

«МОСКОВСКИЙ ФИЗИКО-ТЕХНИЧЕСКИЙ ИНСТИТУТ (национальный
исследовательский университет)
ФИЗТЕХ-ШКОЛА ПРИКЛАДНОЙ МАТЕМАТИКИ И ИНФОРМАТИКИ
КАФЕДРА «ИНТЕЛЛЕКТУАЛЬНЫЕ СИСТЕМЫ»

Исламов Рустем Ильфакович

Распределенные методы второго порядка с быстрой скоростью сходимости и компрессией

03.03.01 — Прикладные математика и физика

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА БАКАЛАВРА

Научный руководитель:

д. ф.-м. н. Стрижов Вадим Викторович

Москва
2021

Аннотация

Данная бакалаврская диссертация основана на статье «Distributed Second Order Methods with Fast Rates and Compressed Communication» [15] за авторством Рустема Исламов, Шун Кяна и Питера Рихтарика.

Мы разработали несколько новых эффективных с точки зрения коммуникации методов второго порядка для распределенной оптимизации. Наш первый метод, NEWTON-STAR, является одним из вариантов метода Ньютона, от которого он наследует свою быструю локальную квадратичную скорость. Однако, в отличие от метода Ньютона, NEWTON-STAR имеет ту же стоимость коммуницирования за одну итерацию, что и градиентный спуск. Хотя этот метод непрактичен, поскольку он опирается на использование некоторых неизвестных параметров, характеризующих Гессиан целевой функции в оптимуме, он служит отправной точкой, которая позволяет нам проектировать практические варианты с доказанными теоретическими гарантиями сходимости. В частности, мы разработали стратегию, основанную на использовании стохастической разреженности для изучения неизвестных параметров итеративным способом, сохраняя эффективность с точки зрения коммуникаций. Применение этой стратегии к NEWTON-STAR приводит к нашему следующему методу, NEWTON-LEARN, для которого мы доказали локальные линейные и сверхлинейные скорости сходимости, не зависящие от числа обусловленности функции. Когда эти методы применимы, они имеют значительно более высокие скорости сходимости по сравнению с современными методами. Наши результаты подтверждаются экспериментальными результатами на реальных наборах данных и показывают улучшение на несколько порядков по сравнению с базовыми и современными методами с точки зрения эффективности коммуницирования.

Abstract

This bachelor thesis is based on paper “Distributed Second Order Methods with Fast Rates and Compressed Communication” [15] written by Rustem Islamov, Xun Qian, and Peter Richtárik.

We develop new communication-efficient second-order method for distributed optimization. Our first method, **NEWTON-STAR**, is a variant of Newton’s method from which it inherits its fast local quadratic rate. However, unlike Newton’s method, **NEWTON-STAR** enjoys the same per iteration communication cost as gradient descent. While this method is impractical as it relies on the use of certain unknown parameters characterizing the Hessian of the objective function at the optimum, it serves as the starting point which enables us design practical variants thereof with strong theoretical guarantees. In particular, we design a stochastic sparsification strategy for learning the unknown parameters in an iterative fashion in a communication efficient manner. Applying this strategy to **NEWTON-STAR** leads to our next method, **NEWTON-LEARN**, for which we prove local linear and superlinear rates independent of the condition number. When applicable, this method can have dramatically superior convergence behavior when compared to state-of-the-art methods. Our results are supported with experimental results on real datasets, and show several orders of magnitude improvement on baseline and state-of-the-art methods in terms of communication complexity.

Contents

1	Introduction	5
1.1	Distributed optimization	5
1.2	The curse of the condition number	6
1.3	Newton’s method to the rescue?	6
1.4	Contributions summary	7
1.5	Related work	8
2	Three Steps Towards an Efficient Distributed Newton Type Method	10
2.1	Naive distributed implementation of Newton’s method	10
2.2	A better implementation taking advantage of the structure of $\mathbf{H}_{ij}(x)$. . .	10
2.3	NEWTON-STAR: Newton’s method with a single Hessian	11
3	NEWTON-LEARN: Learning the Hessian and Local Convergence Theory	13
3.1	The main iteration	13
3.2	Learning the coefficients: the idea	13
3.3	Outline of fast local convergence theory	14
3.4	Compressed learning	14
3.5	NL (learning in the $\lambda > 0$ case)	15
3.5.1	The learning iteration and the NL algorithm	15
3.5.2	Theory	16
4	Experiments	17
4.1	Data sets and parameter settings	17
4.2	Compression operators	18
4.2.1	Random sparsification	18
4.2.2	Random dithering	18
4.2.3	Natural compression	19
4.2.4	Bernoulli compressor	19
4.3	Behavior of NL	19
4.4	Comparison of NL with Newton’s method	19
4.5	Comparison of NL with BFGS	20
4.6	Comparison of NL with ADIANA	20
4.7	Comparison of NL with DINGO	20
A	Proofs for NL (Section 3.5)	23
A.1	Lemma	23
A.2	Proof of Theorem 3.2	24
A.3	Proof of Lemma 3.3	27
	References	29

1 Introduction

The prevalent paradigm for training modern supervised machine learning models is based on (regularized) empirical risk minimization (ERM) [39], and the most commonly used optimization methods deployed for solving ERM problems belong to the class of stochastic first order methods [30, 36]. Since modern training data sets are very large and are becoming larger every year, it is increasingly harder to get by without relying on modern computing architectures which make efficient use of distributed computing. However, in order to develop efficient distributed methods, one has to keep in mind that communication among the different parallel workers (e.g. processors or compute nodes) is typically very slow, and almost invariably forms the main bottleneck in deployed optimization software and systems [3]. For this reason, further advances in the area of communication efficient distributed first order optimization methods for solving ERM problems are highly needed, and research in this area constitutes one of the most important fundamental endeavors in modern machine learning. Indeed, this research field is very active, and numerous advances have been made over the past decade [1, 4, 27, 38, 42, 46, 49].

1.1 Distributed optimization

We consider L2 regularized empirical risk minimization problems of the form

$$\min_{x \in \mathbb{R}^d} \left[P(x) := f(x) + \frac{\lambda}{2} \|x\|^2 \right], \quad (1)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a smooth¹ convex function of the “average of averages” structure

$$f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x), \quad f_i(x) := \frac{1}{m} \sum_{j=1}^m f_{ij}(x), \quad (2)$$

and $\lambda \geq 0$ is a regularization parameter. Here n is the number of parallel workers (nodes), and m is the number of training examples handled by each node². The value $f_{ij}(x)$ denotes the loss of the model parameterized by vector $x \in \mathbb{R}^d$ on the j^{th} example owned by the i^{th} node. This example is denoted as $a_{ij} \in \mathbb{R}^d$, and the corresponding loss function is $\varphi_{ij} : \mathbb{R} \rightarrow \mathbb{R}$, and hence we have

$$f_{ij}(x) := \varphi_{ij}(a_{ij}^\top x). \quad (3)$$

Thus, f represents the average loss/risk over all nm training datapoints, and problem (1) seeks to find the model whose (L2 regularized) empirical risk is minimized. We make the following assumption throughout the paper.

Assumption 1.1. *Problem (1) has at least one optimal solution x^* . For all i and j , the loss function $\varphi_{ij} : \mathbb{R} \rightarrow \mathbb{R}$ is γ -smooth, twice differentiable, and its second derivative $\varphi_{ij}'' : \mathbb{R} \rightarrow \mathbb{R}$ is ν -Lipschitz continuous.*

¹Function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is *smooth* if it is differentiable, and has L_ϕ Lipschitz gradient: $\|\nabla\phi(x) - \nabla\phi(y)\| \leq L_\phi \|x - y\|$ for all $x, y \in \mathbb{R}^d$. We say that L_ϕ is the *smoothness constant* of ϕ .

²All our results can be extended in a straightforward way to the more general case when node i contains m_i training examples. We decided to present the results in the special case $m = m_i$ for all i in order to simplify the notation.

Note that in view of (3), the Hessian of f_{ij} at point x is

$$\mathbf{H}_{ij}(x) := \nabla^2 f_{ij}(x) = h_{ij}(x) a_{ij} a_{ij}^\top, \quad (4)$$

where

$$h_{ij}(x) := \varphi''_{ij}(a_{ij}^\top x). \quad (5)$$

In view of Assumption 1.1, we have $|\varphi''_{ij}(t)| \leq \gamma$ for all $t \in \mathbb{R}$, and

$$|h_{ij}(x) - h_{ij}(y)| \leq \nu |a_{ij}^\top x - a_{ij}^\top y| \leq \nu \|a_{ij}\| \|x - y\| \quad (6)$$

for all $x, y \in \mathbb{R}^d$. Let $R := \max_{ij} \|a_{ij}\|$. The Hessian of f_i is given by

$$\mathbf{H}_i(x) \stackrel{(2)}{=} \frac{1}{m} \sum_{j=1}^m \mathbf{H}_{ij}(x) \stackrel{(4)}{=} \frac{1}{m} \sum_{j=1}^m h_{ij}(x) a_{ij} a_{ij}^\top, \quad (7)$$

and the Hessian of f is given by

$$\mathbf{H}(x) \stackrel{(2)}{=} \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i(x) \stackrel{(7)}{=} \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m h_{ij}(x) a_{ij} a_{ij}^\top. \quad (8)$$

1.2 The curse of the condition number

All first order methods —distributed or not —suffer from a dependence on an appropriately chosen notion of a *condition number*³ —a number that describes the difficulty of solving the problem by the method at hand. A condition number is a function of the goal we are trying to achieve (e.g., minimize the number of iterations vs minimize the number of communications), choice of the loss function, structure of the model we are trying to learn, and last but not least, the size and properties of the training data. In fact, most research in this area is motivated by the desire to design methods that would have a *reduced* dependence on the condition number. This is the case for many of the tricks heavily studied in the literature, including minibatching [45], importance sampling [29, 53], random reshuffling [28], variance reduction [8, 16, 37, 50], momentum [23, 24], adaptivity [26], communication compression [1, 4, 27], and local computation [17, 25, 43]. Research in this area is becoming saturated, and new ideas are needed to make further progress.

1.3 Newton’s method to the rescue?

One of the ideas that undoubtedly crossed everybody’s mind is the trivial observation that there *is* a very old and simple method which does *not* suffer from any conditioning issues: Newton’s method. Indeed, when it works, Newton’s method has a fast *local quadratic convergence rate* which is entirely independent of the condition number of the problem [2]. While this is a very attractive property, developing scalable distributed variants of Newton’s method that could also *provably outperform* gradient based methods remains a largely unsolved problem. To highlight the severity of the issues with extending Newton’s method to stochastic and distributed settings common in machine learning, we note that

³Example: if one wishes to minimize an L -smooth μ -strongly convex function and one cares about the number of gradient type iterations, the appropriate notion of a condition number is $\kappa := \frac{L}{\mu}$.

Table 1: Summary of algorithms proposed and convergence results proved in this paper.

