

Optimal Gradient Compression for Distributed and Federated Learning

Alyazeed Albasyoni Mher Safaryan Laurent Condat Peter Richtárik

King Abdullah University of Science and Technology (KAUST)

September 28, 2020

Abstract

Communicating information, like gradient vectors, between computing nodes in distributed and federated learning is typically an unavoidable burden, resulting in scalability issues. Indeed, communication might be slow and costly. Recent advances in communication-efficient training algorithms have reduced this bottleneck by using compression techniques, in the form of sparsification, quantization, or low-rank approximation. Since compression is a lossy, or inexact, process, the iteration complexity is typically worsened; but the total communication complexity can improve significantly, possibly leading to large computation time savings. In this paper, we investigate the fundamental trade-off between the number of bits needed to encode compressed vectors and the compression error. We perform both worst-case and average-case analysis, providing tight lower bounds. In the worst-case analysis, we introduce an efficient compression operator, *Sparse Dithering*, which is very close to the lower bound. In the average-case analysis, we design a simple compression operator, *Spherical Compression*, which naturally achieves the lower bound. Thus, our new compression schemes significantly outperform the state of the art. We conduct numerical experiments to illustrate this improvement.

Contents

1	Introduction	2
1.1	Related Work	2
1.2	Contributions	3
2	Classes of Compression Operators	4
2.1	Compression operator as composition of encoder and decoder	4
2.2	Two senses of optimality for compression	5
2.3	Dimension-tolerant compression schemes.	5
2.4	Compressed learning algorithms	5
3	Worst-Case Analysis	7
3.1	Asymptotic tightness of the lower bound (1)	7
3.2	New Compressor: Sparse Dithering (SD)	7
3.3	Tighter bounds on minimal communication	8
4	Average-Case Analysis	9
4.1	Lower bound on average communication	9
4.2	Randomized version of Sparse Dithering	9
4.3	New Compressor: Spherical Compression (SC)	10

5 Experiments	10
5.1 Setting	10
5.2 Communication versus Convergence	13
5.3 Total Communication as a Function of α	13
A Discussion on Finite Precision Floats	19
B Proofs for Section 2	19
B.1 Relaxed classes of compression operators	19
B.2 Two senses of optimality for compression: Proof of Proposition 1	20
B.3 Dimension-tolerant compression schemes: Proof of Theorem 1	21
C Proofs for Section 3	21
C.1 Asymptotic tightness of the lower bound (1): Proof of Theorem 2	21
C.2 Deterministic-biased version of SD: Proof of Theorem 3	22
C.3 Tighter bounds on minimal communication: Proofs of Theorems 4 and 5	24
D Proofs for Section 4	25
D.1 Lower bound on average communication: Proof of Theorem 6	25
D.2 Randomized-unbiased version of Sparse Dithering: Proof of Theorem 7	26

1 Introduction

Due to the necessity of huge amounts of data to achieve high-quality machine learning models [27, 35], modern large-scale training procedures are executed in a distributed environment [5, 36]. In such a setup, both storage and computation needs are reduced, as the overall data (potentially too big to fit into a single machine) is partitioned among the nodes and computation is carried out in parallel. However, in order to keep the consensus across the network, compute nodes have to exchange some information about their local progress [31, 17, 4]. The demand of information communication between all machines in a distributed setup is typically a burden, resulting in a scalability issue commonly referred to as *communication bottleneck* [28, 40, 23]. To reduce the amount of information to be transferred, information is passed in a compressed or inexact form. *Information lossy compression* is a common practice, where original information is encoded approximately with essentially fewer bits, while introducing additional controllable distortion into the decoded message.

In the context of *Federated Learning* [21, 24, 15], communication between devices arises naturally, as data is initially decentralized and should remain so, for privacy purposes. Actually, it might be desirable for each unit, or client, to compress/encode/encrypt the information they are going to share, in order to minimize private data disclosures. Another practical scenario where compression methods are useful is when storage capabilities are scarce or there is no need to save complete versions of the data. In such cases, the representation of the data (encoding and decoding schemes) can be optimized and with little or no precision loss, one can allocate significantly less memory space.

1.1 Related Work

Recently, substantial amount of work has been devoted to the advances of communication-efficient training algorithms by utilizing various types of compression mechanisms, such as sparsification [38, 37, 10], quantization [1, 39, 13] and low-rank approximation [36]. Typically, the information communicated by computing nodes consists of local gradients, to which compression operators are applied. For example, one popular example of such compression operator is Top- k [3], which transfers only k coordinates of the gradient with largest magnitudes.

The theoretical foundation of lossy compression has long history and is based on *Rate-Distortion Theory* introduced by Shannon in his seminal papers [29, 30]. Recently, rate-distortion theory has been utilized in the context of model compression [12, 8]. In contrast to this, another line of research is devoted to

the lossless compression methods which is rooted in *Shannon’s source coding theorem* [9]. Both approaches exploit statistical properties of the input messages for analyses, which differs from our setting.

We investigate the problem of lossy compression, namely encoding vectors $x \in \mathbb{R}^d$ without prior knowledge on the distribution, for any $d \geq 1$, into as few bits as possible, while introducing as little distortion as possible. Formally, we measure the distortion of a (possibly randomized) compression operator $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ by its constant $\alpha \in [0, 1]$ such that $\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq \alpha \|x\|^2$ for every x , where the norm is the Euclidean norm (see Definitions 1, 2 and 3 for details). We denote by b the number of bits (in the worst case or in expectation) needed to encode $\mathcal{C}(x)$. Intuitively, b and α cannot be too small at the same time: they are antagonistic and ruled by a fundamental rate-distortion trade-off. As a matter of fact, as shown in [26], the following lower bound, referred to as *uncertainty principle for communication compression*, holds (if omitted, the base of log is assumed to be 2):

$$\alpha 4^{b/d} \geq 1 \quad \text{or, equivalently,} \quad b \geq d \frac{1}{2} \log \frac{1}{\alpha}. \quad (1)$$

In this work, we investigate this trade-off more deeply. We perform two types of analyses: *worst case analysis (WCA)* and *average case analysis (ACA)*. Then, capitalizing on this new knowledge, we design new efficient compression schemes. Note that our derivations deal with real numbers, compressed using a finite number of bits. We should keep in mind that numbers are represented by finite-precision, say 32 bits, floats in computers. We can safely omit this aspect in the derivations; we discuss this point in more details in the Appendix.

1.2 Contributions

Here we summarize our key contributions.

- **(WCA) Tighter bounds on minimal communication.** First, we construct a compression scheme with α distortion and b encoding bits (in the worst case), which satisfies

$$\alpha 4^{b/d} \leq \text{poly}(d)^{1/d} \quad \text{or} \quad b \leq d \frac{1}{2} \log \frac{1}{\alpha} + \mathcal{O}(\log d). \quad (2)$$

This implies the asymptotic tightness of the bound (1) as the dimension d grows (see Theorem 2). Then, we investigate the minimal number of bits (in the worst case) $b^*(\alpha, d)$ as a function of distortion α and dimension d proving that

$$b^*(\alpha, d) = -\log P(\alpha, d) + \log d + \frac{1}{2} \log \log d + e,$$

where $P(\alpha, d) := \frac{1}{2} I_\alpha(\frac{d-1}{2}, \frac{1}{2})$ with I_α being the regularized incomplete beta function, and e is negligible additive error with $|e| \leq \frac{1}{2} \log \log d + \mathcal{O}(1)$ (see Theorem 5), as opposed to $\mathcal{O}(\log d)$ in (1) and (2).

