

OverSketched Newton: Fast Convex Optimization for Serverless Systems

Vipul Gupta¹, Swanand Kadhe¹, Thomas Courtade¹, Michael W. Mahoney², and Kannan Ramchandran¹

¹Department of EECS, University of California, Berkeley

²ICSI and Statistics Department, University of California, Berkeley

Abstract

Motivated by recent developments in serverless systems for large-scale computation as well as improvements in scalable randomized matrix algorithms, we develop OverSketched Newton, a randomized Hessian-based optimization algorithm to solve large-scale convex optimization problems in serverless systems. OverSketched Newton leverages matrix sketching ideas from Randomized Numerical Linear Algebra to compute the Hessian approximately. These sketching methods lead to inbuilt resiliency against stragglers that are a characteristic of serverless architectures. Depending on whether the problem is strongly convex or not, we propose different iteration updates using the approximate Hessian. For both cases, we establish convergence guarantees for OverSketched Newton and empirically validate our results by solving large-scale supervised learning problems on real-world datasets. Experiments demonstrate a reduction of $\sim 50\%$ in total running time on AWS Lambda, compared to state-of-the-art distributed optimization schemes.

1 Introduction

In recent years, there has been tremendous growth in users performing distributed computing operations on the cloud, largely due to extensive and inexpensive commercial offerings like Amazon Web Services (AWS), Google Cloud, Microsoft Azure, etc. Serverless platforms—such as AWS Lambda, Cloud functions and Azure Functions—penetrate a large user base by provisioning and managing the servers on which the computation is performed. These platforms abstract away the need for maintaining servers, since this is done by the cloud provider and is hidden from the user—hence the name *serverless*. Moreover, allocation of these servers is done expeditiously which provides greater elasticity and easy scalability. For example, up to ten thousand machines can be allocated on AWS Lambda in less than ten seconds [1–4].

The use of serverless systems is gaining significant research traction, primarily due to its massive scalability and convenience in operation. It is forecasted that the market share of serverless will grow by USD 9.16 billion during 2019-2023 (at a CAGR of 11%) [5]. Indeed,

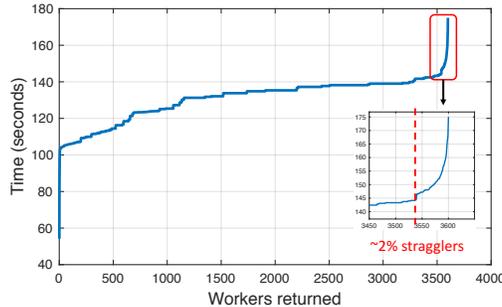


Figure 1: Average job times for 3600 AWS Lambda nodes over 10 trials for distributed matrix multiplication. The median job time is around 135 seconds, and around 2% of the nodes take up to 180 seconds on average.

according to the *Berkeley view on Serverless Computing* [6], serverless systems are expected to dominate the cloud scenario and become the default computing paradigm in the coming years while client-server based computing will witness a considerable decline. For these reasons, using serverless systems for large-scale computation has garnered significant attention from the systems community [3, 4, 7–12].

Due to several crucial differences between the traditional High Performance Computing (HPC) / serverful and serverless architectures, existing distributed algorithms cannot, in general, be extended to serverless computing. First, unlike *serverful* computing, the number of inexpensive workers in serverless platforms is flexible, often scaling into the thousands [3, 4]. This heavy gain in the computation power, however, comes with the disadvantage that the commodity workers in serverless architecture are ephemeral and have low memory.¹ The ephemeral nature of the workers in serverless systems requires that new workers should be invoked every few iterations and data should be communicated to them. Moreover, the workers do not communicate amongst themselves, and instead they read/write data directly from/to a single high-latency data storage entity (e.g., cloud storage like AWS S3 [3]).

Second, unlike HPC/serverful systems, nodes in the serverless systems suffer degradation due to what is known as *system noise*. This can be a result of limited availability of shared resources, hardware failure, network latency, etc. [13, 14]. This results in job time variability, and hence a subset of much slower nodes, often called *stragglers*. These stragglers significantly slow the overall computation time, especially in large or iterative jobs. In Fig. 1, we plot the running times for a distributed matrix multiplication job with 3600 workers on AWS Lambda and demonstrate the effect of stragglers on the total job time. In fact, our experiments consistently demonstrate that at least 2% workers take significantly longer than the median job time, severely degrading the overall efficiency of the system.

Due to these issues, first-order methods, e.g., gradient descent and Nesterov Accelerated Gradient (NAG) methods, tend to perform poorly on distributed serverless architectures [15].

¹For example, serverless nodes in AWS Lambda, Google Cloud Functions and Microsoft Azure Functions have a maximum memory of 3 GB, 2 GB and 1.5 GB, respectively, and a maximum runtime of 900 seconds, 540 seconds and 300 seconds, respectively (these numbers may change over time).

Their slower convergence is made worse on serverless platforms due to persistent stragglers. The straggler effect incurs heavy slowdown due to the accumulation of tail times as a result of a subset of slow workers occurring in each iteration.

Compared to first-order optimization algorithms, second-order methods—which use the gradient as well as Hessian information—enjoy superior convergence rates. For instance, Newton’s method enjoys quadratic convergence for strongly convex and smooth problems, compared to the linear convergence of gradient descent [16]. Moreover, second-order methods do not require step-size tuning and unit step-size provably works for most problems. These methods have a long history in optimization and scientific computing (see, e.g., [16]), but they are less common in machine learning and data science. This is partly since stochastic first order methods suffice for downstream problems [17] and partly since naive implementations of second order methods can perform poorly [18]. However, recent theoretical work has addressed many of these issues [19–23], and recent implementations have shown that high-quality implementations of second order stochastic optimization algorithms can beat state-of-the-art in machine learning applications [24–28] in traditional systems.

1.1 Main Contributions

In this paper, we argue that second-order methods are highly compatible with serverless systems that provide extensive computing power by invoking thousands of workers but are limited by the communication costs and hence the number of iterations; and, to address the challenges of ephemeral workers and stragglers in serverless systems, we propose and analyze a randomized and distributed second-order optimization algorithm, called *OverSketched Newton*. OverSketched Newton uses the technique of matrix sketching from Sub-Sampled Newton (SSN) methods [19–22], which are based on sketching methods from Randomized Numerical Linear Algebra (RandNLA) [29–31], to obtain a good approximation for the Hessian, instead of calculating the full Hessian.

OverSketched Newton has two key components. For straggler-resilient Hessian calculation in serverless systems, we use the sparse sketching based randomized matrix multiplication method from [32]. For straggler mitigation during gradient calculation, we use the recently proposed technique based on error-correcting codes to create redundant computation [33–35]. We prove that, for strongly convex functions, the local convergence rate of OverSketched Newton is linear-quadratic, while its global convergence rate is linear. Then, going beyond the usual strong convexity assumption for second-order methods, we adapt OverSketched Newton using ideas from [22]. For such functions, we prove that a linear convergence rate can be guaranteed with OverSketched Newton. To the best of our knowledge, this is the first work to prove convergence guarantees for weakly-convex problems when the Hessian is computed approximately using ideas from RandNLA.

We extensively evaluate OverSketched Newton on AWS Lambda using several real-world datasets obtained from the LIBSVM repository [36], and we compare OverSketched Newton with several first-order (gradient descent, Nesterov’s method, etc.) and second-order (exact Newton’s method [16], GIANT [24], etc.) baselines for distributed optimization.

We further evaluate and compare different techniques for straggler mitigation, such as speculative execution, coded computing [33, 34], randomization-based sketching [32] and gradient coding [37]. We demonstrate that OverSketched Newton is at least 9x and 2x faster than state-of-the-art first-order and second-order schemes, respectively, in terms of end-to-end training time on AWS Lambda. Moreover, we show that OverSketched Newton on serverless systems outperforms existing distributed optimization algorithms in serverful systems by at least 30%.

1.2 Related Work

Our results tie together three quite different lines of work, each of which we review here briefly.

Existing Straggler Mitigation Schemes: Strategies like speculative execution have been traditionally used to mitigate stragglers in popular distributed computing frameworks like Hadoop MapReduce [38] and Apache Spark [39]. Speculative execution works by detecting workers that are running slower than expected and then allocating their tasks to new workers without shutting down the original straggling task. The worker that finishes first communicates its results. This has several drawbacks, e.g. constant monitoring of tasks is required and late stragglers can still hurt the efficiency.

Recently, many coding-theoretic ideas have been proposed to introduce redundancy into the distributed computation for straggler mitigation (e.g. see [33–35, 37, 40, 41]). The idea of coded computation is to generate redundant copies of the result of distributed computation by encoding the input data using error-correcting-codes. These redundant copies are then used to decode the output of the missing stragglers. Our algorithm to compute gradients in a distributed straggler-resilient manner uses codes to mitigate stragglers, and we compare our performance with speculative execution.

Approximate Second-order Methods: In many machine learning applications, where the data itself is noisy, using the exact Hessian is not necessary. Indeed, using ideas from RandNLA, one can prove convergence guarantees for SSN methods on a single machine, when the Hessian is computed approximately [19–21, 23]. To accomplish this, many sketching schemes can be used (sub-Gaussian, Hadamard, random row sampling, sparse Johnson-Lindenstrauss, etc. [29, 30]), but these methods cannot tolerate stragglers, and thus they do not perform well in serverless environments.

This motivates the use of the *OverSketch* sketch from our recent work in [32]. OverSketch has many nice properties, like subspace embedding, sparsity, input obliviousness, and amenability to distributed implementation. To the best of our knowledge, this is the first work to prove and evaluate convergence guarantees for algorithms based on OverSketch. Our guarantees take into account the amount of communication at each worker and the number of stragglers, both of which are a property of distributed systems.

There has also been a growing research interest in designing and analyzing distributed implementations of stochastic second-order methods [24, 42–46]. However, these implemen-

tations are tailored for serverful distributed systems. Our focus, on the other hand, is on serverless systems.

Distributed Optimization on Serverless Systems: Optimization over the serverless framework has garnered significant interest from the research community. However, these works either evaluate and benchmark existing algorithms (e.g., see [9–11]) or focus on designing new systems frameworks for faster optimization (e.g., see [12]) on serverless. To the best of our knowledge, this is the first work that proposes a large-scale distributed optimization algorithm that specifically caters to *serverless architectures* with *provable convergence guarantees*. We exploit the advantages offered by serverless systems while mitigating the drawbacks such as stragglers and additional overhead per invocation of workers.

2 Newton’s Method: An Overview

We are interested in solving on serverless systems in a distributed and straggler-resilient manner problems of the form:

$$f(\mathbf{w}^*) = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}), \quad (1)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a closed and convex function bounded from below. In the Newton’s method, the update at the $(t+1)$ -th iteration is obtained by minimizing the Taylor’s expansion of the objective function $f(\cdot)$ at \mathbf{w}_t , that is

$$\begin{aligned} \mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^T (\mathbf{w} - \mathbf{w}_t) \right. \\ \left. + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \nabla^2 f(\mathbf{w}_t) (\mathbf{w} - \mathbf{w}_t) \right\}. \end{aligned} \quad (2)$$

For strongly convex $f(\cdot)$, that is, when $\nabla^2 f(\cdot)$ is invertible, Eq. (2) becomes $\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{H}_t^{-1} \nabla f(\mathbf{w}_t)$, where $\mathbf{H}_t = \nabla^2 f(\mathbf{w}_t)$ is the Hessian matrix at the t -th iteration. Given a good initialization and assuming that the Hessian is Lipschitz, the Newton’s method satisfies the update $\|\mathbf{w}_{t+1} - \mathbf{w}^*\|_2 \leq c \|\mathbf{w}_t - \mathbf{w}^*\|_2^2$, for some constant $c > 0$, implying quadratic convergence [16].

One shortcoming for the classical Newton’s method is that it works only for strongly convex objective functions. In particular, if f is weakly-convex², that is, if the Hessian matrix is not positive definite, then the objective function in (2) may be unbounded from below. To address this shortcoming, authors in [22] recently proposed a variant of Newton’s method, called Newton-Minimum-Residual (Newton-MR). Instead of (1), Newton-MR considers the following auxiliary optimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \|\nabla f(\mathbf{w})\|^2.$$

Note that the minimizers of this auxiliary problem and (1) are the same when $f(\cdot)$ is convex. Then, the update direction in the $(t+1)$ -th iteration is obtained by minimizing the Taylor’s

²For the sake of clarity, we call a convex function weakly-convex if it is not strongly convex.

expansion of $\|\nabla f(\mathbf{w}_t + \mathbf{p})\|^2$, that is,

$$\mathbf{p}_t = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\nabla f(\mathbf{w}_t) + \mathbf{H}_t \mathbf{p}\|^2.$$

The general solution of the above problem is given by $\mathbf{p} = -[\mathbf{H}_t]^\dagger \nabla f(\mathbf{w}_t) + (\mathbf{I} - \mathbf{H}_t [\mathbf{H}_t]^\dagger) \mathbf{q}$, $\forall \mathbf{q} \in \mathbb{R}^d$, where $[\cdot]^\dagger$ is the Moore-Penrose inverse. Among these, the minimum norm solution is chosen, which gives the update direction in the t -th iteration as $\mathbf{p}_t = -\mathbf{H}_t^\dagger \nabla f(\mathbf{w}_t)$. Thus, the model update is

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{p}_t = \mathbf{w}_t - [\nabla^2 f(\mathbf{w}_t)]^\dagger \nabla f(\mathbf{w}_t). \quad (3)$$

OverSketched Newton considers both of these variants.

3 OverSketched Newton

In many applications like machine learning where the training data itself is noisy, using the exact Hessian is not necessary. Indeed, many results in the literature prove convergence guarantees for Newton’s method when the Hessian is computed approximately using ideas from RandNLA for a single machine (e.g. [19, 20, 23, 47]). In particular, these methods perform a form of dimensionality reduction for the Hessian using random matrices, called sketching matrices. Many popular sketching schemes have been proposed in the literature, for example, sub-Gaussian, Hadamard, random row sampling, sparse Johnson-Lindenstrauss, etc. [29, 30]. Inspired from these works, we present OverSketched Newton, a stochastic second order algorithm for solving—*on serverless systems, in a distributed, straggler-resilient manner*—problems of the form (1).

Distributed straggler-resilient gradient computation: OverSketched Newton computes the full gradient in each iteration by using tools from error-correcting codes [33, 34]. Our key observation is that, for several commonly encountered optimization problems, gradient computation relies on matrix-vector multiplications (see Sec. 4 for examples). We leverage coded matrix multiplication technique from [34] to perform the large-scale matrix-vector multiplication in a distributed straggler-resilient manner. The idea of coded matrix multiplication is explained in Fig. 2; detailed steps are provided in Algorithm 1.

Distributed straggler-resilient approximate Hessian computation: For several commonly encountered optimization problems, Hessian computation involves matrix-matrix multiplication for a pair of large matrices (see Sec. 4 for several examples). For computing the large-scale matrix-matrix multiplication in parallel in serverless systems, we propose to use a straggler-resilient scheme called *OverSketch* from [32]. OverSketch does *blocked partitioning* of input matrices where each worker works on square blocks of dimension b . Hence, it is more communication efficient than existing coding-based straggler mitigation schemes that do naïve row-column partition of input matrices [40, 48]. We note that it is well known in HPC that blocked partitioning of input matrices can lead to communication-efficient methods for distributed multiplication [32, 49, 50].