Method	Convergence			Rate of the condition number?	Theorem
	result [†]	type	rate		
NEWTON-STAR (12)	$r_{k+1} \leq cr_k^2$	local	quadratic	✓	2.1
NEWTON-LEARN Algorithm 1	$\Phi_1^k \leq \theta_1^k \Phi_1^0$	local	linear	✓	3.2
	$r_{k+1} \leq c\theta_1^k r_k$	local	superlinear	✓	3.2

Quantities for which we prove convergence: (i) distance to solution $r_k := \|x^k - x^*\|$; (ii) Lyapunov function $\Phi_q^k := \|x^k - x^*\|^2 + c_q \sum_{i=1}^n \sum_{j=1}^m (h_{ij}^k - h_{ij}(x^*))^2$ for $q = 1, 2, 3$, where $h_{ij}(x^*) = \varphi''_{ij}(a_{ij}^\top x^*)$ (see (5)); (iii) Function value suboptimality $\Delta_k := P(x^k) - P(x^*)$

[†] constant c is possibly different each time it appears in this table. Refer to the precise statements of the theorems for the exact values.

until recently, we did not even have any Newton-type analogue of SGD that could provably work with small minibatch sizes, let alone minibatch size one [19]. In contrast, SGD with minibatch size one is one of the simplest and well understood variants thereof [29], and much of modern development in the area of SGD methods is much more sophisticated. Most variants of Newton’s method proposed for deployment in machine learning are heuristics, which is to say that they are not supported with any convergence guarantees, or have convergence guarantees without explicit rates, or suffer from rates that are worse than the rates of first order methods.

1.4 Contributions summary

We develop *several new fundamental Newton-type methods* which we hope make a marked step towards the ultimate goal of developing practically useful and *communication efficient distributed second order methods*. Our methods are designed with the explicit goal of supporting *efficient communication* in a distributed setting, and in sharp contrast with most recent work, their design was heavily influenced by our desire to equip them with *strong convergence guarantees* typical for the classical Newton’s method [34, 47] and cubically regularized Newton’s method [12, 31]. Our convergence results are summarized in Table 1.

- **First new method and its local quadratic convergence.** We first show that if we know the Hessian of the objective function at the optimal solution, then we can use it instead of the typical Hessian appearing in Newton’s method, and the resulting algorithm, which we call **NEWTON-STAR (NS)**, inherits local quadratic convergence behavior of Newton’s method (see Theorem 2.1). In a distributed setting with a central orchestrating sever, each compute node only needs to send the local gradient to the server node, and no matrices need to be sent. While this method is not practically useful, it acts as a stepping stone to our next method, in which these deficiencies are removed. This method is described in Section 2.
- **Second new method and its local linear and superlinear convergence.** Motivated by the above result, we propose a *learning scheme which enables us to learn the Hessian at the optimum iteratively in a communication efficient manner*.

This scheme gives rise to our second new method: **NEWTON-LEARN (NL)**. We analyze this method in the case, when all individual loss functions are convex and $\lambda > 0$. Besides the local full gradient, each worker node needs to send additional information to the server node in order to learn the Hessian at the optimum. However, our learning scheme supports *compressed communication* with arbitrary compression level. This level can be chosen so that in each iteration, each node sends an equivalent of a few gradients to the server only. That is, we can achieve $O(d)$ communication complexity in each iteration. We prove local linear convergence for a carefully designed Lyapunov function, and local superlinear convergence for the squared distance to optimum (see Theorems 3.2). Remarkably, all these rates are *independent of the condition number*. The NL method and the associated theory are described in Section 3.

- **Experiments.** Our theory is corroborated with numerical experiments showing the superiority of our method to several state-of-the-art benchmarks, including DCGD [18], DIANA [14, 27], ADIANA [21], BFGS [6, 9, 11, 41], and DINGO [7]. Our method can achieve communication complexity which is *several orders of magnitude better* than competing methods (see Section 4).

1.5 Related work

Several distributed Newton-type methods can be found in recent literature. DANE [40] is a distributed approximate Newton-type method where each worker node needs to solve a subproblem using the full gradient at each iteration, and the new iterate is the average of these subproblem solutions. The linear convergence of DANE was obtained in the strongly convex case. An inexact DANE method in which the subproblem is solved approximately was proposed and studied by Reddi et al. [35]. Moreover, an accelerated version of inexact DANE, called AIDE, was proposed in [35] by a generic acceleration scheme —catalyst [22] — and an optimal communication complexity can be obtained up to logarithmic factors in specific settings. The DiSCO method, which combines inexact damped Newton method and distributed preconditioned conjugate gradient method, was proposed by Zhang and Xiao [52] and analyzed for self-concordant empirical loss. GIANT [48] is a globally improved approximate Newton method which has a better linear convergence rate than first-order methods for quadratic functions, and has local linear-quadratic convergence for strongly convex functions. GIANT and DANE are identical for quadratic programming. The communication cost per iteration of the above methods is $O(d)$. These methods can only achieve linear convergence in the strongly convex case. The comparison of the iteration complexity of the above methods for the ridge regression problem can be found in Table 2 of [48].

Crane and Roosta [7] proposed a distributed Newton-type method called DINGO for solving invex finite-sum problems. Invexity is a special case of non-convexity, which subsumes convexity as a sub-class. A linear convergence rate was obtained for DINGO under certain assumptions using an Armijo-type line search, and at each iteration, several communication rounds are needed assuming two communication rounds for line-search per iteration. The communication cost for each communication round is $O(d)$. The compressed version of DINGO was studied in [10] to reduce the communication cost at each communication round by using the δ -approximate compressor, and the same

Table 2: Comparison of distributed Newton-type methods. Our methods combine the best of both worlds, and are the only methods we know about which do so: we obtain fast rates independent of the condition number, and allow for $O(d)$ communication per communication round.

Method	Convergence rate	Rate independent of the condition number?	Communication cost per iteration	Network structure
DANE [40]	Linear	✗	$O(d)$	Centralized
DiSCO [52]	Linear	✗	$O(d)$	Centralized
AIDE [35]	Linear	✗	$O(d)$	Centralized
GIANT [48]	Linear	✗	$O(d)$	Centralized
DINGO [7]	Linear	✗	$O(d)$	Centralized
DAN [51]	Local quadratic [†]	✓	$O(nd^2)$	Decentralized
DAN-LA [51]	Superlinear	✓	$O(nd)$	Decentralized
NEWTON-STAR this work	Local quadratic	✓	$O(d)$	Centralized
NEWTON-LEARN this work	Local superlinear	✓	$O(d)$	Centralized

[†] DAN converges globally, but the quadratic rate is introduced only after $O(L_2/\mu^2)$ steps, where L_2 is the Lipschitz constant of the Hessian of P , and μ is the strong convexity parameter of P . This is a property it inherits from the recent method of Polyak [32] this method is based on.

rate of convergence as DINGO can be obtained by properly choosing the stepsize and hyper-parameters when δ is a constant. Zhang et al. [51] proposed two decentralized distributed adaptive Newton methods, called DAN and DAN-LA. DAN combines the distributed selective flooding (DSF) algorithm and Polyak’s adaptive Newton method [33], and enters pure Newton method which has quadratic convergence after about $\frac{2M}{\mu^2} \|\nabla P(x^0)\|$ iterations, where M is the Lipschitz constant of the Hessian of P and μ is the strongly convex parameter of P . DAN-LA, which leverages the low-rank approximation method to reduce the communication cost, has global superlinear convergence. At each iteration, both DAN and DAN-LA need $n - 1$ communication rounds, and the communication cost for each communication round is $O(d^2)$ and $O(d)$ respectively.

We compare the convergence rate and per-iteration communication cost with these Newton-type methods in Table 2. Note that the first five methods in the table have rates that depend on the condition number of the problem, and as such, do not have the benefits normally attributed to pure Newton’s method. Note also that the two prior methods which do have rates independent of the condition number have high cost of communication. Our methods combine the best of both worlds, and are the only methods we know about which do so: we obtain fast rates independent of the condition number, and allow for $O(d)$ communication per communication round. We were able to achieve this by a complete redesign of how second order methods should work in the distributed setting. Our methods are not simple extensions of existing schemes, and our proofs use novel arguments and techniques.

2 Three Steps Towards an Efficient Distributed Newton Type Method

In order to better explain the algorithms and results of this paper, we will proceed through several steps in a gradual explanation of the ideas that ultimately lead to our methods. While this is not the process we used to come up with our methods, in retrospect we believe that our methods and results will be understood more easily when seen as having been arrived at in this way. In other words, we have constructed what we believe is a plausible discovery story, one enabling faster and better comprehension. If these ideas seem to follow naturally, it is because we made a conscious effort to make them appear that way. The goal of this paper is to develop communication efficient variants of Newton’s method for solving the distributed optimization problem (1).

2.1 Naive distributed implementation of Newton’s method

Newton’s method applied to problem (1) performs the iteration

$$x^{k+1} = x^k - (\nabla^2 P(x^k))^{-1} \nabla P(x^k) \stackrel{(1)}{=} x^k - (\mathbf{H}(x^k) + \lambda \mathbf{I})^{-1} \nabla P(x^k). \quad (9)$$

A naive way to implement this method in the *parameter server* framework is for each node i to compute the Hessian $\mathbf{H}_i(x^k)$ and gradient $\nabla f_i(x^k)$ and to communicate these objects to the server. The server then averages the local Hessians $\mathbf{H}_i(x^k)$ to produce $\mathbf{H}(x^k)$ via (8), and averages the local gradients $\nabla f_i(x^k)$ to produce $\nabla f(x^k)$. The server then adds $\lambda \mathbf{I}$ to the Hessian, producing $\mathbf{H}(x^k) + \lambda \mathbf{I} = \nabla^2 P(x^k)$, adds λx^k to the gradient, producing $\nabla P(x^k) = \nabla f(x^k) + \lambda x^k$, and subsequently performs the Newton step (9). The resulting vector x^{k+1} is then broadcasted to the nodes and the process is repeated.

This implementation mirrors the way GD and many other first order methods are implemented in the parameter server framework. However, unlike in the case of GD, where only $O(d)$ floats need to be sent and received by each node in each iteration, the upstream communication in Newton’s method requires $O(d^2)$ floats to be communicated by each worker to the server. Since d is typically very large, this is prohibitive in practice. Moreover, computation of the Newton’s step by the parameter server is much more expensive than simple averaging of the gradients performed by gradient type methods. However, in this paper we will not be concerned with the cost of the Newton step itself, as we will assume the server is powerful enough and the network connection is slow enough for this step not to be the main bottleneck of the iteration. Instead, we assume that the communication steps in general, and the $O(d^2)$ communication of the Hessian matrices in particular, is what forms the bottleneck. The $O(d)$ per node communication cost of the local gradients is negligible, and so is the $O(d)$ broadcast of the updated model.

2.2 A better implementation taking advantage of the structure of $\mathbf{H}_{ij}(x)$

The above naive implementation can be improved in the setting when $m < d^2$ by taking advantage of the explicit structure (7) of the local Hessians as a conic combination of

positive semidefinite rank one matrices:

$$\mathbf{H}_i(x) = \frac{1}{m} \sum_{j=1}^m h_{ij}(x) a_{ij} a_{ij}^\top. \quad (10)$$

Indeed, assuming that the server has direct access to all the training data vectors $a_{ij} \in \mathbb{R}^d$ (these vectors can be sent to the server at the start of the process), node i can send the m coefficients $h_{i1}(x), \dots, h_{im}(x)$ to the server instead, and the server is then able to reconstruct the Hessian matrix $\mathbf{H}_i(x)$ from this information. This way, each node sends $O(m + d)$ floats to the server, which is a substantial improvement on the naive implementation in the regime when $m \ll d^2$. However, when $m \gg d$, the upstream communication cost is still substantially larger than the $O(d)$ cost of GD. If the server does not have enough memory to store all vectors a_{ij} , this procedure does not work.

