
Towards a Better Theoretical Understanding of Independent Subnet Training*

Egor Shulgin¹ Peter Richtárik¹

Abstract

Modern advancements in large-scale machine learning would be impossible without the paradigm of data-parallel distributed computing. Since distributed computing with large-scale models imparts excessive pressure on communication channels, significant recent research has been directed toward co-designing communication compression strategies and training algorithms with the goal of reducing communication costs. While pure data parallelism allows better data scaling, it suffers from poor model scaling properties. Indeed, compute nodes are severely limited by memory constraints, preventing further increases in model size. For this reason, the latest achievements in training giant neural network models also rely on some form of model parallelism. In this work, we take a closer theoretical look at Independent Subnetwork Training (IST), which is a recently proposed and highly effective technique for solving the aforementioned problems. We identify fundamental differences between IST and alternative approaches, such as distributed methods with compressed communication, and provide a precise analysis of its optimization performance on a quadratic model.

1. Introduction

A huge part of today’s machine learning success is driven by the possibility of building more and more complex models and training them on increasingly larger datasets. This rapid progress has become feasible due to advancements in distributed optimization, which is necessary for proper scaling when the size of the training data grows (Zinkevich et al., 2010). In a typical scenario, data parallelism is used for efficiency and implies sharding the dataset across computing

devices. This allowed very efficient scaling and acceleration of training moderately sized models by using additional hardware (Goyal et al., 2018). However, this data parallel approach can suffer from communication bottleneck, which has sparked extensive research on distributed optimization with compressed communication of the parameters between nodes (Alistarh et al., 2017; Konečný et al., 2016; Seide et al., 2014).

1.1. The need for model parallelism

Despite its efficiency, data parallelism has some fundamental limitations when it comes to scaling up the model size. As the dimensions of a model increase, the amount of memory required to store and update the parameters also increases, which becomes problematic due to resource constraints on individual devices. This has led to the development of model parallelism (Dean et al., 2012; Richtárik & Takáč, 2016), which splits a large model across multiple nodes, with each node responsible for computations of parts of the model (Farber & Asanovic, 1997; Zhang et al., 1989). However, naive model parallelism also poses challenges because each node can only update its portion of the model based on the data it has access to. This creates a need for very careful management of communication between devices. Thus, a combination of both data and model parallelism is often necessary to achieve efficient and scalable training of huge models.

Independent Subnetwork Training (IST) is a technique that suggests dividing a neural network into smaller subparts, training them in a distributed parallel fashion, and then aggregating the results to update the weights of the whole model. In IST, every subnetwork can operate independently and has fewer parameters than the full model, which not only reduces the load on computing nodes but also results in faster synchronization. This paradigm was pioneered by Yuan et al. (2022) for networks with fully connected layers and was later extended to ResNets (Dun et al., 2022) and Graph architectures (Wolfe et al., 2021). Previous experimental studies have shown that IST is a very promising approach for various applications as it allows to effectively combine data and model parallelism and train larger models with limited compute. In addition, Liao & Kyrillidis (2022) performed theoretical analysis of IST for overparameterized single hidden layer neural networks with ReLU activations.

¹KAUST AI Initiative, King Abdullah University of Science and Technology, Thuwal, Saudi Arabia. Correspondence to: Egor Shulgin <egor.shulgin@kaust.edu.sa>.

*A full version of the paper is available on [arXiv:2306.16484](https://arxiv.org/abs/2306.16484).

The idea of IST was also recently extended to the federated setting via an asynchronous distributed dropout technique (Dun et al., 2023).

Federated Learning. Another important setting when the data is distributed (due to privacy reasons) is Federated Learning (Kairouz et al., 2021; Konečný et al., 2016; McMahan et al., 2017). In this scenario, computing devices are often heterogeneous and more resource-constrained (Caldas et al., 2018) (e.g. mobile phones) in comparison to data-center settings. Such challenges have prompted extensive research efforts into selecting smaller and more efficient submodels for local on-device training (Alam et al., 2022; Charles et al., 2022; Chen et al., 2022; Diao et al., 2021; Horvath et al., 2021; Jiang et al., 2022; Lin et al., 2022; Qiu et al., 2022; Wen et al., 2022; Yang et al., 2022). Many of these works propose approaches to adapt submodels, often tailored to specific neural network architectures, based on the capabilities of individual clients for various machine learning tasks. However, there is a lack of comprehension regarding the theoretical properties of these methods.

1.2. Summary of contributions

After reviewing the literature, we found that a rigorous understanding of IST convergence is virtually non-existent, which motivated this work. The main contributions of this paper include: • A novel approach to analyzing distributed methods that combine data and model parallelism by operating with sparse submodels for a quadratic model. • The first analysis of independent subnetwork training in homogeneous and heterogeneous scenarios without restrictive assumptions on gradient estimators. • Identification of the settings when IST can optimize very efficiently or not converge to the optimal solution but only to an irreducible neighborhood that is also tightly characterized. • Experimental validation of the proposed theory through carefully designed illustrative experiments. The results, together with all the proofs, are given in the Appendix.

2. Formalism and setup

We consider the standard optimization formulation of a distributed/federated learning problem (Wang et al., 2021)

$$\min_{x \in \mathbb{R}^d} \left[f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) \right], \quad (1)$$

where n is the number of clients/workers, and each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}^d$ represents the loss of the model parameterized by vector $x \in \mathbb{R}^d$ on the data of client i .

A typical Stochastic Gradient Descent (SGD)-type method for solving this problem has the form

$$x^{k+1} = x^k - \gamma g^k, \quad g^k = \frac{1}{n} \sum_{i=1}^n g_i^k, \quad (2)$$

where $\gamma > 0$ is the stepsize and g_i^k is a suitably constructed estimator of $\nabla f_i(x^k)$. In the distributed setting, computation of gradient estimators g_i^k is typically performed by clients, and the results are sent to the server, which subsequently performs aggregation via averaging $g^k = \frac{1}{n} \sum_{i=1}^n g_i^k$. The average is then used to update the model x^{k+1} via a gradient-type method (2), and at the next iteration, the model is broadcasted back to the clients. The process is repeated iteratively until a suitable model is found.

One of the main techniques used to accelerate distributed training is lossy *communication compression* (Alistarh et al., 2017; Konečný et al., 2016; Seide et al., 2014), which suggests applying a (possibly randomized) lossy compression mapping \mathcal{C} to a vector/matrix/tensor x before broadcasting. This reduces the bits sent per communication round at the cost of transmitting a less accurate estimate $\mathcal{C}(x)$ of x . Described technique can be formalized in the following definition.

Definition 2.1 (Unbiased compressor). A randomized mapping $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is an **unbiased compression operator** ($\mathcal{C} \in \mathbb{U}(\omega)$ for brevity) if for some $\omega \geq 0$ and $\forall x \in \mathbb{R}^d$

$$\mathbb{E} [\mathcal{C}(x)] = x, \quad \mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq \omega \|x\|^2. \quad (3)$$

A notable example of a mapping from this class is the *random sparsification* ($\text{Rand-}q$ for $q \in \{1, \dots, d\}$) operator defined by

$$\mathbf{C}_{\text{Rand-}q}(x) := \mathbf{C}_q x = \frac{d}{q} \sum_{i \in S} e_i e_i^\top x, \quad (4)$$

where $e_1, \dots, e_d \in \mathbb{R}^d$ are standard unit basis vectors in \mathbb{R}^d , and S is a random subset of $[d] := \{1, \dots, d\}$ sampled from the uniform distribution on the all subsets of $[d]$ with cardinality q . $\text{Rand-}q$ belongs to $\mathbb{U}(d/q - 1)$, which means that the more elements are “dropped” (lower q), the higher the variance ω of the compressor.

In this work, we are mainly interested in a somewhat more general class of operators than mere sparsifiers. In particular, we are interested in compressing via the application of random matrices, i.e., via *sketching*. A sketch $\mathbf{C}_i^k \in \mathbb{R}^{d \times d}$ can be used to represent submodel computations in the following way:

$$g_i^k := \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k), \quad (5)$$

where we require \mathbf{C}_i^k to be a symmetric positive semi-definite matrix. Such gradient estimates correspond to computing the local gradient with respect to a sparse submodel $\mathbf{C}_i^k x^k$, and additionally sketching the resulting gradient with the same matrix \mathbf{C}_i^k to guarantee that the resulting update lies in the lower-dimensional subspace.

Using this notion, IST algorithm (with one local gradient step) can be represented in the following form:

$$x^{k+1} = \frac{1}{n} \sum_{i=1}^n [\mathbf{C}_i^k x^k - \gamma \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k)], \quad (6)$$

which is equivalent to the SGD-type update (2) when the *perfect reconstruction* property holds (with probability one)

$$\mathbf{C}^k := \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k = \mathbf{I},$$

where \mathbf{I} is the identity matrix. This property is inherent for a specific class of compressors that are particularly useful for capturing the concept of an *independent* subnetwork partition.

Definition 2.2 (Permutation sketch). Assume that model size is greater than the number of clients $d \geq n$ and $d = qn$, where $q \geq 1$ is an integer*. Let $\pi = (\pi_1, \dots, \pi_d)$ be a random permutation of $[d]$. Then for all $x \in \mathbb{R}^d$ and each $i \in [n]$, we define $\text{Perm-}q$ operator

$$\mathbf{C}_i := n \cdot \sum_{j=q(i-1)+1}^{qi} e_{\pi_j} e_{\pi_j}^\top. \quad (7)$$

$\text{Perm-}q$ is unbiased and can be conveniently used for representing a structured decomposition of the model, such that every client i is responsible for computations over a submodel $\mathbf{C}_i x^k$.

Our convergence analysis relies on the assumption that was previously used for coordinate descent-type methods.

Assumption 2.3 (Matrix smoothness). A differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is \mathbf{L} -smooth, if there exists a positive semi-definite matrix $\mathbf{L} \in \mathbb{R}^{d \times d}$ such that $\forall x, h \in \mathbb{R}^d$

$$f(x+h) \leq f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle \mathbf{L}h, h \rangle. \quad (8)$$

A standard L -smoothness condition is obtained as a special case of (8) for $\mathbf{L} = L \cdot \mathbf{I}$. Matrix smoothness was previously used for designing data-dependent gradient sparsification to accelerate optimization in communication-constrained settings (Safaryan et al., 2021; Wang et al., 2022).

2.1. Simplifications taken

To conduct a thorough theoretical analysis of methods that combine data with model parallelism, we simplify the algorithm and problem setting to isolate the unique effects of this approach. The following considerations are made:

(1) We assume that every node i computes the true gradient at the submodel $\mathbf{C}_i \nabla f_i(\mathbf{C}_i x^k)$.

(2) A notable difference compared to the original IST is that workers perform a single gradient descent step (or just gradient computation).

*While this condition may look restrictive, it naturally holds for distributed learning in a data-center setting. Permutation sparsifiers were introduced by (Szlendak et al., 2022) and generalized to other scenarios (like $n \geq d$).

(3) Finally, we consider a special case of a quadratic model (9) as a loss function (1).

Condition (1) is mainly for the sake of simplicity and clarity of exposition and can be potentially generalized to stochastic gradient computations. Condition (2) is imposed because local steps did not bring any theoretical efficiency improvements for heterogeneous settings until very recently (Mishchenko et al., 2022), and even then, only with the introduction of additional control variables, which goes against the requirements of resource-constrained device settings. The reason behind (3) is that despite its apparent simplicity, the quadratic problem has been used extensively to study properties of neural networks (Zhang et al., 2019; Zhu et al., 2022). Moreover, it is a non-trivial model, which makes it possible to understand complex optimization algorithms (Arjevani et al., 2020; Cunha et al., 2022; Goujaud et al., 2022). The quadratic problem is suitable for observing complex phenomena and providing theoretical insights, which can also be observed in practical scenarios.