Algorithm 1: Straggler-resilient distributed computation of \mathbf{Ax} using codes

Input : Matrix $\mathbf{A} \in \mathbb{R}^{t \times s}$, vector $x \in \mathbb{R}^s$, and block size parameter b

Result: $\mathbf{y} = \mathbf{Ax}$, where $\mathbf{y} \in \mathbb{R}^s$ is the product of matrix \mathbf{A} and vector \mathbf{x}

- 1 **Initialization:** Divide \mathbf{A} into $T = t/b$ row-blocks, each of dimension $b \times s$
 - 2 **Encoding:** Generate coded \mathbf{A} , say \mathbf{A}_c , in parallel using a 2D product code by arranging the row blocks of \mathbf{A} in a 2D structure of dimension $\sqrt{T} \times \sqrt{T}$ and adding blocks across rows and columns to generate parities; see Fig. 2 in [34] for an illustration
 - 3 **for** $i = 1$ **to** $T + 2\sqrt{T} + 1$ **do**
 - 4 1. Worker W_i receives the i -th row-block of \mathbf{A}_c , say $\mathbf{A}_c(i, :)$, and \mathbf{x} from cloud storage
 - 5 2. W_i computes $\mathbf{y}(i) = \mathbf{A}(i, :)\mathbf{x}$
 - 6 3. Master receives $\mathbf{y}(i)$ from worker W_i
 - 7 **end**
 - 8 **Decoding:** Master checks if it has received results from enough workers to reconstruct \mathbf{y} . Once it does, it decodes \mathbf{y} from available results using the peeling decoder
-

OverSketch uses a sparse sketching matrix based on Count-Sketch [29]. It has similar computational efficiency and accuracy guarantees as that of the Count-Sketch, with two additional properties: it is amenable to distributed implementation; and it is resilient to stragglers. More specifically, the OverSketch matrix is constructed as follows [32].

Recall that the Hessian $\nabla^2 f(\cdot) \in \mathbb{R}^{d \times d}$. First choose the desired sketch dimension m (which depends on d), block-size b (which depends on the memory of the workers), and straggler tolerance $\zeta > 0$ (which depends on the distributed system). Then, define $N = m/b$ and $e = \zeta N$, for some constant $\zeta > 0$. Here ζ is the fraction of stragglers that we want our algorithm to tolerate. Thus, e is the maximum number of stragglers per $N + e$ workers that can be tolerated. The sketch \mathbf{S} is then given by

$$\mathbf{S} = \frac{1}{\sqrt{N}}(\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_{N+e}), \quad (4)$$

where $\mathbf{S}_i \in \mathbb{R}^{n \times b}$, for all $i \in [1, N+e]$, are i.i.d. Count-Sketch matrices³ with sketch dimension b . Note that $\mathbf{S} \in \mathbb{R}^{n \times (m+eb)}$, where $m = Nb$ is the required sketch dimension and e is the over-provisioning parameter to provide resiliency against e stragglers per $N + e$ workers. We leverage the straggler resiliency of OverSketch to obtain the sketched Hessian in a distributed straggler-resilient manner. An illustration of OverSketch is provided in Fig. 3; see Algorithm 2 for details.

Model update: Let $\hat{\mathbf{H}}_t = \mathbf{A}_t^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A}_t$, where \mathbf{A}_t is the square root of the Hessian

³Each of the Count-Sketch matrices \mathbf{S}_i is constructed (independently of others) as follows. First, for every row j , $j \in [n]$, of \mathbf{S}_i , independently choose a column $h(j) \in [b]$. Then, select a uniformly random element from $\{-1, +1\}$, denoted as $\sigma(i)$. Finally, set $\mathbf{S}_i(j, h(j)) = \sigma(i)$ and set $\mathbf{S}_i(j, l) = 0$ for all $l \neq h(j)$. (See [29, 32] for details.)

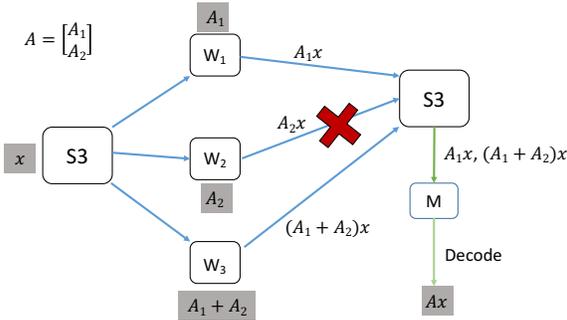


Figure 2: **Coded matrix-vector multiplication:** Matrix \mathbf{A} is divided into 2 row chunks \mathbf{A}_1 and \mathbf{A}_2 . During encoding, redundant chunk $\mathbf{A}_1 + \mathbf{A}_2$ is created. Three workers obtain $\mathbf{A}_1, \mathbf{A}_2$ and $\mathbf{A}_1 + \mathbf{A}_2$ from the cloud storage S3, respectively, and then multiply by \mathbf{x} and write back the result to the cloud. The master M can decode $\mathbf{A}\mathbf{x}$ from the results of any two workers, thus being resilient to one straggler (W_2 in this case).

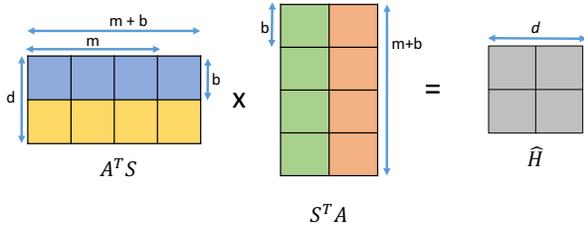


Figure 3: **OverSketch-based approximate Hessian computation:** First, the matrix \mathbf{A} —satisfying $\mathbf{A}^T \mathbf{A} = \nabla^2 f(\mathbf{w}_t)$ —is sketched in parallel using the sketch in (4). Then, each worker receives a block of each of the sketched matrices $\mathbf{A}^T \mathbf{S}$ and $\mathbf{S}^T \mathbf{A}$, multiplies them, and communicates back its results for reduction. During reduction, stragglers can be ignored by the virtue of “over” sketching. For example, here the desired sketch dimension m is increased by block-size b for obtaining resiliency against one straggler for each block of $\hat{\mathbf{H}}$.

$\nabla^2 f(\mathbf{w}_t)$, and \mathbf{S}_t is an independent realization of (4) at the t -th iteration. For strongly-convex functions, the update direction is $\mathbf{p}_t = -\hat{\mathbf{H}}_t^{-1} \nabla f(\mathbf{w}_t)$. We use line-search to choose the step-size, that is, find

$$\alpha_t = \max_{\alpha \leq 1} \alpha \quad \text{such that} \quad f(\mathbf{w}_t + \alpha \mathbf{p}_t) \leq f(\mathbf{w}_t) + \alpha \beta \mathbf{p}_t^T \nabla f(\mathbf{w}_t), \quad (5)$$

for some constant $\beta \in (0, 1/2]$. For weakly-convex functions, the update direction (inspired by Newton-MR [22]) is $\mathbf{p}_t = -\hat{\mathbf{H}}_t^\dagger \nabla f(\mathbf{w}_t)$, where $\hat{\mathbf{H}}_t^\dagger$ is the Moore-Penrose inverse of $\hat{\mathbf{H}}_t$. To find the update \mathbf{w}_{t+1} , we find the right step-size α_t using line-search in (5), but with $f(\cdot)$ replaced by $\|\nabla f(\cdot)\|^2$ and $\nabla f(\mathbf{w}_t)$ replaced by $2\hat{\mathbf{H}}_t \nabla f(\mathbf{w}_t)$, according to the objective in $\|\nabla f(\cdot)\|^2$. More specifically, for some constant $\beta \in (0, 1/2]$,

$$\alpha_t = \max_{\alpha \leq 1} \alpha \quad \text{such that} \quad \|\nabla f(\mathbf{w}_t + \alpha \mathbf{p}_t)\|^2 \leq \|\nabla f(\mathbf{w}_t)\|^2 + 2\alpha \beta \mathbf{p}_t^T \hat{\mathbf{H}}_t \nabla f(\mathbf{w}_t). \quad (6)$$

Note that for OverSketched Newton, we use $\hat{\mathbf{H}}_t$ in the line-search since the exact Hessian is not available. The update in the t -th iteration in both cases is given by

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{p}_t.$$

Note that (5) line-search can be solved approximately in single machine systems using Armijo backtracking line search [16, 51]. OverSketched Newton is concisely described in Algorithm 3. In Section 3.2, we describe how to implement distributed line-search in serverless systems when the data is stored in the cloud. Next, we prove convergence guarantees for OverSketched Newton that uses the sketch matrix in (4) and full gradient for approximate Hessian computation.

Algorithm 2: Approximate Hessian calculation on serverless systems using OverSketch

Input : Matrices $\mathbf{A} \in \mathbb{R}^{n \times d}$, required sketch dimension m , straggler tolerance e , block-size b . Define $N = m/b$

Result: $\hat{\mathbf{H}} \approx \mathbf{A}^T \times \mathbf{A}$

- 1 **Sketching:** Use sketch in Eq. (4) to obtain $\tilde{\mathbf{A}} = \mathbf{S}^T \mathbf{A}$ distributedly (see Algorithm 5 in [32] for details)
 - 2 **Block partitioning:** Divide $\tilde{\mathbf{A}}$ into $(N + e) \times d/b$ matrix of $b \times b$ blocks
 - 3 **Computation phase:** Each worker takes a block of $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{A}}^T$ each and multiplies them. This step invokes $(N + e)d^2/b^2$ workers, where $N + e$ workers compute one block of $\hat{\mathbf{H}}$
 - 4 **Termination:** Stop computation when any N out of $N + e$ workers return their results for each block of $\hat{\mathbf{H}}$
 - 5 **Reduction phase:** Invoke d^2/b^2 workers to aggregate results during the computation phase, where each worker will calculate one block of $\hat{\mathbf{H}}$
-

3.1 Convergence Guarantees

First, we focus our attention to strongly convex functions. We consider the following assumptions. We note that these assumptions are standard for analyzing approximate Newton methods, (e.g., see [19, 20, 23]).

Assumptions:

1. f is twice-differentiable;
2. f is k -strongly convex ($k > 0$), that is,

$$\nabla^2 f(\mathbf{w}) \succeq k\mathbf{I};$$

3. f is M -smooth ($k \leq M < \infty$), that is,

$$\nabla^2 f(\mathbf{w}) \preceq M\mathbf{I};$$

4. the Hessian is L -Lipschitz continuous, that is, for any $\Delta \in \mathbb{R}^d$

$$\|\nabla^2 f(\mathbf{w} + \Delta) - \nabla^2 f(\mathbf{w})\|_2 \leq L\|\Delta\|_2,$$

where $\|\cdot\|_2$ is the spectral norm for matrices.

We first prove the following “global” convergence guarantee which shows that OverSketched Newton would converge from any random initialization of $\mathbf{w}_0 \in \mathbb{R}^d$ with high probability.

Theorem 3.1 (Global convergence for strongly-convex f). *Consider Assumptions 1, 2, and 3 and step-size α_t given by Eq. (5). Let \mathbf{w}^* be the optimal solution of (1). Let ϵ and μ be positive constants. Then, using the sketch in (4) with a sketch dimension $Nb + eb = \Omega(\frac{d^{1+\mu}}{\epsilon^2})$*

Algorithm 3: OverSketched Newton in a nutshell

Input : Convex function f ; Initial iterate $\mathbf{w}_0 \in \mathbb{R}^d$; Line search parameter $0 < \beta \leq 1/2$;
Number of iterations T

- 1 **for** $t = 1$ to T **do**
- 2 Compute full gradient \mathbf{g}_t in a distributed fashion using Algorithm 1
- 3 Compute sketched Hessian matrix $\hat{\mathbf{H}}_t$ in a distributed fashion using Algorithm 2
- 4 **if** f is strongly-convex **then**
- 5 Compute the update direction at the master as: $\mathbf{p}_t = -[\hat{\mathbf{H}}_t]^{-1}\nabla f(\mathbf{w}_t)$
- 6 Compute step-size α_t satisfying the line-search condition (5) in a distributed fashion
- 7 **else**
- 8 Compute the update direction at the master as: $\mathbf{p}_t = -[\hat{\mathbf{H}}_t]^\dagger\nabla f(\mathbf{w}_t)$
- 9 Find step-size α_t satisfying the line-search condition (6) in a distributed fashion
- 10 **end**
- 11 Compute the model update $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t\mathbf{p}_t$ at the master
- 12 **end**

and the number of column-blocks $N + e = \Theta_\mu(1/\epsilon)$, the updates for OverSketched Newton, for any $\mathbf{w}_t \in \mathbb{R}^d$, satisfy

$$f(\mathbf{w}_{t+1}) - f(\mathbf{w}^*) \leq (1 - \rho)(f(\mathbf{w}_t) - f(\mathbf{w}^*)),$$

with probability at least $1 - 1/d^\tau$, where $\rho = \frac{2\alpha_t\beta k}{M(1+\epsilon)}$ and $\tau > 0$ is a constant depending on μ and constants in $\Omega(\cdot)$ and $\Theta(\cdot)$. Moreover, α_t satisfies $\alpha_t \geq \frac{2(1-\beta)(1-\epsilon)k}{M}$.

Proof. See Section 6.1. □

Theorem 3.1 guarantees the global convergence of OverSketched Newton starting with any initial estimate $\mathbf{w}_0 \in \mathbb{R}^d$ to the optimal solution \mathbf{w}^* with at least a linear rate.

Next, we can also prove an additional “local” convergence guarantee for OverSketched Newton, under the assumption that \mathbf{w}_0 is sufficiently close to \mathbf{w}^* .

Theorem 3.2 (Local convergence for strongly-convex f). Consider Assumptions 1, 2, and 4 and step-size $\alpha_t = 1$. Let \mathbf{w}^* be the optimal solution of (1) and γ and β be the minimum and maximum eigenvalues of $\nabla^2 f(\mathbf{w}^*)$, respectively. Let $\epsilon \in (0, \gamma/(8\beta)]$ and $\mu > 0$. Then, using the sketch in (4) with a sketch dimension $Nb + eb = \Omega(\frac{d^{1+\mu}}{\epsilon^2})$ and the number of column-blocks $N + e = \Theta_\mu(1/\epsilon)$, the updates for OverSketched Newton, with initialization \mathbf{w}_0 such that $\|\mathbf{w}_0 - \mathbf{w}^*\|_2 \leq \frac{\gamma}{8L}$, follow

$$\|\mathbf{w}_{t+1} - \mathbf{w}^*\|_2 \leq \frac{25L}{8\gamma}\|\mathbf{w}_t - \mathbf{w}^*\|_2^2 + \frac{5\epsilon\beta}{\gamma}\|\mathbf{w}_t - \mathbf{w}^*\|_2 \quad \text{for } t = 1, 2, \dots, T,$$

with probability at least $1 - T/d^\tau$, where $\tau > 0$ is a constant depending on μ and constants in $\Omega(\cdot)$ and $\Theta(\cdot)$.

Proof. See Section 6.2. □

Theorem 3.2 implies that the convergence is linear-quadratic in error $\Delta_t = \mathbf{w}_t - \mathbf{w}^*$. Initially, when $\|\Delta_t\|_2$ is large, the first term of the RHS will dominate and the convergence will be quadratic, that is, $\|\Delta_{t+1}\|_2 \lesssim \frac{25L}{8\gamma} \|\Delta_t\|_2^2$. In later stages, when $\|\mathbf{w}_t - \mathbf{w}^*\|_2$ becomes sufficiently small, the second term of RHS will start to dominate and the convergence will be linear, that is, $\|\Delta_{t+1}\|_2 \lesssim \frac{5\epsilon\beta}{\gamma} \|\Delta_t\|_2$. At this stage, the sketch dimension can be increased to reduce ϵ to diminish the effect of the linear term and improve the convergence rate in practice. Note that, for second order methods, the number of iterations T is in the order of tens in general while the number of features d is typically in thousands. Hence, the probability of failure is generally small (and can be made negligible by choosing τ appropriately).