- **(WCA) Near optimal and practical compressor.** Motivated by these lower bounds we turn to the construction of a compression method which would be optimal and implementable in high dimensions. The example compression schemes in Theorem 2 ensuring (2) or in Theorem 5 are optimal but impractical, due to the exponential computation time to compress a vector. To make the scheme efficient, we slightly depart from the optimal boundary and propose a new efficient compression method—*Sparse Dithering (SD)*. Both deterministic (biased) and randomized (unbiased) versions of SD are analyzed, and comparisons with existing methods are made, showing that we outperform the state of the art. In the special case, the encoding of deterministic SD with $\alpha = 1/10$ distortion requires at most $30 + \log d + 3.35d$ bits, which is optimal within $1.69d$ additional bits (see Theorem 3).

- **(ACA) Lower bound on average communication.** Switching to the average case analysis, we establish a lower bound $-\log P(\alpha, d) \leq B$ on the expected number of bits B needed to encode a compression operator from $\mathbb{C}(\alpha)$ (see Definition 3).

- **(ACA) Compressor with optimal average communication.** As an attempt to reach the lower bound obtained in the average case analysis, we first analyze the randomized (and unbiased) version of SD. We prove that with variance $\omega > 0$ it requires at most

$$30 + \log d + \left(\log 3 + \frac{1}{2\sqrt{\omega}} \right) d$$

bits in expectation (see Theorem 7). In the special case of $\omega = 1/4$, it provides $\approx 9.9\times$ bandwidth savings. However, this scheme is suboptimal with respect to the lower bound. We finally present a simple compression operator—*Spherical Compression*—which attains the lower bound with less than 3 extra bits, namely it communicates $B < -\log P(\alpha, d) + 3$ bits in expectation (see Theorem 8).

2 Classes of Compression Operators

Here we formally define and perform preliminary analysis for three general classes of compression operators, that will be considered throughout the paper. We start with the most common and well studied class of unbiased compressors [2, 39, 18, 38].

Definition 1 (ω -compressors). *We denote by $\mathbb{U}(\omega)$ the class of unbiased compression operators $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ with variance $\omega \geq 0$; that is, $\mathbb{E}[\mathcal{C}(x)] = x$ and*

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

Another broad class of compressions operators, for which compressed learning algorithms have been successfully analysed [16, 32, 42, 6], is the class of biased operators, which are contractive in expectation.

Definition 2 (α -contractive operators). *We denote by $\mathbb{B}(\alpha)$ the class of (possibly biased and randomized) compression operators $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ with $\alpha \in [0, 1]$ -contractive property; that is,*

$$\mathbb{E}[\|\mathcal{C}(x) - x\|^2] \leq \alpha \|x\|^2, \quad \forall x \in \mathbb{R}^d. \quad (4)$$

Analogous to parameter ω for the variance, the parameter α is referred to as normalized variance or distortion threshold¹. It has been shown, that the class $\mathbb{U}(\omega)$ can be embedded into $\mathbb{B}(\alpha)$. Specifically, if $\mathcal{C} \in \mathbb{U}(\omega)$ then $\frac{1}{\omega+1}\mathcal{C} \in \mathbb{B}(\frac{\omega}{\omega+1})$ (see e.g. Lemma 1 in [26]). We will also consider the subclass of strictly contractive operators which, compared to operators from $\mathbb{B}(\alpha)$, are contractive for all realizations rather than in expectation:

Definition 3 (Strictly α -contractive operators). *We denote by $\mathbb{C}(\alpha)$ the class of (possibly biased and randomized) compression operators $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ with $\alpha \in [0, 1]$ -strictly contractive property; that is,*

$$\|\mathcal{C}(x) - x\|^2 \leq \alpha \|x\|^2, \quad \forall x \in \mathbb{R}^d. \quad (5)$$

2.1 Compression operator as composition of encoder and decoder

Generally speaking, compression is a two-sided notion, in the sense that one end encodes the message, while the other end decodes it to estimate the original information. An encoder is any mapping $E: \mathbb{R}^d \rightarrow \{0, 1\}^*$ which maps a given vector $x \in \mathbb{R}^d$ to some finite word from the set of all finite words $\{0, 1\}^*$ with the binary alphabet $\{0, 1\}$. A decoder, on the other hand, is a mapping $D: \{0, 1\}^* \rightarrow \mathbb{R}^d$ which aims to reconstruct the initial vector $x \in \mathbb{R}^d$ from the finite binary codeword $E(x)$. Thus, a compression operator $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ can be decomposed into an encoder and decoder so that $\mathcal{C}(x) = D(E(x))$. The number of bits needed to transfer a compressed version of $x \in \mathbb{R}^d$ is the length $|E(x)|$ of the binary word $E(x)$. In the worst case analysis we are interested in the length of the longest codeword $\sup_{x \in \mathbb{R}^d} |E(x)|$, while in average case analysis we investigate the size of the longest expected codeword $\sup_x \mathbb{E}_{\mathcal{C}} [|E(x)|]$.

Notice that a compression operator from any of the three classes requires countably many bits in order to encode points near $x = 0$ and $x = \infty$. We address this issue in the Appendix by considering relaxed classes of compression operators capturing finite representation of a single float in machines. From now on, we exclude trivial cases $\omega = 0$, $\alpha \in \{0, 1\}$ and assume $\omega > 0$, $\alpha \in (0, 1)$.

¹note that the definition of distortion in rate-distortion theory is slightly different than what we define.

2.2 Two senses of optimality for compression

It is worth distinguishing between optimality within a class in a single step of communication and optimality of total communication throughout the optimization process leading to ϵ -accuracy, e.g. $\frac{\|x^t - x^*\|^2}{\|x^0 - x^*\|^2} \leq \epsilon$ for a prescribed ϵ , where t is the iteration counter. Our theoretical contributions mainly deal with the first sense of optimality. Regarding the second view of optimality, the following proposition shows that Compressed Gradient Descent (CGD) can converge at significantly different speeds for different operators from $\mathbb{B}(\alpha)$.

Proposition 1. *If $\mathcal{C} \in \mathbb{B}(\alpha)$, the iteration complexity of CGD is $\frac{1}{1-\alpha}$ times bigger than for GD; that is CGD needs $\frac{1}{1-\alpha}$ times more iterations than GD to obtain the same ϵ -accuracy. Moreover, if \mathcal{C} is additionally unbiased, then only $1 + \alpha$ times more iterations are sufficient.*

Thus, if we aim to minimize the total communication complexity ensuring convergence to ϵ -accuracy, then the optimal operator \mathcal{C}^* should be either unbiased, or it will need to satisfy not only the direct condition, $\mathbb{E}[|E_{\mathcal{C}^*}(x)|] \leq \mathbb{E}[|E_{\mathcal{C}}(x)|]$ for all operators $\mathcal{C} \in \mathbb{B}(\alpha)$, but also the additional condition $\mathbb{E}[|E_{\mathcal{C}^*}(x)|] \leq (1 - \alpha^2)\mathbb{E}[|E_{\mathcal{U}}(x)|]$ for all unbiased operators $\mathcal{U} \in \mathbb{B}(\alpha)$. It is important to see that when α is close to 1, then this additional constraint is hard to satisfy when \mathcal{C}^* is not unbiased. For $\alpha < 1$, we show that this is indeed the case by obtaining an optimal biased operator $\mathcal{C}^* \in \mathbb{B}(\alpha)$, which we call *Spherical Compression*, and another unbiased one, which we call *Sparse Dithering*. We show that the latter is more suitable in practice due to its unbiasedness, and hence, convergence occurs in much fewer iterations, and that this is most pronounced when α is close to 1. In addition to being computationally efficient, we show that Sparse Dithering can guarantee reducing the total training communication by $\approx 9.9\times$ compared to full precision gradient communication of 32-bits floats.