Having said that, we consider a special case of problem (1) for symmetric matrices \mathbf{L}_i

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad f_i(x) \equiv \frac{1}{2} x^\top \mathbf{L}_i x - x^\top \mathbf{b}_i. \quad (9)$$

In this case, $f(x)$ is $\bar{\mathbf{L}}$ -smooth, and $\nabla f(x) = \bar{\mathbf{L}}x - \bar{\mathbf{b}}$, where $\bar{\mathbf{L}} = \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i$ and $\bar{\mathbf{b}} := \frac{1}{n} \sum_{i=1}^n \mathbf{b}_i$.

3. Results in the interpolation case

First, let us examine the case of $\mathbf{b}_i \equiv 0$, which we call interpolation for quadratics, and perform the analysis for general sketches \mathbf{C}_i^k . In this case, the gradient estimator (2) takes the form

$$g^k = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k) = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{L}_i \mathbf{C}_i^k x^k = \bar{\mathbf{B}}^k x^k \quad (10)$$

where $\bar{\mathbf{B}}^k := \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{L}_i \mathbf{C}_i^k$. We prove the following result for a method with such an estimator.

Theorem 3.1. Consider the method (2) with estimator (10) for a quadratic problem (9) with $\bar{\mathbf{L}} \succ 0$ and $\mathbf{b}_i \equiv 0$. Then if $\bar{\mathbf{W}} := \frac{1}{2} \mathbb{E} [\bar{\mathbf{L}} \bar{\mathbf{B}}^k + \bar{\mathbf{B}}^k \bar{\mathbf{L}}] \succeq 0$ and there exists a constant $\theta > 0$:

$$\mathbb{E} [\bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k] \preceq \theta \bar{\mathbf{W}}, \quad (11)$$

and the step size is chosen as $0 < \gamma \leq \frac{1}{\theta}$, the iterates satisfy

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}}^2] \leq \frac{2(f(x^0) - \mathbb{E}[f(x^K)])}{\gamma K}. \quad (12)$$

This theorem establishes an $\mathcal{O}(1/K)$ convergence rate with a constant step size up to a stationary point. Note that we employ weighted norms in our analysis, as the considered class of loss functions satisfies the matrix $\bar{\mathbf{L}}$ -smoothness Assumption 2.3. The use of standard Euclidean distance may result in loose bounds that do not recover correct rates for special cases like gradient descent.

It is important to highlight that the inequality (11) may not hold (for any $\theta > 0$) in the general case as the matrix $\bar{\mathbf{W}}$ is not guaranteed to be positive (semi-)definite in the case of general sampling. The intuition behind this issue is that arbitrary sketches \mathbf{C}_i^k can result in the gradient estimator g^k , which is misaligned with the true gradient $\nabla f(x^k)$. Specifically, the inner product $\langle \nabla f(x^k), g^k \rangle$ can be negative, and there is no expected descent after one step.

Next, we give examples of samplings for which the inequality (11) can be satisfied.

1. Identity. Consider $\mathbf{C}_i \equiv \mathbf{I}$. Then $\bar{\mathbf{B}}^k = \bar{\mathbf{L}}, \bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k = \bar{\mathbf{L}}^3, \bar{\mathbf{W}} = \bar{\mathbf{L}}^2 \succ 0$ and hence (11) is satisfied for $\theta = \lambda_{\max}(\bar{\mathbf{L}})$. So, (12) says that if we choose $\gamma = 1/\theta$, then

$$\frac{1}{K} \sum_{k=0}^{K-1} \|\nabla f(x^k)\|_{\mathbf{I}}^2 \leq \frac{2\lambda_{\max}(\bar{\mathbf{L}})(f(x^0) - f(x^K))}{K},$$

which exactly matches the rate of gradient descent in the non-convex setting.

2. Permutation. Assume[†] $n = d$ and the use of Perm-1 sketch $\mathbf{C}_i^k = n e_{\pi_i^k} e_{\pi_i^k}^\top$, where $\pi^k = (\pi_1^k, \dots, \pi_n^k)$ is a random permutation of $[n]$. Then

$$\mathbb{E} [\bar{\mathbf{B}}^k] = \frac{1}{n} \sum_{i=1}^n n^2 \mathbb{E} [\mathbf{C}_i^k \mathbf{L}_i \mathbf{C}_i^k] = \sum_{i=1}^n \mathbf{D}_i = n \bar{\mathbf{D}},$$

where $\bar{\mathbf{D}} := \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i, \mathbf{D}_i := \text{Diag}(\mathbf{L}_i)$. Then inequality (11) leads to

$$n \bar{\mathbf{D}} \bar{\mathbf{L}} \bar{\mathbf{D}} \preceq \frac{\theta}{2} (\bar{\mathbf{L}} \bar{\mathbf{D}} + \bar{\mathbf{D}} \bar{\mathbf{L}}), \quad (13)$$

which may not always hold as $\bar{\mathbf{L}} \bar{\mathbf{D}} + \bar{\mathbf{D}} \bar{\mathbf{L}}$ is not guaranteed to be positive-definite—even in the case of $\bar{\mathbf{L}} \succ 0$. However, such a condition can be enforced via a slight modification of the permutation sketches, which is done in Section 3.2. The limitation of such an approach is that the resulting compressors are no longer unbiased.

Next, we focus on the particular case of **permutation** sketches, which are the most suitable for model partitioning according to Independent Subnetwork Training (IST). In the rest of this section, we discuss how the condition (11) can be enforced via a specially designed preconditioning of the problem (9) or modification of the sketch mechanism (7).

[†]This is mainly done to simplify the presentation. Results can be generalized to the case of $n \neq d$ in a similar manner as in (Szlendak et al., 2022), which can be found in the Appendix.

3.1. Homogeneous problem preconditioning

To start, consider a homogeneous setting $f_i(x) = \frac{1}{2} x^\top \mathbf{L} x$, so $\mathbf{L}_i \equiv \mathbf{L}$. Now define $\mathbf{D} = \text{Diag}(\mathbf{L})$ – a diagonal matrix with elements equal to the diagonal of \mathbf{L} . Then, the problem can be converted to

$$f_i(\mathbf{D}^{-\frac{1}{2}} x) = \frac{1}{2} (\mathbf{D}^{-\frac{1}{2}} x)^\top \mathbf{L} (\mathbf{D}^{-\frac{1}{2}} x) = \frac{1}{2} x^\top \tilde{\mathbf{L}} x,$$

where $\tilde{\mathbf{L}} := \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$. It is equivalent to the original problem after changing the variables $\tilde{x} := \mathbf{D}^{-\frac{1}{2}} x$. Note that $\mathbf{D} = \text{Diag}(\mathbf{L})$ is positive-definite as $\mathbf{L} \succ 0$, and therefore $\tilde{\mathbf{L}} \succ 0$. Moreover, the preconditioned matrix $\tilde{\mathbf{L}}$ has all ones on the diagonal: $\text{Diag}(\tilde{\mathbf{L}}) = \mathbf{I}$. If we now combine (14) with Perm-1 sketches

$$\mathbb{E} [\bar{\mathbf{B}}^k] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}_i \tilde{\mathbf{L}} \mathbf{C}_i \right] = n \text{Diag}(\tilde{\mathbf{L}}) = n \mathbf{I}.$$

Therefore, inequality (11) takes the form $\bar{\mathbf{W}} = n \tilde{\mathbf{L}} \succeq \frac{1}{\theta} n^2 \tilde{\mathbf{L}}$, which holds for $\theta \geq n$, and the left-hand side of (12) can be transformed (for an accurate comparison to standard methods) in the following way:

$$\|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1} \bar{\mathbf{W}} \tilde{\mathbf{L}}^{-1}}^2 \geq n \lambda_{\min}(\tilde{\mathbf{L}}^{-1}) \|\nabla f(x^k)\|_{\mathbf{I}}^2.$$

The resulting convergence guarantee

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{I}}^2 \right] \leq \frac{2\lambda_{\max}(\tilde{\mathbf{L}})(f(x^0) - \mathbb{E}[f(x^K)])}{K},$$

which matches classical gradient descent.

3.2. Heterogeneous sketch preconditioning

In contrast to the homogeneous case, the heterogeneous problem $f_i(x) = \frac{1}{2} x^\top \mathbf{L}_i x$ cannot be so easily preconditioned by a simple change of variables $\tilde{x} := \mathbf{D}^{-\frac{1}{2}} x$, as every client i has its own matrix \mathbf{L}_i . However, this problem can be fixed via the following modification of Perm-1 , which scales the output according to the diagonal elements of the local smoothness matrix \mathbf{L}_i :

$$\tilde{\mathbf{C}}_i := \sqrt{n / [\mathbf{L}_i]_{\pi_i, \pi_i}} e_{\pi_i} e_{\pi_i}^\top. \quad (14)$$

In this case, $\mathbb{E} [\tilde{\mathbf{C}}_i \mathbf{L}_i \tilde{\mathbf{C}}_i] = \mathbf{I}, \mathbb{E} [\bar{\mathbf{B}}^k] = \mathbf{I}$, and $\bar{\mathbf{W}} = \bar{\mathbf{L}}$. Then inequality (11) is satisfied for $\theta \geq 1$.

If one inputs these results into (12), such convergence guarantee can be obtained

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{I}}^2 \right] \leq \frac{2\lambda_{\max}(\bar{\mathbf{L}})(f(x^0) - \mathbb{E}[f(x^K)])}{K},$$

which matches the gradient descent result as well. Thus, we can conclude that heterogeneity does not bring such a fundamental challenge in this scenario. In addition, the method with Perm-1 is significantly better in terms of computational and communication complexity, as it requires calculation of the local gradients with respect to much smaller submodels and transmits only sparse updates.

4. Irreducible bias in the general case

Now we look at the most general heterogeneous case with different matrices and linear terms $f_i(x) \equiv \frac{1}{2}x^\top \mathbf{L}_i x - x^\top \mathbf{b}_i$. In this instance, the gradient estimator (2) takes the form

$$g^k = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k (\mathbf{L}_i \mathbf{C}_i^k x^k - \mathbf{b}_i) = \overline{\mathbf{B}}^k x^k - \overline{\mathbf{C}}\mathbf{b}, \quad (15)$$

where $\overline{\mathbf{C}}\mathbf{b} = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{b}_i$. Herewith let us use a heterogeneous permutation sketch preconditioner (14), as in Section 3.2. Then $\mathbb{E}[\overline{\mathbf{B}}^k] = \mathbf{I}$ and $\mathbb{E}[\overline{\mathbf{C}}\mathbf{b}] = \frac{1}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b}$, where $\widetilde{\mathbf{D}}\mathbf{b} := \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i$. Furthermore, the expected gradient estimator (15) results in $\mathbb{E}[g^k] = x^k - \frac{1}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b}$ and can be transformed in the following manner:

$$\begin{aligned} \mathbb{E}[g^k] &= \overline{\mathbf{L}}^{-1} \overline{\mathbf{L}} x^k \pm \overline{\mathbf{L}}^{-1} \overline{\mathbf{b}} - \frac{1}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b} \\ &= \overline{\mathbf{L}}^{-1} \nabla f(x^k) + h, \end{aligned}$$

where $h := \overline{\mathbf{L}}^{-1} \overline{\mathbf{b}} - \frac{1}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b}$. Obtained formula reflects the decomposition of the estimator into the optimally preconditioned true gradient and a bias, depending on terms \mathbf{b}_i .