Though the works [19, 20, 23, 47, 52] also prove convergence guarantees for approximate Hessian-based optimization, no convergence results exist for the OverSketch matrix in Eq. (4) to the best of our knowledge. OverSketch has many nice properties like sparsity, input obliviousness, and amenability to distributed implementation, and our convergence guarantees take into account the block-size b (that captures the amount of communication at each worker) and the number of stragglers e , both of which are a property of the distributed system. On the other hand, algorithms in [19, 20, 23, 47, 52] are tailored to run on a single machine.

Next, we consider the case of weakly-convex functions. For this case, we consider two more assumptions on the Hessian matrix, similar to [22]. These assumptions are a relaxation of the strongly-convex case.

Assumptions:

5. There exists some $\eta > 0$ such that, $\forall \mathbf{w} \in \mathbb{R}^d$,

$$\|(\nabla^2 f(\mathbf{w}))^\dagger\|_2 \leq 1/\eta.$$

This assumption establishes regularity on the pseudo-inverse of $\nabla^2 f(\mathbf{x})$. It also implies that $\|\nabla^2 f(\mathbf{w})\mathbf{p}\| \geq \eta\|\mathbf{p}\| \forall p \in \text{Range}(\nabla^2 f(\mathbf{w}))$, that is, the minimum ‘non-zero’ eigenvalue of $\nabla^2 f(\mathbf{w})$ is lower bounded by η ; just as in the k -strongly convex case, the smallest eigenvalue is greater than k .

6. Let $\mathbf{U} \in \mathbb{R}^{d \times d}$ be any arbitrary orthogonal basis for $\text{Range}(\nabla^2 f(\mathbf{w}))$, there exists $0 < \nu \leq 1$, such that,

$$\|\mathbf{U}^T \nabla \mathbf{f}(\mathbf{w})\|^2 \geq \nu \|\nabla \mathbf{f}(\mathbf{w})\|^2 \quad \forall \mathbf{w} \in \mathbb{R}^d.$$

This assumption ensures that there is always a non-zero component of the gradient in the subspace spanned by the Hessian, and, thus, ensures that the model update $-\hat{\mathbf{H}}_t^\dagger \nabla f(\mathbf{w}_t)$ will not be zero.

Note that the above assumptions are always satisfied by strongly-convex functions. Next, we prove global convergence of OverSketched Newton when the objective is weakly-convex.

Theorem 3.3 (Global convergence for weakly-convex f). *Consider Assumptions 1,3,4,5 and 6 and step-size α_t given by Eq. (6). Let $\epsilon \in \left(0, \frac{(1-\beta)\nu\eta}{2M}\right]$ and $\mu > 0$. Then, using an OverSketch matrix with a sketch dimension $Nb + eb = \Omega\left(\frac{d^{1+\mu}}{\epsilon^2}\right)$ and the number of column-blocks $N + e = \Theta_\mu(1/\epsilon)$, the updates for OverSketched Newton, for any $\mathbf{w}_t \in \mathbb{R}^d$, satisfy*

$$\|\nabla f(\mathbf{w}_{t+1})\|^2 \leq \left(1 - 2\beta\alpha\nu\frac{(1-\epsilon)\eta}{M(1+\epsilon)}\right)\|\nabla f(\mathbf{w}_t)\|^2,$$

with probability at least $1 - 1/d^\tau$, where $\alpha = \frac{\eta}{2Q}[(1-\beta)\nu\eta - 2\epsilon M]$, $Q = (L\|\nabla f(\mathbf{w}_0)\| + M^2)$, \mathbf{w}_0 is the initial iterate of the algorithm and $\tau > 0$ is a constant depending on μ and constants in $\Omega(\cdot)$ and $\Theta(\cdot)$.

Proof. See Section 6.3. □

Even though we present the above guarantees for the sketch matrix in Eq. (4), our analysis is valid for any sketch that satisfies the *subspace embedding* property (Lemma 6.1; see [29] for details on subspace embedding property of sketches). To the best of our knowledge, this is the first work to prove the convergence guarantees for weakly-convex functions when the Hessian is calculated approximately using sketching techniques. Later, authors in [53] extended the analysis to the case of general Hessian perturbations with additional assumptions on the type of perturbation.

3.2 Distributed Line Search

Here, we describe a line-search procedure for distributed serverless optimization, which is inspired by the line-search method from [24] for serverful systems. To solve for the step-size α_t as described in the optimization problem in (5), we set $\beta = 0.1$ and choose a candidate set $\mathcal{S} = \{4^0, 4^1, \dots, 4^{-5}\}$. After the master calculates the descent direction \mathbf{p}_t in the t -th iteration, the i -th worker calculates $f_i(\mathbf{w}_t + \alpha\mathbf{p}_t)$ for all values of α in the candidate set \mathcal{S} , where $f_i(\cdot)$ depends on the local data available at the i -th worker and $f(\cdot) = \sum_i f_i(\cdot)$ ⁴.

The master then sums the results from workers to obtain $f(\mathbf{w}_t + \alpha\mathbf{p}_t)$ for all values of α in \mathcal{S} and finds the largest α that satisfies the Armijo condition in (5)⁵. Note that line search requires an additional round of communication where the master communicates \mathbf{p}_t to the workers through cloud and the workers send back the function values $f_i(\cdot)$. Finally, the master finds the best step-size from set \mathcal{S} and finds the model estimate \mathbf{w}_{t+1} .

⁴For the weakly-convex case, the workers calculate $\nabla f_i(\cdot)$ instead of $f_i(\cdot)$, and the master calculates $\|\nabla f(\cdot)\|^2$ instead of $f(\cdot)$.

⁵Note that codes can be used to mitigate stragglers during distributed line-search in a manner similar to the gradient computation phase.

4 OverSketched Newton on Serverless Systems: Examples

Here, we describe several examples where our general approach can be applied.

4.1 Logistic Regression using OverSketched Newton

The optimization problem for supervised learning using Logistic Regression takes the form

$$\min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right\}. \quad (7)$$

Here, $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^{d \times 1}$ and $y_1, \dots, y_n \in \mathbb{R}$ are training sample vectors and labels, respectively. The goal is to learn the feature vector $\mathbf{w}^* \in \mathbb{R}^{d \times 1}$. Let $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ and $\mathbf{y} = [y_1, \dots, y_n] \in \mathbb{R}^{n \times 1}$ be the example and label matrices, respectively. The gradient for the problem in (7) is given by

$$\nabla f(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \frac{-y_i \mathbf{x}_i}{1 + e^{y_i \mathbf{w}_i^T \mathbf{x}_i}} + \lambda \mathbf{w}.$$

Calculation of $\nabla f(\mathbf{w})$ involves two matrix-vector products, $\boldsymbol{\alpha} = \mathbf{X}^T \mathbf{w}$ and $\nabla f(\mathbf{w}) = \frac{1}{n} \mathbf{X} \boldsymbol{\beta} + \lambda \mathbf{w}$, where $\beta_i = \frac{-y_i}{1 + e^{y_i \alpha_i}} \forall i \in [1, \dots, n]$. When the example matrix is large, these matrix-vector products are performed distributedly using codes. Faster convergence is obtained by second-order methods which will additionally compute the Hessian $\mathbf{H} = \frac{1}{n} \mathbf{X} \boldsymbol{\Lambda} \mathbf{X}^T + \lambda \mathbf{I}_d$, where $\boldsymbol{\Lambda}$ is a diagonal matrix with entries given by $\Lambda(i, i) = \frac{e^{y_i \alpha_i}}{(1 + e^{y_i \alpha_i})^2}$. The product $\mathbf{X} \boldsymbol{\Lambda} \mathbf{X}^T$ is computed approximately in a distributed straggler-resilient manner using the sketch matrix in (4). Using the result of distributed multiplication, the Hessian matrix \mathbf{H} is calculated at the master and the model is updated as $\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{H}^{-1} \nabla f(\mathbf{w}_t)$. In practice, efficient algorithm like conjugate gradient, that provide a good estimate in a small number of iterations, can be used locally at the master to solve for \mathbf{w}_{t+1} [54].⁶

We provide a detailed description of OverSketched Newton for large-scale logistic regression for serverless systems in Algorithm 4. Steps 4, 8, and 14 of the algorithm are computed in parallel on AWS Lambda. All other steps are simple vector operations that can be performed locally at the master, for instance, the user’s laptop. Steps 4 and 8 are executed in a straggler-resilient fashion using the coding scheme in [34], as illustrated in Fig. 1 and described in detail in Algorithm 1.

We use the coding scheme in [34] since the encoding can be implemented in parallel and requires less communication per worker compared to the other schemes, for example schemes in [33, 40], that use Maximum Distance Separable (MDS) codes. Moreover, the

⁶Note that here we have assumed that the number of features is small enough to perform the model update locally at the master. This is not necessary, and straggler resilient schemes, such as in [35], can be used to perform distributed conjugate gradient in serverless systems.

Algorithm 4: OverSketched Newton: Logistic Regression for Serverless Computing

```
1 Input Data (stored in cloud storage): Example Matrix  $\mathbf{X} \in \mathbb{R}^{d \times n}$  and vector
    $\mathbf{y} \in \mathbb{R}^{n \times 1}$  (stored in cloud storage), regularization parameter  $\lambda$ , number of iterations
    $T$ , Sketch  $\mathbf{S}$  as defined in Eq. (4)
2 Initialization: Define  $\mathbf{w}^1 = \mathbf{0}^{d \times 1}, \boldsymbol{\beta} = \mathbf{0}^{n \times 1}, \boldsymbol{\gamma} = \mathbf{0}^{n \times 1}$ , Encode  $\mathbf{X}$  and  $\mathbf{X}^T$  as
   described in Algorithm 1
3 for  $t = 1$  to  $T$  do
4    $\boldsymbol{\alpha} = \mathbf{X}\mathbf{w}^t$  ; // Compute in parallel using Algorithm 1
5   for  $i = 1$  to  $n$  do
6      $\beta_i = \frac{-y_i}{1+e^{y_i\alpha_i}}$  ;
7   end
8    $\mathbf{g} = \mathbf{X}^T\boldsymbol{\beta}$  ; // Compute in parallel using Algorithm 1
9    $\nabla f(\mathbf{w}^t) = \mathbf{g} + \lambda\mathbf{w}^t$  ;
10  for  $i = 1$  to  $n$  do
11     $\gamma(i) = \frac{e^{y_i\alpha_i}}{(1+e^{y_i\alpha_i})^2}$  ;
12  end
13   $\mathbf{A} = \sqrt{\text{diag}(\boldsymbol{\gamma})}\mathbf{X}^T$ 
14   $\hat{\mathbf{H}} = \mathbf{A}^T\mathbf{S}\mathbf{S}^T\mathbf{A}$  ; // Compute in parallel using Algorithm 2
15   $\mathbf{H} = \frac{1}{n}\hat{\mathbf{H}} + \lambda\mathbf{I}_d$  ;
16   $\mathbf{w}^{t+1} = \mathbf{w}^t - \mathbf{H}^{-1}\nabla f(\mathbf{w}^t)$  ;
17 end
Result:  $\mathbf{w}^* = \mathbf{w}_{T+1}$ 
```

decoding scheme takes linear time and is applicable on real-valued matrices. Note that since the example matrix \mathbf{X} is constant in this example, the encoding of \mathbf{X} is done only once before starting the optimization algorithm. Thus, the encoding cost can be amortized over iterations. Moreover, decoding over the resultant product vector requires negligible time and space, even when n is scaling into the millions.

The same is, however, not true for the matrix multiplication for Hessian calculation (step 14 of Algorithm 4), as the matrix \mathbf{A} changes in each iteration, thus encoding costs will be incurred in every iteration if error-correcting codes are used. Moreover, encoding and decoding a huge matrix stored in the cloud incurs heavy communication cost and becomes prohibitive. Motivated by this, we use OverSketch in step 14, as described in Algorithm 2, to calculate an approximate matrix multiplication, and hence the Hessian, efficiently in serverless systems with inbuilt straggler resiliency.⁷

⁷We also evaluate the exact Hessian-based algorithm with speculative execution, i.e., recomputing the straggling jobs, and compare it with OverSketched Newton in Sec. 5.

4.2 Softmax Regression using OverSketched Newton

We take unregularized softmax regression as an illustrative example for the weakly convex case. The goal is to find the weight matrix $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K]$ that fit the training data $\mathbf{X} \in \mathbb{R}^{d \times N}$ and $\mathbf{y} \in \mathbb{R}^{K \times N}$. Here $\mathbf{w}_i \in \mathbb{R}^d$ represents the weight vector for the i -th class for all $i \in [1, K]$ and K is the total number of classes. Hence, the resultant feature dimension for softmax regression is dK . The optimization problem is of the form

$$f(\mathbf{W}) = \sum_{n=1}^N \left[\sum_{k=1}^K y_{kn} \mathbf{w}_k^T \mathbf{x}_n - \log \sum_{l=1}^K \exp(\mathbf{w}_l^T \mathbf{x}_n) \right]. \quad (8)$$

The gradient vector for the i -th class is given by

$$\nabla f_i(\mathbf{W}) = \sum_{n=1}^N \left[\frac{\exp(\mathbf{w}_i^T \mathbf{x}_n)}{\sum_{l=1}^K \exp(\mathbf{w}_l^T \mathbf{x}_n)} - y_{in} \right] \mathbf{x}_n \quad \forall i \in [1, K], \quad (9)$$

which can be written as matrix products $\boldsymbol{\alpha}_i = \mathbf{X}^T \mathbf{w}_i$ and $\nabla f_i(\mathbf{W}) = \mathbf{X} \boldsymbol{\beta}_i$, where the entries of $\boldsymbol{\beta}_i \in \mathbb{R}^N$ are given by $\beta_{in} = \left(\frac{\exp(\alpha_{in})}{\sum_{l=1}^K \exp(\alpha_{ln})} - y_{in} \right)$. Thus, the full gradient matrix is given by $\nabla f(\mathbf{W}) = \mathbf{X} \boldsymbol{\beta}$ where the entries of $\boldsymbol{\beta} \in \mathbb{R}^{N \times K}$ are dependent on $\boldsymbol{\alpha} \in \mathbb{R}^{N \times K}$ as above and the matrix $\boldsymbol{\alpha}$ is given by $\boldsymbol{\alpha} = \mathbf{X}^T \mathbf{W}$. We assume that the number of classes K is small enough such that tall matrices $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are small enough for the master to do local calculations on them.

Since the effective number of features is $d \times K$, the Hessian matrix is of dimension $dK \times dK$. The (i, j) -th component of the Hessian, say \mathbf{H}_{ij} , is

$$\mathbf{H}_{ij}(\mathbf{W}) = \frac{d}{d\mathbf{w}_j} \nabla f_i(\mathbf{W}) = \frac{d}{d\mathbf{w}_j} \mathbf{X} \boldsymbol{\beta}_i = \mathbf{X} \frac{d}{d\mathbf{w}_j} \boldsymbol{\beta}_i = \mathbf{X} \mathbf{Z}_{ij} \mathbf{X}^T \quad (10)$$

where $\mathbf{Z}_{ij} \in \mathbb{R}^{N \times N}$ is a diagonal matrix whose n -th diagonal entry is

$$Z_{ij}(n) = \frac{\exp(\alpha_{in})}{\sum_{l=1}^K \exp(\alpha_{ln})} \left(\mathbb{I}(i=j) - \frac{\exp(\alpha_{jn})}{\sum_{l=1}^K \exp(\alpha_{ln})} \right) \quad \forall n \in [1, N], \quad (11)$$

where $\mathbb{I}(\cdot)$ is the indicator function and $\boldsymbol{\alpha} = \mathbf{X} \mathbf{W}$ was defined above. The full Hessian matrix is obtained by putting together all such \mathbf{H}_{ij} 's in a $dK \times dK$ matrix and can be expressed in a matrix-matrix multiplication form as

$$\nabla^2 f(\mathbf{W}) = \begin{bmatrix} \mathbf{H}_{11} & \cdots & \mathbf{H}_{1K} \\ \vdots & \ddots & \vdots \\ \mathbf{H}_{K1} & \cdots & \mathbf{H}_{KK} \end{bmatrix} = \begin{bmatrix} \mathbf{X} \mathbf{Z}_{11} \mathbf{X}^T & \cdots & \mathbf{X} \mathbf{Z}_{1K} \mathbf{X}^T \\ \vdots & \ddots & \vdots \\ \mathbf{X} \mathbf{Z}_{K1} \mathbf{X}^T & \cdots & \mathbf{X} \mathbf{Z}_{KK} \mathbf{X}^T \end{bmatrix} = \bar{\mathbf{X}} \bar{\mathbf{Z}} \bar{\mathbf{X}}^T, \quad (12)$$

where $\bar{\mathbf{X}} \in \mathbb{R}^{dK \times NK}$ is a block diagonal matrix that contains \mathbf{X} in the diagonal blocks and $\bar{\mathbf{Z}} \in \mathbb{R}^{NK \times NK}$ is formed by stacking all the \mathbf{Z}_{ij} 's for $i, j \in [1, K]$. In OverSketched Newton, we compute this multiplication using sketching in serverless systems for efficiency and resiliency to stragglers. Assuming $d \times K$ is small enough, the master can then calculate the update \mathbf{p}_t using efficient algorithms such the minimum-residual method [22, 55].