2.3 Dimension-tolerant compression schemes.

By dimension-tolerant compression, we mean a collection of operators $\bar{\mathcal{C}} = (\mathcal{C}_d)_{d \geq d_0}$ that can be used to compress vectors $x \in \mathbb{R}^d$ for any $d \geq d_0$ and there exists a non-trivial fixed upper bound ($\bar{\omega} < \infty$ or $\bar{\alpha} < 1$) for variances (ω_d or α_d), i.e. $\omega_d \leq \bar{\omega} < \infty$ or $\alpha_d \leq \bar{\alpha} < 1$ for any $d \geq d_0$.

Below we show that for such collection of compression schemes, it is necessary and sufficient to use at least a constant amount of bits per dimension on average and this constant can be arbitrarily small.

Theorem 1. *The following holds:*

- (i) *If $\bar{\mathcal{C}} = (\mathcal{C}_d)_{d \geq d_0}$ is a dimension-tolerant compression composed of operators from $\mathbb{U}(\omega)$ ($\mathbb{B}(\alpha)$ or $\mathbb{C}(\alpha)$), then there exists a positive constant $c > 0$ (independent of d) such that for any $d \geq d_0$ at least cd bits are required in the worst case to encode $\mathcal{C}_d(x) \in \mathbb{R}^d$ for any $x \in \mathbb{R}^d$.*
- (ii) *Let $c > 0$ be a fixed positive constant. Then there exists a dimension-tolerant compression $\bar{\mathcal{C}} = (\mathcal{C}_d)_{d \geq 3}$ composed of operators from $\mathbb{U}(\omega)$ with $\omega = \mathcal{O}(1/c)$ ($\mathbb{B}(\alpha)$ or $\mathbb{C}(\alpha)$ with $\alpha = \frac{1}{1+\Omega(c)}$) such that $\mathcal{C}_d(x) \in \mathbb{R}^d$ can be encoded with cd bits for any $x \in \mathbb{R}^d$.*

Thus, $\Theta(d)$ bits need to be transmitted in order to bound the variance by a constant. The same asymptotic bound, $\Theta(d)$ bits per node, on total communication holds for distributed mean estimation [41, 20, 33].

2.4 Compressed learning algorithms

To highlight the importance of investigating the communication-variance trade-off of compression operators, we present how these operators affect the performance of compressed learning algorithms. For the sake of simplicity, consider distributed *Compressed Gradient Descent (CGD)* with compression operator $\mathcal{C} \in \mathbb{U}(\omega)$ solving the following smooth non-convex optimization problem

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

COMPRESSED LEARNING ALGORITHM	OBJECTIVE FUNCTION	ITERATION COMPLEXITY
COMPRESSED GD (CGD) [6, 18]	L -SMOOTH, μ -CONVEX	$\tilde{\mathcal{O}}\left(\frac{\kappa}{1-\alpha}\right), \tilde{\mathcal{O}}((\omega+1)\kappa)$
ACCELERATED CGD [22]	L -SMOOTH, μ -CONVEX	$\tilde{\mathcal{O}}((\omega+1)\sqrt{\kappa})$
ACCELERATED CGD [22]	L -SMOOTH, CONVEX	$\mathcal{O}\left((\omega+1)\sqrt{L/\varepsilon}\right)$
DISTRIBUTED CGD-DIANA [25, 14]	L -SMOOTH, μ -CONVEX	$\tilde{\mathcal{O}}\left(\omega + \frac{\omega\kappa}{n} + \kappa\right)$
DISTRIBUTED ACGD-DIANA [22]	L -SMOOTH, μ -CONVEX	$\tilde{\mathcal{O}}\left(\omega + \sqrt{(\omega/n + \sqrt{\omega/n})\omega\kappa} + \sqrt{\kappa}\right)$
QUANTIZED SGD (QSGD) [2]	L -SMOOTH, CONVEX	$\mathcal{O}\left(\frac{\omega}{n} \frac{1}{\varepsilon^2} + \frac{L}{\varepsilon}\right)$
DISTRIBUTED COMPRESSED SGD [13]	L -SMOOTH, NON-CONVEX	$\mathcal{O}\left((\omega+1)(\omega/n+1)\frac{L}{\varepsilon^2}\right)$
COMPRESSED SGD WITH ERROR FEEDBACK (EF-SGD) [32, 6]	L -SMOOTH, μ -CONVEX	$\tilde{\mathcal{O}}\left(\frac{\kappa}{1-\alpha} + \frac{1}{\mu\varepsilon}\right)$
COMPRESSED EF-SGD [16]	L -SMOOTH, NON-CONVEX	$\mathcal{O}\left(\frac{L^2}{\varepsilon}\left(\frac{1}{\varepsilon} + \frac{1}{(1-\alpha)^2}\right)\right)$
DOUBLSQUEEZE [34]	SMOOTH, NON-CONVEX	$\mathcal{O}\left(\frac{1}{n\varepsilon^2} + \frac{1}{1-\alpha} \frac{1}{\varepsilon^{1.5}} + \frac{1}{\varepsilon}\right)$

Table 1: Iteration complexities of various compressed learning algorithms with respect to the variance (ω or α) of the compression operator. For smooth and strongly convex (μ -convex with $\mu > 0$) objectives $\kappa = L/\mu$ indicates the condition number. $\tilde{\mathcal{O}}$ hides logarithmic factor $\log 1/\varepsilon$, n denotes the number of nodes, ε is the desired convergence accuracy.

where n is the number of nodes or machines available and $f_i(x)$ is the loss function corresponding to the data stored at node i . Hence, CGD algorithm iteratively performs the updates $x^{t+1} = x^t - \gamma_t g^t$ with unbiased gradient estimator

$$g^t = \frac{1}{n} \sum_{i=1}^n g_i^t := \frac{1}{n} \sum_{i=1}^n \mathcal{C}(\nabla f_i(x^t)).$$

Using smoothness of the loss function $f(x)$, the expected loss is upper bounded as follows:

$$\begin{aligned} \mathbb{E}[f(x^{t+1})|x^t] &= \mathbb{E}_t[f(x^t - \gamma_t g^t)] \\ &\stackrel{L\text{-smoothness}}{\leq} f(x^t) - \gamma_t \|\nabla f(x^t)\|^2 + \frac{L\gamma_t^2}{2} \mathbb{E}_t[\|g^t\|^2] \\ &= f(x^t) - \frac{2\gamma_t - L\gamma_t^2}{2} \|\nabla f(x^t)\|^2 + \frac{L\gamma_t^2}{2} \mathbb{E}_t[\|g^t - \nabla f(x^t)\|^2], \end{aligned}$$

where $L > 0$ is the smoothness parameter. Now, the term that is affected by compression and slowing down the convergence is the last one, namely the variance of estimator g^t , which can be transform into

$$\begin{aligned} \mathbb{E}_t[\|g^t - \nabla f(x^t)\|^2] &= \mathbb{E}_t\left[\left\|\frac{1}{n} \sum_{i=1}^n (g_i^t - \nabla f_m(x^t))\right\|^2\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}_t[\|g_i^t - \nabla f_i(x^t)\|^2] \stackrel{(3)}{\leq} \frac{\omega}{n^2} \sum_{i=1}^n \|\nabla f_i(x^t)\|^2. \end{aligned}$$

Clearly, in case of no compression ($\omega = 0$), this term vanishes. Thus, the slowdown caused by the compression operator $\mathcal{C} \in \mathbb{U}(\omega)$ is controlled by its parameter ω .