Estimator (16) can be directly plugged (with proper conditioning) into the general SGD update (2)

$$\mathbb{E}[x^{k+1}] = (1 - \gamma)^{k+1} x^0 + \frac{\gamma}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b} \sum_{j=0}^k (1 - \gamma)^j. \quad (16)$$

The resulting recursion (16) is exact, and its asymptotic limit can be analyzed. Thus, for constant $\gamma < 1$, by using the formula for the sum of the first k terms of a geometric series, one gets

$$\mathbb{E}[x^k] = (1 - \gamma)^k x^0 + \frac{1 - (1 - \gamma)^k}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b} \xrightarrow[k \rightarrow \infty]{} \frac{1}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b},$$

which shows that in the limit, the first initialization term (with x^0) vanishes while the second converges to $\frac{1}{\sqrt{n}} \widetilde{\mathbf{D}}\mathbf{b}$. This reasoning shows that the method does not converge to the exact solution

$$\mathbb{E}[x^k] \rightarrow x^\infty \neq x^* \in \arg \min_{x \in \mathbb{R}^d} \left\{ \frac{1}{2} x^\top \overline{\mathbf{L}} x - x^\top \overline{\mathbf{b}} \right\},$$

which for the positive-definite $\overline{\mathbf{L}}$ can be defined as $x^* = \overline{\mathbf{L}}^{-1} \overline{\mathbf{b}}$, while $x^\infty = \frac{1}{n\sqrt{n}} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i$. So, in general, there is an unavoidable bias. However, in the limit case: $n = d \rightarrow \infty$, the bias diminishes.

Theorem 4.1. *Consider the method (2) with the estimator (15) for the quadratic problem (9) with the positive-definite matrix $\overline{\mathbf{L}} \succ 0$. Assume that for every $\mathbf{D}_i := \text{Diag}(\mathbf{L}_i)$ matrices $\mathbf{D}_i^{-\frac{1}{2}}$ exist, scaled permutation sketches (14) are used,*

and heterogeneity is bounded as $\mathbb{E}[\|g^k - \mathbb{E}[g^k]\|_{\overline{\mathbf{L}}}^2] \leq \sigma^2$. Then, for the step size chosen as follows:

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1/2 - \beta}{\beta + 1/2}, \quad (17)$$

where $\gamma_{c,\beta} \in (0, 1]$ for $\beta \in (0, 1/2)$, the iterates satisfy

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla f(x^k)\|_{\overline{\mathbf{L}}^{-1}}^2] &\leq \frac{2(f(x^0) - \mathbb{E}[f(x^K)])}{\gamma K} \\ &+ \left(\frac{1 - \gamma}{0.5\beta} + \gamma \right) \|h\|_{\overline{\mathbf{L}}}^2 + \gamma \sigma^2, \end{aligned} \quad (18)$$

where $h = \overline{\mathbf{L}}^{-1} \overline{\mathbf{b}} - \frac{1}{n^{3/2}} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i$.

Note that the derived convergence upper bound has a neighborhood proportional to the bias of the gradient estimator h and level of heterogeneity σ^2 . Some of these terms with factor γ can be eliminated by decreasing the learning rate (e.g., $\sim 1/\sqrt{k}$). However, such a strategy does not diminish the term with a multiplier $2\beta^{-1}(1 - \gamma)$, making the neighborhood irreducible. Moreover, this term can be eliminated for $\gamma = 1$, which also minimizes the first term that decreases as $1/K$. However, this step size choice maximizes the terms with factor γ . Thus, there exists an inherent trade-off between convergence speed and the size of the neighborhood.

In addition, convergence to the stationary point is measured by the weighted $\overline{\mathbf{L}}^{-1}$ squared norm of the gradient. At the same time, the neighborhood term depends on the weighted by $\overline{\mathbf{L}}$ norm of h . This fine-grained decoupling is achieved by carefully applying the Fenchel-Young inequality and provides a tighter characterization of the convergence compared to using standard Euclidean distances.

Homogeneous case. In this scenario, every worker has access to all data $f_i(x) \equiv \frac{1}{2}x^\top \mathbf{L} x - x^\top \mathbf{b}$. Then diagonal preconditioning of the problem can be used, as in the previous Section 3.1. This results in a gradient $\nabla f(x) = \tilde{\mathbf{L}} x - \tilde{\mathbf{b}}$ for $\tilde{\mathbf{L}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$ and $\tilde{\mathbf{b}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{b}$. If this expression is further combined with a permutation sketch scaled by $1/\sqrt{n}$ $\mathbf{C}_i^k := \sqrt{n} e_{\pi_i} e_{\pi_i}^\top$, the resulting gradient estimator is:

$$g^k = x^k - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}} = \tilde{\mathbf{L}}^{-1} \nabla f(x^k) + \tilde{h}, \quad (19)$$

for $\tilde{h} = \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}$. In this case, the heterogeneity term σ^2 from the upper bound (18) disappears as $\mathbb{E}[\|g^k - \mathbb{E}[g^k]\|_{\overline{\mathbf{L}}}^2] = 0$, which can decrease the neighborhood size. However, the bias term depending on \tilde{h} still remains, as the method does not converge to the exact solution $x^k \rightarrow x^\infty \neq x^* = \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}}$ for positive-definite $\tilde{\mathbf{L}}$. Nevertheless the method's fixed point $x^\infty = \tilde{\mathbf{b}}/\sqrt{n}$ and solution x^* can coincide when $\tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} = \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}$, which means that $\tilde{\mathbf{b}}$ is the right eigenvector of matrix $\tilde{\mathbf{L}}^{-1}$ with eigenvalue $\frac{1}{\sqrt{n}}$.

References

- Ajalloeian, A. and Stich, S. U. On the convergence of SGD with biased gradients. *arXiv preprint arXiv:2008.00051*, 2020.
- Alam, S., Liu, L., Yan, M., and Zhang, M. FedRolex: Model-heterogeneous federated learning with rolling sub-model extraction. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=OtxyysUdBE>.
- Alistarh, D., Grubic, D., Li, J., Tomioka, R., and Vojnovic, M. QSGD: Communication-efficient SGD via gradient quantization and encoding. *Advances in Neural Information Processing Systems*, 30, 2017.
- Arjevani, Y., Shamir, O., and Srebro, N. A tight convergence analysis for stochastic gradient descent with delayed updates. In *Algorithmic Learning Theory*, pp. 111–132. PMLR, 2020.
- Caldas, S., Konečný, J., McMahan, H. B., and Talwalkar, A. Expanding the reach of federated learning by reducing client resource requirements. *arXiv preprint arXiv:1812.07210*, 2018.
- Charles, Z., Bonawitz, K., Chiknavaryan, S., McMahan, B., et al. Federated select: A primitive for communication- and memory-efficient federated learning. *arXiv preprint arXiv:2208.09432*, 2022.
- Chayti, E. M. and Karimireddy, S. P. Optimization with access to auxiliary information. *arXiv preprint arXiv:2206.00395*, 2022.
- Chen, Y., Chen, Z., Wu, P., and Yu, H. Fedobd: Opportunistic block dropout for efficiently training large-scale neural networks through federated learning. *arXiv preprint arXiv:2208.05174*, 2022.
- Chraïbi, S., Khaled, A., Kovalev, D., Richtárik, P., Salim, A., and Takáč, M. Distributed fixed point methods with compressed iterates. *arXiv preprint arXiv:2102.07245*, 2019.
- Cunha, L., Gidel, G., Pedregosa, F., Scieur, D., and Paquette, C. Only tails matter: Average-case universality and robustness in the convex regime. In *International Conference on Machine Learning*, pp. 4474–4491. PMLR, 2022.
- Dean, J., Corrado, G., Monga, R., Chen, K., Devin, M., Mao, M., Ranzato, M., Senior, A., Tucker, P., Yang, K., et al. Large scale distributed deep networks. *Advances in Neural Information Processing Systems*, 25, 2012.
- Diao, E., Ding, J., and Tarokh, V. HeteroFL: Computation and communication efficient federated learning for heterogeneous clients. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=TNkPBBYFkXg>.
- Dun, C., Wolfe, C. R., Jermaine, C. M., and Kyrillidis, A. ResIST: Layer-wise decomposition of resnets for distributed training. In *Uncertainty in Artificial Intelligence*, pp. 610–620. PMLR, 2022.
- Dun, C., Hipolito, M., Jermaine, C., Dimitriadis, D., and Kyrillidis, A. Efficient and light-weight federated learning via asynchronous distributed dropout. In *International Conference on Artificial Intelligence and Statistics*, pp. 6630–6660. PMLR, 2023.
- Farber, P. and Asanovic, K. Parallel neural network training on multi-spert. In *Proceedings of 3rd International Conference on Algorithms and Architectures for Parallel Processing*, pp. 659–666. IEEE, 1997.
- Gorbunov, E., Hanzely, F., and Richtárik, P. A unified theory of SGD: Variance reduction, sampling, quantization and coordinate descent. In *International Conference on Artificial Intelligence and Statistics*, pp. 680–690. PMLR, 2020.
- Goujoud, B., Scieur, D., Dieuleveut, A., Taylor, A. B., and Pedregosa, F. Super-acceleration with cyclical step-sizes. In *International Conference on Artificial Intelligence and Statistics*, pp. 3028–3065. PMLR, 2022.
- Gower, R. M., Loizou, N., Qian, X., Sailanbayev, A., Shulgin, E., and Richtárik, P. SGD: General analysis and improved rates. *Proceedings of the 36th International Conference on Machine Learning, Long Beach, California*, 2019.
- Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., Tulloch, A., Jia, Y., and He, K. Accurate, large minibatch SGD: Training ImageNet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2018.
- Horvath, S., Laskaridis, S., Almeida, M., Leontiadis, I., Venieris, S., and Lane, N. FjORD: Fair and accurate federated learning under heterogeneous targets with ordered dropout. *Advances in Neural Information Processing Systems*, 34:12876–12889, 2021.
- Jiang, Y., Wang, S., Valls, V., Ko, B. J., Lee, W.-H., Leung, K. K., and Tassiulas, L. Model pruning enables efficient federated learning on edge devices. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K. A., Charles, Z., Cormode, G., Cummings, R., D’Oliveira, R. G. L., Eichner,

-
- H., Rouayheb, S. E., Evans, D., Gardner, J., Garrett, Z., Gascón, A., Ghazi, B., Gibbons, P. B., Gruteser, M., Harchaoui, Z., He, C., He, L., Huo, Z., Hutchinson, B., Hsu, J., Jaggi, M., Javidi, T., Joshi, G., Khodak, M., Konečný, J., Korolova, A., Koushanfar, F., Koyejo, S., Lepoint, T., Liu, Y., Mittal, P., Mohri, M., Nock, R., Özgür, A., Pagh, R., Qi, H., Ramage, D., Raskar, R., Raykova, M., Song, D., Song, W., Stich, S. U., Sun, Z., Suresh, A. T., Tramèr, F., Vepakomma, P., Wang, J., Xiong, L., Xu, Z., Yang, Q., Yu, F. X., Yu, H., and Zhao, S. Advances and open problems in federated learning. *Found. Trends Mach. Learn.*, 14(1-2):1–210, 2021. doi: 10.1561/22000000083. URL <https://doi.org/10.1561/22000000083>.
- Khaled, A. and Richtárik, P. Gradient descent with compressed iterates. *arXiv preprint arXiv:1909.04716*, 2019.
- Khaled, A. and Richtárik, P. Better theory for SGD in the nonconvex world. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=AU4qHN2VkS>. Survey Certification.
- Khaled, A., Mishchenko, K., and Richtárik, P. Tighter theory for local SGD on identical and heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, pp. 4519–4529. PMLR, 2020.
- Khairat, S., Feyzmahdavian, H. R., and Johansson, M. Distributed learning with compressed gradients. *arXiv preprint arXiv:1806.06573*, 2018.
- Konečný, J., McMahan, H. B., Yu, F. X., Richtárik, P., Suresh, A. T., and Bacon, D. Federated learning: Strategies for improving communication efficiency. *NIPS Private Multi-Party Machine Learning Workshop*, 2016.
- Liao, F. and Kyrillidis, A. On the convergence of shallow neural network training with randomly masked neurons. *Transactions on Machine Learning Research*, 2022. URL <https://openreview.net/forum?id=e7mYYMSyZH>.
- Lin, R., Xiao, Y., Yang, T.-J., Zhao, D., Xiong, L., Motta, G., and Beaufays, F. Federated pruning: Improving neural network efficiency with federated learning. *arXiv preprint arXiv:2209.06359*, 2022.
- Lin, T., Stich, S. U., Barba, L., Dmitriev, D., and Jaggi, M. Dynamic model pruning with feedback. In *International Conference on Learning Representations*, 2019.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
- Mishchenko, K., Malinovsky, G., Stich, S., and Richtárik, P. ProxSkip: Yes! Local gradient steps provably lead to communication acceleration! Finally! In *International Conference on Machine Learning*, pp. 15750–15769. PMLR, 2022.
- Mohtashami, A., Jaggi, M., and Stich, S. Masked training of neural networks with partial gradients. In *International Conference on Artificial Intelligence and Statistics*, pp. 5876–5890. PMLR, 2022.
- Nesterov, Y. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- Qiu, X., Fernandez-Marques, J., Gusmao, P. P., Gao, Y., Parcollet, T., and Lane, N. D. ZeroFL: Efficient on-device training for federated learning with local sparsity. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=2sDQwC_hmnM.
- Richtárik, P. and Takáč, M. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2014.
- Richtárik, P. and Takáč, M. Distributed coordinate descent method for learning with big data. *Journal of Machine Learning Research*, 17(75):1–25, 2016.
- Safaryan, M., Hanzely, F., and Richtárik, P. Smoothness matrices beat smoothness constants: Better communication compression techniques for distributed optimization. *Advances in Neural Information Processing Systems*, 34: 25688–25702, 2021.
- Seide, F., Fu, H., Droppo, J., Li, G., and Yu, D. 1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs. In *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- Shulgin, E. and Richtárik, P. Shifted compression framework: Generalizations and improvements. In *The 38th Conference on Uncertainty in Artificial Intelligence*, 2022.
- Szlendak, R., Tyurin, A., and Richtárik, P. Permutation compressors for provably faster distributed nonconvex optimization. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=GugZ5DzzAu>.
- Wang, B., Safaryan, M., and Richtárik, P. Theoretically better and numerically faster distributed optimization with smoothness-aware quantization techniques. *Advances in Neural Information Processing Systems*, 35:9841–9852, 2022.