4.3 Other Example Problems

In this section, we describe several other commonly encountered optimization problems that can be solved using OverSketched Newton.

Ridge Regularized Linear Regression: The optimization problem is

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2n} \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2. \quad (13)$$

The gradient in this case can be written as $\frac{1}{n} \mathbf{X}(\boldsymbol{\beta} - \mathbf{y}) + \lambda \mathbf{w}$, where $\boldsymbol{\beta} = \mathbf{X}^T \mathbf{w}$, where the training matrix \mathbf{X} and label vector \mathbf{y} were defined previously. The Hessian is given by $\nabla^2 f(\mathbf{w}) = \mathbf{X}\mathbf{X}^T + \lambda \mathbf{I}$. For $n \gg d$, this can be computed approximately using the sketch matrix in (4).

Linear programming via interior point methods: The following linear program can be solved using OverSketched Newton

$$\underset{\mathbf{Ax} \leq \mathbf{b}}{\text{minimize}} \mathbf{c}^T \mathbf{x}, \quad (14)$$

where $\mathbf{x} \in \mathbb{R}^{m \times 1}$, $\mathbf{c} \in \mathbb{R}^{m \times 1}$, $\mathbf{b} \in \mathbb{R}^{n \times 1}$ and $\mathbf{A} \in \mathbb{R}^{n \times m}$ is the constraint matrix with $n > m$. In algorithms based on interior point methods, the following sequence of problems using Newton's method

$$\min_{\mathbf{x} \in \mathbb{R}^m} f(\mathbf{x}) = \min_{\mathbf{x} \in \mathbb{R}^m} \left(\tau \mathbf{c}^T \mathbf{x} - \sum_{i=1}^n \log(b_i - \mathbf{a}_i \mathbf{x}) \right), \quad (15)$$

where \mathbf{a}_i is the i -th row of \mathbf{A} , τ is increased geometrically such that when τ is very large, the logarithmic term does not affect the objective value and serves its purpose of keeping all intermediates solution inside the constraint region. The update in the t -th iteration is given by $\mathbf{x}_{t+1} = \mathbf{x}_t - (\nabla^2 f(\mathbf{x}_t))^{-1} \nabla f(\mathbf{x}_t)$, where \mathbf{x}_t is the estimate of the solution in the t -th iteration. The gradient can be written as $\nabla f(\mathbf{x}) = \tau \mathbf{c} + \mathbf{A}^T \boldsymbol{\beta}$ where $\beta_i = 1/(b_i - \alpha_i)$ and $\boldsymbol{\alpha} = \mathbf{A}\mathbf{x}$.

The Hessian for the objective in (15) is given by

$$\nabla^2 f(\mathbf{x}) = \mathbf{A}^T \text{diag} \frac{1}{(b_i - \alpha_i)^2} \mathbf{A}. \quad (16)$$

The square root of the Hessian is given by $\nabla^2 f(\mathbf{x})^{1/2} = \text{diag} \frac{1}{|b_i - \alpha_i|} \mathbf{A}$. The computation of Hessian requires $O(nm^2)$ time and is the bottleneck in each iteration. Thus, we can use sketching to mitigate stragglers while evaluating the Hessian efficiently, i.e. $\nabla^2 f(\mathbf{x}) \approx (\mathbf{S} \nabla^2 f(\mathbf{x})^{1/2})^T \times (\mathbf{S} \nabla^2 f(\mathbf{x})^{1/2})$, where \mathbf{S} is the OverSketch matrix defined in (4).

Lasso Regularized Linear Regression: The optimization problem takes the following form

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_1, \quad (17)$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$ is the measurement matrix, the vector $\mathbf{y} \in \mathbb{R}^n$ contains the measurements, $\lambda \geq 0$ and $d \gg n$. To solve (17), we consider its dual variation

$$\min_{\|\mathbf{X}^T \mathbf{z}\|_\infty \leq \lambda, \mathbf{z} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|_2^2,$$

which is amenable to interior point methods and can be solved by optimizing the following sequence of problems where τ is increased geometrically

$$\min_{\mathbf{z}} f(\mathbf{z}) = \min_{\mathbf{z}} \left(\frac{\tau}{2} \|\mathbf{y} - \mathbf{z}\|_2^2 - \sum_{j=1}^d \log(\lambda - \mathbf{x}_j^T \mathbf{z}) - \sum_{j=1}^d (\lambda + \mathbf{x}_j^T \mathbf{z}) \right),$$

where \mathbf{x}_j is the j -th column of \mathbf{X} . The gradient can be expressed in few matrix-vector multiplications as $\nabla f(\mathbf{z}) = \tau(\mathbf{z} - \mathbf{y}) + \mathbf{X}(\boldsymbol{\beta} - \boldsymbol{\gamma})$, where $\beta_i = 1/(\lambda - \alpha_i)$, $\gamma_i = 1/(\lambda + \alpha_i)$, and $\boldsymbol{\alpha} = \mathbf{X}^T \mathbf{z}$. Similarly, the Hessian can be written as $\nabla^2 f(\mathbf{z}) = \tau \mathbf{I} + \mathbf{X} \boldsymbol{\Lambda} \mathbf{X}^T$, where $\boldsymbol{\Lambda}$ is a diagonal matrix whose entries are given by $\Lambda_{ii} = 1/(\lambda - \alpha_i)^2 + 1/(\lambda + \alpha_i)^2 \forall i \in [1, n]$.

Other common problems where OverSketched Newton is applicable include Linear Regression, Support Vector Machines (SVMs), Semidefinite programs, etc.

5 Experimental Results

In this section, we evaluate OverSketched Newton on AWS Lambda using real-world and synthetic datasets, and we compare it with state-of-the-art distributed optimization algorithms⁸. We use the serverless computing framework, Pywren [3]. Our experiments are focused on logistic and softmax regression, which are popular supervised learning problems, but they can be reproduced for other problems described in Section 4. We present experiments on the following datasets:

Dataset	Training Samples	Features	Testing samples
Synthetic	300,000	3000	100,000
EPSILON	400,000	2000	100,000
WEBPAGE	48,000	300	15,000
a9a	32,000	123	16,000
EMNIST	240,000	7840	40,000

For comparison of OverSketched Newton with existing distributed optimization schemes, we choose recently-proposed Globally Improved Approximate Newton Direction (GIANT) [24]. The reason is that GIANT boasts a better convergence rate than many existing distributed second-order methods for linear and logistic regression, when $n \gg d$. In GIANT, and other similar distributed second-order algorithms, the training data is evenly divided among

⁸A working implementation of OverSketched Newton is available at <https://github.com/vvipgupta/OverSketchedNewton>

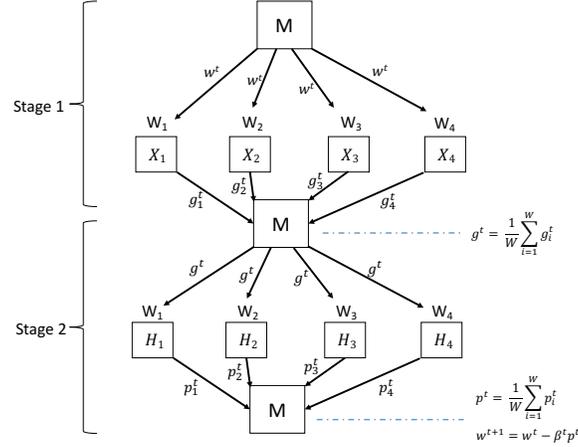
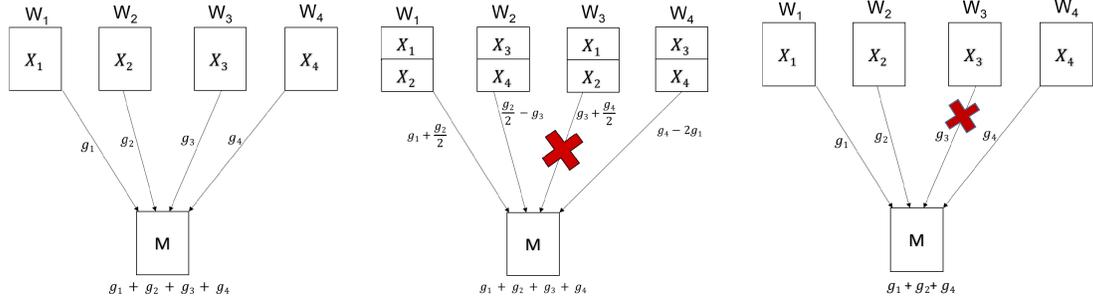


Figure 4: GIANT: The two stage second order distributed optimization scheme with four workers. First, master calculates the full gradient by aggregating local gradients from workers. Second, the master calculates approximate Hessian using local second-order updates from workers.



(a) Simple Gradient Descent where each worker stores one-fourth fraction of the whole data and sends back a partial gradient corresponding to its own data to the master
 (b) Gradient Coding described in [37] with W_3 straggling. To get the global gradient, master would compute $g_1 + g_2 + g_3 + g_4 = 3(g_1 + \frac{g_2}{2}) - (\frac{g_2}{2} - g_3) + (g_4 - 2g_1)$
 (c) Mini-batch gradient descent, where the stragglers are ignored during gradient aggregation and the gradient is later scaled according to the size of mini-batch

Figure 5: Different gradient descent schemes in serverful systems in presence of stragglers

workers, and the algorithms proceed in two stages. First, the workers compute partial gradients using local training data, which is then aggregated by the master to compute the exact gradient. Second, the workers receive the full gradient to calculate their local second-order estimate, which is then averaged by the master. An illustration is shown in Fig. 4.

For straggler mitigation in such serverful systems based algorithms, [37] proposes a scheme for coding gradient updates called *gradient coding*, where the data at each worker

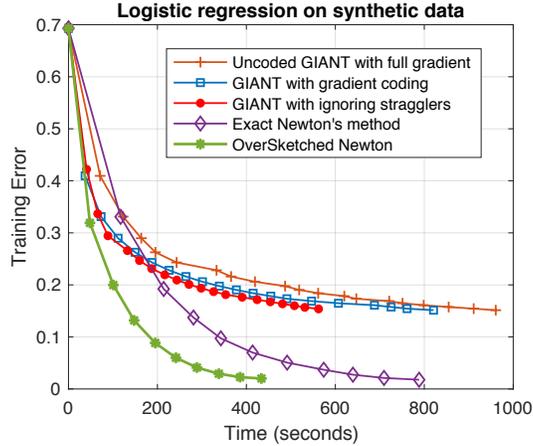


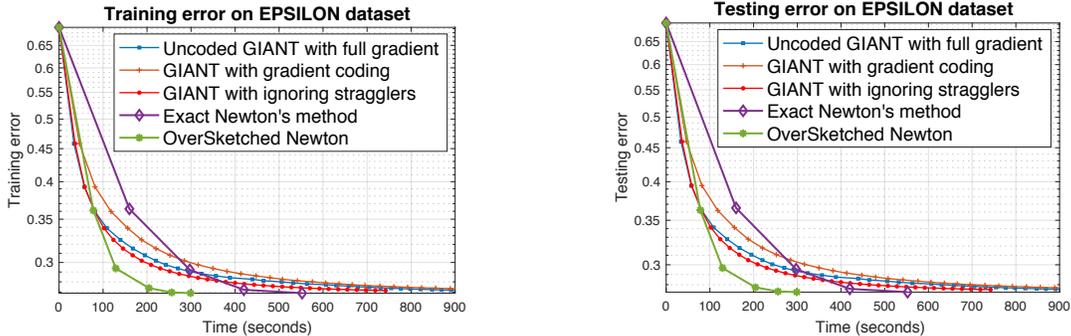
Figure 6: Convergence comparison of GIANT (employed with different straggler mitigation methods), exact Newton’s method and OverSketched Newton for Logistic regression on AWS Lambda. The synthetic dataset considered has 300,000 examples and 3000 features.

is repeated multiple times to compute redundant copies of the gradient. See Figure 5b for illustration. Figure 5a illustrates the scheme that waits for all workers and Figure 5c illustrates the ignoring stragglers approach. We use the three schemes for dealing with stragglers illustrated in Figure 5 during the two stages of GIANT, and we compare their convergence with OverSketched Newton. We further evaluate and compare the convergence exact Newton’s method (employed with speculative execution, that is, reassigning and recomputing the work for straggling workers).

5.1 Comparisons with Existing Second-Order Methods on AWS Lambda

In Figure 6, we present our results on a synthetic dataset with $n = 300,000$ and $d = 3000$ for logistic regression on AWS Lambda. Each column $\mathbf{x}_i \in \mathbb{R}^d$, for all $i \in [1, n]$, is sampled uniformly randomly from the cube $[-1, 1]^d$. The labels y_i are sampled from the logistic model, that is, $\mathbb{P}[y_i = 1] = 1/(1 + \exp(\mathbf{x}_i \mathbf{w} + b))$, where the weight vector \mathbf{w} and bias b are generated randomly from the normal distribution. The vector \mathbf{w} and bias b are generated randomly from the normal distribution.

The orange, blue and red curves demonstrate the convergence for GIANT with the full gradient (that waits for all the workers), gradient coding and mini-batch gradient (that ignores the stragglers while calculating gradient and second-order updates) schemes, respectively. The purple and green curves depict the convergence for the exact Newton’s method and OverSketched Newton, respectively. The gradient coding scheme is applied for one straggler, that is the data is repeated twice at each worker. We use 60 Lambda workers for executing GIANT in parallel. Similarly, for Newton’s method, we use 60 workers for matrix-vector multiplication in steps 4 and 8 of Algorithm 4, 3600 workers for exact Hessian computation and 600 workers for sketched Hessian computation with a sketch dimension of $10d = 30,000$



(a) Training error for logistic regression on EP-SILON dataset

(b) Testing error for logistic regression on EP-SILON dataset

Figure 7: Comparison of training and testing errors for logistic regression on EPSILON dataset with several Newton based schemes on AWS Lambda. OverSketched Newton outperforms others by at least 46%. Testing error closely follows training error.