Similarly, for compression operators from $\mathbb{B}(\alpha)$ or $\mathbb{C}(\alpha)$, the parameter α controls the slowdown. Table 1 summarizes iteration complexities of various learning algorithms exploiting compressed communication and exposes the dependence of the variance (ω and α) of compression operator. The conclusion from this discussion and from Table 1 is that to facilitate fast and communication-efficient training process, one needs to design compression operators minimizing both variance and number of encoding bits. This is the

motivation of our work. Therefore, compression operators developed in this paper can be incorporated in any compressed learning algorithm, including all the ones in Table 1.

3 Worst-Case Analysis

We start our analysis of compression operators with respect to the number of encoding bits in the worst case. First, we show that the lower bound (1) for the class $\mathbb{B}(\alpha)$ is asymptotically tight for any $\alpha \in (0, 1)$. Then, we design an efficient compression operator from $\mathbb{C}(\alpha)$, *Sparse Dithering*, which is within a small constant factor of being optimal. Finally, we derive asymptotically tighter lower and upper bounds.

3.1 Asymptotic tightness of the lower bound (1)

First, we show that for any fixed $\alpha \in (0, 1)$, the constant 1 in the lower bound (1) is not improvable. We denote by $\mathbb{S}^d = \{x \in \mathbb{R}^d : \|x\| = 1\}$ the unit sphere of \mathbb{R}^d .

Theorem 2. *For any given $\alpha \in (0, 1)$ and $d \geq 3$ there exists an α -contractive compression operator $\mathcal{C}: \mathbb{S}^d \rightarrow \mathbb{R}^d$, such that*

$$\alpha 4^{b/d} \leq (1600d^2 \log d)^{2/d}, \quad (6)$$

where b is the number of bits (in the worst case) needed to encode $\mathcal{C}(x) \in \mathbb{R}^d$ for any unit vector $x \in \mathbb{S}^d$. In particular, for any $\alpha \in (0, 1)$ and $\epsilon > 0$ one can choose d large enough such that compression operator \mathcal{C} satisfies

$$\alpha 4^{b/d} < 1 + \epsilon. \quad (7)$$

Remark 1. *Using covering results from [11] (see Theorem 1), the constant 1600 in (6) can be reduced up to 2. Using tighter inequalities for the Γ function, the term d^2 can be improved as well. However, these will not improve the inequality (7). Notice that the right hand side of (6) approaches 1 quickly; for $d = 10^3$ it is ≈ 1.047 .*

Note that the compression operator in this theorem acts on \mathbb{S}^d , not \mathbb{R}^d . However, allocating an additional constant amount of bits for the norm $\|x\|$ (say 31 bits in *float32* format), we can extend the domain of compression operators without hurting the asymptotic tightness. Thus, the lower bound (1) is asymptotically tight for the class $\mathbb{B}(\alpha)$.

Although the construction of this theorem yields an optimal contractive operator, it is infeasible to apply in high dimensions.

3.2 New Compressor: Sparse Dithering (SD)

With the aim of constructing both optimal and efficient compression operators, we introduce a new compression scheme—*Sparse Dithering (SD)*—which is efficient in high dimension and nearly optimal. In some sense, SD can be viewed as an effective combination of Top- k sparsification [3] and random dithering with uniform levels [2]. The essential novelty is the encoding scheme and better upper bound on the number of communicated bits. In this section, we present a deterministic and hence biased version of SD.

Construction and variance bound. To compress a given nonzero vector $x \in \mathbb{R}^d$, we first compress the normalized vector $u = x/\|x\| \in \mathbb{S}^d$ and then rescale it. To quantize the coordinates of the unit vector u , we apply dithering with levels $2k_i h$, $k_i \geq 0$, where $h = \sqrt{\nu/d}$ is the half-step and $\nu > 0$ is a free parameter. For each coordinate u_i , $i \in [d]$ we choose the nearest level so that $\|u_i - 2k_i h\| \leq h$. Letting $\hat{u}_i = \text{sign}(u_i) 2k_i h$ we have $|u_i - \hat{u}_i| \leq h$ for all $i \in [d]$. Therefore,

$$\|u - \hat{u}\|^2 = \sum_{i=1}^d (u_i - \hat{u}_i)^2 \leq dh^2 = \nu.$$

Note that, after applying the scaling factor $\|x\|$, this gives a compression with variance at most ν . However, $\|x\|$ is not always the best option. Specifically, we can choose the scaling factor $\gamma > 0$ so to minimize the variance $\|x - \gamma\hat{u}\|^2$, which yields the optimal factor $\gamma^* = \frac{\langle x, \hat{u} \rangle}{\|\hat{u}\|^2}$ with the optimal variance of $\|x - \gamma^*\hat{u}\|^2 = \sin^2 \varphi \|x\|^2$, where $\varphi \in [0, \pi/2]$ is the angle² between x and \hat{u} . Hence, defining the compression operator as $\mathcal{C}(x) = \gamma^*\hat{u}$, we have the following bound on the variance:

$$\|\mathcal{C}(x) - x\|^2 \leq \min(\nu, \sin^2 \varphi) \|x\|^2.$$

Encoding scheme. We now describe the corresponding encoding scheme into a sequence of bits. With the following notations:

$$\gamma := 2h\gamma^* \in \mathbb{R}_+, \quad k := (k_i)_{i=1}^d \in \mathbb{N}_+^d, \quad s := (\text{sign}(u_i k_i))_{i=1}^d \in \{-1, 0, 1\}^d,$$

the compression operator can be written as $\mathcal{C}(x) = \gamma^*\hat{u} = 2h\gamma^* \text{sign}(u) k = \gamma s k$. So, we need to encode the triple (γ, s, k) . As $\gamma \in \mathbb{R}_+$, we need only 31 bits for the scaling factor. Next we encode s . Let

$$n_0 := |\{i \in [d] : s_i = 0\}| = |\{i \in [d] : k_i = 0\}|$$

be the number of coordinates u_i that are compressed to 0. To communicate s , we first send the locations of those n_0 coordinates and then $d - n_0$ bits for the values ± 1 . Sending n_0 positions can be done by sending $\log d$ bits representing the number n_0 , afterwards sending $\log \binom{d}{n_0}$ bits for the positions. Finally, it remains to encode k , for which we only need to send nonzero entries, since the positions of $k_i = 0$ are already encoded. We encode $k_i \geq 1$ with k_i bits: $k_i - 1$ ones followed by a zero. Hence, encoding k requires $\sum k_i$ bits.

A theoretical upper bound on the total number of bits for any choice of parameter $\nu > 0$ is given in the Appendix. Below, we highlight one special case of $\nu = 1/10$.

Theorem 3. *Deterministic SD compression operator with parameter $\nu = 1/10$ belongs to $\mathbb{C}(1/10)$ communicating $30 + \log d + 3.35d$ bits at most. In addition, ignoring $30 + \log d$ negligible bits, SD is within a factor of*

$$\log_4(\alpha 4^{b/d}) = \log_4\left(\frac{1}{10} 4^{3.35}\right) \approx 1.69$$

of optimality; that is, at most 1.69d more bits are sent in comparison to optimal compression with the same normalized variance $1/10$.

3.3 Tighter bounds on minimal communication

We first look into the tightness of (1) when the normalized variance α approaches 1. In particular, for $\alpha = 1 - \frac{1}{d}$ the lower bound (1) implies that the number of bits b is lower bounded by some constant. However, the following holds:

Theorem 4. *For any compression operator from $\mathbb{B}(\alpha)$, with $\alpha \in (0, 1)$, at least $\log d$ bits are needed.*

As briefly mentioned before, the lower bound (1) is tight up to a $\mathcal{O}(\log d)$ additive error term. Here we perform a deeper analysis of the same lower bound.