2.3 NEWTON-STAR: Newton’s method with a single Hessian

We now introduce a simple idea which, surprisingly, enables us to *remove the need to iteratively communicate any coefficients altogether*. Assume, for the sake of argument, that we know the values $h_{ij}(x^*)$ for all i, j . That is, assume the server has access to coefficients $h_{ij}(x^*)$ for all i, j , and that each node i has access to coefficients $h_{ij}(x^*)$ for $j = 1, \dots, m$, i.e., to the vector

$$h_i(x) := (h_{i1}(x), \dots, h_{im}(x)) \in \mathbb{R}^m \quad (11)$$

for $x = x^*$. Next, consider the following new Newton-like method which we call **NEWTON-STAR (NS)**, where the “star” points to the method’s reliance on the knowledge of the optimal solution x^* :

$$\begin{aligned} x^{k+1} &= x^k - (\nabla^2 P(x^*))^{-1} \nabla P(x^k) \\ &\stackrel{(1)}{=} x^k - (\mathbf{H}(x^*) + \lambda \mathbf{I})^{-1} \left(\frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) + \lambda x^k \right). \end{aligned} \quad (12)$$

Since the server knows $\mathbf{H}(x^*)$, all that the nodes need to communicate are the local gradients $\nabla f_i(x^k)$, which costs $O(d)$ per node. The server then computes x^{k+1} , broadcasts it back to the nodes, and the process is repeated. This method has the same per-iteration $O(d)$ communication complexity as GD. However, as we show next, the number of iterations (which is the same as the number of communications) of **NEWTON-STAR** does not depend on the condition number — a property it borrows from the classical Newton’s method. The following theorem says that **NEWTON-STAR** enjoys *local quadratic convergence*.

Theorem 2.1 (Local quadratic convergence). *Let Assumption 1.1 hold, and assume that $\mathbf{H}(x^*) \succeq \mu^* \mathbf{I}$ for some $\mu^* \geq 0$ (for instance, this holds if f is μ^* -strongly convex) and that $\mu^* + \lambda > 0$. Then for any starting point $x^0 \in \mathbb{R}^d$, the iterates of **NEWTON-STAR** for solving problem (1) satisfy the following inequality:*

$$\|x^{k+1} - x^*\| \leq \frac{\nu}{2(\mu^* + \lambda)} \cdot \left(\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|a_{ij}\|^3 \right) \cdot \|x^k - x^*\|^2. \quad (13)$$

Proof. By the first order optimality conditions, we have

$$\nabla f(x^*) + \lambda x^* = 0. \quad (14)$$

Let $\mathbf{H}_* := \mathbf{H}(x^*)$. Since $\mathbf{H}_* \succeq \mu^* \mathbf{I}$, we have $\mathbf{H}_* + \lambda \mathbf{I} \succeq (\mu^* + \lambda) \mathbf{I}$, and hence

$$\|(\mathbf{H}_* + \lambda \mathbf{I})^{-1}\| \leq \frac{1}{\mu^* + \lambda}. \quad (15)$$

Using (14) and (15) and subsequently applying Jensen's inequality to the function $x \mapsto \|x\|$, we get

$$\begin{aligned} \|x^{k+1} - x^*\| &= \left\| x^k - x^* - (\mathbf{H}_* + \lambda \mathbf{I})^{-1} \nabla P(x^k) \right\| \\ &\stackrel{(14)}{=} \left\| (\mathbf{H}_* + \lambda \mathbf{I})^{-1} [(\mathbf{H}_* + \lambda \mathbf{I})(x^k - x^*) - (\nabla f(x^k) - \nabla f(x^*) + \lambda(x^k - x^*))] \right\| \\ &\stackrel{(15)}{\leq} \frac{1}{\mu^* + \lambda} \|(\mathbf{H}_* + \lambda \mathbf{I})(x^k - x^*) - (\nabla f(x^k) - \nabla f(x^*)) - \lambda(x^k - x^*)\| \\ &= \frac{1}{\mu^* + \lambda} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i(x^*)(x^k - x^*) - \frac{1}{n} \sum_{i=1}^n (\nabla f_i(x^k) - \nabla f_i(x^*)) \right\| \\ &\leq \frac{1}{n(\mu^* + \lambda)} \sum_{i=1}^n \|\mathbf{H}_i(x^*)(x^k - x^*) - (\nabla f_i(x^k) - \nabla f_i(x^*))\| \\ &\stackrel{(7)}{=} \frac{1}{n(\mu^* + \lambda)} \sum_{i=1}^n \left\| \frac{1}{m} \sum_{j=1}^m [h_{ij}(x^*) a_{ij} a_{ij}^\top (x^k - x^*) - (\nabla f_{ij}(x^k) - \nabla f_{ij}(x^*))] \right\|. \quad (16) \end{aligned}$$

We now use the fundamental theorem of calculus to express difference of gradients $\nabla f_{ij}(x^k) - \nabla f_{ij}(x^*)$ in an integral, obtaining

$$\nabla f_{ij}(x^k) - \nabla f_{ij}(x^*) = \int_0^1 \nabla^2 f_{ij}(x^* + \tau(x^k - x^*))(x^k - x^*) d\tau. \quad (17)$$

Plugging this representation into (16) and noting that $\nabla^2 f_{ij}(x) \equiv \mathbf{H}_{ij}(x)$ (see (4)), we can continue:

$$\begin{aligned} \|x^{k+1} - x^*\| &\stackrel{(16)+(17)}{\leq} \frac{1}{n(\mu^* + \lambda)} \sum_{i=1}^n \left\| \frac{1}{m} \sum_{j=1}^m (h_{ij}(x^*) a_{ij} a_{ij}^\top (x^k - x^*) \right. \\ &\quad \left. - \int_0^1 \mathbf{H}_{ij}(x^* + \tau(x^k - x^*))(x^k - x^*) d\tau) \right\| \\ &\stackrel{(4)}{=} \frac{1}{n(\mu^* + \lambda)} \sum_{i=1}^n \left\| \frac{1}{m} \sum_{j=1}^m (h_{ij}(x^*) a_{ij} a_{ij}^\top (x^k - x^*) \right. \\ &\quad \left. - \int_0^1 h_{ij}(x^* + \tau(x^k - x^*)) a_{ij} a_{ij}^\top (x^k - x^*) d\tau) \right\| \\ &= \frac{1}{n(\mu^* + \lambda)} \sum_{i=1}^n \left\| \frac{1}{m} \sum_{j=1}^m a_{ij} a_{ij}^\top (x^k - x^*) \left(h_{ij}(x^*) - \int_0^1 h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right) \right\| \\ &\leq \frac{\|x^k - x^*\|}{(\mu^* + \lambda)} \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|a_{ij}\|^2 \left| \int_0^1 h_{ij}(x^*) - h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right|. \quad (18) \end{aligned}$$

In the last step we have again used Jensen's inequality applied to the function $x \mapsto \|x\|$, followed by inequalities of the form $\|\mathbf{A}_{ij} x t_{ij}\| \leq \|\mathbf{A}_{ij}\| \|x\| |t_{ij}|$ for $\mathbf{A}_{ij} = a_{ij} a_{ij}^\top$, $x = x^k - x^*$ and $t_{ij} \in \mathbb{R}$.

From (6) we obtain $|h_{ij}(x^*) - h_{ij}(x^* + \tau(x^k - x^*))| \leq \nu \tau \|a_{ij}\| \cdot \|x^k - x^*\|$, which implies that

$$\left| \int_0^1 h_{ij}(x^*) - h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right| \leq \int_0^1 \nu \tau \|a_{ij}\| \cdot \|x^k - x^*\| d\tau = \frac{\nu \|a_{ij}\|}{2} \cdot \|x^k - x^*\|.$$

Plugging this into (18), we finally arrive at (13). \square

Note that we do not need to assume f to be convex or strongly convex. All we need to assume is positive definiteness of the Hessian at the optimum. This implies local strong convexity, and since our convergence result is local, that is all we need.

3 NEWTON-LEARN: Learning the Hessian and Local Convergence Theory

In Sections 2.1, 2.2 and 2.3 we have gone through three steps in our story, with the first true innovation and contribution of this paper being the NEWTON-STAR method and its rate. We have now sufficiently prepared the ground to motivate our first *key* contribution: the NEWTON-LEARN method. We only outline the basic insights behind this method here; the details are included in Section 3.

3.1 The main iteration

In NEWTON-LEARN we maintain a sequence of vectors

$$h_i^k = (h_{i1}^k, \dots, h_{im}^k) \in \mathbb{R}^m,$$

for all $i = 1, \dots, n$ throughout the iterations $k \geq 0$ with the goal of *learning* the values $h_{ij}(x^*)$ for all i, j . That is, we construct the sequence with the explicit intention to enforce

$$h_{ij}^k \rightarrow h_{ij}(x^*) \quad \text{as} \quad k \rightarrow \infty. \quad (19)$$

Using $h_{ij}^k \approx h_{ij}(x^*)$ we estimate the Hessian $\mathbf{H}(x^*)$ via

$$\mathbf{H}(x^*) \approx \mathbf{H}^k := \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m h_{ij}^k a_{ij} a_{ij}^\top, \quad (20)$$

and then perform a similar iteration to (12):

$$x^{k+1} = x^k - (\mathbf{H}^k + \lambda \mathbf{I})^{-1} \nabla P(x^k). \quad (21)$$

3.2 Learning the coefficients: the idea

To complete the description of the method, we need to explain how the vectors h_i^{k+1} are updated. This is also the place where we can force the method to be communication efficient. Indeed, if we can design a rule that would enforce the update vectors $h_i^{k+1} - h_i^k$ to be *sparse*, say

$$\|h_i^{k+1} - h_i^k\|_0 \leq s \quad (22)$$

for some $1 \leq s \leq m$ and all i and k , then the upstream communication by each node in each iteration would be of the order $O(s + d)$ only (provided the server has access to all vectors a_{ij})! That is, each node i only needs to communicate s entries of the update vector $h_i^{k+1} - h_i^k$ as the rest are equal to zero, and each node also needs to communicate the d dimensional gradient $\nabla f_i(x^k)$. Note that $O(s + d)$ can be interpreted as

an *interpolation* of the $O(m + d)$ per-iteration communication complexity of the structure-aware implementation of Newton’s method from Section 2.2, and of the $O(d)$ per-iteration communication complexity of **NEWTON-STAR** described in Section 2.3.

In the more realistic regime when the server does *not* have access to the data $\{a_{ij}\}$, we ask each worker i to additionally send the corresponding s vectors a_{ij} , which costs extra $O(sd)$ in communication per node. However, when $s = O(1)$, this is the same per-iteration communication effort as that of GD.

We develop the update rule defining the evolution of the vectors h_1^k, \dots, h_n^k . This rule (see (25)) applies to the $\lambda > 0$ case and leads to **NEWTON-LEARN** which we call **NL** (see Algorithm 1). This rule and the method are described in Section 3.5.