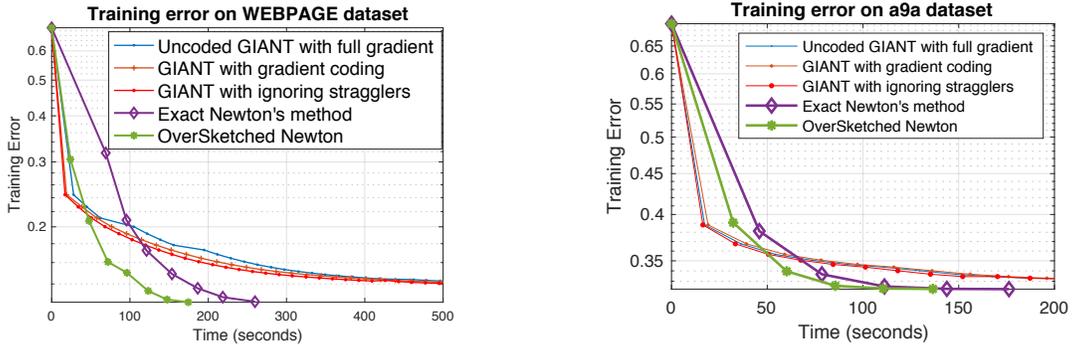
in step 14 of Algorithm 4. In all cases, unit step-size was used to update the model⁹

Remark 1. *In our experiments, we choose the number of workers in such a way that each worker receives approximately the same amount of data to work with, regardless of the algorithm. This is motivated by the fact that the memory at each worker is the bottleneck in serverless systems (e.g., in AWS Lambda, the memory at each worker can be as low as 128 MB). Note that this is unlike serverful/HPC systems, where the number of workers is the bottleneck.*

An important point to note from Fig. 6 is that the uncoded scheme (that is, the one that waits for all stragglers) has the worst performance. The implication is that good straggler/fault mitigation algorithms are essential for computing in the serverless setting. Secondly, the mini-batch scheme outperforms the gradient coding scheme by 25%. This is because gradient coding requires additional communication of data to serverless workers (twice when coding for one straggler, see [37] for details) at each invocation to AWS Lambda. On the other hand, the exact Newton’s method converges much faster than GIANT, even though it requires more time per iteration.

The number of iterations needed for convergence for OverSketched Newton and exact Newton (that exactly computes the Hessian) is similar, but OverSketched Newton converges in almost half the time due to an efficient computation of (approximate) Hessian (which is the computational bottleneck and thus reduces time per iteration).

⁹Line-search in Section 3 was mainly introduced to prove theoretical guarantees. In our experiments, we observe that constant step-size works well for OverSketched Newton.



(a) Logistic regression on WEBPAGE dataset

(b) Logistic regression on a9a dataset

Figure 8: Logistic regression on WEBPAGE and a9a datasets with several Newton based schemes on AWS Lambda. OverSketched Newton outperforms others by at least 25%.

5.1.1 Logistic Regression on EPSILON, WEBPAGE and a9a Datasets

In Figure 7, we repeat the above experiment with EPSILON classification dataset obtained from [36], with $n = 0.4$ million and $d = 2000$. We plot training and testing errors for logistic regression for the schemes described in the previous section. Here, we use 100 workers for GIANT, and 100 workers for matrix-vector multiplications for gradient calculation in OverSketched Newton. We use gradient coding designed for three stragglers in GIANT. This scheme performs worse than uncoded GIANT that waits for all the stragglers due to the repetition of training data at workers. Hence, one can conclude that the communication costs dominate the stragglers costs. In fact, it can be observed that the mini-batch gradient scheme that ignores the stragglers outperforms the gradient coding and uncoded schemes for GIANT.

During exact Hessian computation, we use 10,000 serverless workers with speculative execution to mitigate stragglers (i.e., recomputing the stragglers jobs) compared to OverSketched Newton that uses 1500 workers with a sketch dimension of $15d = 30,000$. OverSketched Newton requires a significantly smaller number of workers, as once the square root of Hessian is sketched in a distributed fashion, it can be copied into local memory of the master due to dimension reduction, and the Hessian can be calculated locally. Testing error follows training error closely, and important conclusions remain the same as in Figure 6. OverSketched Newton outperforms GIANT and exact Newton-based optimization by at least 46% in terms of running time.

We repeated the above experiments for classification on the WEBPAGE ($n = 49,749$ and $d = 300$) and a9a ($n = 32,561$ and $d = 123$) datasets [36]. For both datasets, we used 30 workers for each iteration in GIANT and any matrix-vector multiplications. Exact hessian calculation invokes 900 workers as opposed to 300 workers for OverSketched Newton, where the sketch dimension was $10d = 3000$. The results for training loss on logistic regression are shown in Figure 8. Testing error closely follows the training error in both cases. OverSketched

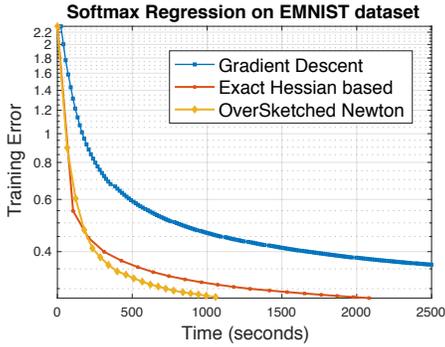


Figure 9: Convergence comparison of gradient descent, exact Newton’s method and OverSketched Newton for Softmax regression on AWS Lambda.

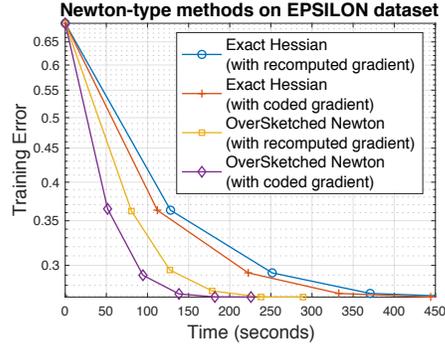


Figure 10: Convergence comparison of speculative execution and coded computing for gradient and Hessian computing with logistic regression on AWS Lambda.

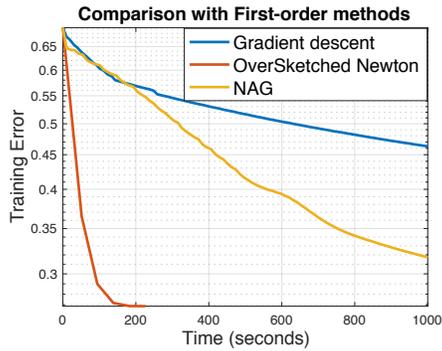


Figure 11: Convergence comparison of gradient descent, NAG and OverSketched Newton on AWS Lambda.

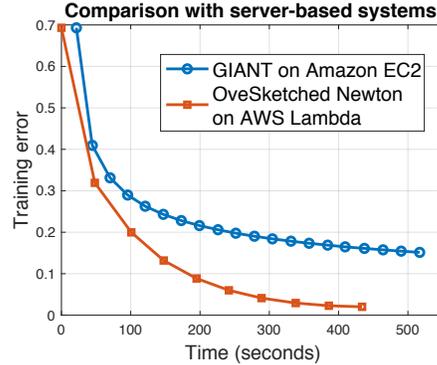


Figure 12: Convergence comparison of GIANT on AWS EC2 and OverSketched Newton on AWS Lambda.

Newton outperforms exact Newton and GIANT by at least $\sim 25\%$ and $\sim 75\%$, respectively, which is similar to the trends witnessed heretofore.

Remark 2. Note that conventional distributed second-order methods for serverful systems—which distribute training examples evenly across workers (such as [24, 42–46])—typically find a “localized approximation” (localized to each machine) of second-order update at each worker and then aggregate it. OverSketched Newton, on the other hand, uses the massive storage and compute power in serverless systems to find a more “globalized approximation” (globalized in the sense of across machine). Thus, it performs better in practice.

5.2 Softmax Regression on EMIST

In Fig. 9, we solve unregularized softmax regression, which is weakly convex (see Sec. 4.2 for details). We use the Extended MNIST (EMNIST) dataset [56] with $N = 240,000$ training examples, $d = 784$ features and $K = 10$ classes. Note that GIANT cannot be applied here as the objective function is not strongly convex. We compare the convergence rate of OverSketched Newton, exact Hessian and gradient descent based schemes.

For gradient computation in all three schemes, we use 60 workers. However, exact Newton scheme requires 3600 workers to calculate the $dK \times dK$ Hessian and recomputes the straggling jobs, while OverSketched Newton requires only 360 workers to calculate the sketch in parallel with sketch dimension $6dK = 47,040$. The approximate Hessian is then computed locally at the master using its sketched square root, where the sketch dimension is $6dK = 47,040$. The step-size is fixed and is determined by hyperparameter tuning before the start of the algorithm. Even for the weakly-convex case, second-order methods tend to perform better. Moreover, the runtime of OverSketched Newton outperforms both gradient descent and Exact Newton based methods by $\sim 75\%$ and $\sim 50\%$, respectively.

5.3 Coded computing versus Speculative Execution

In Figure 10, we compare the effect of straggler mitigation schemes, namely speculative execution, that is, restarting the jobs with straggling workers, and coded computing on the convergence rate during training and testing. We regard OverSketch based matrix multiplication as a coding scheme in which some redundancy is introduced during “over” sketching for matrix multiplication. There are four different cases, corresponding to gradient and hessian calculation using either speculative execution or coded computing. For speculative execution, we wait for at least 90% of the workers to return (this works well as the number of stragglers is generally less than 10%) and restart the jobs that did not return till this point.

For both exact Hessian and OverSketched Newton, using codes for distributed gradient computation outperforms speculative execution based straggler mitigation. Moreover, computing the Hessian using OverSketch is significantly better than exact computation in terms of running time as calculating the Hessian is the computational bottleneck in each iteration.

5.4 Comparison with First-Order Methods on AWS Lambda

In Figure 11, we compare gradient descent and Nesterov Accelerated Gradient (NAG) (while ignoring the stragglers) with OverSketched Newton for logistic regression on EPSILON dataset. We observed that for first-order methods, there is only a slight difference in convergence for a mini-batch gradient when the batch size is 95%. Hence, for gradient

descent and NAG, we use 100 workers in each iteration while ignoring the stragglers.¹⁰ These first-order methods were given the additional advantage of backtracking line-search, which determined the optimal amount to move in given a descent direction.¹¹ Overall, OverSketched Newton with unit step-size significantly outperforms gradient descent and NAG with backtracking line-search.

5.5 Comparison with Serverful Optimization

In Fig. 12, we compare OverSketched Newton on AWS Lambda with existing distributed optimization algorithm GIANT in serverful systems (AWS EC2). The results are plotted on synthetically generated data for logistic regression. For serverful programming, we use Message Passing Interface (MPI) with one `c3.8xlarge` master and 60 `t2.medium` workers in AWS EC2. In [4], the authors observed that many large-scale linear algebra operations on serverless systems take at least 30% more time compared to MPI-based computation on serverful systems. However, as shown in Fig. 12, we observe a slightly surprising trend that OverSketched Newton outperforms MPI-based optimization (that uses existing state-of-the-art optimization algorithm). This is because OverSketched Newton exploits the flexibility and massive scale at disposal in serverless, and thus produces a better approximation of the second-order update than GIANT.¹²

6 Proofs

To complete the proofs in this section, we will need the following lemma.

Lemma 6.1. *Let $\hat{\mathbf{H}}_t = \mathbf{A}_t^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A}_t$ where \mathbf{S}_t is the sparse sketch matrix in (4) with sketch dimension $m = \Omega(d^{1+\mu}/\epsilon^2)$ and $N = \Theta_\mu(1/\epsilon)$. Then, the following holds*

$$\lambda_{\min}(\hat{\mathbf{H}}_t) \geq (1 - \epsilon)\lambda_{\min}(\nabla^2 f(\mathbf{w}_t)), \quad (18)$$

$$\lambda_{\max}(\hat{\mathbf{H}}_t) \leq (1 + \epsilon)\lambda_{\max}(\nabla^2 f(\mathbf{w}_t)) \quad (19)$$

with probability at least $1 - \frac{1}{d^\tau}$, where $\tau > 0$ is a constant depending on μ and the constants in $\Theta(\cdot)$ and $\Omega(\cdot)$, and $\lambda_{\max}(\cdot)$ and $\lambda_{\min}(\cdot)$ denote the maximum and minimum eigenvalues, respectively. In general,

$$\lambda_i(\nabla^2 f(\mathbf{w}_t)) - \epsilon\lambda_{\max}(\nabla^2 f(\mathbf{w}_t)) \leq \lambda_i(\hat{\mathbf{H}}_t) \leq \lambda_i(\nabla^2 f(\mathbf{w}_t)) + \epsilon\lambda_{\max}(\nabla^2 f(\mathbf{w}_t)),$$

¹⁰We note that stochastic methods such as SGD perform worse than gradient descent since their update quality is poor, requiring more iterations (hence, more communication) to converge while not using the massive compute power of serverless. For example, 20% minibatch SGD in the setup of Fig. 11 requires 1.9× more time than gradient descent with same number of workers.

¹¹We remark that backtracking line-search required $\sim 13\%$ of the total time for NAG. Hence, as can be seen from Fig. 11, any well-tuned step-size method would still be significantly slower than OverSketched Newton.

¹²We do not compare with exact Newton in serverful systems since the data is large and stored in the cloud. Computing the exact Hessian would require a large number of workers (e.g., we use 10,000 workers for exact Newton in EPSILON dataset) which is infeasible in existing serverful systems.

where $\lambda_i(\cdot)$ is the i -th eigenvalue.

Proof. We note than N is the number of non-zero elements per row in the sketch \mathbf{S}_t in (4) after ignoring stragglers. We use Theorem 8 in [57] to bound the singular values for the sparse sketch \mathbf{S}_t in (4) with sketch dimension $m = \Omega(d^{1+\mu}/\epsilon^2)$ and $N = \Theta(1/\epsilon)$. It says that $\mathbb{P}(\forall \mathbf{x} \in \mathbb{R}^n, \|\mathbf{S}_t \mathbf{x}\|_2 \in (1 \pm \epsilon/3)\|\mathbf{x}\|_2) > 1 - 1/d^\tau$, where $\tau > 0$ depends on μ and the constants in $\Theta(\cdot)$ and $\Omega(\cdot)$. Thus, $\|\mathbf{S}_t \mathbf{x}\|_2 \in (1 \pm \epsilon/3)\|\mathbf{x}\|_2$, which implies that

$$\|\mathbf{S}_t \mathbf{x}\|_2^2 \in (1 + \epsilon^2/9 \pm 2\epsilon/3)\|\mathbf{x}\|_2^2,$$

with probability at least $1 - 1/d^\tau$. For $\epsilon \leq 1/2$, this leads to the following inequality

$$\|\mathbf{S}_t \mathbf{x}\|_2^2 \in (1 \pm \epsilon)\|\mathbf{x}\|_2^2 \Rightarrow |\mathbf{x}^T(\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I})\mathbf{x}| \leq \epsilon\|\mathbf{x}\|_2^2 \quad \forall x \in \mathbb{R}^n \quad (20)$$

with probability at least $1 - 1/d^\tau$. Also, since $(1 - \epsilon)\mathbf{x}^T \mathbf{x} \leq \mathbf{x}^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{x} \leq (1 + \epsilon)\mathbf{x}^T \mathbf{x} \quad \forall x \in \mathbb{R}^n$ by the inequality above, replacing \mathbf{x} by $\mathbf{A}\mathbf{y}$, we get

$$(1 - \epsilon)\mathbf{y}^T \mathbf{A}^T \mathbf{A} \mathbf{y} \leq \mathbf{y}^T \mathbf{A}^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A} \mathbf{y} \leq (1 + \epsilon)\mathbf{y}^T \mathbf{A}^T \mathbf{A} \mathbf{y} \quad (21)$$

with probability at least $1 - 1/d^\tau$. Let \mathbf{y}_1 be the unit norm eigenvector corresponding to the minimum eigenvalue of $\hat{\mathbf{H}}_t = \mathbf{A}_t^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A}_t$. Since the above inequality is true for all \mathbf{y} , we have

$$\begin{aligned} \mathbf{y}_1^T \mathbf{A}_t^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A}_t \mathbf{y}_1 &\geq (1 - \epsilon)\mathbf{y}_1^T \mathbf{A}_t^T \mathbf{A}_t \mathbf{y}_1 \geq (1 - \epsilon)\lambda_{\min}(\mathbf{A}_t^T \mathbf{A}_t) = (1 - \epsilon)\lambda_{\min}(\nabla^2 f(\mathbf{w}_t)) \\ &\Rightarrow \lambda_{\min}(\hat{\mathbf{H}}_t) \geq (1 - \epsilon)\lambda_{\min}(\nabla^2 f(\mathbf{w}_t)) \end{aligned}$$

with probability at least $1 - 1/d^\tau$. Along similar lines, we can prove that $\lambda_{\max}(\hat{\mathbf{H}}_t) \leq (1 + \epsilon)\lambda_{\max}(\nabla^2 f(\mathbf{w}_t))$ with probability at least $1 - 1/d^\tau$ using the right hand inequality in (21). Together, these prove the first result.