Definition 4. *For a fixed $\alpha \in (0, 1)$ and dimension d , consider compression operators $\mathcal{C} \in \mathbb{B}(\alpha)$ with underlying encoder E , decoder D and define $b^*(\alpha, d)$ as the minimum number of bits in the worst case:*

$$b^*(\alpha, d) = \min_{\mathcal{C} \in \mathbb{B}(\alpha)} \max_{\|x\|=1} |E(x)|.$$

In other words, for any compression operator from $\mathbb{B}(\alpha)$ there exists a unit vector that cannot be encoded into less than $b^(\alpha, d)$ bits and it shows the least amount of bits with such property.*

²in case of $\mathcal{C}(x) = 0$ we let $\varphi = \pi/2$.

Combining lower bound (1) with (6) of Theorem 2, yields

$$b^*(\alpha, d) = \frac{1}{2} \log \frac{1}{\alpha} + e \quad \text{with error term } 0 \leq e = \mathcal{O}(\log d).$$

Denoting $P(\alpha, d) := \frac{1}{2} I_\alpha \left(\frac{d-1}{2}, \frac{1}{2} \right) \in (0, \frac{1}{2})$, where I_α is the regularized incomplete beta function, we show tighter asymptotic behavior:

Theorem 5. *With error term $|e| \leq \frac{1}{2} \log \log d + \mathcal{O}(1)$,*

$$b^*(\alpha, d) = -\log P(\alpha, d) + \log d + \frac{1}{2} \log \log d + e.$$

4 Average-Case Analysis

Now we switch to the average-case analysis for the class $\mathbb{C}(\alpha)$. First, we prove a lower bound for communicated bits in expectation. Then we analyze the randomized version of Sparse Dithering, which, having better theoretical guarantees than random dithering, is suboptimal in this analysis. Finally, we present a new compression operator from $\mathbb{C}(\alpha)$, *Spherical Compression*, which is provably optimal.

4.1 Lower bound on average communication

In this section, we consider compression operators from $\mathbb{C}(\alpha)$ and investigate the trade-off between normalized variance α and expected number of bits

$$B = \sup_{\|x\|=1} \mathbb{E}_{\mathcal{C}} [|E(x)|].$$

In other words, we study the trade-off for strictly α -contractive operators that encode any unit vector with no more than B bits in expectation. In such a setting, we show the following lower bound on B .

Theorem 6. *Let $\mathcal{C} \in \mathbb{C}(\alpha)$ be a compression operator such that $\mathcal{C}(x) \in \mathbb{R}^d$ can be transferred with B bits in expectation for any unit vector $x \in \mathbb{S}^d$. Then $-\log P(\alpha, d) \leq B$.*

4.2 Randomized version of Sparse Dithering

Here we randomize Sparse Dithering to make it unbiased and estimate the number of encoding bits it needs in expectation. First we decompose the to-be-compressed vector $x \in \mathbb{R}^d$ into the magnitude and unit direction $u = x/\|x\|$ as before. To randomize the scheme, each coordinate u_i gets rounded to one of the two nearest neighbors, so as to preserve unbiasedness; that is, if $2k_i h \leq |u_i| \leq 2(k_i + 1)h$ for some $k_i \geq 0$, then $\hat{u}_i = \text{sign}(u_i) 2\hat{k}_i h$ where

$$\hat{k}_i = \begin{cases} k_i & \text{with prob. } \frac{2(k_i+1)h - |u_i|}{2h} \\ k_i + 1 & \text{with prob. } \frac{|u_i| - 2k_i h}{2h}. \end{cases}$$

Clearly, $\mathbb{E}[\hat{u}] = u$ and defining $\mathcal{C}(x) = \|x\| \hat{u}$, we maintain unbiasedness $\mathbb{E}[\mathcal{C}(x)] = x$. The encoding scheme is the same as in the deterministic case. Upper bounding the expected number of bits and the variance, we obtain:

Theorem 7. *Randomized SD compression with parameter $\nu = \omega$ belongs to $\mathbb{U}(\omega)$, communicating at most*

$$30 + \log d + \left(\log 3 + \frac{1}{2\sqrt{\omega}} \right) d$$

bits in expectation. In particular, with $\omega = 1/4$ variance (ignoring $30 + \log d$ negligible factors), it uses $(1 + \log 3)d \approx 2.6d$ bits in each iteration (about 12 times less than full precision case) and forces up to $1 + \omega = 5/4$ times more iterations, leading to ≈ 9.9 times bandwidth savings.

As mentioned earlier, SD is similar to random dithering with uniform levels, namely with \sqrt{d} levels. However, with a different parametrization ν and better encoding strategy, SD provides better theoretical guarantees. Indeed, random dithering with \sqrt{d} levels communicates $\approx 2.8d$ bits in expectation and requires $1 + \omega = 2$ times more iterations, resulting in a factor of ≈ 5.7 in bandwidth saving (see Theorem 3.2 and Corollary 3.3 of [2]). For more comparisons on bandwidth savings see Table 2 in the Appendix.

4.3 New Compressor: Spherical Compression (SC)

It can be shown that randomized SD compression discussed in the previous section is suboptimal with respect to the lower bound of Theorem 6. Here we provide a simple compression operator—*Spherical Compression (SC)*—that achieves this lower bound with less than 3 overhead bits.

Construction and variance bound. As before, we transmit the magnitude and direction separately. For a given unit vector $x \in \mathbb{S}^d$, SC generates a sequence $(x^t)_{t=1}^T$ of i.i.d. points with $\|x^t\|^2 = 1 - \alpha$, and terminates once $\|x^T - x\|^2 \leq \alpha$ for some $T \geq 1$. The last generated point x^T is the compressed version of x we need to communicate, that is $\mathcal{C}(x) = x^T$. It follows directly from this construction that $\mathcal{C} \in \mathbb{C}(\alpha)$.

Encoding scheme. The crucial part of the encoding scheme is that it is enough to communicate only T . Indeed, the communication process is the following. Importantly, the emitter and receiver have agreed on using the same random seed for generating i.i.d. points (x^t) , before the compression of any vector is performed.

Then, upon receiving the number of trials T , the decoder can reproduce the same sequence x^1, x^2, \dots, x^T and recover x^T . Consequently, it remains to encode the random integer T into a binary code.

Upper bound on B . First we show that T follows a geometric distribution with parameter $p = P(\alpha, d)$. Indeed, T can be viewed as the number of trials before the first success happens after a series of failures.