3.3 Outline of fast local convergence theory

We show in Theorem 3.2 (covering **NL**) and Theorem ?? (covering **NL2**) that **NEWTON-LEARN** enjoys a local linear rate wrt a certain Lyapunov function which involves the term $\|x^k - x^*\|^2$ and also all terms of the form $\|h_i^k - h_i(x^*)\|^2$. This means that i) the main iteration (21) works, i.e., x^k converges to x^* at a local linear rate, and that ii) the learning procedure works, and the desired convergence described in (19) occurs at a local linear rate. In addition, we also establish a local superlinear rate of $\|x^k - x^*\|^2$. Remarkably, these rates are *independent of any condition number, which is in sharp contrast with virtually all results on distributed Newton-type methods we are aware of*.

Moreover, we wish to remark that second order methods are not typically analyzed using a Lyapunov style analysis. Indeed, we only know of a couple works that do so. First, Kovalev et al. [19] develop stochastic Newton and cubic Newton methods of a different structure and scope from ours. They do not consider distributed optimization nor communication compression. Second, Kovalev et al. [20] develop a stochastic BFGS method. Again, their method and scope is very different from ours. Hence, *our analysis may be of independent interest as it adds to the arsenal of theoretical tools which could be used in a more precise analysis of other second order methods*.

3.4 Compressed learning

Instead of merely relying on sparse updates for the vectors h_i^k (see (22)), we provide a more general communication compression strategy which includes sparsification as a special case [1]. We do so via the use of a *random compression* operator. We say that a randomized map $\mathcal{C} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is a *compression operator (compressor)* if there exists a constant $\omega \geq 0$ such that the following relations hold for all $x \in \mathbb{R}^m$:

$$\mathbb{E}[\mathcal{C}(x)] = x \tag{23}$$

$$\mathbb{E}\|\mathcal{C}(x)\|^2 \leq (\omega + 1)\|x\|^2. \tag{24}$$

The identity compressor $\mathcal{C}(x) \equiv x$ satisfies these relations with $\omega = 0$. The larger the variance parameter ω is allowed to be, the easier it can be to construct a compressor \mathcal{C} for which the value $\mathcal{C}(x)$ can be encoded using a small number of bits only. We refer the reader to [5] for a list of several compressors and their properties.

3.5 NL (learning in the $\lambda > 0$ case)

We now consider the case where all loss functions φ_{ij} are convex and $\lambda > 0$.

Assumption 3.1. *Each φ_{ij} is convex, $\lambda > 0$.*

When combined with Assumption 1.1, Assumption 3.1 implies that $\varphi''_{ij}(t) \geq 0$ for all t , hence $h_{ij}(x) = \varphi''_{ij}(a_i^\top x) \geq 0$ for all $x \in \mathbb{R}^d$. In particular, $h_{ij}(x^*) \geq 0$ for all i, j . Since we wish to construct a sequence of vectors $h_i^k = (h_{i1}^k, \dots, h_{im}^k) \in \mathbb{R}^m$ satisfying $h_{ij}^k \rightarrow h_{ij}(x^*)$, it makes sense to try to enforce all vectors in this sequence to have *nonnegative entries*:

$$h_{ij}^k \geq 0.$$

Since \mathbf{H}^k arises as a linear combination of the rank-one matrices $a_{ij}a_{ij}^\top$ (see (20)), this makes \mathbf{H}^k positive semidefinite, which in turn means that the matrix $\mathbf{H}^k + \lambda \mathbf{I}$ appearing in the main iteration (21) of NEWTON-LEARN is invertible, and hence the iteration is *well defined*.⁴

3.5.1 The learning iteration and the NL algorithm

In particular, in NEWTON-LEARN each node i computes the vector $h_i(x^k) \in \mathbb{R}^m$ of second derivatives defined in (11), and then performs the update

$$h_i^{k+1} = [h_i^k + \eta \mathcal{C}_i^k(h_i(x^k) - h_i^k)]_+, \quad (25)$$

where $\eta > 0$ is a learning rate, \mathcal{C}_i^k is a freshly sampled compressor by node i at iteration k . By $[\cdot]_+$ we denote the positive part function applied element-wise, defined for scalars as follows: $[t]_+ = t$ if $t \geq 0$ and $[t]_+ = 0$ otherwise.

We remark that it is possible to interpret the learning procedure (25) as one step of projected stochastic gradient descent (SGD) applied to a certain quadratic optimization problem whose unique solution is the vector $h_i(x^k)$.

The NL algorithm (Algorithm 1) arises as the combination of the Newton-like update (21) (adjusted to take account of the explicit regularizer) and the learning procedure (25). It is easy to see that the update rule for \mathbf{H}^k in NL is designed to ensure that \mathbf{H}^k remains of the form $\mathbf{H}^k = \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i^k$, where $\mathbf{H}_i^k = \frac{1}{m} \sum_{j=1}^m h_{ij}^k a_{ij} a_{ij}^\top$. The update rule for x^k , performed by the server, is identical to (21), with an extra provision for the regularizer. The vector x^{k+1} is broadcasted to all workers. Let us comment on how the key communication step is implemented. If the server does not have direct access to the training data vectors $\{a_{ij}\}$, we choose Option 1, otherwise we choose Option 2. A key property of NL is that the server is able to maintain copies of the learning vectors h_i^k *without the need for these vectors to be communicated by the workers to the server*. Indeed, provided the workers and the server agree on the same set of initial vectors h_1^0, \dots, h_n^0 , update (25) can be independently computed by the server as well from its memory state h_i^k and the compressed message $\mathcal{C}_i^k(h_i(x^k) - h_i^k)$ received from node i . This strategy is reminiscent of the way the key step in the first-order method DIANA [14, 27] is executed. In this sense, NL can be seen as arising from a successful marriage of Newton's method and the DIANA trick.

⁴Positive definiteness of Hessian estimates is enforced in several popular quasi-Newton methods as well; for instance, in the BFGS method [6, 9, 11, 41]. However, quasi-Newton methods operate in a markedly different manner, and the way in which positive definiteness is enforced there is also different.

Algorithm 1 NL: NEWTON-LEARN ($\lambda > 0$ case)

Parameters: learning rate $\eta > 0$

Initialization: $x^0 \in \mathbb{R}^d$; $h_1^0, \dots, h_n^0 \in \mathbb{R}_+^m$; $\mathbf{H}^0 = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m h_{ij}^0 a_{ij} a_{ij}^\top \in \mathbb{R}^{d \times d}$

for $k = 0, 1, 2, \dots$ **do**

 Broadcast x^k to all workers

for each node $i = 1, \dots, n$ **do**

 Compute local gradient $\nabla f_i(x^k)$

$h_i^{k+1} = [h_i^k + \eta \mathcal{C}_i^k(h_i(x^k) - h_i^k)]_+$

 Send $\nabla f_i(x^k)$ and $\mathcal{C}_i^k(h_i(x^k) - h_i^k)$ to server

Option 1: Send $\{a_{ij} : h_{ij}^{k+1} - h_{ij}^k \neq 0\}$ to server

Option 2: Do nothing if server knows $\{a_{ij} : \forall j\}$

end for

$x^{k+1} = x^k - (\mathbf{H}^k + \lambda \mathbf{I})^{-1} \left(\frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) + \lambda x^k \right)$

$\mathbf{H}^{k+1} = \mathbf{H}^k + \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m (h_{ij}^{k+1} - h_{ij}^k) a_{ij} a_{ij}^\top$

end for

3.5.2 Theory

In our theoretical results we rely on the Lyapunov function

$$\Phi_1^k := \|x^k - x^*\|^2 + \frac{1}{3mn\eta\nu^2 R^2} \mathcal{H}^k, \quad \mathcal{H}^k := \sum_{i=1}^n \|h_i^k - h_i(x^*)\|^2.$$

Our main theorem follows.

Theorem 3.2 (Convergence of NL). *Let Assumptions 1.1 and 3.1 hold. Let $\eta \leq \frac{1}{\omega+1}$ and assume that $\|x^k - x^*\|^2 \leq \frac{\lambda^2}{12\nu^2 R^6}$ for all $k \geq 0$. Then for Algorithm 1 we have the inequalities*

$$\begin{aligned} \mathbb{E}[\Phi_1^k] &\leq \theta_1^k \Phi_1^0, \\ \mathbb{E} \left[\frac{\|x^{k+1} - x^*\|^2}{\|x^k - x^*\|^2} \right] &\leq \theta_1^k \left(6\eta + \frac{1}{2} \right) \frac{\nu^2 R^6}{\lambda^2} \Phi_1^0, \end{aligned}$$

where $\theta_1 := 1 - \min \left\{ \frac{\eta}{2}, \frac{5}{8} \right\}$.

Since the stepsize bound $\eta \leq \frac{1}{\omega+1}$ is independent of the condition number, the linear convergence rates of $\mathbb{E}[\Phi_1^k]$ and $\mathbb{E} \left[\frac{\|x^{k+1} - x^*\|^2}{\|x^k - x^*\|^2} \right]$ are both *independent of the condition number*. Next, we explore under what conditions we can guarantee for all the iterates to stay in a small neighborhood.

Lemma 3.3. *Let Assumptions 1.1 and 3.1 hold. Assume h_{ij}^k is a convex combination of $\{h_{ij}(x^0), h_{ij}(x^1), \dots, h_{ij}(x^k)\}$ for all i, j and k . Assume $\|x^0 - x^*\|^2 \leq \frac{\lambda^2}{12\nu^2 R^6}$. Then*

$$\|x^k - x^*\|^2 \leq \frac{\lambda^2}{12\nu^2 R^6} \quad \text{for all } k \geq 0.$$

It is easy to verify that if we choose $h_{ij}^0 = h_{ij}(x^0)$ and use the random sparsification compressor and $\eta \leq \frac{1}{\omega+1}$, then h_{ij}^k is always a convex combination of $\{h_{ij}(x^0), h_{ij}(x^1), \dots, h_{ij}(x^k)\}$ for $k \geq 0$. Thus, from Lemma 3.3 we can guarantee that all the iterates stay in the small neighborhood assumed in Theorem 3.2 as long as the initial point x^0 is in it.

4 Experiments

We now study the empirical performance of our second order method **NL** and compare it with relevant benchmarks and with state-of-the-art methods. We test on the regularized logistic regression problem

$$\min_{x \in \mathbb{R}^d} \left\{ \frac{1}{n} \sum_{i=1}^n \frac{1}{m} \sum_{j=1}^m \log(1 + \exp(-b_{ij} a_{ij}^\top x)) + \frac{\lambda}{2} \|x\|^2 \right\},$$

where $\{a_{ij}, b_{ij}\}_{j \in [m]}$ are data samples at the i -th node.

4.1 Data sets and parameter settings

In our experiments we use five standard datasets from the LIBSVM library: **a2a**, **a7a**, **a9a**, **w8a** and **phishing**. Besides, we generated an artificial dataset **artificial** as follows: each of the d elements of the data vector $a_{ij} \in \mathbb{R}^d$ was sampled from the normal distribution $\mathcal{N}(10, 10)$. The corresponding label b_{ij} was sampled uniformly at random from $\{-1, 1\}$. We partitioned each dataset across several nodes (selection of n) in order to capture a variety of scenarios. See Table 3 for more details on all the datasets and the choice of n .

