^2 C_3} \leq \delta\varepsilon, \quad (25)$$

where $\delta \in (0, 1)$ and ε is an expected sub-optimality, then the lower bound on the number of iterations follows.

E. PROOF OF COROLLARY 4.1

To prove the result, we need to quantify ρ . This can be done by using the following lemma:

Lemma 5. [14] *If for all $x_k \in \mathbb{R}^n$ it holds with probability 1 that $\text{Null}(S_k^T H(x_k) S_k) = \text{Null}(S_k)$ and $\text{Range}(H(x_k)) \subset \text{Range}(\mathbb{E}_{S_k \sim \mathcal{D}}[S_k^T S_k])$, then $\rho(x_k) = \lambda_{\min}^+(\mathbb{E}_{S_k \sim \mathcal{D}}[\hat{P}(x_k)])$ which is positive, where*

$$\hat{P}(x_k) := H^{1/2}(x_k) S_k (S_k^T H(x_k) S_k)^{\dagger} S_k^T H^{1/2}(x_k). \quad (26)$$

From this lemma, we can quantify ρ by considering the cases when the columns of S_k are chosen randomly without replacement from a basis of orthonormal eigenvectors of $H(x_k)$. Let $\tilde{\Lambda}_{S_k}$ be an $m \times m$ diagonal matrix such that the eigenvalue corresponding to column i of S_k is the i^{th} element on the diagonal of $\tilde{\Lambda}_{S_k}$ and let its square root be $\tilde{\Lambda}_{S_k}^{\frac{1}{2}}$. Then,

$$\begin{aligned} \hat{P}(x_k) & = H^{1/2}(x_k) S_k (S_k^T H(x_k) S_k)^{-1} S_k^T H^{1/2}(x_k) \\ & = S_k \tilde{\Lambda}_{S_k}^{\frac{1}{2}} (\tilde{\Lambda}_{S_k}^{\frac{1}{2}} S_k^T S_k \tilde{\Lambda}_{S_k}^{\frac{1}{2}})^{-1} \tilde{\Lambda}_{S_k}^{\frac{1}{2}} S_k^T = S_k S_k^T. \end{aligned}$$

The eigenvectors of all realizations of $\hat{P}(x_k)$ are the eigenvectors of $H(x_k)$, with eigenvalues 1 for each vector that is

among the columns of S_k and eigenvalues 0 for the other vectors. Therefore, the orthonormal eigenvectors of $H(x_k)$ are also the eigenvectors of $\mathbb{E}_{S_k \sim \mathcal{D}}[\hat{P}(x_k)]$. Since the probability that this eigenvector is among the columns of S_k is m/n , $v^T \mathbb{E}_{S_k \sim \mathcal{D}}[\hat{P}(x_k)]v = m/n$ for any eigenvector v . Thus, we can prove that $\rho = m/n$.

From Theorem 1, the iteration complexity bound can be approximated in Eq.(8). If we choose $\gamma = 1/\hat{L}$ and some $\sigma \in (0, 1)$ such that $\alpha = (\sqrt{C_1^2/4 + (1 - \sigma)\rho\hat{\mu}\gamma} - C_1/2)/C_2$, then

$$\rho\hat{\mu}\gamma - \alpha C_1 - \alpha^2 C_3 = \sigma m\hat{\mu}/(n\hat{L}).$$

Plugging this expression into Eq.(8), we complete the proof.