In general, Eq. (20) implies that the eigenvalues of $(\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I})$ are in the set $[-\epsilon, \epsilon]$. Thus, all the eigenvalues of $\mathbf{A}_t^T (\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I}) \mathbf{A}_t$ are in the set $[-\epsilon\lambda_{\max}(\nabla^2 f(\mathbf{w}_t)), \epsilon\lambda_{\max}(\nabla^2 f(\mathbf{w}_t))]$. Also, we can write

$$\hat{\mathbf{H}}_t = \mathbf{A}_t^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A}_t = \mathbf{A}_t^T \mathbf{A}_t + \mathbf{A}_t^T (\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I}) \mathbf{A}_t.$$

Now, applying Weyl's inequality (see [58], Section 1.3) on symmetric matrices $\hat{\mathbf{H}}_t = \mathbf{A}_t^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A}_t$, $\nabla^2 f(\mathbf{w}_t) = \mathbf{A}_t^T \mathbf{A}_t$ and $\mathbf{A}_t^T (\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I}) \mathbf{A}_t$, we get

$$\lambda_i(\nabla^2 f(\mathbf{w}_t)) - \epsilon\lambda_{\max}(\nabla^2 f(\mathbf{w}_t)) \leq \lambda_i(\hat{\mathbf{H}}_t) \leq \lambda_i(\nabla^2 f(\mathbf{w}_t)) + \epsilon\lambda_{\max}(\nabla^2 f(\mathbf{w}_t)),$$

which proves the second result. □

6.1 Proof of Theorem 3.1

Let's define $\mathbf{w}_\tau = \mathbf{w}_t + \tau \mathbf{p}_t$, where the descent direction \mathbf{p}_t is given by $\mathbf{p}_t = -\hat{\mathbf{H}}_t^{-1} \nabla f(\mathbf{w}_t)$. Also, from Lemma 6.1, we have

$$\lambda_{\min}(\hat{\mathbf{H}}_t) \geq (1 - \epsilon) \lambda_{\min}(\nabla^2 f(\mathbf{w}_t)) \text{ and } \lambda_{\max}(\hat{\mathbf{H}}_t) \leq (1 + \epsilon) \lambda_{\max}(\nabla^2 f(\mathbf{w}_t)),$$

with probability at least $1 - 1/d^r$. Using the above inequalities and the fact that $f(\cdot)$ is k -strongly convex and M -smooth, we get

$$(1 - \epsilon)k\mathbf{I} \preceq \hat{\mathbf{H}}_t \preceq (1 + \epsilon)M\mathbf{I}, \quad (22)$$

with probability at least $1 - 1/d^r$.

Next, we show that there exists an $\alpha > 0$ such that the Armijo line search condition in (5) is satisfied. From the smoothness of $f(\cdot)$, we get (see [59], Theorem 2.1.5)

$$\begin{aligned} f(\mathbf{w}_\alpha) - f(\mathbf{w}_t) &\leq (\mathbf{w}_\alpha - \mathbf{w}_t)^T \nabla f(\mathbf{w}_t) + \frac{M}{2} \|\mathbf{w}_\alpha - \mathbf{w}_t\|^2, \\ &= \alpha \mathbf{p}_t^T \nabla f(\mathbf{w}_t) + \alpha^2 \frac{M}{2} \|\mathbf{p}_t\|^2. \end{aligned}$$

Now, for \mathbf{w}_α to satisfy the Armijo rule, α should satisfy

$$\begin{aligned} \alpha \mathbf{p}_t^T \nabla f(\mathbf{w}_t) + \alpha^2 \frac{M}{2} \|\mathbf{p}_t\|^2 &\leq \alpha \beta \mathbf{p}_t^T \nabla f(\mathbf{w}_t) \\ \Rightarrow \alpha \frac{M}{2} \|\mathbf{p}_t\|^2 &\leq (\beta - 1) \mathbf{p}_t^T \nabla f(\mathbf{w}_t) \\ \Rightarrow \alpha \frac{M}{2} \|\mathbf{p}_t\|^2 &\leq (1 - \beta) \mathbf{p}_t^T \hat{\mathbf{H}}_t \mathbf{p}_t, \end{aligned}$$

where the last inequality follows from the definition of \mathbf{p}_t . Now, using the lower bound from (22), \mathbf{w}_α satisfies Armijo rule for all

$$\alpha \leq \frac{2(1 - \beta)(1 - \epsilon)k}{M}.$$

Hence, we can always find an $\alpha_t \geq \frac{2(1 - \beta)(1 - \epsilon)k}{M}$ using backtracking line search such that \mathbf{w}_{t+1} satisfies the Armijo condition, that is,

$$\begin{aligned} f(\mathbf{w}_{t+1}) - f(\mathbf{w}_t) &\leq \alpha_t \beta \mathbf{p}_t^T \nabla f(\mathbf{w}_t) \\ &= -\alpha_t \beta \nabla f(\mathbf{w}_t)^T \hat{\mathbf{H}}_t^{-1} \nabla f(\mathbf{w}_t) \\ &\leq -\frac{\alpha_t \beta}{\lambda_{\max}(\hat{\mathbf{H}}_t)} \|\nabla f(\mathbf{w}_t)\|^2 \end{aligned}$$

which in turn implies

$$f(\mathbf{w}_t) - f(\mathbf{w}_{t+1}) \geq \frac{\alpha_t \beta}{M(1 + \epsilon)} \|\nabla f(\mathbf{w}_t)\|^2 \quad (23)$$

with probability at least $1 - 1/d^r$. Here the last inequality follows from the bound in (22). Moreover, k -strong convexity of $f(\cdot)$ implies (see [59], Theorem 2.1.10)

$$f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq \frac{1}{2k} \|\nabla f(\mathbf{w}_t)\|^2.$$

Using the inequality from (23) in the above inequality, we get

$$\begin{aligned} f(\mathbf{w}_t) - f(\mathbf{w}_{t+1}) &\geq \frac{2\alpha_t\beta k}{M(1+\epsilon)} (f(\mathbf{w}_t) - f(\mathbf{w}^*)) \\ &\geq \rho (f(\mathbf{w}_t) - f(\mathbf{w}^*)), \end{aligned}$$

where $\rho = \frac{2\alpha_t\beta k}{M(1+\epsilon)}$. Rearranging, we get

$$f(\mathbf{w}_{t+1}) - f(\mathbf{w}^*) \leq (1 - \rho)(f(\mathbf{w}_t) - f(\mathbf{w}^*))$$

with probability at least $1 - 1/d^r$, which proves the desired result.

6.2 Proof of Theorem 3.2

According to OverSketched Newton update, \mathbf{w}_{t+1} is obtained by solving

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^T (\mathbf{w} - \mathbf{w}_t) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \hat{\mathbf{H}}_t (\mathbf{w} - \mathbf{w}_t) \right\}.$$

Thus, we have, for any $\mathbf{w} \in \mathbb{R}^d$,

$$\begin{aligned} &f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^T (\mathbf{w} - \mathbf{w}_t) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \hat{\mathbf{H}}_t (\mathbf{w} - \mathbf{w}_t), \\ &\geq f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^T (\mathbf{w}_{t+1} - \mathbf{w}_t) + \frac{1}{2} (\mathbf{w}_{t+1} - \mathbf{w}_t)^T \hat{\mathbf{H}}_t (\mathbf{w}_{t+1} - \mathbf{w}_t), \\ &\Rightarrow \nabla f(\mathbf{w}_t)^T (\mathbf{w} - \mathbf{w}_{t+1}) + \frac{1}{2} (\mathbf{w} - \mathbf{w}_t)^T \hat{\mathbf{H}}_t (\mathbf{w} - \mathbf{w}_t) - \frac{1}{2} (\mathbf{w}_{t+1} - \mathbf{w}_t)^T \hat{\mathbf{H}}_t (\mathbf{w}_{t+1} - \mathbf{w}_t) \geq 0, \\ &\Rightarrow \nabla f(\mathbf{w}_t)^T (\mathbf{w} - \mathbf{w}_{t+1}) + \frac{1}{2} \left[(\mathbf{w} - \mathbf{w}_t)^T \hat{\mathbf{H}}_t (\mathbf{w} - \mathbf{w}_{t+1}) + (\mathbf{w} - \mathbf{w}_{t+1})^T \hat{\mathbf{H}}_t (\mathbf{w}_{t+1} - \mathbf{w}_t) \right] \geq 0. \end{aligned}$$

Substituting \mathbf{w} by \mathbf{w}^* in the above expression and calling $\Delta_t = \mathbf{w}^* - \mathbf{w}_t$, we get

$$\begin{aligned} &-\nabla f(\mathbf{w}_t)^T \Delta_{t+1} + \frac{1}{2} \left[\Delta_{t+1}^T \hat{\mathbf{H}}_t (2\Delta_t - \Delta_{t+1}) \right] \geq 0, \\ &\Rightarrow \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_t - \nabla f(\mathbf{w}_t)^T \Delta_{t+1} \geq \frac{1}{2} \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_{t+1}. \end{aligned}$$

Now, due to the optimality of \mathbf{w}^* , we have $\nabla f(\mathbf{w}^*)^T \Delta_{t+1} \geq 0$. Hence, we can write

$$\Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_t - (\nabla f(\mathbf{w}_t) - \nabla f(\mathbf{w}^*))^T \Delta_{t+1} \geq \frac{1}{2} \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_{t+1}.$$

Next, substituting $\nabla f(\mathbf{w}_t) - \nabla f(\mathbf{w}^*) = \left(\int_0^1 \nabla^2 f(\mathbf{w}^* + p(\mathbf{w}_t - \mathbf{w}^*)) dp \right) (\mathbf{w}_t - \mathbf{w}^*)$ in the above inequality, we get

$$\Delta_{t+1}^T (\hat{\mathbf{H}}_t - \nabla^2 f(\mathbf{w}_t)) \Delta_t + \Delta_{t+1}^T \left(\nabla^2 f(\mathbf{w}_t) - \int_0^1 \nabla^2 f(\mathbf{w}^* + p(\mathbf{w}_t - \mathbf{w}^*)) dp \right) \Delta_t \geq \frac{1}{2} \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_{t+1}.$$

Using Cauchy-Schwartz inequality in the LHS above, we get

$$\|\Delta_{t+1}\|_2 \|\Delta_t\|_2 \left(\|\hat{\mathbf{H}}_t - \nabla^2 f(\mathbf{w}_t)\|_2 + \int_0^1 \|\nabla^2 f(\mathbf{w}_t) - \nabla^2 f(\mathbf{w}^* + p(\mathbf{w}_t - \mathbf{w}^*))\|_2 dp \right) \geq \frac{1}{2} \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_{t+1}.$$

Now, using the L -Lipschitzness of $\nabla^2 f(\cdot)$ in the inequality above, we get

$$\begin{aligned} \frac{1}{2} \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_{t+1} &\leq \|\Delta_{t+1}\|_2 \|\Delta_t\|_2 \|\hat{\mathbf{H}}_t - \nabla^2 f(\mathbf{w}_t)\|_2 + \frac{L}{2} \|\Delta_{t+1}\|_2 \|\Delta_t\|_2^2 \int_0^1 (1-p) dp, \\ \Rightarrow \frac{1}{2} \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_{t+1} &\leq \|\Delta_{t+1}\|_2 \left(\|\Delta_t\|_2 \|\hat{\mathbf{H}}_t - \nabla^2 f(\mathbf{w}_t)\|_2 + \frac{L}{2} \|\Delta_t\|_2^2 \right). \end{aligned} \quad (24)$$

Note that for the positive definite matrix $\nabla^2 f(\mathbf{w}_t) = \mathbf{A}_t^T \mathbf{A}_t$, we have $\|\mathbf{A}_t\|_2^2 = \|\nabla^2 f(\mathbf{w}_t)\|_2$. Moreover,

$$\|\hat{\mathbf{H}}_t - \nabla^2 f(\mathbf{w}_t)\|_2 = \|\mathbf{A}_t^T (\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I}) \mathbf{A}_t\|_2 \leq \|\mathbf{A}_t\|_2^2 \|\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I}\|_2$$

Now, using Equation 20 from the proof of Lemma 6.1, we get $\|\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I}\|_2 = \lambda_{\max}(\mathbf{S}_t \mathbf{S}_t^T - \mathbf{I}) \leq \epsilon$. Using this to bound the RHS of (24), we have, with probability at least $1 - 1/d^T$,

$$\begin{aligned} \frac{1}{2} \Delta_{t+1}^T \hat{\mathbf{H}}_t \Delta_{t+1} &\leq \|\Delta_{t+1}\|_2 \left(\epsilon \|\nabla^2 f(\mathbf{w}_t)\|_2 \|\Delta_t\|_2 + \frac{L}{2} \|\Delta_t\|_2^2 \right) \\ \frac{1}{2} \|\mathbf{S}_t \mathbf{A}_t \Delta_{t+1}\|_2^2 &\leq \|\Delta_{t+1}\|_2 \left(\epsilon \|\nabla^2 f(\mathbf{w}_t)\|_2 \|\Delta_t\|_2 + \frac{L}{2} \|\Delta_t\|_2^2 \right), \end{aligned}$$

where the last inequality follows from $\hat{\mathbf{H}}_t = \mathbf{A}_t^T \mathbf{S}_t^T \mathbf{S}_t \mathbf{A}_t$. Now, since the sketch dimension $m = \Omega(d^{1+\mu}/\epsilon^2)$, using Eq. (20) from the proof of Lemma 1 in above inequality, we get, with probability at least $1 - 1/d^T$,

$$\begin{aligned} \frac{1}{2} (1 - \epsilon) \|\mathbf{A}_t \Delta_{t+1}\|_2^2 &\leq \|\Delta_{t+1}\|_2 \left(\epsilon \|\nabla^2 f(\mathbf{w}_t)\|_2 \|\Delta_t\|_2 + \frac{L}{2} \|\Delta_t\|_2^2 \right), \\ \Rightarrow \frac{1}{2} (1 - \epsilon) \Delta_{t+1}^T \nabla^2 f(\mathbf{w}_t) \Delta_{t+1} &\leq \|\Delta_{t+1}\|_2 \left(\epsilon \|\nabla^2 f(\mathbf{w}_t)\|_2 \|\Delta_t\|_2 + \frac{L}{2} \|\Delta_t\|_2^2 \right). \end{aligned}$$

Now, since γ and β are the minimum and maximum eigenvalues of $\nabla^2 f(\mathbf{w}^*)$, we get

$$\frac{1}{2} (1 - \epsilon) \|\Delta_{t+1}\|_2 (\gamma - L \|\Delta_t\|_2) \leq \epsilon (\beta + L \|\Delta_t\|_2) \|\Delta_t\|_2 + \frac{L}{2} \|\Delta_t\|_2^2$$

by the Lipschitzness of $\nabla^2 f(\mathbf{w})$, that is, $|\Delta_{t+1}^T (\nabla^2 f(\mathbf{w}_t) - \nabla^2 f(\mathbf{w}^*)) \Delta_{t+1}| \leq L \|\Delta_t\|_2 \|\Delta_{t+1}\|_2^2$. Rearranging for $\epsilon \leq \gamma/(8\beta) < 1/2$, we get

$$\|\Delta_{t+1}\|_2 \leq \frac{4\epsilon\beta}{\gamma - L \|\Delta_t\|_2} \|\Delta_t\|_2 + \frac{5L}{2(\gamma - L \|\Delta_t\|_2)} \|\Delta_t\|_2^2, \quad (25)$$

with probability at least $1 - 1/d^T$.