In our case, trials correspond to generating i.i.d. points x^t and success means $x \in C^d(x^t, \sqrt{\alpha})$ which happens with probability $P(\alpha, d) \in (0, 1/2)$. Therefore, the expected number of points x^t we need to generate until we get into α -vicinity of the initial point x is $\mathbb{E}[T] = 1/p = 1/P(\alpha, d) > 2$. Next, we encode T with the Golomb–Rice coding scheme, which is known to be optimal for geometric distributions. Define integer $m \geq 0$ from $1/2p \leq 2^m < 1/p$ and decompose T as $T = 2^m q + r$ with $q \geq 0$, $0 \leq r < 2^m$. The quotient q is encoded with unary coding as a string of q zeros followed by a 1. The remainder r is communicated with exactly m bits using truncated binary coding. There is no need to send the value of m as it can be computed from p , which depends only on α and d . Hence, the total number of bits to encode T is no more than $q + m + 1$. Note that $m < \log 1/p = -\log p$ is fixed, while q depends on T and $q \leq 2^{-m}T$. Hence,

$$B = \mathbb{E}[q + m + 1] < 2^{-m} \mathbb{E}[T] - \log p + 1 = \frac{1}{2^m p} - \log p + 1 \leq -\log p + 3,$$

which implies:

Theorem 8. *In the average-case analysis, Spherical Compression is optimal up to 3 extra bits; that is, it communicates $B < -\log P(\alpha, d) + 3$ bits in expectation.*

Remark 2. *It is worth mentioning that the above compression operator satisfies $\|\mathcal{C}(x) - x\|^2 \leq \alpha \|x\|^2$ in the worst case, not in expectation. Moreover, because of the symmetry of spheres and caps $C^d(x^t, \sqrt{\alpha})$, it can be seen from the construction that $\mathbb{E}[\mathcal{C}(x)]$ points to the same direction as the initial vector x . Thus, with an appropriate (fixed) scaling factor, it can be made unbiased as well.*

5 Experiments

5.1 Setting

We consider both l_2 -regularized logistic regression and ridge regression. In both cases, we use regularizing coefficient $\lambda = \frac{1}{n}$. We run this on multiple datasets, and show that our compression methods provide

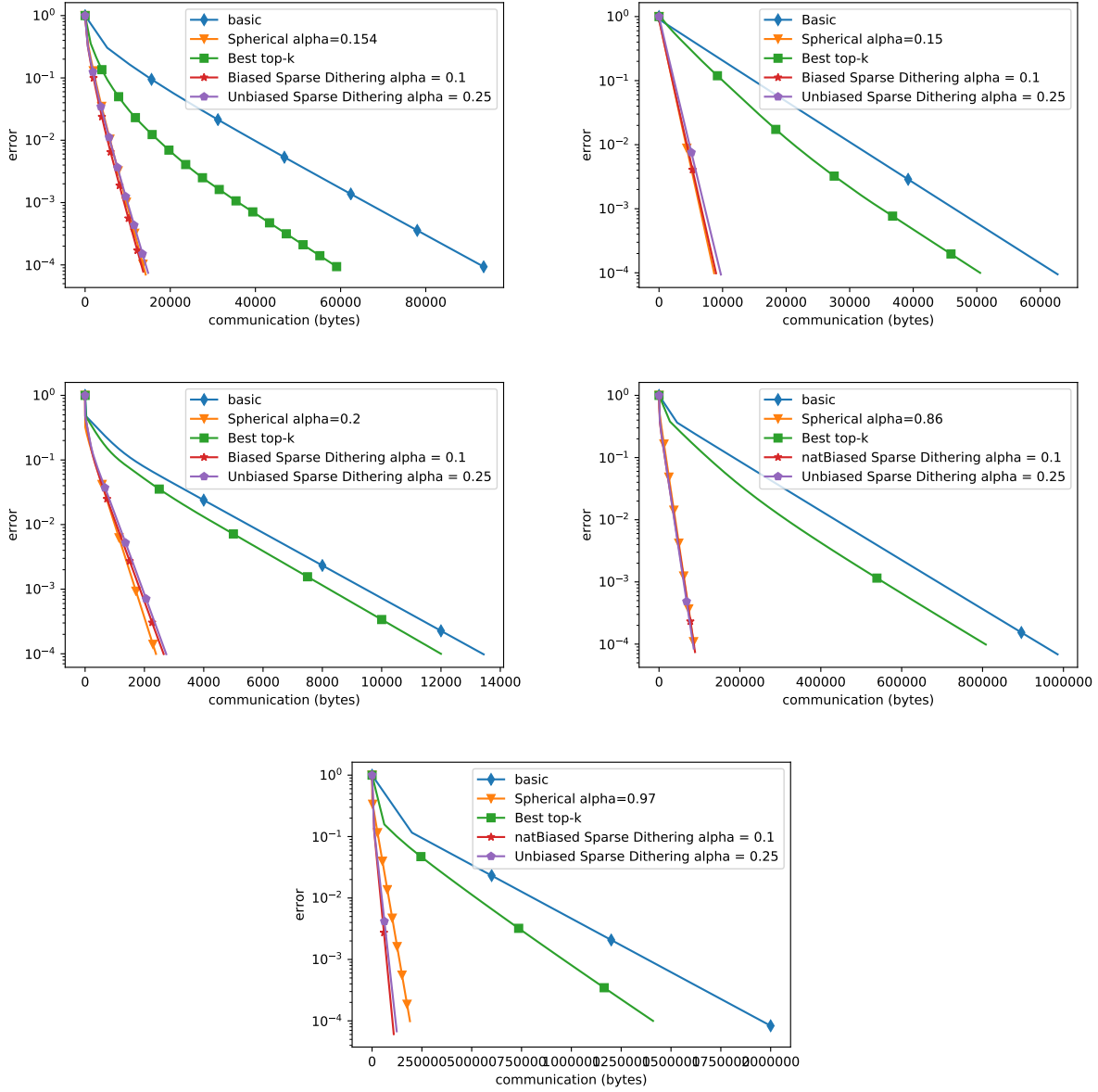


Figure 1: The first two plots correspond to ridge regression (Housing, Bodyfat datasets), while the next three plots correspond to regularized logistic regression (Breast Cancer, Madelon, Mushrooms datasets). This shows convergence as a function of total communication (in bytes), for various selected compression operators.

significant savings in communication (measured in bytes). The algorithm we used is Compressed Gradient Descent, which consists in iterating