In all experiments we use the theoretical parameters (e.g., stepsizes) for all the three algorithms: vanilla Distributed Compressed Gradient Descent (DCGD) [18], DIANA [27], and ADIANA [21].

Table 3: Datasets used in the experiments, and the number of worker nodes n used in each case.

	# workers n	# data points ($= nm$)	# features d
a2a	15	2 265	123
a7a	100	16 100	123
a9a	80	32 560	123
w8a	142	49 700	300
phishing	100	11 000	68
artificial	100	1 000	200

As the initial approximation of the Hessian in BFGS [6, 9, 11, 41], we use $\mathbf{H}^0 = \nabla^2 P(x^0)$, and the stepsize is 1. We set the same constants in DINGO [7] as they did: $\theta = 10^{-4}$, $\phi = 10^{-6}$, $\rho = 10^{-4}$, and use backtracking line search for DINGO to select the largest stepsize in $\{1, 2^{-1}, 2^{-2}, 2^{-4}, \dots, 2^{-10}\}$. We conduct experiments for three values of the regularization parameter λ : $10^{-3}, 10^{-4}, 10^{-5}$. In the figures we plot the

relation of the optimality gap $P(x^k) - P(x^*)$ and the number of accumulated transmitted bits⁵ or the number of iterations. The optimal value $P(x^*)$ in each case is the function value at the 20-th iterate of standard Newton’s method. We adopt the realistic setting where the server does not have access to the local data (**Option 1**).

4.2 Compression operators

For the first order methods we use three compression operators: random sparsification [44], random dithering [1], and natural compression [13] (all defined below). For random- r sparsification, the number of communicated bits per iteration is $32r + \log_2 \binom{d}{r}$, and we choose $r = d/4$. For random dithering, we choose $s = \sqrt{d}$, which means the number of communicated bits per iteration is $2.8d + 32$. For natural compression, the number of communicated bits per iteration is $9d$ bits.

For NL we use the random- r sparsification operator with a selection of values of r . In addition, we use the random sparsification operator \mathcal{C}_p induced by the random- r compressor. This compressor is also defined below.

4.2.1 Random sparsification

The random sparsification compressor [44], denoted random- r , is a randomized mapping $\mathcal{C} : \mathbb{R}^m \rightarrow \mathbb{R}$ defined as

$$\mathcal{C}(x) := \frac{m}{r} \cdot \xi \circ x$$

where $\xi \in \mathbb{R}^m$ is a random vector distributed uniformly at random on the discrete set $\{y \in \{0, 1\}^m : \|y\|_0 = r\}$, where $\|y\|_0 := \{i \mid y_i \neq 0\}$ and \circ is the Hadamard product. The variance parameter associate with this compressor is $\omega = \frac{m}{r} - 1$.

4.2.2 Random dithering

The random dithering compressor [1, 14] with s levels is defined via

$$\mathcal{C}(x) := \text{sign}(x) \cdot \|x\|_q \cdot \frac{1}{s} \cdot \xi_s,$$

where $\|x\|_q := (\sum_i |x_i|^q)^{1/q}$ and $\xi_s \in \mathbb{R}^m$ is a random vector with i -th element defined as

$$\xi_s(i) := \begin{cases} l + 1 & \text{with probability } \frac{|x_i|}{\|x\|_q} s - l \\ l & \text{otherwise} \end{cases}.$$

Here, l satisfies $\frac{|x_i|}{\|x\|_q} \in [\frac{l}{s}, \frac{l+1}{s}]$ and $s \in \mathbb{N}_+$ denotes the levels of the rounding. The variance parameter of this compressor is $\omega \leq 2 + \frac{m^{1/2} + m^{1/q}}{s}$ [14]. For $q = 2$, one can get the improved bound $\omega \leq \min\{\frac{m}{s^2}, \frac{\sqrt{m}}{s}\}$ [1].

⁵In all plots, “communicated bits” refers to the total number of bits sent by all nodes to the server.

4.2.3 Natural compression

The natural compression [13] operator $\mathcal{C}_{\text{nat}} : \mathbb{R}^m \rightarrow \mathbb{R}$ is obtained by applying the random mapping $\mathcal{C} : \mathbb{R} \rightarrow \mathbb{R}$, defined next, to each coordinate of x independently. We define $\mathcal{C}(0) = 0$ and for $t \neq 0$, we let

$$\mathcal{C}(t) := \begin{cases} \text{sign}(t) \cdot 2^{\lfloor \log_2 |t| \rfloor} & \text{with probability } p(t) := \frac{2^{\lceil \log_2 |t| \rceil} - |t|}{2^{\lceil \log_2 |t| \rceil}} \\ \text{sign}(t) \cdot 2^{\lceil \log_2 |t| \rceil} & \text{with probability } 1 - p(t) \end{cases}$$

The variance parameter of natural compression is $\omega = \frac{1}{8}$.

4.2.4 Bernoulli compressor

A variant of any compression operator $\mathcal{C} : \mathbb{R}^m \rightarrow \mathbb{R}$ can be constructed as follows:

$$\mathcal{C}_p(x) := \begin{cases} \frac{1}{p}\mathcal{C}(x) & \text{with probability } p \\ 0 & \text{otherwise} \end{cases}, \quad (26)$$

where $p \in (0, 1]$ is a probability parameter. It is easy to verify that \mathcal{C}_p is still a compression operator with variance parameter $\omega_p := \frac{\omega+1}{p} - 1$, where ω is the variance parameter of the underlying compressor \mathcal{C} .

4.3 Behavior of NL

Before we compare our method NL with competing baselines, we investigate how is their performance affected by the choice of the sparsification parameter r defining the random- r sparsification operator \mathcal{C} . Likewise, we vary the probability parameter p defining the induced Bernoulli compressor \mathcal{C}_p .

According to the results summarized in Figure 1, the best performance of NL is obtained for $r = 1$ and $p = 1$. We will use these parameter settings for NL in our subsequent experiments where we compare our method with several baselines and state-of-the-art methods.

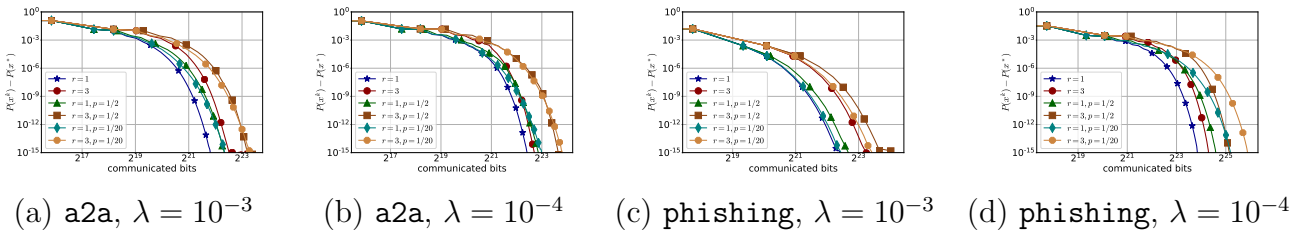


Figure 1: Performance of NL across a few values of r defining the random- r compressor, and a few values of p defining the induced Bernoulli compressor \mathcal{C}_p .

4.4 Comparison of NL with Newton's method

In our next experiment we compare NL using different values of r for random- r compression, with Newton's method; see Figure 2. We clearly see that Newton's method performs

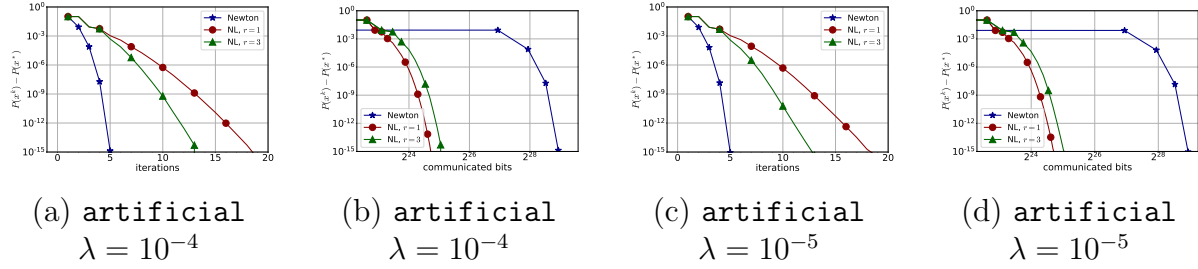


Figure 2: Comparison of NL with Newton’s method in terms of iteration complexity for (a), (c); in terms of communication complexity for (b), (d).

better than NL in terms of iteration complexity, as expected. However, our methods have better communication efficiency than Newton’s method, by *several orders of magnitude*. Moreover, we see that the smaller r is, the better NL performs in terms of communication complexity. In Figure 3 we perform a similar comparison for several more datasets, but focus on communication complexity only. The conclusions are unchanged: our methods NL have superior performance.

4.5 Comparison of NL with BFGS

In our next test, we compare NL with BFGS in Figure 4. As we can see, our methods have better communication efficiency than BFGS, by *several orders of magnitude*.

4.6 Comparison of NL with ADIANA

Next, we compare NL with ADIANA [21] using three different compression operators: natural compression (DIANA-NC), random sparsification (DIANA-RS, $r = d/4$) and random dithering (DIANA-RD, $s = \sqrt{d}$); see Figure 5. Based on the experimental results, we can conclude that NL outperforms all three versions ADIANA for all types of compression, often by *several degrees of magnitude*.

4.7 Comparison of NL with DINGO

In our next experiment, we compare NL with DINGO [7]. The results, presented in Figure 6, show that our method are more communication efficient than DINGO by *many orders of magnitude*. This is true for all experiments.

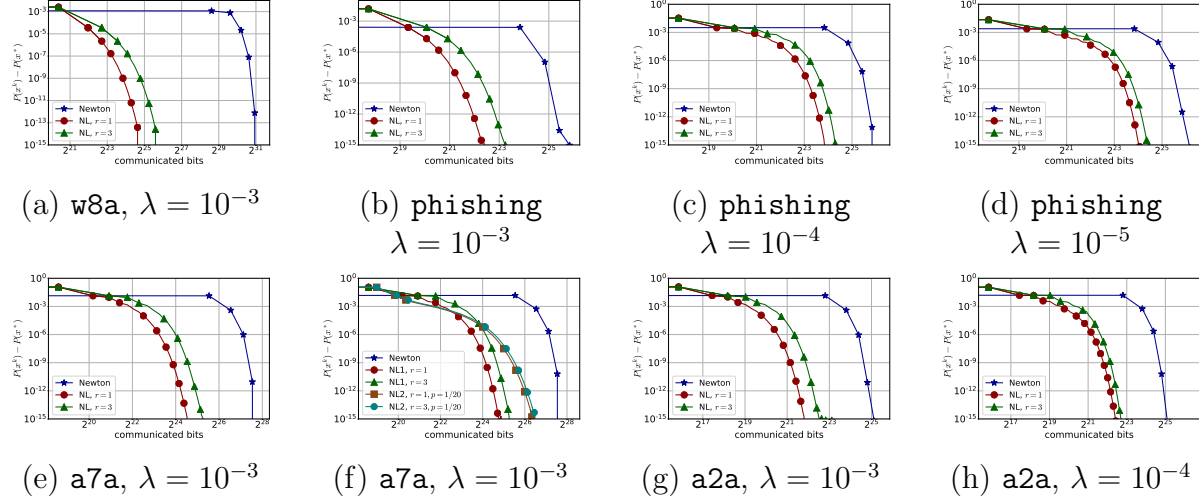


Figure 3: Comparison of NL with Newton's method in terms of communication complexity.