-
- Wang, J., Charles, Z., Xu, Z., Joshi, G., McMahan, H. B., Al-Shedivat, M., Andrew, G., Avestimehr, S., Daly, K., Data, D., et al. A field guide to federated optimization. *arXiv preprint arXiv:2107.06917*, 2021.
- Wen, D., Jeon, K.-J., and Huang, K. Federated dropout—a simple approach for enabling federated learning on resource constrained devices. *IEEE Wireless Communications Letters*, 11(5):923–927, 2022.
- Wolfe, C. R., Yang, J., Chowdhury, A., Dun, C., Bayer, A., Segarra, S., and Kyrillidis, A. GIST: Distributed training for large-scale graph convolutional networks. *arXiv preprint arXiv:2102.10424*, 2021.
- Yang, T.-J., Guliani, D., Beaufays, F., and Motta, G. Partial variable training for efficient on-device federated learning. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4348–4352. IEEE, 2022.
- Yuan, B., Wolfe, C. R., Dun, C., Tang, Y., Kyrillidis, A., and Jermaine, C. Distributed learning of fully connected neural networks using independent subnet training. *Proceedings of the VLDB Endowment*, 15(8):1581–1590, 2022.
- Zhang, G., Li, L., Nado, Z., Martens, J., Sachdeva, S., Dahl, G., Shallue, C., and Grosse, R. B. Which algorithmic choices matter at which batch sizes? insights from a noisy quadratic model. *Advances in neural information processing systems*, 32, 2019.
- Zhang, X., Mckenna, M., Mesirov, J., and Waltz, D. An efficient implementation of the back-propagation algorithm on the connection machine CM-2. *Advances in neural information processing systems*, 2, 1989.
- Zhou, H., Lan, T., Venkataramani, G., and Ding, W. On the convergence of heterogeneous federated learning with arbitrary adaptive online model pruning. *arXiv preprint arXiv:2201.11803*, 2022. URL <https://openreview.net/forum?id=p3EhUXVMeyn>.
- Zhu, L., Liu, C., Radhakrishnan, A., and Belkin, M. Quadratic models for understanding neural network dynamics. *arXiv preprint arXiv:2205.11787*, 2022.
- Zinkevich, M., Weimer, M., Li, L., and Smola, A. Parallelized stochastic gradient descent. *Advances in neural information processing systems*, 23, 2010.

Appendix

Contents

1	Introduction	1
1.1	The need for model parallelism	1
1.2	Summary of contributions	2
2	Formalism and setup	2
2.1	Simplifications taken	3
3	Results in the interpolation case	3
3.1	Homogeneous problem preconditioning	4
3.2	Heterogeneous sketch preconditioning	4
4	Irreducible bias in the general case	5
A	Basic and auxiliary facts	10
B	Proofs	10
B.1	Permutation sketch computations	10
B.1.1	Heterogeneous sketch preconditioning.	10
B.2	Interpolation case: proof of Theorem 3.1	11
B.3	Non-zero solution	12
B.3.1	Generic convergence analysis for heterogeneous case: proof of Theorem 4.1.	13
B.3.2	Homogeneous case	14
B.4	Generalization to $n \neq d$ case.	16
C	Discussion	17
C.1	Issues with existing approaches	18
C.2	Comparison to SGD-type methods	18
C.3	Improvements over previous analysis	19
D	Comparison to previous related works	19
E	Experiments	21
F	Conclusions and Future Work	22

A. Basic and auxiliary facts

L-matrix smoothness:

$$f(x+h) \leq f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} \langle \mathbf{L}h, h \rangle, \quad \forall x, h \in \mathbb{R}^d. \quad (20)$$

Basic Inequalities. For all vectors $a, b \in \mathbb{R}^d$ and random vector $X \in \mathbb{R}^d$:

$$2\langle a, b \rangle = \|a\|^2 + \|b\|^2 - \|a - b\|^2, \quad (21)$$

$$\mathbf{E} \|X - a\|^2 = \mathbf{E} \|X - \mathbf{E} X\|^2 + \|\mathbf{E} X - a\|^2. \quad (22)$$

Lemma A.1 (Fenchel–Young inequality). *For any function f and its convex conjugate f^* , Fenchel’s inequality (also known as the Fenchel–Young inequality) holds for every $x, y \in \mathbb{R}^d$*

$$\langle x, y \rangle \leq f(x) + f^*(y).$$

The proof follows from the definition of conjugate: $f^*(y) := \sup_{x'} \{\langle y, x' \rangle - f(x')\} \geq \langle y, x \rangle - f(x)$.

In the case of a quadratic function $f(x) = \beta \|x\|_{\mathbf{L}}^2$, we can compute $f^*(y) = \frac{1}{4}\beta^{-1} \|y\|_{\mathbf{L}^{-1}}^2$. Thus

$$\langle x, y \rangle \leq \beta \|x\|_{\mathbf{L}}^2 + \frac{1}{4}\beta^{-1} \|y\|_{\mathbf{L}^{-1}}^2. \quad (23)$$

B. Proofs

B.1. Permutation sketch computations

All derivations in this section are performed for the $n = d$ case.

Classical Permutation Sketching. Perm-1: $\mathbf{C}_i = n e_{\pi_i} e_{\pi_i}^\top$, where $\pi = (\pi_1, \dots, \pi_n)$ is a random permutation of $[n]$. For the homogeneous problem $\mathbf{L}_i \equiv \mathbf{L}$:

$$\mathbb{E} [\overline{\mathbf{B}}^k] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}_i \mathbf{L} \mathbf{C}_i \right] = n \text{Diag}(\mathbf{L}) \quad (24)$$

Then

$$2 \overline{\mathbf{W}} = \mathbb{E} [\mathbf{L} \overline{\mathbf{B}}^k + \overline{\mathbf{B}}^k \mathbf{L}] = n (\mathbf{L} \text{Diag}(\mathbf{L}) + \text{Diag}(\mathbf{L}) \mathbf{L}) \quad (25)$$

and

$$\mathbb{E} [\overline{\mathbf{B}}^k \mathbf{L} \overline{\mathbf{B}}^k] = n^2 \text{Diag}(\mathbf{L}) \text{L} \text{Diag}(\mathbf{L}). \quad (26)$$

By repeating basically the same calculations for $\mathbf{C}'_i = \sqrt{n} e_{\pi_i} e_{\pi_i}^\top$ we have that

$$\mathbb{E} [\overline{\mathbf{B}}^k] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}'_i \mathbf{L} \mathbf{C}'_i \right] = \text{Diag}(\mathbf{L}), \quad (27)$$

and $\mathbb{E} [\overline{\mathbf{B}}^k \mathbf{L} \overline{\mathbf{B}}^k] = \text{Diag}(\mathbf{L}) \text{L} \text{Diag}(\mathbf{L})$, $2 \overline{\mathbf{W}} = \mathbb{E} [\mathbf{L} \overline{\mathbf{B}}^k + \overline{\mathbf{B}}^k \mathbf{L}] = \text{L} \text{Diag}(\mathbf{L}) + \text{Diag}(\mathbf{L}) \mathbf{L}$.

B.1.1. HETEROGENEOUS SKETCH PRECONDITIONING.

We recall the following modification of Perm-1:

$$\tilde{\mathbf{C}}_i := \sqrt{n / [\mathbf{L}_i]_{\pi_i, \pi_i}} e_{\pi_i} e_{\pi_i}^\top. \quad (28)$$

Then

$$\mathbb{E} [\tilde{\mathbf{C}}_i \mathbf{L}_i \tilde{\mathbf{C}}_i] = \mathbb{E} [n [\mathbf{L}_i]_{\pi_i, \pi_i}^{-1} e_{\pi_i} e_{\pi_i}^\top \mathbf{L}_i e_{\pi_i} e_{\pi_i}^\top] = \frac{1}{n} \sum_{j=1}^n n e_j \mathbf{I}_{j,j} e_j^\top = \mathbf{I}. \quad (29)$$

and

$$\begin{aligned}
\mathbb{E}[\bar{\mathbf{B}}^k] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{C}}_i \mathbf{L}_i \tilde{\mathbf{C}}_i\right] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E}\left[n[\mathbf{L}_i]_{\pi_i, \pi_i}^{-1} e_{\pi_i} e_{\pi_i}^\top \mathbf{L}_i e_{\pi_i} e_{\pi_i}^\top\right] \\
&= \frac{1}{n} \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n n[\mathbf{L}_i]_{j,j}^{-1} e_j [\mathbf{L}_i]_{j,j} e_j^\top \\
&= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n e_j e_j^\top \\
&= \mathbf{I}.
\end{aligned}$$

Thus $\bar{\mathbf{W}} = \frac{1}{2} \mathbb{E}[\bar{\mathbf{L}} \bar{\mathbf{B}}^k + \bar{\mathbf{B}}^k \bar{\mathbf{L}}] = \bar{\mathbf{L}}$. On the left hand side of inequality (11), we have

$$\begin{aligned}
\mathbb{E}[\bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{C}}_i \mathbf{L}_i \tilde{\mathbf{C}}_i \bar{\mathbf{L}} \frac{1}{n} \sum_{i=j}^n \tilde{\mathbf{C}}_j \mathbf{L}_j \tilde{\mathbf{C}}_j\right] \\
&= \frac{1}{n^2} \sum_{i,j=1}^n \mathbb{E}\left[\tilde{\mathbf{C}}_i \mathbf{L}_i \tilde{\mathbf{C}}_i \bar{\mathbf{L}} \tilde{\mathbf{C}}_j \mathbf{L}_j \tilde{\mathbf{C}}_j\right] \\
&= \sum_{i,j=1}^n e_i e_i^\top \bar{\mathbf{L}} e_j e_j^\top \\
&= \mathbf{I} \bar{\mathbf{L}} \mathbf{I} \\
&= \bar{\mathbf{L}}.
\end{aligned}$$

B.2. Interpolation case: proof of Theorem 3.1

In the quadratic interpolation regime, the linear term is zero $f_i(x) = \frac{1}{2} x^\top \mathbf{L}_i x$, and the gradient estimator has the form

$$g^k = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k) = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{L}_i \mathbf{C}_i^k x^k = \bar{\mathbf{B}}^k x^k. \quad (30)$$

Proof. First, we prove the **stationary point** convergence result (12).