$$x^{t+1} = x^t - \frac{1}{L} \mathcal{C}(\nabla f(x^t)),$$

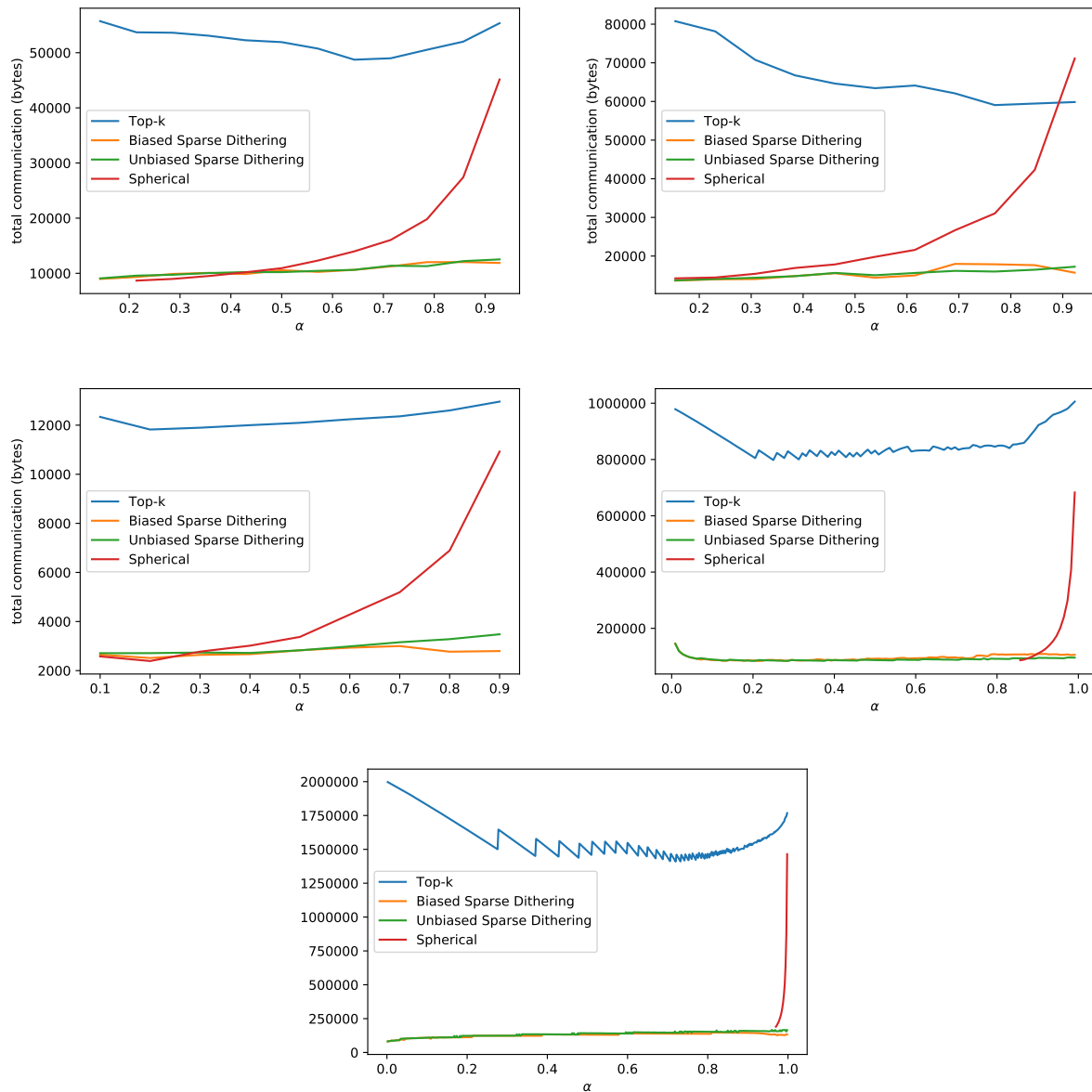


Figure 2: The first two plots correspond to ridge regression (Housing, Bodyfat datasets), while the next three plots correspond to regularized logistic regression (Breast Cancer, Madelon, Mushrooms datasets). This shows total communication needed to achieve $\epsilon = 10^{-4}$ as a function of α for various operators. For Top- k , we set $\alpha = 1 - \frac{k}{d}$, as predicted by the theory.

where f is the loss function, L is the smoothness constant of f computed explicitly. We stop the process whenever $\frac{\|x^t - x^*\|^2}{\|x^0 - x^*\|^2} \leq 10^{-4}$, where x^* is the minimizer of f and is computed beforehand for all problems.

5.2 Communication versus Convergence

In this experiment, we look at convergence, measured as $\frac{\|x^t - x^*\|^2}{\|x^0 - x^*\|^2}$ with respect to the number of bits communicated, for various compression operators. As shown in Figure 1, our compression operators significantly outperform all other operators. It is important that we compare our methods with the benchmark ‘Basic’, which sends a 32-bits float for every element in the gradient, sending a total of $32d$ bits at every iteration. In addition, we run these experiments with Top- k for all $1 \leq k \leq d$ and pick the best representative in the comparison, naming it ‘Best Top- k ’.

5.3 Total Communication as a Function of α

Here, we let $\alpha = 1 - k/d$ vary and in Figure 2 we show the total number of bits communicated before converging to ϵ -accuracy, with $\epsilon = 10^{-4}$. This clearly shows the superiority of our methods. It is important to note that Sparse Dithering can beat the optimal Spherical Compression, because it is unbiased, so it requires significantly less iterations. It’s important to note that for $1 \leq k \leq d$, we plot Top- k at $\alpha = 1 - k/d$.

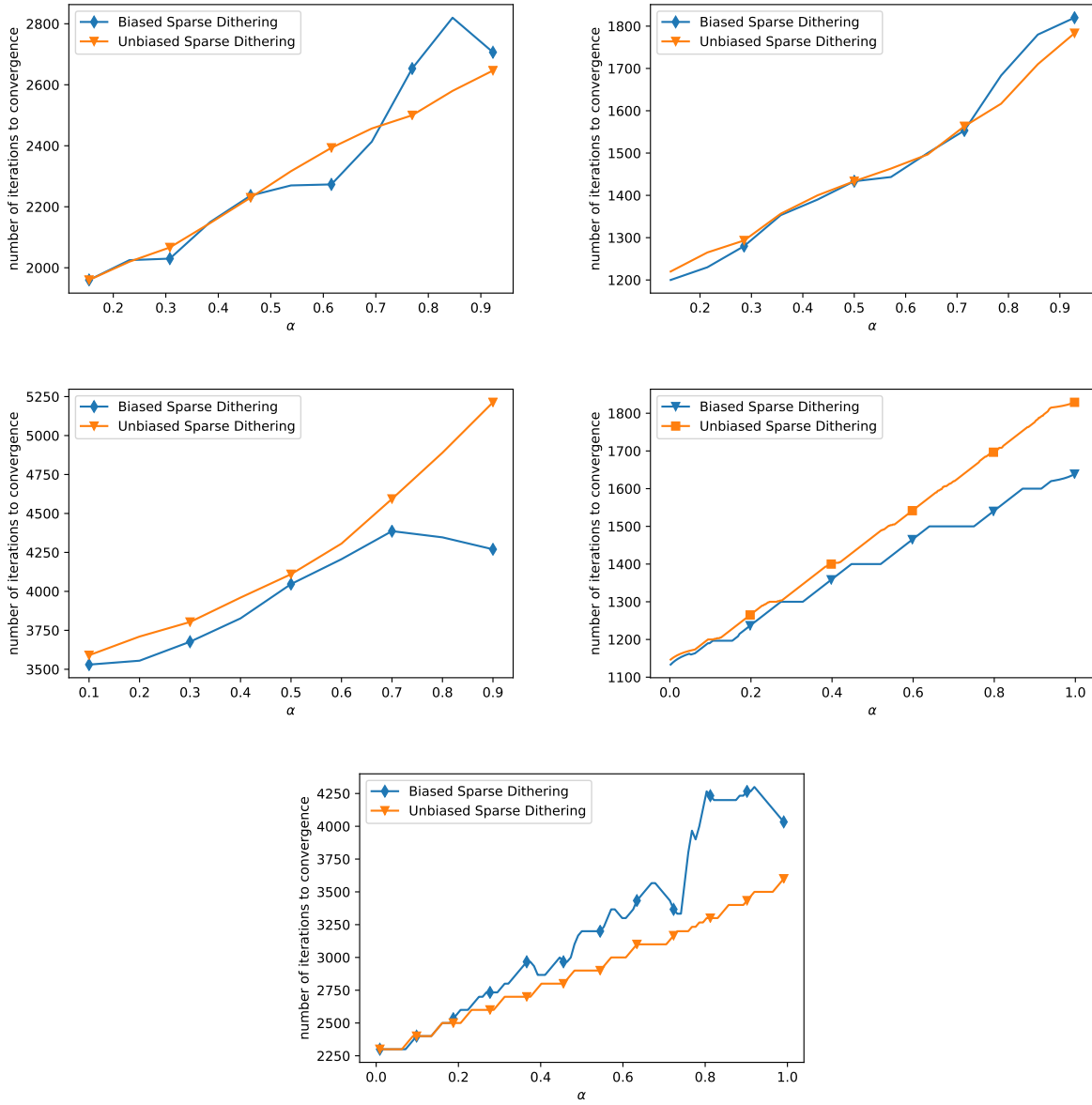


Figure 3: The first two plots correspond to ridge regression (Housing, Bodyfat datasets), while the next three plots correspond to regularized logistic regression (Breast Cancer, Madelon, Mushrooms datasets). This shows the number of iterations as a function of α for both sparse dithering methods. Both curves look like $Y = 1 + X$, as predicted by the theory.