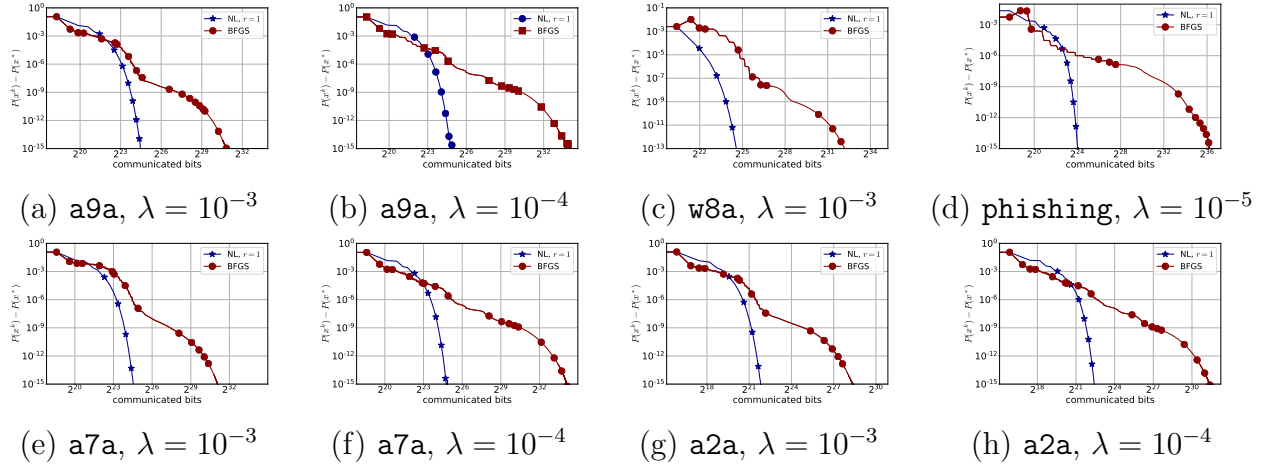


Figure 4: Comparison of NL and BFGS in terms of communication complexity.

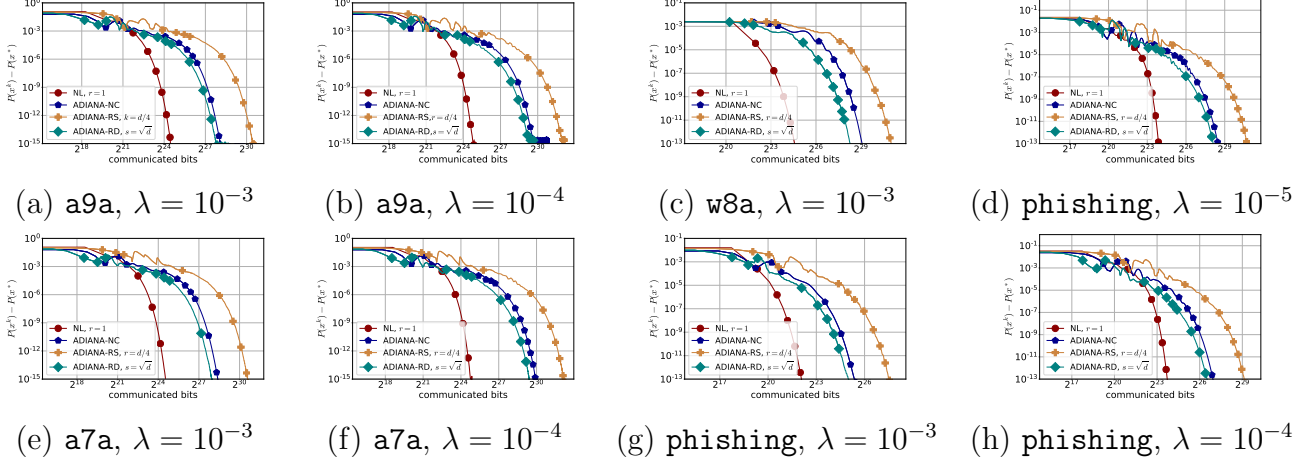


Figure 5: Comparison of NL with ADIANA in terms of communication complexity.

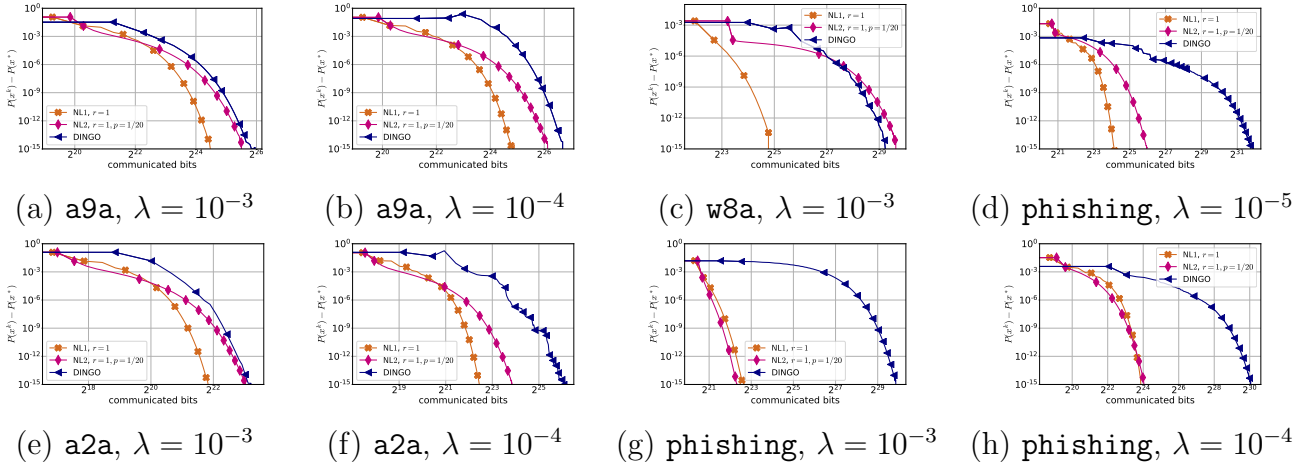


Figure 6: Comparison of NL with DINGO in terms of communication complexity.

Appendix

A Proofs for NL (Section 3.5)

Let $W^k := (h_1^k, \dots, h_n^k, x^k)$.

A.1 Lemma

We first establish a lemma.

Lemma A.1. *Let $S^2 := \sum_{i=1}^n \sum_{j=1}^m \|a_{ij}\|^2$. For $\eta \leq \frac{1}{\omega+1}$ and all $k \geq 0$, we have*

$$\mathbb{E} [\mathcal{H}^{k+1} \mid W^k] \leq (1 - \eta) \mathcal{H}^k + \eta \nu^2 S^2 \|x^k - x^*\|^2.$$

Proof. First, recall that

$$\mathcal{H}^k = \sum_{i=1}^n \|h_i^k - h_i(x^*)\|^2. \quad (27)$$

Since $h_i(x^*) \in \mathbb{R}_+^m$ due to convexity of f_{ij} , we have

$$\|[x]_+ - h_i(x^*)\|^2 \leq \|x - h_i(x^*)\|^2, \quad \text{for all } x \in \mathbb{R}^m. \quad (28)$$

Then as long as $\eta \leq \frac{1}{\omega+1}$, we have

$$\begin{aligned} \mathbb{E} [\mathcal{H}^{k+1} \mid W^k] &= \mathbb{E} \left[\sum_{i=1}^n \|h_i^{k+1} - h_i(x^*)\|^2 \mid W^k \right] \\ &= \sum_{i=1}^n \mathbb{E} \left[\|[h_i^k + \eta \mathcal{C}_i^k(h_i(x^k) - h_i^k)]_+ - h_i(x^*)\|^2 \mid W^k \right] \\ &\stackrel{(28)}{\leq} \sum_{i=1}^n \mathbb{E} \left[\|h_i^k + \eta \mathcal{C}_i^k(h_i(x^k) - h_i^k) - h_i(x^*)\|^2 \mid W^k \right] \\ &= \sum_{i=1}^n \|h_i^k - h_i(x^*)\|^2 + 2\eta \sum_{i=1}^n \mathbb{E} [\langle \mathcal{C}_i^k(h_i(x^k) - h_i^k), h_i^k - h_i(x^*) \rangle \mid W^k] \\ &\quad + \eta^2 \sum_{i=1}^n \mathbb{E} [\|\mathcal{C}_i^k(h_i(x^k) - h_i^k)\|^2 \mid W^k] \\ &\stackrel{(27)+(23)+(24)}{\leq} \mathcal{H}^k + 2\eta \sum_{i=1}^n \langle h_i(x^k) - h_i^k, h_i^k - h_i(x^*) \rangle \\ &\quad + \eta^2 \sum_{i=1}^n (\omega + 1) \|h_i(x^k) - h_i^k\|^2. \end{aligned} \quad (29)$$

Using the stepsize restriction $\eta \leq \frac{1}{\omega+1}$, we can bound $\eta^2(\omega + 1) \leq \eta$. Plugging this back in

(29), we get

$$\begin{aligned}
\mathbb{E} [\mathcal{H}^{k+1} \mid W^k] &\leq \mathcal{H}^k + \eta \sum_{i=1}^n \langle h_i(x^k) - h_i^k, h_i(x^k) + h_i^k - 2h_i(x^*) \rangle \\
&= \mathcal{H}^k + \eta \sum_{i=1}^n \left(\|h_i(x^k) - h_i(x^*)\|^2 - \|h_i^k - h_i(x^*)\|^2 \right) \\
&= (1 - \eta) \mathcal{H}^k + \eta \sum_{i=1}^n \|h_i(x^k) - h_i(x^*)\|^2 \\
&\stackrel{(11)}{=} (1 - \eta) \mathcal{H}^k + \eta \sum_{i=1}^n \sum_{j=1}^m (h_{ij}(x^k) - h_{ij}(x^*))^2 \\
&\stackrel{(6)}{\leq} (1 - \eta) \mathcal{H}^k + \eta \sum_{i=1}^n \sum_{j=1}^m \nu^2 \|a_{ij}\|^2 \|x^k - x^*\|^2 \\
&\leq (1 - \eta) \mathcal{H}^k + \eta \nu^2 S^2 \|x^k - x^*\|^2.
\end{aligned}$$

□

A.2 Proof of Theorem 3.2

It is easy to see that

$$\mathbf{H}^k = \frac{1}{n} \sum_{i=1}^n \mathbf{H}_i^k, \quad \mathbf{H}_i^k = \frac{1}{m} \sum_{j=1}^m h_{ij}^k a_{ij} a_{ij}^\top. \quad (30)$$