Using the $\bar{\mathbf{L}}$ -smoothness of function f , we get

$$\begin{aligned}
f(x^{k+1}) &\stackrel{(2)}{=} f(x^k - \gamma g^k) \stackrel{(8)}{\leq} f(x^k) - \langle \nabla f(x^k), \gamma g^k \rangle + \frac{\gamma^2}{2} \|g^k\|_{\bar{\mathbf{L}}}^2 \\
&\stackrel{(10)}{=} f(x^k) - \gamma \langle \bar{\mathbf{L}} x^k, \bar{\mathbf{B}}^k x^k \rangle + \frac{\gamma^2}{2} \|\bar{\mathbf{B}}^k x^k\|_{\bar{\mathbf{L}}}^2 \\
&= f(x^k) - \gamma (x^k)^\top \bar{\mathbf{L}} \bar{\mathbf{B}}^k x^k + \frac{\gamma^2}{2} (x^k)^\top \bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k x^k.
\end{aligned}$$

After applying conditional expectation, using its linearity, and the fact that

$$x^\top \mathbf{A} x = \frac{1}{2} x^\top (\mathbf{A} + \mathbf{A}^\top) x$$

we get

$$\begin{aligned}
\mathbb{E} [f(x^{k+1}) | x^k] &\leq f(x^k) - \gamma(x^k)^\top \mathbb{E} [\bar{\mathbf{L}} \bar{\mathbf{B}}^k] x^k + \frac{\gamma^2}{2} (x^k)^\top \mathbb{E} [\bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k] x^k \\
&= f(x^k) - \gamma(x^k)^\top \bar{\mathbf{W}} x^k + \frac{\gamma^2}{2} (x^k)^\top \mathbb{E} [\bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k] x^k \\
&= f(x^k) - \gamma(\nabla f(x^k))^\top \bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1} \nabla f(x^k) \\
&\quad + \frac{\gamma^2}{2} (\nabla f(x^k))^\top \bar{\mathbf{L}}^{-1} \mathbb{E} [\bar{\mathbf{B}}^k \bar{\mathbf{L}} \bar{\mathbf{B}}^k] \bar{\mathbf{L}}^{-1} \nabla f(x^k) \\
&\stackrel{(11)}{\leq} f(x^k) - \gamma \|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}}^2 + \frac{\theta \gamma^2}{2} \|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}}^2 \\
&= f(x^k) - \gamma(1 - \theta \gamma/2) \|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}}^2 \\
&\leq f(x^k) - \frac{\gamma}{2} \|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}}^2,
\end{aligned}$$

where the last inequality holds for the stepsize $\gamma \leq \frac{1}{\theta}$.

Rearranging gives

$$\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}}^2 \leq \frac{2}{\gamma} (f(x^k) - \mathbb{E} [f(x^{k+1}) | x^k]),$$

which after averaging gives the desired result

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1} \bar{\mathbf{W}} \bar{\mathbf{L}}^{-1}}^2 \right] \leq \frac{2}{\gamma K} \sum_{k=0}^{K-1} (f(x^k) - \mathbb{E} [f(x^{k+1})]) = \frac{2(f(x^0) - \mathbb{E} [f(x^K)])}{\gamma K}. \quad (31)$$

□

B.3. Non-zero solution

As a reminder, in the most general case, the problem has the form

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad f_i(x) \equiv \frac{1}{2} x^\top \mathbf{L}_i x - x^\top \mathbf{b}_i.$$

with the gradient estimator

$$g^k = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \nabla f_i(\mathbf{C}_i^k x^k) = \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k (\mathbf{L}_i \mathbf{C}_i^k x^k - \mathbf{b}_i) = \bar{\mathbf{B}}^k x^k - \frac{1}{n} \sum_{i=1}^n \mathbf{C}_i^k \mathbf{b}_i. \quad (32)$$

General calculations for estimator (15). In the heterogeneous case, the following sketch preconditioner is used

$$\tilde{\mathbf{C}}_i := \sqrt{n / [\mathbf{L}_i]_{\pi_i, \pi_i}} e_{\pi_i} e_{\pi_i}^\top.$$

Then $\mathbb{E} [\bar{\mathbf{B}}^k] = \mathbf{I}$ (calculation was done as in Section B.1.1) and

$$\begin{aligned}
\mathbb{E} [\bar{\mathbf{C}} \mathbf{b}] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\tilde{\mathbf{C}}_i^k \mathbf{b}_i] \\
&= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\sqrt{n} [\mathbf{L}_i]_{\pi_i, \pi_i}^{-\frac{1}{2}} e_{\pi_i} e_{\pi_i}^\top \mathbf{b}_i \right] \\
&= \frac{1}{n} \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n \sqrt{n} [\mathbf{L}_i]_{j,j}^{-\frac{1}{2}} e_j [\mathbf{b}_i]_j
\end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{i=1}^n \frac{1}{n} \sqrt{n} \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i \\
&= \frac{1}{\sqrt{n}} \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i \\
&= \frac{1}{\sqrt{n}} \underbrace{\overline{\mathbf{D}}^{-\frac{1}{2}} \bar{\mathbf{b}}}_{\mathbf{D} \bar{\mathbf{b}}}
\end{aligned}$$

B.3.1. GENERIC CONVERGENCE ANALYSIS FOR HETEROGENEOUS CASE: PROOF OF THEOREM 4.1.

Here we formulate and further prove a more general version of Theorem 4.1, which is obtained as a special case of the next result for $c = 1/2$.

Theorem B.1. *Consider the method (2) with estimator (15) for a quadratic problem (9) with positive-definite matrix $\bar{\mathbf{L}} \succ 0$. Then, if for every $\mathbf{D}_i := \text{Diag}(\mathbf{L}_i)$ matrices $\mathbf{D}_i^{-\frac{1}{2}}$ exist, scaled permutation sketches $\mathbf{C}_i := \sqrt{n}[\mathbf{L}_i^{-\frac{1}{2}}]_{\pi_i, \pi_i} e_{\pi_i} e_{\pi_i}^\top$ are used and heterogeneity is bounded as $\mathbb{E} [\|g^k - \mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2] \leq \sigma^2$. Then, the step size is chosen as*

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1 - c - \beta}{\beta + 1/2}, \quad (33)$$

where $\gamma_{c,\beta} \in (0, 1]$ for $\beta + c < 1$, the iterates satisfy

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1}}^2] \leq \frac{f(x^0) - \mathbb{E}[f(x^K)]}{c\gamma K} + \left(\frac{1 - \gamma}{c\beta} + \frac{\gamma}{2c} \right) \|h\|_{\bar{\mathbf{L}}}^2 + \frac{\gamma}{2c} \sigma^2. \quad (34)$$

where $\bar{\mathbf{L}} = \frac{1}{n} \sum_{i=1}^n \mathbf{L}_i$, $h = \bar{\mathbf{L}}^{-1} \bar{\mathbf{b}} - \frac{1}{\sqrt{n}} \frac{1}{n} \sum_{i=1}^n \mathbf{D}_i^{-\frac{1}{2}} \mathbf{b}_i$ and $\bar{\mathbf{b}} = \frac{1}{n} \sum_{i=1}^n \mathbf{b}_i$.

Proof. By using \mathbf{L} -smoothness

$$\begin{aligned}
\mathbb{E} [f(x^{k+1}) | x^k] &\stackrel{(8)}{\leq} f(x^k) - \gamma \langle \nabla f(x^k), \mathbb{E}[g^k] \rangle + \frac{\gamma^2}{2} \mathbb{E} [\|g^k\|_{\bar{\mathbf{L}}}^2] \\
&\stackrel{(16),(22)}{=} f(x^k) - \gamma \langle \nabla f(x^k), \bar{\mathbf{L}}^{-1} \nabla f(x^k) + h \rangle \\
&\quad + \frac{\gamma^2}{2} \left(\|\mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2 + \mathbb{E} [\|g^k - \mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2] \right) \\
&\stackrel{(16)}{=} f(x^k) - \gamma \left(\langle \nabla f(x^k), \bar{\mathbf{L}}^{-1} \nabla f(x^k) \rangle + \langle \nabla f(x^k), h \rangle \right) \\
&\quad + \frac{\gamma^2}{2} \left(\|\bar{\mathbf{L}}^{-1} \nabla f(x^k) + h\|_{\bar{\mathbf{L}}}^2 + \mathbb{E} [\|g^k - \mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2] \right) \\
&\stackrel{(21)}{=} f(x^k) - \gamma \left(\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1}}^2 + \langle \nabla f(x^k), h \rangle \right) + \frac{\gamma^2}{2} \mathbb{E} [\|g^k - \mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2] \\
&\quad + \frac{\gamma^2}{2} \left(\|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1}}^2 + 2 \langle \nabla f(x^k), h \rangle + \|h\|_{\bar{\mathbf{L}}}^2 \right) \\
&\leq f(x^k) - \gamma (1 - \gamma/2) \|\nabla f(x^k)\|_{\bar{\mathbf{L}}^{-1}}^2 + \frac{\gamma^2}{2} \sigma^2 \\
&\quad - \gamma (1 - \gamma) \langle \nabla f(x^k), h \rangle + \frac{\gamma^2}{2} \|h\|_{\bar{\mathbf{L}}}^2,
\end{aligned}$$

where the last inequality follows from the grouping of similar terms and bounded heterogeneity

$$\mathbb{E} [\|g^k - \mathbb{E}[g^k]\|_{\bar{\mathbf{L}}}^2] = \mathbb{E} \left[\left\| g^k - \left(\bar{\mathbf{L}}^{-1} \nabla f(x^k) + h \right) \right\|_{\bar{\mathbf{L}}}^2 \right] \quad (35)$$

$$= \mathbb{E} \left[\left\| \bar{\mathbf{B}}^k x^k - \bar{\mathbf{C}} \bar{\mathbf{b}} - \left(x^k - \frac{1}{\sqrt{n}} \widetilde{\mathbf{D}} \bar{\mathbf{b}} \right) \right\|_{\bar{\mathbf{L}}}^2 \right] \leq \sigma^2. \quad (36)$$

Next, using a Fenchel-Young inequality (23) for $\langle \nabla f(x^k), -h \rangle$ and $1 - \gamma \geq 0$

$$\begin{aligned} \mathbb{E} [f(x^{k+1}) | x^k] &\leq f(x^k) - \gamma(1 - \gamma/2) \|\nabla f(x^k)\|_{\mathbf{L}^{-1}}^2 + \frac{\gamma^2}{2} (\|h\|_{\mathbf{L}}^2 + \sigma^2) \\ &\quad + \gamma(1 - \gamma) \left[\beta \|\nabla f(x^k)\|_{\mathbf{L}^{-1}}^2 + 0.25\beta^{-1} \|h\|_{\mathbf{L}}^2 \right] \\ &\leq f(x^k) - \gamma(1 - \gamma/2 - \beta(1 - \gamma)) \|\nabla f(x^k)\|_{\mathbf{L}^{-1}}^2 \\ &\quad + \gamma \left\{ \left(\beta^{-1}(1 - \gamma) + \frac{\gamma}{2} \right) \|h\|_{\mathbf{L}}^2 + \frac{\gamma}{2} \sigma^2 \right\}, \end{aligned} \quad (37)$$

where in the last inequality we grouped similar terms and used the fact that $0.25 < 1$.