Let ξ_T be the event that the above inequality (in (25)) is true for $t = 0, 1, \dots, T$. Thus,

$$\mathbb{P}(\xi_T) \geq \left(1 - \frac{1}{d^T}\right)^T \geq 1 - \frac{T}{d^T},$$

where the second inequality follows from Bernoulli's inequality. Next, assuming that the event ξ_T holds, we prove that $\|\Delta_t\|_2 \leq \gamma/5L$ using induction. We can verify the base case using the initialization condition, i.e. $\|\Delta_0\|_2 \leq \gamma/8L$. Now, assuming that $\|\Delta_{t-1}\|_2 \leq \gamma/5L$ and using it in the inequality (25), we get

$$\begin{aligned} \|\Delta_t\|_2 &\leq \frac{4\epsilon\beta}{\gamma} \times \frac{\gamma}{5L} + \frac{5L}{2\gamma} \times \frac{\gamma^2}{25L^2} \\ &= \frac{4\epsilon\beta}{5L} + \frac{\gamma}{10L} \\ &\leq \frac{\gamma}{L} \left(\frac{1}{10} + \frac{1}{10} \right) \leq \frac{\gamma}{5L}, \end{aligned}$$

where the last inequality uses the fact that $\epsilon \leq \gamma/(8\beta)$. Thus, by induction,

$$\|\Delta_t\|_2 \leq \gamma/(5L) \quad \forall t \geq 0 \quad \text{with probability at least } 1 - T/d^T.$$

Using this in (25), we get the desired result, that is,

$$\|\Delta_{t+1}\|_2 \leq \frac{5\epsilon\beta}{\gamma} \|\Delta_t\|_2 + \frac{25L}{8\gamma} \|\Delta_t\|_2^2,$$

with probability at least $1 - T/d^T$.

6.3 Proof of Theorem 3.3

Let us define a few short notations for convenience. Say $\mathbf{g}_t = \nabla f(\mathbf{w}_t)$ and $\mathbf{H}_t = \nabla^2 f(\mathbf{w}_t) = \mathbf{A}_t^T \mathbf{A}_t$, and we know that $\hat{\mathbf{H}}_t = \mathbf{A}_t^T \mathbf{S}_t \mathbf{S}_t^T \mathbf{A}_t$. Moreover, all the results with approximate Hessian $\hat{\mathbf{H}}_t$ hold with probability $1 - 1/d^T$. We skip its mention in most of the proof for brevity. The following lemmas will assist us in the proof.

Lemma 6.2. *M-smoothness of $f(\cdot)$ and L-Lipchitzness of $\nabla^2 f(\cdot)$ imply*

$$\|\nabla^2 f(\mathbf{y})\nabla f(\mathbf{y}) - \nabla^2 f(\mathbf{x})\nabla f(\mathbf{x})\| \leq Q\|\mathbf{y} - \mathbf{x}\| \quad (26)$$

for all $\mathbf{x} \in \mathbb{R}^d$, $Q = (L\delta + M^2)$, where $\mathbf{y} \in \mathcal{Y}$, where $\mathcal{Y} = \{\mathbf{y} \in \mathbb{R}^d \mid \|\nabla f(\mathbf{y})\| \leq \delta\}$ and $\delta > 0$ is some constant.

Proof. We have

$$\begin{aligned} LHS &= \|\nabla^2 f(\mathbf{y})\nabla f(\mathbf{y}) - \nabla^2 f(\mathbf{x})\nabla f(\mathbf{x})\| \\ &= \|\nabla^2 f(\mathbf{y}) - \nabla^2 f(\mathbf{x})\|\nabla f(\mathbf{y}) + \nabla^2 f(\mathbf{x})(\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}))\| \end{aligned}$$

By applying triangle inequality and Cauchy-Schwarz to above equation, we get

$$LHS \leq \|\nabla^2 f(\mathbf{y}) - \nabla^2 f(\mathbf{x})\|_2 \|\nabla f(\mathbf{y})\| + \|\nabla^2 f(\mathbf{x})\|_2 \|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\|$$

From the smoothness of $f(\cdot)$, that is, Lipschitzness of gradient, we get $\|\nabla^2 f(\mathbf{x})\|_2 \leq M \forall x \in \mathbb{R}^d$. Additionally, using Lipschitzness of Hessian, we get

$$\begin{aligned} LHS &\leq (L\|\nabla f(\mathbf{y})\| + M^2)\|\mathbf{y} - \mathbf{x}\| \\ &\leq (L\delta + M^2)\|\mathbf{y} - \mathbf{x}\| \end{aligned}$$

for $\mathbf{y} \in \mathcal{Y}$. This proves the desired result. \square

Lemma 6.3. Let $\mathbf{A}^T = \mathbf{U}\sqrt{\Sigma}\mathbf{V}^T$ and $\mathbf{A}^T\mathbf{S}_t = \hat{\mathbf{U}}\sqrt{\hat{\Sigma}}\hat{\mathbf{V}}^T$ be the truncated Singular Value Decompositions (SVD) of \mathbf{A}^T and $\mathbf{A}^T\mathbf{S}_t$, respectively. Thus, $\mathbf{H}_t = \mathbf{U}\Sigma\mathbf{U}^T$ and $\hat{\mathbf{H}}_t = \hat{\mathbf{U}}\hat{\Sigma}\hat{\mathbf{U}}^T$. Then, for all $\mathbf{g} \in \mathbb{R}^d$, we have

$$\|\hat{\mathbf{U}}^T \mathbf{g}\|^2 \geq \frac{(1-\epsilon)\eta}{M(1+\epsilon)} \|\mathbf{U}^T \mathbf{g}\|^2, \quad (27)$$

where η is defined in Assumption (5).

Proof. For all $\mathbf{g} \in \mathbb{R}^d$, using the fact that $\mathbf{A} = \mathbf{V}\sqrt{\Sigma}\mathbf{U}^T$, we get

$$\begin{aligned} \|\mathbf{A}\mathbf{g}\|^2 &= (\mathbf{U}^T \mathbf{g})^T \Sigma (\mathbf{U}^T \mathbf{g}) \\ &\geq \lambda_{\min}(\Sigma) \|\mathbf{U}^T \mathbf{g}\|^2 \\ &\geq \eta \|\mathbf{U}^T \mathbf{g}\|^2, \end{aligned} \quad (28)$$

where the last inequality uses Assumption (5). In a similar fashion, we can obtain

$$\begin{aligned} \|\mathbf{S}_t^T \mathbf{A}\mathbf{g}\|^2 &= (\hat{\mathbf{U}}^T \mathbf{g})^T \hat{\Sigma} (\hat{\mathbf{U}}^T \mathbf{g}) \\ &\leq \lambda_{\max}(\hat{\Sigma}) \|\hat{\mathbf{U}}^T \mathbf{g}\|^2 \\ &\leq M(1+\epsilon) \|\hat{\mathbf{U}}^T \mathbf{g}\|^2, \end{aligned} \quad (29)$$

where the last inequality uses M -smoothness of $f(\cdot)$ and Lemma 6.1. Also, from the subspace embedding property of \mathbf{S}_t (see Lemma 6.1), we have

$$\|\mathbf{S}_t^T \mathbf{A}\mathbf{g}\|^2 \geq (1-\epsilon) \|\mathbf{A}\mathbf{g}\|^2.$$

Now, using the above inequality and Eqs. (28) and (29), we get

$$\|\hat{\mathbf{U}}^T \mathbf{g}\|^2 \geq \frac{(1-\epsilon)\eta}{M(1+\epsilon)} \|\mathbf{U}^T \mathbf{g}\|^2, \quad (30)$$

which is the desired result. \square

Now we are ready to prove Theorem 3.3. Let $\mathbf{H}_t = \mathbf{U}\Sigma\mathbf{U}^T$ and $\hat{\mathbf{H}}_t = \hat{\mathbf{U}}\hat{\Sigma}\hat{\mathbf{U}}^T$ be the truncated SVDs of \mathbf{H}_t and $\hat{\mathbf{H}}_t$, respectively. Also, let α_t be the step-size obtained using line-search in (6) in the t -th iteration. Thus, Eq. (6) with the update direction $\mathbf{p}_t = -\hat{\mathbf{H}}_t^\dagger \mathbf{g}_t$ implies

$$\begin{aligned} \|\mathbf{g}_{t+1}\|^2 &\leq \|\mathbf{g}_t\|^2 - 2\beta\alpha_t \langle \hat{\mathbf{H}}_t \mathbf{g}_t, \hat{\mathbf{H}}_t^\dagger \mathbf{g}_t \rangle \\ &= \|\mathbf{g}_t\|^2 - 2\beta\alpha_t \|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2, \end{aligned} \quad (31)$$

where the last equality uses the fact that $\hat{\mathbf{H}}_t^\dagger$ can be expressed as $\hat{\mathbf{H}}_t^\dagger = \hat{\mathbf{U}}\hat{\Sigma}^{-1}\hat{\mathbf{U}}^T$. Note that Lemma 6.2 implies that the function $\|\nabla f(\mathbf{y})\|^2/2$ is smooth for all $\mathbf{y} \in \mathcal{Y}$, where $\mathcal{Y} = \{\mathbf{y} \in \mathbb{R}^d \mid \|\nabla f(\mathbf{y})\| \leq \delta\}$. Smoothness in turn implies the following property (see [59], Theorem 2.1.10)

$$\frac{1}{2}\|\nabla f(\mathbf{y})\|^2 \leq \frac{1}{2}\|\nabla f(\mathbf{x})\|^2 + \langle \nabla^2 f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2}Q\|\mathbf{y} - \mathbf{x}\|^2 \quad \forall \mathbf{x}, \mathbf{y} \in (Y), \quad (32)$$

where $Q = L\delta + M^2$. We take $\delta = \|\nabla f(\mathbf{w}_0)\|$ where \mathbf{w}_0 is the initial point of our algorithm. Due to line-search condition in (6), it holds that $\|\nabla f(\mathbf{w}_t)\| \leq \|\nabla f(\mathbf{w}_0)\| \quad \forall t > 0$. Thus, substituting $\mathbf{x} = \mathbf{w}_t$ and $\mathbf{y} = \mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{p}_t$, we get

$$\begin{aligned} \frac{1}{2}\|\mathbf{g}_{t+1}\|^2 &\leq \frac{1}{2}\|\mathbf{g}_t\|^2 + \langle \mathbf{H}_t \mathbf{g}_t, \alpha_t \mathbf{p}_t \rangle + \frac{1}{2}Q\alpha^2\|\mathbf{p}_t\|^2 \\ \Rightarrow \|\mathbf{g}_{t+1}\|^2 &\leq \|\mathbf{g}_t\|^2 + \langle 2\mathbf{H}_t \mathbf{g}_t, \alpha_t \mathbf{p}_t \rangle + Q\alpha^2\|\mathbf{p}_t\|^2, \end{aligned} \quad (33)$$

where

$$Q = L\|\nabla f(\mathbf{w}_0)\| + M^2.$$

Also, since the minimum non-zero eigenvalue of $\mathbf{H}_t \geq \eta$ from Assumption (5), the minimum non-zero eigenvalue of $\hat{\mathbf{H}}_t$ is at least $\eta - \epsilon M$ from Lemma 6.1. Thus,

$$\|\hat{\mathbf{H}}_t^\dagger\|_2 \leq 1/(\eta - \epsilon M). \quad (34)$$

Moreover,

$$\|\mathbf{p}_t\| = \|\hat{\mathbf{H}}_t^\dagger \mathbf{g}_t\| \leq \|\hat{\mathbf{H}}_t^\dagger\|_2 \|\mathbf{g}_t\| \leq \frac{\|\mathbf{g}_t\|}{(\eta - \epsilon M)}. \quad (35)$$

Using this in (33), we get

$$\|\mathbf{g}_{t+1}\|^2 \leq \|\mathbf{g}_t\|^2 - 2\alpha_t \langle \mathbf{H}_t \mathbf{g}_t, \hat{\mathbf{H}}_t^\dagger \mathbf{g}_t \rangle + Q\alpha^2 \frac{\|\mathbf{g}_t\|^2}{(\eta - \epsilon M)^2}. \quad (36)$$

Now,

$$\begin{aligned} -\langle \mathbf{H}_t \mathbf{g}_t, \hat{\mathbf{H}}_t^\dagger \mathbf{g}_t \rangle &= -\langle \hat{\mathbf{H}}_t \mathbf{g}_t, \hat{\mathbf{H}}_t^\dagger \mathbf{g}_t \rangle + \langle (\hat{\mathbf{H}}_t - \mathbf{H}_t) \mathbf{g}_t, \hat{\mathbf{H}}_t^\dagger \mathbf{g}_t \rangle \\ \Rightarrow -\langle \mathbf{H}_t \mathbf{g}_t, \hat{\mathbf{H}}_t^\dagger \mathbf{g}_t \rangle &\leq -\|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 + \|\mathbf{g}_t\|^2 \|\hat{\mathbf{H}}_t - \mathbf{H}_t\|_2 \|\hat{\mathbf{H}}_t^\dagger\|_2, \end{aligned}$$

where the last inequality is obtained by applying the triangle inequality and Cauchy-Schwartz inequality. This can be further simplified using Lemma 6.1 and Eq. (34) as

$$-\langle \mathbf{H}_t \mathbf{g}_t, \hat{\mathbf{H}}_t^\dagger \mathbf{g}_t \rangle \leq -\|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 + \frac{\epsilon M}{(\eta - M\epsilon)} \|\mathbf{g}_t\|^2$$

Using the above in Eq. (36), we get

$$\|\mathbf{g}_{t+1}\|^2 \leq \|\mathbf{g}_t\|^2 + 2\alpha_t(-\|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 + \frac{\epsilon M}{(\eta - M\epsilon)} \|\mathbf{g}_t\|^2) + Q\alpha_t^2 \frac{\|\mathbf{g}_t\|^2}{(\eta - \epsilon M)^2} \quad (37)$$

Note that the upper bound in Eq. (37) always holds. Also, we want the inequality in (31) to hold for some $\alpha_t > 0$. Therefore, we want α_t to satisfy the following (and hope that it is always satisfied for some $\alpha_t > 0$)

$$\begin{aligned} \|\mathbf{g}_t\|^2 + 2\alpha_t(-\|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 + \frac{\epsilon M}{(\eta - M\epsilon)} \|\mathbf{g}_t\|^2) + Q\alpha_t^2 \frac{\|\mathbf{g}_t\|^2}{(\eta - \epsilon M)^2} &\leq \|\mathbf{g}_t\|^2 - 2\beta\alpha_t \|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 \\ \Rightarrow Q\alpha_t^2 \frac{\|\mathbf{g}_t\|^2}{(\eta - \epsilon M)^2} &\leq 2\alpha_t \left[(1 - \beta) \|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 - \frac{\epsilon M}{(\eta - M\epsilon)} \|\mathbf{g}_t\|^2 \right] \\ \Rightarrow \alpha_t &\leq \frac{2(\eta - \epsilon M)^2}{Q} \left[(1 - \beta) \frac{\|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2}{\|\mathbf{g}_t\|^2} - \frac{\epsilon M}{(\eta - M\epsilon)} \right]. \end{aligned} \quad (38)$$

Thus, any α_t satisfying the above inequality would definitely satisfy the line-search termination condition in

Now, using Lemma 6.3 and Assumption (6), we have

$$\|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 \geq \frac{(1 - \epsilon)\eta}{M(1 + \epsilon)} \|\mathbf{U}_t^T \mathbf{g}_t\|^2 \geq \frac{(1 - \epsilon)\eta}{M(1 + \epsilon)} \nu \|\mathbf{g}_t\|^2. \quad (39)$$

Using the above in Eq. (38) to find an iteration independent bound on α_t , we get

$$\alpha_t \leq \frac{2(\eta - \epsilon M)^2}{Q} \left[(1 - \beta)\nu - \frac{\epsilon M}{(\eta - M\epsilon)} \right]. \quad (40)$$

Hence, line-search will always terminate for all α_t that satisfy the above inequality. This can be further simplified by assuming that ϵ is small enough such that $\epsilon < \eta/2M$. Thus, $\eta - M\epsilon > \eta/2$, and the sufficient condition on α_t in (40) becomes

$$\alpha_t \leq \frac{\eta}{2Q} [(1 - \beta)\nu\eta - 2\epsilon M]. \quad (41)$$

For a positive α_t to always exist, we require ϵ to further satisfy

$$\epsilon \leq \frac{(1 - \beta)\nu\eta}{2M}, \quad (42)$$

which is tighter than the initial upper bound on ϵ . Now, Eqs. (31) and (39) proves the desired result, that is

$$\|\mathbf{g}_{t+1}\|^2 \leq \|\mathbf{g}_t\|^2 - 2\beta\alpha_t \|\hat{\mathbf{U}}_t^T \mathbf{g}_t\|^2 \leq \left(1 - 2\beta\alpha_t \nu \frac{(1 - \epsilon)\eta}{M(1 + \epsilon)} \right) \|\mathbf{g}_t\|^2.$$

Thus, OverSketched Newton for the weakly-convex case enjoys a uniform linear convergence rate of decrease in $\|\nabla f(\mathbf{w})\|^2$.