By the first order necessary optimality conditions, we have

$$\nabla f(x^*) + \lambda x^* = \nabla P(x^*) = 0. \quad (31)$$

Furthermore, since we maintain $\mathbf{H}^k \succeq 0$ for all k , we have $\mathbf{H}^k + \lambda \mathbf{I} \succeq \lambda \mathbf{I}$, whence

$$\|(\mathbf{H}^k + \lambda \mathbf{I})^{-1}\| \leq \frac{1}{\lambda}. \quad (32)$$

Then we can write

$$\begin{aligned}
\|x^{k+1} - x^*\| &\stackrel{(21)}{=} \|x^k - x^* - (\mathbf{H}^k + \lambda \mathbf{I})^{-1}(\nabla f(x^k) + \lambda x^k)\| \\
&= \|(\mathbf{H}^k + \lambda \mathbf{I})^{-1} ((\mathbf{H}^k + \lambda \mathbf{I})(x^k - x^*) - \nabla f(x^k) - \lambda x^k)\| \\
&\stackrel{(32)}{\leq} \frac{1}{\lambda} \|(\mathbf{H}^k + \lambda \mathbf{I})(x^k - x^*) - \nabla f(x^k) - \lambda x^k\| \\
&\stackrel{(31)}{=} \frac{1}{\lambda} \|\mathbf{H}^k(x^k - x^*) - (\nabla f(x^k) - \nabla f(x^*))\| \\
&= \frac{1}{\lambda} \left\| \frac{1}{n} \sum_{i=1}^n [\mathbf{H}_i^k(x^k - x^*) - (\nabla f_i(x^k) - \nabla f_i(x^*))] \right\| \\
&\leq \frac{1}{n\lambda} \sum_{i=1}^n \|\mathbf{H}_i^k(x^k - x^*) - (\nabla f_i(x^k) - \nabla f_i(x^*))\|.
\end{aligned}$$

where the last step follows from applying Jensen's inequality to the convex function $x \mapsto \|x\|$. Using (30) and (2) we get

$$\|x^{k+1} - x^*\| \leq \frac{1}{n\lambda} \sum_{i=1}^n \left\| \frac{1}{m} \sum_{j=1}^m [h_{ij}^k a_{ij} a_{ij}^\top (x^k - x^*) - (\nabla f_{ij}(x^k) - \nabla f_{ij}(x^*))] \right\|. \quad (33)$$

We can now express the difference of the gradients in integral form via the fundamental theorem of calculus, obtaining

$$\begin{aligned} \nabla f_{ij}(x^k) - \nabla f_{ij}(x^*) &= \int_0^1 \mathbf{H}_{ij}(x^* + \tau(x^k - x^*))(x^k - x^*) d\tau \\ &\stackrel{(4)}{=} \int_0^1 h_{ij}(x^* + \tau(x^k - x^*)) a_{ij} a_{ij}^\top (x^k - x^*) d\tau. \end{aligned}$$

We can plug this back into (33), which gives

$$\begin{aligned} &\|x^{k+1} - x^*\| \\ &\leq \frac{1}{n\lambda} \sum_{i=1}^n \frac{1}{m} \left\| \sum_{j=1}^m \left(h_{ij}^k a_{ij} a_{ij}^\top (x^k - x^*) - \int_0^1 h_{ij}(x^* + \tau(x^k - x^*)) a_{ij} a_{ij}^\top (x^k - x^*) d\tau \right) \right\| \\ &= \frac{1}{n\lambda} \sum_{i=1}^n \frac{1}{m} \left\| \sum_{j=1}^m a_{ij} a_{ij}^\top (x^k - x^*) \left(h_{ij}^k - \int_0^1 h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right) \right\| \\ &\leq \frac{\|x^k - x^*\|}{\lambda} \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|a_{ij}\|^2 \left| \int_0^1 h_{ij}^k - h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right|. \quad (34) \end{aligned}$$

From (6), we have

$$\begin{aligned} &\left| \int_0^1 h_{ij}^k - h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right| \\ &\leq \int_0^1 |h_{ij}^k - h_{ij}(x^* + \tau(x^k - x^*))| d\tau \\ &\leq |h_{ij}^k - h_{ij}(x^*)| + \int_0^1 |h_{ij}(x^*) - h_{ij}(x^* + \tau(x^k - x^*))| d\tau \\ &\stackrel{(6)}{\leq} |h_{ij}^k - h_{ij}(x^*)| + \int_0^1 \tau \nu \|a_{ij}\| \cdot \|x^k - x^*\| d\tau \\ &= |h_{ij}^k - h_{ij}(x^*)| + \frac{\nu \|a_{ij}\|}{2} \|x^k - x^*\|. \quad (35) \end{aligned}$$

By squaring both sides of (34), applying Jensen's inequality to the function $t \mapsto t^2$ in the

form $\left(\frac{1}{nm} \sum_i \sum_j t_{ij}\right)^2 \leq \frac{1}{nm} \sum_i \sum_j t_{ij}^2$, and plugging in the bound (35), we get

$$\begin{aligned}
& \|x^{k+1} - x^*\|^2 \\
& \leq \frac{\|x^k - x^*\|^2}{\lambda^2} \left(\frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|a_{ij}\|^2 \left| \int_0^1 h_{ij}^k - h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right| \right)^2 \\
& \leq \frac{\|x^k - x^*\|^2}{\lambda^2} \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \left(\|a_{ij}\|^2 \left| \int_0^1 h_{ij}^k - h_{ij}(x^* + \tau(x^k - x^*)) d\tau \right| \right)^2 \\
& \stackrel{(35)}{\leq} \frac{\|x^k - x^*\|^2}{\lambda^2} \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|a_{ij}\|^4 \left(|h_{ij}^k - h_{ij}(x^*)| + \frac{\nu \|a_{ij}\|}{2} \|x^k - x^*\| \right)^2 \\
& \leq \frac{\|x^k - x^*\|^2}{\lambda^2} \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \|a_{ij}\|^4 \left(2|h_{ij}^k - h_{ij}(x^*)|^2 + \frac{\nu^2 \|a_{ij}\|^2}{2} \|x^k - x^*\|^2 \right) \\
& \leq \frac{\|x^k - x^*\|^2}{\lambda^2} \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m R^4 \left(2|h_{ij}^k - h_{ij}(x^*)|^2 + \frac{\nu^2 R^2}{2} \|x^k - x^*\|^2 \right).
\end{aligned}$$

In the last step we have used Young's inequality.⁶

Since $R := \max_{i,j} \{\|a_{ij}\|\}$, we can further bound

$$\|x^{k+1} - x^*\|^2 \leq \frac{2R^4}{nm\lambda^2} \mathcal{H}^k \|x^k - x^*\|^2 + \frac{\nu^2 R^6}{2\lambda^2} \|x^k - x^*\|^4. \quad (36)$$

Assume $\|x^k - x^*\|^2 \leq \frac{\lambda^2}{12\nu^2 R^6}$ for all $k \geq 0$. Then from (36), we have

$$\begin{aligned}
\|x^{k+1} - x^*\|^2 & \leq \frac{\lambda^2}{12\nu^2 R^6} \cdot \frac{2R^4}{nm\lambda^2} \mathcal{H}^k + \frac{\lambda^2}{12\nu^2 R^6} \cdot \frac{\nu^2 R^6 \|x^k - x^*\|^2}{2\lambda^2} \\
& \leq \frac{1}{6nm\nu^2 R^2} \mathcal{H}^k + \frac{1}{24} \|x^k - x^*\|^2,
\end{aligned}$$

and by taking expectation, we have

$$\mathbb{E} [\|x^{k+1} - x^*\|^2 \mid W^k] \leq \frac{1}{6nm\nu^2 R^2} \mathcal{H}^k + \frac{1}{24} \|x^k - x^*\|^2. \quad (37)$$

Next, Lemma A.1 implies that

$$\mathbb{E} [\mathcal{H}^{k+1} \mid W^k] \leq (1 - \eta) \mathcal{H}^k + \eta nm\nu^2 R^2 \|x^k - x^*\|^2. \quad (38)$$

Recall that $\Phi_1^{k+1} := \|x^{k+1} - x^*\|^2 + \frac{1}{3\eta nm\nu^2 R^2} \mathcal{H}^{k+1}$. Combining (37) and (38), we get

$$\begin{aligned}
\mathbb{E} [\Phi_1^{k+1} \mid W^k] &= \mathbb{E} [\|x^{k+1} - x^*\|^2 \mid W^k] + \frac{1}{3\eta nm\nu^2 R^2} \mathbb{E} [\mathcal{H}^{k+1} \mid W^k] \\
&\stackrel{(37)}{\leq} \frac{1}{3\eta nm\nu^2 R^2} \left(1 - \eta + \frac{\eta}{2} \right) \mathcal{H}^k + \left(\frac{1}{24} + \frac{1}{3} \right) \|x^k - x^*\|^2 \\
&\leq \left(1 - \min \left\{ \frac{\eta}{2}, \frac{5}{8} \right\} \right) \Phi_1^k = \theta_1^k \Phi_1^k.
\end{aligned} \quad (39)$$

⁶ $(a+b)^2 \leq 2a^2 + 2b^2$.

By applying the tower property, we get

$$\mathbb{E} [\Phi_1^{k+1}] = \mathbb{E} [\mathbb{E} [\Phi_1^{k+1} | W^k]] \stackrel{(39)}{\leq} \theta_1 \mathbb{E} [\Phi_1^k].$$

Unrolling the recursion, we get $\mathbb{E} [\Phi_1^k] \leq \theta_1^k \Phi_1^0$, and the first claim is proved.

We further have $\mathbb{E} [\|x^k - x^*\|^2] \leq \theta_1^k \Phi_1^0$ and $\mathbb{E} [\mathcal{H}^k] \leq \theta_1^k 3\eta nm \nu^2 R^2 \Phi_1^0$. Assume $x^k \neq x^*$ for all k . Then from (36), we have

$$\frac{\|x^{k+1} - x^*\|^2}{\|x^k - x^*\|^2} \leq \frac{2R^4}{nm\lambda^2} \mathcal{H}^k + \frac{\nu^2 R^6}{2\lambda^2} \|x^k - x^*\|^2,$$

and by taking expectation, we obtain

$$\begin{aligned} \mathbb{E} \left[\frac{\|x^{k+1} - x^*\|^2}{\|x^k - x^*\|^2} \right] &\leq \frac{2R^4}{mn\lambda^2} \mathbb{E} [\mathcal{H}^k] + \frac{\nu^2 R^6}{2\lambda^2} \mathbb{E} [\|x^k - x^*\|^2] \\ &\leq \theta_1^k \left(6\eta + \frac{1}{2} \right) \frac{\nu^2 R^6}{\lambda^2} \Phi_1^0. \end{aligned}$$

A.3 Proof of Lemma 3.3

We prove this by induction. First, we have $\|x^0 - x^*\|^2 \leq \frac{\lambda^2}{12\nu^2 R^6}$ by the assumption. We assume $\|x^k - x^*\|^2 \leq \frac{\lambda^2}{12\nu^2 R^6}$ holds for all $k \leq K$. For $k \leq K$, since h_{ij}^k is a convex combination of $\{h_{ij}(x^0), h_{ij}(x^1), \dots, h_{ij}(x^k)\}$ for all i, j , we have

$$h_{ij}^k = \sum_{t=0}^k \rho_t h_{ij}(x^t), \quad \text{with} \quad \sum_{t=0}^k \rho_t = 1, \quad \text{and} \quad \rho_t \geq 0,$$

which implies that

$$\begin{aligned} \sum_{i=1}^n \|h_i^k - h_i(x^*)\|^2 &= \sum_{i=1}^n \sum_{j=1}^m |h_{ij}^k - h_{ij}(x^*)|^2 \\ &= \sum_{i=1}^n \sum_{j=1}^m \left| \sum_{t=0}^k \rho_t (h_{ij}(x^t) - h_{ij}(x^*)) \right|^2 \\ &\leq \sum_{i=1}^n \sum_{j=1}^m \sum_{t=0}^k \rho_t |h_{ij}(x^t) - h_{ij}(x^*)|^2 \\ &\stackrel{(6)}{\leq} \sum_{i=1}^n \sum_{j=1}^m \sum_{t=0}^k \rho_t \nu^2 \|a_{ij}\|^2 \|x^t - x^*\|^2 \\ &\stackrel{\text{As. 3.1}}{\leq} \nu^2 R^2 \sum_{i=1}^n \sum_{j=1}^m \sum_{t=0}^k \rho_t \cdot \frac{\lambda^2}{12\nu^2 R^6} \\ &= \frac{mn\lambda^2}{12R^4}, \end{aligned}$$

for $k \leq K$.