Now to guarantee that $1 - \gamma/2 - \beta(1 - \gamma) \geq c > 0$, we choose the step size using

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1 - c - \beta}{\beta + 1/2}, \quad (38)$$

where $\gamma_{c,\beta} > 0$ for $\beta + c < 1$. This means that β can not arbitrarily grow to diminish β^{-1} .

Then, after standard manipulations and unrolling the recursion

$$\gamma c \|\nabla f(x^k)\|_{\mathbf{L}^{-1}}^2 \leq f(x^k) - \mathbb{E} [f(x^{k+1}) | x^k] + \gamma(\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\mathbf{L}}^2 + \frac{\gamma^2}{2} \sigma^2 \quad (39)$$

we obtain

$$\frac{c}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{L}^{-1}}^2 \right] \leq \frac{f(x^0) - \mathbb{E} [f(x^K)]}{\gamma K} + (\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\mathbf{L}}^2 + \frac{\gamma}{2} \sigma^2. \quad (40)$$

□

B.3.2. HOMOGENEOUS CASE

The main difference compared to the result in the previous subsection is that the gradient estimator expression (19) holds deterministically (without expectation \mathbb{E}). That is why $g^k = \mathbb{E} [g^k]$ and heterogeneity term σ^2 equals to 0.

We provide the full statement and proof for the homogeneous result discussed in Section 4.

Theorem B.2. *Consider the method (2) with estimator (19) for a homogeneous quadratic problem (9) with positive-definite matrix $\mathbf{L}_i \equiv \mathbf{L} \succ 0$. Then if exists $\mathbf{D}^{-\frac{1}{2}}$ for $\mathbf{D} := \text{Diag}(\mathbf{L})$, scaled permutation sketch $\mathbf{C}'_i = \sqrt{n} e_{\pi_i} e_{\pi_i}^\top$ is used and the step size is chosen as*

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1 - c - \beta}{\beta + 1/2}, \quad (41)$$

where $\gamma_{c,\beta} > 0$ for $\beta + c < 1$. Then the iterates satisfy

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\mathbf{L}^{-1}}^2 \right] \leq \frac{f(x^0) - \mathbb{E} [f(x^K)]}{c\gamma K} + \left(\frac{1 - \gamma}{c\beta} + \frac{\gamma}{2c} \right) \|h\|_{\mathbf{L}}^2, \quad (42)$$

where $\tilde{\mathbf{L}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$, $h = \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}$ and $\tilde{\mathbf{b}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{b}$.

Proof. By using \mathbf{L} -smoothness

$$\begin{aligned} \mathbb{E} [f(x^k - \gamma g^k) | x^k] &\stackrel{(8)}{\leq} f(x^k) - \langle \nabla f(x^k), \gamma \mathbb{E} [g^k] \rangle + \frac{\gamma^2}{2} \mathbb{E} \left[\|g^k\|_{\mathbf{L}}^2 \right] \\ &\leq f(x^k) - \gamma \left\langle \nabla f(x^k), \tilde{\mathbf{L}}^{-1} \nabla f(x^k) + h \right\rangle + \frac{\gamma^2}{2} \left\| \tilde{\mathbf{L}}^{-1} \nabla f(x^k) + h \right\|_{\tilde{\mathbf{L}}}^2 \\ &\stackrel{(21)}{=} f(x^k) - \gamma \left(\left\langle \nabla f(x^k), \tilde{\mathbf{L}}^{-1} \nabla f(x^k) \right\rangle + \langle \nabla f(x^k), h \rangle \right) \\ &\quad + \frac{\gamma^2}{2} \left(\|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 + 2 \langle \nabla f(x^k), h \rangle + \|h\|_{\tilde{\mathbf{L}}}^2 \right) \\ &= f(x^k) - \gamma(1 - \gamma/2) \|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 + \frac{\gamma^2}{2} \|h\|_{\tilde{\mathbf{L}}}^2 - \gamma(1 - \gamma) \langle \nabla f(x^k), h \rangle \end{aligned}$$

Next by using a Fenchel-Young inequality (23) for $\langle \nabla f(x^k), -h \rangle$ and $1 - \gamma \geq 0$

$$\begin{aligned} \mathbb{E} [f(x^{k+1}) | x^k] &\leq f(x^k) - \gamma(1 - \gamma/2) \|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 + \frac{\gamma^2}{2} \|h\|_{\tilde{\mathbf{L}}}^2 \\ &\quad + \gamma(1 - \gamma) \left[\beta \|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 + 0.25\beta^{-1} \|h\|_{\tilde{\mathbf{L}}}^2 \right] \\ &= f(x^k) - \gamma(1 - \gamma/2 - \beta(1 - \gamma)) \|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 \\ &\quad + \gamma(\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2. \end{aligned}$$

Now to guarantee that $1 - \gamma/2 - \beta(1 - \gamma) \geq c > 0$ we choose the step size as

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1 - c - \beta}{\beta + 1/2}, \quad (43)$$

where $\gamma_{c,\beta} \geq 0$ for $\beta + c < 1$.

Then after standard manipulations and unrolling the recursion

$$\gamma c \|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 \leq f(x^k) - \mathbb{E} [f(x^{k+1}) | x^k] + \gamma(\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2 \quad (44)$$

we obtain the formulated result

$$\frac{c}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 \right] \leq \frac{f(x^0) - \mathbb{E} [f(x^K)]}{\gamma K} + (\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2. \quad (45)$$

□

Remark B.3. 1) The first term in the convergence upper bound (42) is minimized by maximizing product $c \cdot \gamma$, which motivates to choose $c > 0$ and $\gamma \leq 1$ as large as possible. Although due to the constraint on the step size (and $\beta > 0$)

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1 - c - \beta}{\beta + 1/2}, \quad (46)$$

constant $c \in (0, 1)$. So, by maximizing c the value $\gamma_{c,\beta}$ becomes smaller, thus there is a trade-off.

2) The second term or the neighborhood size (multiplier in front of $\|h\|_{\tilde{\mathbf{L}}}^2$)

$$\Psi(\beta, \gamma) := \frac{\beta^{-1}(1 - \gamma) + \gamma/2}{c} = \frac{\beta^{-1}(1 - \gamma) + \gamma/2}{1 - \gamma/2 - \beta(1 - \gamma)} \quad (47)$$

can be numerically minimized (e.g. by using WolframAlpha) with constraints $\gamma \in (0, 1]$ and $\beta > 0$. The solution of such optimization problem is $\gamma^* \approx 1$ and $\beta^* \approx \xi \in \{3.992, 2.606, 2.613\}$. In fact, $\Psi(\beta^*, \gamma^*) \approx 0.5$.

Functional gap convergence. Note that for the quadratic optimization problem (9)

$$\|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 = \left\langle \tilde{\mathbf{L}} x^k - \tilde{\mathbf{b}}, \tilde{\mathbf{L}}^{-1} (\tilde{\mathbf{L}} x^k - \tilde{\mathbf{b}}) \right\rangle = 2(f(x^k) - f(x^*)). \quad (48)$$

Then by rearranging and subtracting $f^* := f(x^*)$ from both sides of inequality (44) we obtain

$$\begin{aligned} \mathbb{E} [f(x^{k+1}) | x^k] - f^* &\leq f(x^k) - f^* - \gamma c \|\nabla f(x^k)\|_{\tilde{\mathbf{L}}^{-1}}^2 + \gamma(\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2 \\ &\stackrel{(48)}{=} (f(x^k) - f^*) - \gamma c \cdot 2(f(x^k) - f^*) + \gamma(\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2 \\ &= (1 - 2\gamma c)(f(x^k) - f^*) + \gamma(\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2. \end{aligned}$$

After unrolling the recursion

$$\begin{aligned} \mathbb{E} [f(x^{k+1}) | x^k] - f^* &\leq (1 - 2\gamma c)^k (f(x^0) - f^*) + \gamma(\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2 \sum_{i=0}^k (1 - 2\gamma c)^i \\ &\leq (1 - 2\gamma c)^k (f(x^0) - f^*) + \frac{1}{2c} (\beta^{-1}(1 - \gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2. \end{aligned}$$

This result is formalized in the following Theorem.

Theorem B.4. Consider the method (2) with estimator (19) for a homogeneous quadratic problem (9) with positive-definite matrix $\mathbf{L}_i \equiv \mathbf{L} \succ 0$. Then if exists $\mathbf{D}^{-\frac{1}{2}}$ for $\mathbf{D} := \text{Diag}(\mathbf{L})$, scaled permutation sketch $\mathbf{C}'_i = \sqrt{n}e_{\pi_i}e_{\pi_i}^\top$ is used and the step size is chosen as

$$0 < \gamma \leq \gamma_{c,\beta} := \frac{1-c-\beta}{\beta+1/2}, \quad (49)$$

where $\gamma_{c,\beta} > 0$ for $\beta + c < 1$. Then the iterates satisfy

$$\mathbb{E} [f(x^k)] - f^* \leq (1-2\gamma c)^k (f(x^0) - f^*) + \frac{1}{2c} (\beta^{-1}(1-\gamma) + \gamma/2) \|h\|_{\tilde{\mathbf{L}}}^2, \quad (50)$$

where $h = \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}$ and $\tilde{\mathbf{L}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$, $\tilde{\mathbf{b}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{b}$.

This result shows that for a proper choice of the step size $\gamma = 1$ and constant $c = 1/2$, the functional gap can converge in basically one iteration to the neighborhood of size

$$\|h\|_{\tilde{\mathbf{L}}}^2 = \left\langle \tilde{\mathbf{L}} \left(\tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}} \right), \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} - \frac{1}{\sqrt{n}} \tilde{\mathbf{b}} \right\rangle,$$

which equals zero if $\tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}} = \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}$. This condition is the same as the condition we obtained in Section 4 with asymptotic analysis of the iterates in the homogeneous case.

Discussion of the trace. Consider a positive-definite $\mathbf{L} \succ 0$ such that $\exists \mathbf{D}^{-\frac{1}{2}}$. Thus $\tilde{\mathbf{L}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{L} \mathbf{D}^{-\frac{1}{2}}$ has only ones on the diagonal and $\text{tr}(\tilde{\mathbf{L}}) = n$. Then

$$n \cdot \text{tr}(\tilde{\mathbf{L}}^{-1}) = \text{tr}(\tilde{\mathbf{L}}) \text{tr}(\tilde{\mathbf{L}}^{-1}) = (\lambda_1 + \dots + \lambda_n) \left(\frac{1}{\lambda_1} + \dots + \frac{1}{\lambda_n} \right) \geq n^2,$$

where the last inequality is due to the relation between harmonic and arithmetic means. Therefore $\text{tr}(\tilde{\mathbf{L}}^{-1}) = \lambda_1^{-1} + \dots + \lambda_n^{-1} \geq n$ and sum of $\tilde{\mathbf{L}}^{-1}$ eigenvalues has to be greater than n .

B.4. Generalization to $n \neq d$ case.

Our results can be generalized in a similar way as in (Szlendak et al., 2022).