7 Conclusions

We proposed OverSketched Newton, a straggler-resilient distributed optimization algorithm for serverless systems. It uses the idea of matrix sketching from RandNLA to find an approximate second-order update in each iteration. We proved that OverSketched Newton has a local linear-quadratic convergence rate for the strongly-convex case, where the dependence on the linear term can be made to diminish by increasing the sketch dimension. Moreover, it has a linear global convergence rate for weakly-convex functions. By exploiting the massive scalability of serverless systems, OverSketched Newton produces a global approximation of the second-order update. Empirically, this translates into faster convergence than state-of-the-art distributed optimization algorithms on AWS Lambda.

Acknowledgments

This work was partially supported by NSF grants CCF-1748585 and CNS-1748692 to SK, and NSF grants CCF-1704967 and CCF-0939370 (Center for Science of Information) to TC, and ARO, DARPA, NSF, and ONR grants to MWM, and NSF Grant CCF-1703678 to KR. The authors would like to additionally thank Fred-Roosta and Yang Liu for helpful discussions regarding our proof techniques and AWS for providing promotional cloud credits for research.

References

- [1] I. Baldini, P. C. Castro, K. S.-P. Chang, P. Cheng, S. J. Fink, V. Ishakian, N. Mitchell, V. Muthusamy, R. M. Rabbah, A. Slominski, and P. Suter, “Serverless computing: Current trends and open problems,” *CoRR*, vol. abs/1706.03178, 2017.
- [2] J. Spillner, C. Mateos, and D. A. Monge, “Faaster, better, cheaper: The prospect of serverless scientific computing and HPC,” in *Latin American High Performance Computing Conference*, pp. 154–168, Springer, 2017.
- [3] E. Jonas, Q. Pu, S. Venkataraman, I. Stoica, and B. Recht, “Occupy the cloud: distributed computing for the 99%,” in *Proceedings of the 2017 Symposium on Cloud Computing*, pp. 445–451, ACM, 2017.
- [4] V. Shankar, K. Krauth, Q. Pu, E. Jonas, S. Venkataraman, I. Stoica, B. Recht, and J. Ragan-Kelley, “numpywren: serverless linear algebra,” *ArXiv e-prints*, Oct. 2018.
- [5] Technavio, “Serverless architecture market by end-users and geography - global forecast 2019-2023.” <https://www.technavio.com/report/serverless-architecture-market-industry-analysis>.

- [6] E. Jonas, J. Schleier-Smith, V. Sreekanti, C.-C. Tsai, A. Khandelwal, Q. Pu, V. Shankar, J. Carreira, K. Krauth, N. Yadwadkar, *et al.*, “Cloud programming simplified: A berkeley view on serverless computing,” *arXiv preprint arXiv:1902.03383*, 2019.
- [7] L. Feng, P. Kudva, D. D. Silva, and J. Hu, “Exploring serverless computing for neural network training,” in *2018 IEEE 11th International Conference on Cloud Computing (CLOUD)*, vol. 00, pp. 334–341, Jul 2018.
- [8] V. Ishakian, V. Muthusamy, and A. Slominski, “Serving deep learning models in a serverless platform,” *arXiv e-prints*, p. arXiv:1710.08460, Oct. 2017.
- [9] A. Aytekin and M. Johansson, “Harnessing the Power of Serverless Runtimes for Large-Scale Optimization,” *arXiv e-prints*, p. arXiv:1901.03161, Jan. 2019.
- [10] H. Wang, D. Niu, and B. Li, “Distributed machine learning with a serverless architecture,” in *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, pp. 1288–1296, IEEE, 2019.
- [11] L. Feng, P. Kudva, D. Da Silva, and J. Hu, “Exploring serverless computing for neural network training,” in *2018 IEEE 11th International Conference on Cloud Computing (CLOUD)*, pp. 334–341, IEEE, 2018.
- [12] J. Carreira, P. Fonseca, A. Tumanov, A. Zhang, and R. Katz, “Cirrus: a serverless framework for end-to-end ml workflows,” in *Proceedings of the ACM Symposium on Cloud Computing*, pp. 13–24, 2019.
- [13] J. Dean and L. A. Barroso, “The tail at scale,” *Commun. ACM*, vol. 56, pp. 74–80, Feb. 2013.
- [14] T. Hoefler, T. Schneider, and A. Lumsdaine, “Characterizing the influence of system noise on large-scale applications by simulation,” in *Proc. of the ACM/IEEE Int. Conf. for High Perf. Comp., Networking, Storage and Analysis*, pp. 1–11, 2010.
- [15] J. M. Hellerstein, J. Faleiro, J. E. Gonzalez, J. Schleier-Smith, V. Sreekanti, A. Tumanov, and C. Wu, “Serverless computing: One step forward, two steps back,” *arXiv preprint arXiv:1812.03651*, 2018.
- [16] J. Nocedal and S. Wright, *Numerical optimization*. Springer Science & Business Media, 2006.
- [17] L. Bottou, F. E. Curtis, and J. Nocedal, “Optimization methods for large-scale machine learning,” *Siam Review*, vol. 60, no. 2, pp. 223–311, 2018.
- [18] N. S. Wadia, D. Duckworth, S. S. Schoenholz, E. Dyer, and J. Sohl-Dickstein, “Whitening and second order optimization both destroy information about the dataset, and can make generalization impossible,” *arXiv e-prints*, p. arXiv:2008.07545, Aug. 2020.
- [19] F. Roosta-Khorasani and M. W. Mahoney, “Sub-Sampled Newton Methods I: Globally Convergent Algorithms,” *arXiv e-prints*, p. arXiv:1601.04737, Jan. 2016.

- [20] F. Roosta-Khorasani and M. W. Mahoney, “Sub-Sampled Newton Methods II: Local Convergence Rates,” *arXiv e-prints*, p. arXiv:1601.04738, Jan. 2016.
- [21] P. Xu, F. Roosta, and M. W. Mahoney, “Newton-type methods for non-convex optimization under inexact hessian information,” 2017.
- [22] F. Roosta, Y. Liu, P. Xu, and M. W. Mahoney, “Newton-MR: Newton’s method without smoothness or convexity,” *arXiv preprint arXiv:1810.00303*, 2018.
- [23] M. Pilanci and M. J. Wainwright, “Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence,” *SIAM Jour. on Opt.*, vol. 27, pp. 205–245, 2017.
- [24] S. Wang, F. Roosta-Khorasani, P. Xu, and M. W. Mahoney, “GIANT: Globally improved approximate Newton method for distributed optimization,” in *Advances in Neural Information Processing Systems*, pp. 2332–2342, 2018.
- [25] C.-H. Fang, S. B. Kylasa, F. Roosta-Khorasani, M. W. Mahoney, and A. Grama, “Distributed Second-order Convex Optimization,” *ArXiv e-prints*, July 2018.
- [26] Z. Yao, A. Gholami, K. Keutzer, and M. Mahoney, “Pyhessian: Neural networks through the lens of the hessian,” *arXiv preprint arXiv:1912.07145*, 2019.
- [27] Z. Yao, A. Gholami, S. Shen, K. Keutzer, and M. W. Mahoney, “Adahessian: An adaptive second order optimizer for machine learning,” *arXiv preprint arXiv:2006.00719*, 2020.
- [28] R. Anil, V. Gupta, T. Koren, K. Regan, and Y. Singer, “Second order optimization made practical,” *arXiv preprint arXiv:2002.09018*, 2020.
- [29] D. P. Woodruff, “Sketching as a tool for numerical linear algebra,” *Found. Trends Theor. Comput. Sci.*, vol. 10, pp. 1–157, 2014.
- [30] M. W. Mahoney, *Randomized algorithms for matrices and data*. Foundations and Trends in Machine Learning, Boston: NOW Publishers, 2011.
- [31] A. Gittens, A. Devarakonda, E. Racah, M. Ringenbun, L. Gerhardt, J. Kottalam, J. Liu, K. Maschhoff, S. Canon, J. Chhugani, *et al.*, “Matrix factorizations at scale: A comparison of scientific data analytics in spark and c+ mpi using three case studies,” in *2016 IEEE International Conference on Big Data (Big Data)*, pp. 204–213, IEEE, 2016.
- [32] V. Gupta, S. Wang, T. Courtade, and K. Ramchandran, “Oversketch: Approximate matrix multiplication for the cloud,” *IEEE International Conference on Big Data, Seattle, WA, USA*, 2018.
- [33] K. Lee, M. Lam, R. Pedarsani, D. Papailiopoulos, and K. Ramchandran, “Speeding up distributed machine learning using codes,” *IEEE Transactions on Information Theory*, vol. 64, no. 3, pp. 1514–1529, 2018.

- [34] T. Baharav, K. Lee, O. Ocal, and K. Ramchandran, “Straggler-proofing massive-scale distributed matrix multiplication with d-dimensional product codes,” in *IEEE Int. Sym. on Information Theory (ISIT)*, IEEE, 2018.
- [35] V. Gupta, D. Carrano, Y. Yang, V. Shankar, T. Courtade, and K. Ramchandran, “Serverless straggler mitigation using local error-correcting codes,” *IEEE International Conference on Distributed Computing and Systems (ICDCS)*, Singapore, 2020.
- [36] C.-C. Chang and C.-J. Lin, “Libsvm: a library for support vector machines,” *ACM transactions on intelligent systems and technology (TIST)*, vol. 2, no. 3, p. 27, 2011.
- [37] R. Tandon, Q. Lei, A. G. Dimakis, and N. Karampatziakis, “Gradient coding: Avoiding stragglers in distributed learning,” in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70, pp. 3368–3376, PMLR, 2017.
- [38] J. Dean and S. Ghemawat, “Mapreduce: Simplified data processing on large clusters,” *Commun. ACM*, vol. 51, pp. 107–113, Jan. 2008.
- [39] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica, “Spark: Cluster computing with working sets,” in *Proceedings of the 2Nd USENIX Conference on Hot Topics in Cloud Computing*, pp. 10–10, 2010.
- [40] Q. Yu, M. Maddah-Ali, and S. Avestimehr, “Polynomial codes: an optimal design for high-dimensional coded matrix multiplication,” in *Advances in Neural Inf. Processing Systems 30*, pp. 4403–4413, 2017.
- [41] Y. Yang, P. Grover, and S. Kar, “Coded distributed computing for inverse problems,” in *Advances in Neural Information Processing Systems 30*, pp. 709–719, Curran Associates, Inc., 2017.
- [42] O. Shamir, N. Srebro, and T. Zhang, “Communication-efficient distributed optimization using an approximate Newton-type method,” in *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32, ICML’14*, pp. II–1000–II–1008, JMLR.org, 2014.
- [43] Y. Zhang and X. Lin, “Disco: Distributed optimization for self-concordant empirical loss,” in *Proceedings of the 32nd International Conference on Machine Learning* (F. Bach and D. Blei, eds.), vol. 37 of *Proceedings of Machine Learning Research*, (Lille, France), pp. 362–370, PMLR, 07–09 Jul 2015.
- [44] S. J. Reddi, A. Hefny, S. Sra, B. Pöczos, and A. Smola, “On variance reduction in stochastic gradient descent and its asynchronous variants,” in *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2, NIPS’15*, (Cambridge, MA, USA), pp. 2647–2655, MIT Press, 2015.
- [45] C. Duenner, A. Lucchi, M. Gargiani, A. Bian, T. Hofmann, and M. Jaggi, “A distributed second-order algorithm you can trust,” in *Proceedings of the 35th International Conference on Machine Learning*, vol. 80, pp. 1358–1366, PMLR, 10–15 Jul 2018.

- [46] V. Smith, S. Forte, C. Ma, M. Takác, M. I. Jordan, and M. Jaggi, “Cocoa: A general framework for communication-efficient distributed optimization,” *arXiv preprint arXiv:1611.02189*, 2016.
- [47] R. Bollapragada, R. Byrd, and J. Nocedal, “Exact and Inexact Subsampled Newton Methods for Optimization,” *arXiv e-prints*, p. arXiv:1609.08502, Sep 2016.
- [48] K. Lee, C. Suh, and K. Ramchandran, “High-dimensional coded matrix multiplication,” in *IEEE Int. Sym. on Information Theory (ISIT), 2017*, pp. 2418–2422, IEEE, 2017.
- [49] E. Solomonik and J. Demmel, “Communication-optimal parallel 2.5D matrix multiplication and LU factorization algorithms,” in *Proceedings of the 17th International Conference on Parallel Processing*, pp. 90–109, 2011.
- [50] R. A. van de Geijn and J. Watts, “Summa: Scalable universal matrix multiplication algorithm,” tech. rep., 1995.
- [51] S. Boyd and L. Vandenberghe, *Convex Optimization*. New York, NY, USA: Cambridge University Press, 2004.
- [52] A. S. Berahas, R. Bollapragada, and J. Nocedal, “An Investigation of Newton-Sketch and Subsampled Newton Methods,” *arXiv e-prints*, p. arXiv:1705.06211, May 2017.
- [53] Y. Liu and F. Roosta, “Stability analysis of Newton-MR under hessian perturbations,” *arXiv preprint arXiv:1909.06224*, 2019.
- [54] J. R. Shewchuk *et al.*, “An introduction to the conjugate gradient method without the agonizing pain,” 1994.
- [55] J. Levin, “Note on convergence of minres,” *Multivariate behavioral research*, vol. 23, no. 3, pp. 413–417, 1988.
- [56] G. Cohen, S. Afshar, J. Tapson, and A. van Schaik, “Emnist: an extension of mnist to handwritten letters,” *arXiv preprint arXiv:1702.05373*, 2017.
- [57] J. Nelson and H. L. Nguyen, “Osnap: Faster numerical linear algebra algorithms via sparser subspace embeddings,” in *2013 IEEE 54th Annual Symposium on Foundations of Computer Science*, pp. 117–126, Oct 2013.
- [58] T. Tao, *Topics in random matrix theory*, vol. 132. American Mathematical Soc., 2012.
- [59] Y. Nesterov, *Introductory Lectures on Convex Optimization: A Basic Course*. Springer Publishing Company, Incorporated, 1 ed., 2014.