Combining the above inequality and (36), we arrive at

$$\begin{aligned}
\|x^{K+1} - x^*\|^2 &\leq \frac{2\|x^K - x^*\|^2 R^4}{mn\lambda^2} \sum_{i=1}^n \|h_i^K - h_i(x^*)\|^2 + \frac{\nu^2 R^6 \|x^K - x^*\|^4}{2\lambda^2} \\
&\leq \frac{2\|x^K - x^*\|^2 R^4}{mn\lambda^2} \cdot \frac{mn\lambda^2}{12R^4} + \frac{\nu^2 R^6 \|x^K - x^*\|^4}{2\lambda^2} \\
&\leq \frac{1}{6} \|x^K - x^*\|^2 + \frac{\nu^2 R^6 \|x^K - x^*\|^4}{2\lambda^2} \\
&\leq \frac{1}{6} \cdot \frac{\lambda^2}{12\nu^2 R^6} + \frac{\nu^2 R^6}{2\lambda^2} \cdot \left(\frac{\lambda^2}{12\nu^2 R^6} \right)^2 \\
&\leq \frac{\lambda^2}{12\nu^2 R^6}.
\end{aligned}$$

References

- [1] D. Alistarh, D. Grubic, J. Li, R. Tomioka, and M. Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in Neural Information Processing Systems*, pages 1709–1720, 2017.
- [2] Amir Beck. *Introduction to Nonlinear Optimization: Theory, Algorithms, and Applications with MATLAB*. Society for Industrial and Applied Mathematics, USA, 2014. ISBN 1611973643.
- [3] Ron Bekkerman, Mikhail Bilenko, and John Langford. *Scaling up machine learning: Parallel and distributed approaches*. Cambridge University Press, 2011.
- [4] J. Bernstein, Y. X. Wang, K. Azizzadenesheli, and A. Anandkumar. SignSGD: Compressed optimisation for non-convex problems. *The 35th International Conference on Machine Learning*, pages 560–569, 2018.
- [5] Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compression for distributed learning. *arXiv:2002.12410*, 2020.
- [6] Charles G Broyden. Quasi-newton methods and their application to function minimisation. *Mathematics of Computation*, 21(99):368–381, 1967.
- [7] Rixon Crane and Fred Roosta. DINGO: Distributed Newton-type method for gradient-norm optimization. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32, pages 9498–9508. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/9718db12cae6be37f7349779007ee589-Paper.pdf>.
- [8] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 1646–1654. Curran Associates, Inc., 2014.
- [9] Rodger Fletcher. A new approach to variable metric algorithms. *The Computer Journal*, 13(3):317–323, 1970.
- [10] Avishek Ghosh, Raj Kumar Maity, Arya Mazumdar, and Kannan Ramchandran. Communication efficient distributed approximate Newton method. In *IEEE International Symposium on Information Theory (ISIT)*, 2020. doi: 10.1109/ISIT44484.2020.9174216.
- [11] Donald Goldfarb. A family of variable-metric methods derived by variational means. *Mathematics of Computation*, 24(109):23–26, 1970.
- [12] Andreas Griewank. The modification of Newton’s method for unconstrained optimization by bounding cubic terms. Technical report, Department of Applied Mathematics and Theoretical Physics, University of Cambridge, 1981. Technical Report NA/12.

- [13] Samuel Horváth, Chen-Yu Ho, Ľudovít Horváth, Atal Narayan Sahu, Marco Canini, and Peter Richtárik. Natural compression for distributed deep learning. *arXiv preprint arXiv:1905.10988*, 2019.
- [14] Samuel Horváth, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. Stochastic distributed learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*, 2019.
- [15] Rustem Islamov, Xun Qian, and Peter Richtárik. Distributed second order methods with fast rates and compressed communication. *arXiv preprint arXiv:2102.07158*, 2021.
- [16] Rie Johnson and Tong Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *NIPS*, pages 315–323, 2013.
- [17] Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local SGD on identical and heterogeneous data. In *The 23rd International Conference on Artificial Intelligence and Statistics (AISTATS 2020)*, 2020.
- [18] Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with compressed gradients. In *arXiv preprint arXiv:1806.06573*, 2018.
- [19] Dmitry Kovalev, Konstantin Mishchenko, and Peter Richtárik. Stochastic Newton and cubic Newton methods with simple local linear-quadratic rates. In *NeurIPS Beyond First Order Methods Workshop*, 2019.
- [20] Dmitry Kovalev, Robert M. Gower, Peter Richtárik, and Alexander Rogozin. Fast linear convergence of randomized BFGS. *arXiv:2002.11337*, 2020.
- [21] Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for compressed gradient descent in distributed and federated optimization. In *International Conference on Machine Learning*, 2020.
- [22] Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. A universal catalyst for first-order optimization. *arXiv preprint arXiv:1506.02186*, 2015.
- [23] Nicolas Loizou and Peter Richtárik. Linearly convergent stochastic heavy ball method for minimizing generalization error. In *NIPS Workshop on Optimization for Machine Learning*, 2017.
- [24] Nicolas Loizou and Peter Richtárik. Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. *arXiv:1712.09677*, 2017.
- [25] Chenxin Ma, Jakub Konečný, Martin Jaggi, Virginia Smith, Michael I. Jordan, Peter Richtárik, and Martin Takáč. Distributed optimization with arbitrary local solvers. *Optimization Methods and Software*, 32(4):813–848, 2017.
- [26] Yura Malitsky and Konstantin Mishchenko. Adaptive gradient descent without descent. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 6702–6712. PMLR, 13–18 Jul 2019.

- [27] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*, 2019.
- [28] Konstantin Mishchenko, Ahmed Khaled, and Peter Richtárik. Random reshuffling: Simple analysis with vast improvements. In *34th Conference on Neural Information Processing Systems (NeurIPS 2020)*, 2020.
- [29] Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic gradient descent, weighted sampling, and the randomized Kaczmarz algorithm. *Mathematical Programming*, 155(1–2):549–573, 2015.
- [30] Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on Optimization*, 19(4):1574–1609, 2009.
- [31] Yurii Nesterov and Boris T. Polyak. Cubic regularization of Newton method and its global performance. *Mathematical Programming*, 108(1):177–205, 2006.
- [32] Boris Polyak and Andrey Tremba. New versions of Newton method: step-size choice, convergence domain and under-determined equations. *arXiv preprint arXiv:1703.07810*, 2019.
- [33] Boris Polyak and Andrey Tremba. New versions of newton method: step-size choice, convergence domain and under-determined equations. *Optimization Methods and Software*, 35(6):1272–1303, 2020.
- [34] Joseph Raphson. Analysis aequationum universalis seu ad aequationes algebraicas resolvendas methodus generalis, & expedita, ex nova infinitarum serierum methodo, deducta ac demonstrata. *Oxford: Richard Davis*, 1697.
- [35] Sashank J. Reddi, Jakub Konečný, Peter Richtárik, Barnabás Póczos, and Alex Smola. AIDE: fast and communication efficient distributed optimization. *arXiv:1608.06879*, 2016.
- [36] H. Robbins and S. Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22:400–407, 1951.
- [37] M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *Math. Program.*, 162(1-2):83–112, 2017.
- [38] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu. 1-bit stochastic gradient descent and its application to data- parallel distributed training of speech DNNs. *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [39] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: from theory to algorithms*. Cambridge University Press, 2014.
- [40] Ohad Shamir, Nati Srebro, and Tong Zhang. Communication-efficient distributed optimization using an approximate Newton-type method. In *Proceedings of the 31st International Conference on Machine Learning, PMLR*, volume 32, pages 1000–1008, 2014.

- [41] David F Shanno. Conditioning of quasi-Newton methods for function minimization. *Mathematics of computation*, 24(111):647–656, 1970.
- [42] S. U. Stich and S. P. Karimireddy. The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication. *arXiv: 1909.05350*, 2019.
- [43] Sebastian U. Stich. Local SGD converges fast and communicates little. In *International Conference on Learning Representations*, 2020.
- [44] Sebastian U. Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with memory. In *Advances in Neural Information Processing Systems*, pages 4447–4458, 2018.
- [45] Martin Takáč, Avleen Bijral, Peter Richtárik, and Nathan Srebro. Mini-batch primal and dual methods for SVMs. In *30th International Conference on Machine Learning*, pages 537–552, 2013.
- [46] H. Tang, X. Lian, T. Zhang, and J. Liu. DoubleSqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. In *Proceedings of the 36th International Conference on Machine Learning*, pages 6155–6165, 2019.
- [47] John Wallis. A treatise of algebra, both historical and practical. *Philosophical Transactions of the Royal Society of London*, 15(173):1095–1106, 1685. doi: 10.1098/rstl.1685.0053. URL <https://royalsocietypublishing.org/doi/abs/10.1098/rstl.1685.0053>.
- [48] Shusen Wang, Fred Roosta, Peng Xu, and Michael W Mahoney. Giant: Globally improved approximate newton method for distributed optimization. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages 2332–2342. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/dabd8d2ce74e782c65a973ef76fd540b-Paper.pdf>.
- [49] W. Wen, C. Xu, F. Yan, C. Wu, Y. Wang, and H. Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. *Advances in Neural Information Processing Systems*, pages 1509–1519, 2017.
- [50] Lin Xiao and Tong Zhang. A proximal stochastic gradient method with progressive variance reduction. *SIAM Journal on Optimization*, 24(4):2057–2075, 2014.
- [51] Jiaqi Zhang, Keyou You, and Tamer Başar. Distributed adaptive Newton methods with globally superlinear convergence. *arXiv preprint arXiv:2002.07378*, 2020.
- [52] Yuchen Zhang and Lin Xiao. DiSCO: Distributed optimization for self-concordant empirical loss. In *Proceedings of the 32nd International Conference on Machine Learning*, PMLR, volume 37, pages 362–370, 2015.
- [53] Peilin Zhao and Tong Zhang. Stochastic optimization with importance sampling. *The 32nd International Conference on Machine Learning*, 37:1–9, 2015.