1) $d = qn$, for integer $q \geq 1$. Let $\pi = (\pi_1, \dots, \pi_d)$ be a random permutation of $\{1, \dots, d\}$. Then for each $i \in \{1, \dots, n\}$ define

$$\mathbf{C}'_i := \sqrt{n} \cdot \sum_{j=q(i-1)+1}^{qi} e_{\pi_j} e_{\pi_j}^\top. \quad (51)$$

Matrix $\mathbb{E} [\overline{\mathbf{B}}^k]$ for the homogeneous preconditioned case can be computed as follows:

$$\begin{aligned} \mathbb{E} [\overline{\mathbf{B}}^k] &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}'_i \tilde{\mathbf{L}} \mathbf{C}'_i \right] \\ &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\sum_{j=q(i-1)+1}^{qi} n e_{\pi_j} e_{\pi_j}^\top \tilde{\mathbf{L}} e_{\pi_j} e_{\pi_j}^\top \right] \\ &= \sum_{i=1}^n \sum_{j=q(i-1)+1}^{qi} \mathbb{E} \left[e_{\pi_j} e_{\pi_j}^\top \tilde{\mathbf{L}} e_{\pi_j} e_{\pi_j}^\top \right] \\ &= \sum_{i=1}^n \sum_{j=q(i-1)+1}^{qi} \frac{1}{d} \sum_{l=1}^d e_l e_l^\top \tilde{\mathbf{L}} e_l e_l^\top \\ &= \sum_{i=1}^n \sum_{j=q(i-1)+1}^{qi} \frac{1}{d} \text{Diag}(\tilde{\mathbf{L}}) \end{aligned}$$

$$\begin{aligned}
&= n \frac{q}{d} \text{Diag}(\tilde{\mathbf{L}}) \\
&= \text{Diag}(\tilde{\mathbf{L}}) \\
&= \mathbf{I}.
\end{aligned}$$

As for the linear term

$$\begin{aligned}
\mathbb{E}[\mathbf{C}' \mathbf{b}] &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}'_i \tilde{\mathbf{b}} \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\sum_{j=q(i-1)+1}^{qi} \sqrt{n} e_{\pi_j} e_{\pi_j}^\top \tilde{\mathbf{b}} \right] \\
&= \frac{1}{\sqrt{n}} \sum_{i=1}^n \sum_{j=q(i-1)+1}^{qi} \frac{1}{d} \mathbf{I} \tilde{\mathbf{b}} = \frac{\sqrt{n}q}{d} \mathbf{I} \tilde{\mathbf{b}} = \frac{1}{\sqrt{n}} \tilde{\mathbf{b}}.
\end{aligned}$$

2) $n = qd$, for integer $q \geq 1$. Define the multiset $S := \{1, \dots, 1, 2, \dots, 2, \dots, d, \dots, d\}$, where each number occurs precisely q times. Let $\pi = (\pi_1, \dots, \pi_n)$ be a random permutation of S . Then for each $i \in \{1, \dots, n\}$ define

$$\mathbf{C}'_i := \sqrt{d} \cdot e_{\pi_i} e_{\pi_i}^\top. \quad (52)$$

$$\begin{aligned}
\mathbb{E}[\overline{\mathbf{B}}^k] &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}'_i \tilde{\mathbf{L}} \mathbf{C}'_i \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[d e_{\pi_i} e_{\pi_i}^\top \tilde{\mathbf{L}} e_{\pi_i} e_{\pi_i}^\top \right] \\
&= \frac{1}{n} \sum_{i=1}^n \frac{1}{d} \sum_{j=1}^d d e_j e_j^\top \tilde{\mathbf{L}} e_j e_j^\top = \frac{1}{n} \sum_{i=1}^n \text{Diag}(\tilde{\mathbf{L}}) = \mathbf{I}.
\end{aligned}$$

The linear term

$$\mathbb{E}[\mathbf{C}' \mathbf{b}] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}'_i \tilde{\mathbf{b}} \right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[\sqrt{d} e_{\pi_i} e_{\pi_i}^\top \tilde{\mathbf{b}} \right] = \frac{\sqrt{d}}{n} \sum_{i=1}^n \frac{1}{d} \mathbf{I} \tilde{\mathbf{b}} = \frac{1}{\sqrt{d}} \tilde{\mathbf{b}}.$$

To sum up both cases, in a homogeneous preconditioned setting $\mathbb{E}[\overline{\mathbf{B}}^k] = \mathbf{I}$ and

$$\mathbb{E}[\mathbf{C}' \mathbf{b}] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbf{C}'_i \mathbf{b} \right] = \tilde{\mathbf{b}} / \sqrt{\min(n, d)}.$$

Similar modifications and calculations can be performed for heterogeneous scenarios. The case when n does not divide d and vice versa is generalized using constructions from (Szlendak et al., 2022).

C. Discussion

A generalized analog of IST is formalized as an iterative procedure in Algorithm 1.

Remark C.1. Matrix $\overline{\mathbf{W}}$ in case of permutation sketches may not be positive-definite. Consider the following homogeneous ($\mathbf{L}_i \equiv \mathbf{L}$) two-dimensional problem example

$$\mathbf{L} = \begin{bmatrix} a & c \\ c & b \end{bmatrix}. \quad (53)$$

Then

$$\overline{\mathbf{W}} = \frac{1}{2} [\overline{\mathbf{L}} \overline{\mathbf{D}} + \overline{\mathbf{D}} \overline{\mathbf{L}}] = \begin{bmatrix} a^2 & c(a+b)/2 \\ c(a+b)/2 & b^2 \end{bmatrix}, \quad (54)$$

which for $c > \frac{2ab}{a+b}$ has $\det(\overline{\mathbf{W}}) < 0$, and thus $\overline{\mathbf{W}} \not\prec 0$ according to Sylvester's criterion.

Algorithm 1 Distributed Submodel (Stochastic) Gradient Descent

```
1: Parameters: learning rate  $\gamma > 0$ ; sketches  $\mathbf{C}_1, \dots, \mathbf{C}_n$ ; initial model  $x^0 \in \mathbb{R}^d$ 
2: for  $k = 0, 1, 2 \dots$  do
3:   Select submodels  $w_i^k = \mathbf{C}_i^k x^k$  for  $i \in [n]$  and broadcast to all computing nodes
4:   for  $i = 1, \dots, n$  in parallel do
5:     Compute local (stochastic) gradient w.r.t. submodel:  $\mathbf{C}_i^k \nabla f_i(w_i^k)$ 
6:     Take (maybe multiple) gradient descent step  $z_i^+ = w_i^k - \gamma \mathbf{C}_i^k \nabla f_i(w_i^k)$ 
7:     Send  $z_i^+$  to the server
8:   end for
9:   Aggregate/merge received submodels:  $x^{k+1} = \frac{1}{n} \sum_{i=1}^n z_i^+$ 
10: end for
```

C.1. Issues with existing approaches

Consider the simplest gradient type method with compressed model in the single node setting

$$x^{k+1} = x^k - \gamma \nabla f(\mathcal{C}(x^k)). \quad (55)$$

Algorithms belonging to this family require a different analysis in comparison to SGD (Gorbunov et al., 2020; Gower et al., 2019), Distributed Compressed Gradient Descent (Alistarh et al., 2017; Khirirat et al., 2018) and Randomized Coordinate Descent (Nesterov, 2012; Richtárik & Takáč, 2014) type methods because the gradient estimator is no longer unbiased

$$\mathbb{E}[\nabla f(\mathcal{C}(x))] \neq \nabla f(x) = \mathbb{E}[\mathcal{C}(\nabla f(x))]. \quad (56)$$

That is why such kind of algorithms are harder to analyze. So, prior results for *unbiased* SGD (Khaled & Richtárik, 2023) can not be directly reused. Furthermore, the nature of the bias in this type of gradient estimator does not exhibit additive (zero-mean) noise, thereby preventing the application of previous analyses for biased SGD (Ajalloeian & Stich, 2020).

An assumption like bounded stochastic gradient norm extensively used in previous works (Lin et al., 2019; Zhou et al., 2022) hinders an accurate understanding of such methods. This assumption hides the fundamental difficulty of analyzing biased gradient estimator:

$$\mathbb{E}[\|\nabla f(\mathcal{C}(x))\|^2] \leq G \quad (57)$$

and may not hold even for quadratic functions $f(x) = x^\top \mathbf{A}x$. In addition, in the distributed setting such condition can result in vacuous bounds (Khaled et al., 2020) as it does not allow to accurately capture heterogeneity.

C.2. Comparison to SGD-type methods

Let us contrast obtained result (18) with non-convex rate of SGD (Khaled & Richtárik, 2023) with constant step size γ for L -smooth and lower-bounded f

$$\min_{k \in \{0, \dots, K-1\}} \|\nabla f(x^k)\|^2 \leq \frac{6(f(x^0) - \inf f)}{\gamma K} + \gamma LC, \quad (58)$$

where constant C depends, for example, on the variance of stochastic gradient estimates. Observe that the first term in the compared upper bounds (58) and (18) is almost identical and decreases with speed $1/K$. But unlike (18) the neighborhood for SGD can be completely eliminated by reducing the step size γ . This highlights a fundamental difference of our results to unbiased methods.

The intuition behind this issue is that for SGD-type methods like Compressed Gradient Descent

$$x^{k+1} = x^k - \mathcal{C}(\nabla f(x^k)) \quad (59)$$

the gradient estimate is unbiased and enjoys the property that variance

$$\mathbb{E}[\|\mathcal{C}(\nabla f(x^k)) - \nabla f(x^k)\|^2] \leq \omega \|\nabla f(x^k)\|^2 \quad (60)$$

goes down to zero as the method progresses because $\nabla f(x^k) \rightarrow \nabla f(x^*) = 0$ in the unconstrained case. In addition, any stationary point x^* ceases to be a fixed point of the iterative procedure as

$$x^* \neq x^* - \nabla f(\mathcal{C}(x^*)), \quad (61)$$

in the general case, unlike for Compressed Gradient Descent with both biased and unbiased compressors \mathcal{C} . So, even if the method (computing gradient at sparse model) is initialized from the *solution* after one gradient step, it may get away from there.

C.3. Improvements over previous analysis

Independent Subnetwork Training (Yuan et al., 2022). There are several improvements over the previous works that tried to theoretically analyze the convergence of Distributed IST.

The first difference is that our results allow for an almost arbitrary level of model sparsification, i.e., work for any $\omega \geq 0$ as permutation sketches can be viewed as a special case of compression operators (2.1). This improves significantly over the work of (Yuan et al., 2022), which demands[‡] $\omega \lesssim \mu^2/L^2$. Such a requirement is very restrictive as the condition number L/μ of the loss function f is typically very large for any non-trivial optimization problem. Thus, the sparsifier’s (4) variance $\omega = d/q - 1$ has to be very close to 0 and $q \approx d$. So, the previous theory allows almost no compression (sparsification) because it is based on the analysis of Gradient Descent with Compressed Iterates (Khaled & Richtárik, 2019).

The second distinction is that the original IST work (Yuan et al., 2022) considered a single node setting and thus their convergence bounds did not capture the effect of heterogeneity, which we believe is of crucial importance for distributed setting (Chraïbi et al., 2019; Shulgin & Richtárik, 2022). Besides, they consider Lipschitz continuity of the loss function f , which is not satisfied for a simple quadratic model. A more detailed comparison including additional assumptions on the gradient estimator made in (Yuan et al., 2022) is presented in the Appendix.

FL with Model Pruning. In a recent work (Zhou et al., 2022) made an attempt to analyze a variant of the FedAvg algorithm with sparse local initialization and compressed gradient training (pruned local models). They considered a case of L -smooth loss and sparsification operator satisfying a similar condition to (2.1). However, they also assumed that the squared norm of stochastic gradient is uniformly bounded (57), which is “pathological” (Khaled et al., 2020) especially in the case of local methods as it does not allow to capture the very important effect of heterogeneity and can result in vacuous bounds.

In the next section we show some limitations of other relevant previous approaches to training with compressed models: too restrictive assumptions on the algorithm (Mohtashami et al., 2022) or not applicability in our problem setting (Chayti & Karimireddy, 2022).

D. Comparison to previous related works

Overview of theory provided in the original IST work (Yuan et al., 2022). The authors consider the following method

$$x^{k+1} = \mathcal{C}(x^k) - \gamma \nabla f_{i_k}(\mathcal{C}(x^k)), \quad (62)$$

where $[\mathcal{C}(x)]_i = x_i \cdot \text{Be}(p)$ [§] is a Bernoulli sparsifier and i_k is sampled uniformly at random from $[n]$.

The analysis in (Yuan et al., 2022) relies on the assumptions

1. L_i -smoothness of individual losses f_i ;
2. Q -Lipschitz continuity of f : $|f(x) - f(y)| \leq Q\|x - y\|$;
3. Error bound (or PL-condition): $\|\nabla f(x)\| \geq \mu\|x^* - x\|$, where x^* is the global optimum;

[‡] μ refers to constant from Polyak-Łojasiewicz (or strong convexity) condition. In case of a quadratic problem with positive-definite matrix \mathbf{A} : $\mu = \lambda_{\min}(\mathbf{A})$

[§] $\mathcal{B}_p(x) := \begin{cases} x/p & \text{with probability } p \\ 0 & \text{with probability } 1 - p \end{cases}$

-
4. Stochastic gradient variance: $\mathbb{E} \left[\|\nabla f_{i_k}(x)\|^2 \right] \leq M + M_f \|\nabla f(x)\|^2$;
 5. $\mathbb{E} [\nabla f_{i_k}(\mathcal{C}(x^k)) | x^k] = \nabla f(x^k) + \varepsilon, \quad \|\varepsilon\| \leq B$.

Convergence result from Theorem 1 (Yuan et al., 2022) for step size $\gamma = 1/(2L_{\max})$:

$$\min_{k \in \{1, \dots, K\}} \mathbb{E} \left[\|\nabla f(x^k)\|^2 \right] \leq \frac{f(x^0) - f(x^*)}{\alpha(K+1)} + \frac{1}{\alpha} \cdot \left(\frac{BQ}{2L_{\max}} + \frac{5L_{\max}\omega}{2} \|x^*\|^2 + \frac{M}{4L_{\max}} \right), \quad (63)$$

where $\alpha := \frac{1}{2L_{\max}} \left(1 - \frac{M_f}{2} \right) - \frac{5\omega L_{\max}}{2\mu^2}$, $\omega := \frac{1}{p} - 1 < \frac{\mu^2}{10L_{\max}^2}$, and $L_{\max} := \max_i L_i$.

If Lipschitzness and Assumption 5 are replaced with *norm condition*:

$$\|\mathbb{E} [\nabla f_{i_k}(\mathcal{C}(x^k)) | x^k] - \nabla f(x^k)\| \leq \theta \|\nabla f(x^k)\| \quad (64)$$

they obtain the following (for step size $\gamma = 1/2L_{\max}$)

$$\min_{k \in \{1, \dots, K\}} \mathbb{E} \left[\|\nabla f(x^k)\|^2 \right] \leq \frac{f(x^0) - f(x^*)}{\alpha(K+1)} + \frac{1}{\alpha} \cdot \left(\frac{5L_{\max}\omega}{2} \|x^*\|^2 + \frac{M}{4L_{\max}} \right), \quad (65)$$

where $\alpha = \frac{1}{2L_{\max}} \left(\frac{1}{2} - \theta - \frac{M_f}{2} \right) - \frac{5\omega L_{\max}}{2\mu^2}$ and $\omega = \frac{1}{p} - 1 < \frac{\mu^2}{5L_{\max}^2 \left(\frac{1}{2} - \theta - \frac{M_f}{2} \right)}$.

Remark D.1. The original method (62) does not incorporate gradient sparsification, which can create a significant disparity between theory and practice. This is because the gradient computed at the compressed model, denoted as $\nabla f(\mathcal{C}(x))$, is not guaranteed to be sparse and representative of the submodel computations. Such modification of the method also significantly simplifies theoretical analysis, as using a single sketch (instead of CLC) allows for an unbiased gradient estimator.

Through our analysis of the IST gradient estimator in Equation (19), we discover that conditions—such as Assumption 5 and Inequality (64)—are not satisfied, even in the homogeneous setting for a simple quadratic problem. Furthermore, it is evident that such conditions are also not met for logistic loss. At the same time, in general, it is expected that insightful theory for general (non-)convex functions should yield appropriate results for quadratic problems. Additionally, it remains unclear whether the norm condition (64) is satisfied in practical scenarios. The situation is not straightforward—even for quadratic problems—as we show in the expression for σ^2 in Equation (35).

Masked training (Mohtashami et al., 2022). The authors consider the following “Partial SGD” method

$$\begin{aligned} \hat{x}^k &= x^k + \delta x^k = x^k - (1-p) \odot x^k \\ x^{k+1} &= x^k - \gamma p \odot \nabla f(\hat{x}^k, \xi^k), \end{aligned} \quad (66)$$

where $\nabla f(x, \xi)$ is an unbiased stochastic gradient estimator of a L -smooth loss function f , \odot is an element-wise product, and p is a binary sparsification mask.

Mohtashami et al. (Mohtashami et al., 2022) make the following “bounded perturbation” assumption

$$\max_k \frac{\|\delta x^k\|}{\max \{ \|p^k \odot \nabla f(x^k)\|, \|p^k \odot \nabla f(\hat{x}^k)\| \}} \leq \frac{1}{2L}. \quad (67)$$

This inequality may not hold for a simple convex case. Consider a function $f(x) = \frac{1}{2} x^\top A x$, for

$$A = \begin{pmatrix} a & 0 \\ 0 & c \end{pmatrix}, \quad x^0 = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad p^0 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}. \quad (68)$$

Then condition (67) (at iteration $k = 0$) will be equivalent to

$$\frac{x_1}{cx_2} \leq \frac{1}{2a} \Leftrightarrow 2 \leq \frac{2a}{c} \leq \frac{x_2}{x_1},$$

which clearly does not hold for an arbitrary initialization x^0 .

In addition, convergence bound in Theorem 1 of (Mohtashami et al., 2022) suggests choosing the step size as $\gamma_0 \alpha^k$, where

$$\alpha^k = \min \left\{ 1, \frac{\langle p^k \odot \nabla f(x^k), p^k \odot \nabla f(\hat{x}^k) \rangle}{\|p^k \odot \nabla f(\hat{x}^k)\|^2} \right\} \quad (69)$$

is not guaranteed to be positive to the inner product $\langle p^k \odot \nabla f(x^k), p^k \odot \nabla f(\hat{x}^k) \rangle$, which may lead to non-convergence of the method.

Optimization with access to auxiliary information framework (Chayti & Karimireddy, 2022) suggests modeling training with compressed models via performing gradient steps with respect to function $h(x) := \mathbb{E}_{\mathcal{M}} [f(1_{\mathcal{M}} \odot x)]$. This function allows access to a sparse/low-rank version of the original model $f(x)$. They impose the following bounded Hessian dissimilarity assumption on h and f

$$\|\nabla^2 f(x) - \mathbb{E}_{\mathcal{M}} [\mathbf{D}_{\mathcal{M}} \nabla^2 f(1_{\mathcal{M}} \odot x) \mathbf{D}_{\mathcal{M}}]\|_2 \leq \delta, \quad (70)$$

where $1_{\mathcal{M}}$ and $\mathbf{D}_{\mathcal{M}} = \text{Diag}(1_{\mathcal{M}})$ refer to a binary vector and matrix sparsification masks.

This approach relies on variance-reduction and requires gradient computations on the full model x , and thus it is not suitable for our problem setting.

E. Experiments

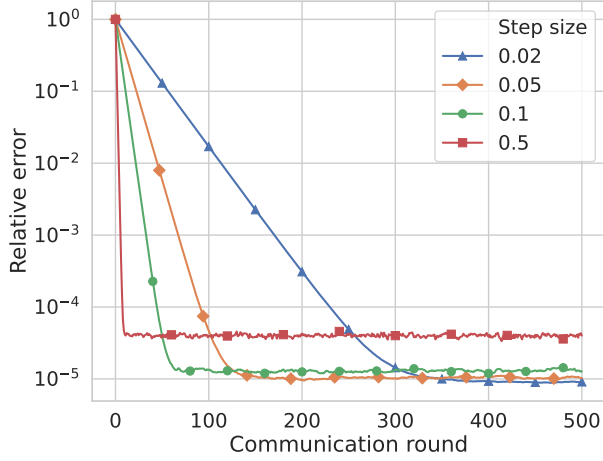
To empirically validate our theoretical framework and its implications, we focus on carefully controlled settings that satisfy the assumptions of our work. Specifically, we consider a quadratic problem defined in (9). As a reminder, the local loss function is defined as

$$f_i(x) = \frac{1}{2} x^\top \mathbf{L}_i x - x^\top \mathbf{b}_i,$$

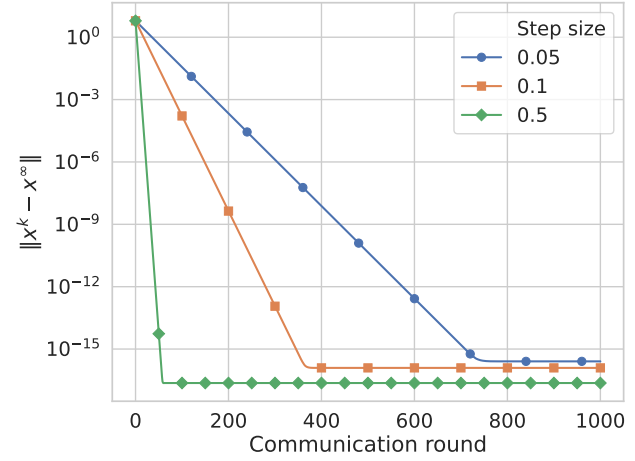
where $\mathbf{L}_i = \mathbf{B}_i^\top \mathbf{B}_i$. Entries of the matrices $\mathbf{B}_i \in \mathbb{R}^{d \times d}$, vectors $\mathbf{b}_i \in \mathbb{R}^d$, and initialization $x^0 \in \mathbb{R}^d$ are generated from a standard Gaussian distribution $\mathcal{N}(0, 1)$.

Heterogeneous setting. In Figure 1(a), we present the performance of the simplified Independent Subnetwork Training (IST) algorithm (update (2) with estimator (15)) for a heterogeneous problem. We fix the dimension d to 1000 and the number of computing nodes n to 10. We evaluate the logarithm of a relative functional error $\log(f(x^k) - f(x^*)) / (f(x^0) - f(x^*))$, while the horizontal axis denotes the number of communication rounds required to achieve a certain error tolerance. According to our theory (50), the method converges to a neighborhood of the solution, which depends on the chosen step size. Specifically, a larger step size allows for faster convergence but results in a larger neighborhood.

Homogeneous setting. In Figure 1(b), we demonstrate the convergence of the iterates x^k for a homogeneous problem with $d = n = 50$. The results are in close agreement with our theoretical predictions for the estimator (19). We observe that the distance to the method's expected fixed point $x^\infty = \tilde{\mathbf{b}} / \sqrt{n}$ decreases linearly for different step size values. This confirms that IST may not converge to the optimal solution $x^* = \tilde{\mathbf{L}}^{-1} \tilde{\mathbf{b}}$ of the original problem (9) in general (no interpolation) cases. In addition, there are no visible oscillations in comparison to the heterogeneous case.



(a) Function convergence for heterogeneous case.



(b) Iterates convergence for homogeneous case.

Figure 1. Performance of simplified IST on quadratic problem for varying step size values.

Simulations were performed on a machine with 24 Intel(R) Xeon(R) Gold 6246 CPU @ 3.30 GHz.

F. Conclusions and Future Work

In this study, we introduced a novel approach to understanding training with combined model and data parallelism for a quadratic model. Our framework sheds light on distributed submodel optimization, which reveals the advantages and limitations of Independent Subnetwork Training (IST). Moreover, we accurately characterized the behavior of the considered method in both homogeneous and heterogeneous scenarios without imposing restrictive assumptions on the gradient estimators.

In future research, it would be valuable to explore extensions of our findings to settings that are closer to scenarios, such as cross-device federated learning. This could involve investigating partial participation support, leveraging local training benefits, and ensuring robustness against stragglers. Additionally, it would be interesting to generalize our results to non-quadratic scenarios without relying on pathological assumptions. Another potential promising research direction is algorithmic modifications of the original IST to solve the fundamental problems highlighted in this work and acceleration of training.