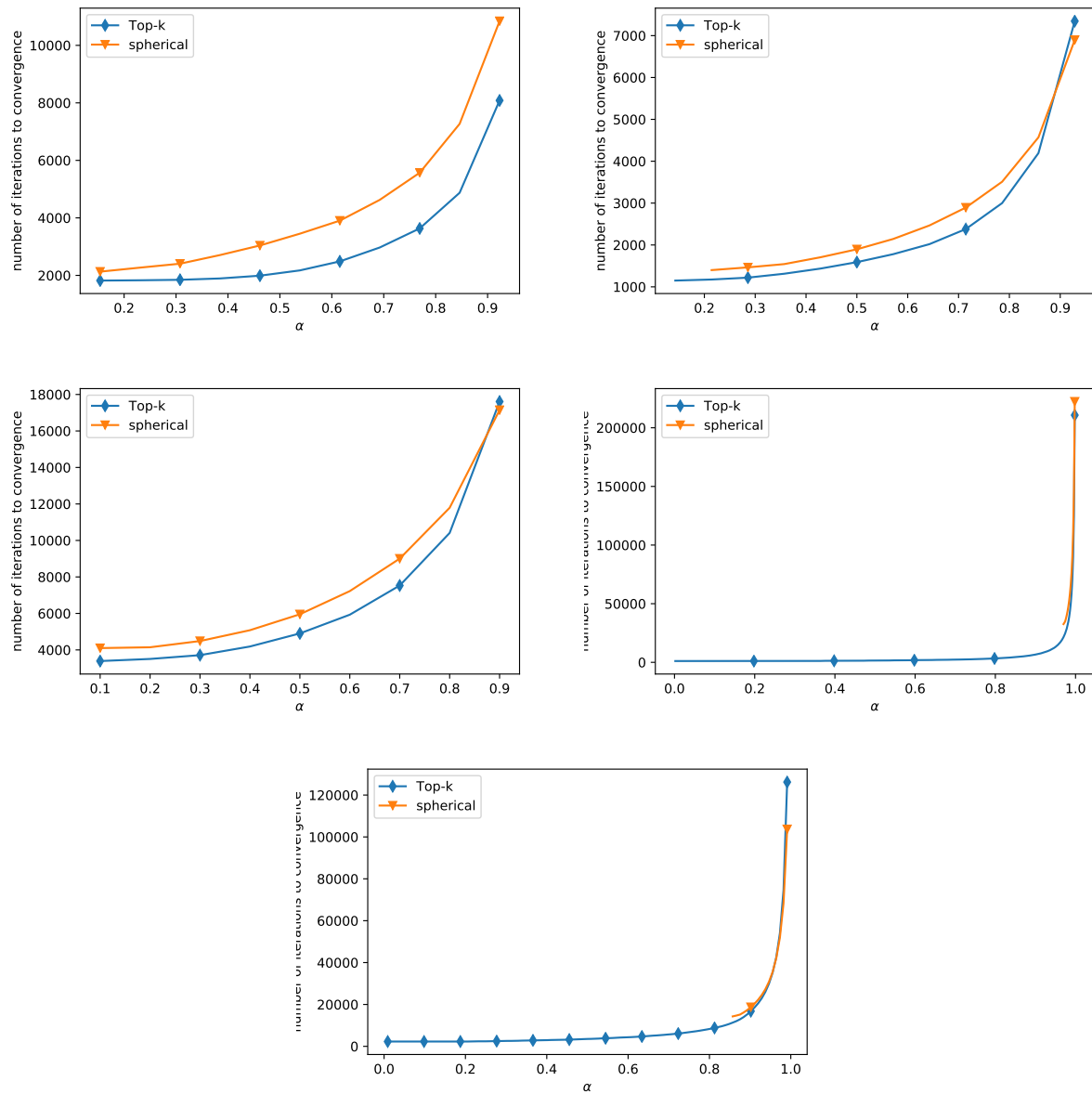


Figure 4: The first two plots correspond to ridge regression (Housing, Bodyfat datasets), while the next three plots correspond to regularized logistic regression (Breast Cancer, Madelon, Mushrooms datasets). This shows the number of iterations as a function of α for both Top- k and Spherical compressions. Both curves look like $Y = \frac{1}{1-X}$, as predicted by the theory.

References

- [1] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient SGD via gradient quantization and encoding. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2017.
- [2] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-efficient sgd via gradient quantization and encoding. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2017.
- [3] Dan Alistarh, Torsten Hoefler, Mikael Johansson, Sarit Khirirat, Nikola Konstantinov, and Cédric Renggli. The convergence of sparsified gradient methods. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2018.
- [4] Debraj Basu, Deepesh Data, Can Karakus, and Suhas N. Diggavi. Qsparse-local-SGD: Distributed SGD with quantization, sparsification, and local computations. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2019.
- [5] Ron Bekkerman, Mikhail Bilenko, and John Langford. *Scaling up machine learning: Parallel and distributed approaches*. Cambridge University Press, 2011.
- [6] Alexandre Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On biased compression for distributed learning. preprint arXiv:2002.12410, 2020.
- [7] Károly Böröczky and Gergely Wintsche. Covering the sphere by equal spherical balls. In Boris Aronov, Saugata Basu, János Pach, and Micha Sharir, editors, *Discrete and Computational Geometry: The Goodman-Pollack Festschrift*, pages 235–251. Springer, Berlin, Heidelberg, 2003.
- [8] Yuheng Bu, Weihao Gao, Shaofeng Zou, and Venugopal V. Veeravalli. Information-theoretic understanding of population risk improvement with model compression. In *AAAI Conference on Artificial Intelligence*, pages 3300–3307, February 2020.
- [9] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, USA, 2006.
- [10] N. Dryden, T. Moon, S. A. Jacobs, and B. V. Essen. Communication quantization for data-parallel training of deep neural networks. In *2nd Workshop on Machine Learning in HPC Environments (MLHPC)*, pages 1–8, Nov 2016.
- [11] Ilya Dumer. Covering spheres with spheres. *Discrete & Computational Geometry*, 38:665–679, 2007.
- [12] Weihao Gao, Yu-Han Liu, Chong Wang, and Sewoong Oh. Rate distortion for model Compression: From theory to practice. In *Int. Conf. Machine Learning (ICML)*, volume PMLR 97, pages 2102–2111, 2019.
- [13] Samuel Horváth, Chen-Yu Ho, Ľudovít Horváth, Atal Narayan Sahu, Marco Canini, and Peter Richtárik. Natural compression for distributed deep learning. 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. preprint arXiv:1904.05115, 2019.
- [15] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank J. Reddi, Sebastian U. Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for on-device federated learning. preprint arXiv:1910.06378, 2019.
- [16] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U Stich, and Martin Jaggi. Error feedback fixes SignSGD and other gradient compression schemes. preprint arXiv:1901.09847, 2019.

- [17] Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter theory for local SGD on identical and heterogeneous data. In *Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, 2020.
- [18] Sarit Khirirat, Hamid Reza Feyzmahdavian, and Mikael Johansson. Distributed learning with compressed gradients. preprint arXiv:1806.06573, 2018.
- [19] Martin Kochol. Constructive approximation of a ball by polytopes. *Mathematica Slovaca*, 44(1):99–105, 1994.
- [20] Jakub Konečný and Peter Richtárik. Randomized distributed mean estimation: accuracy vs communication. *Frontiers in Applied Mathematics and Statistics*, 4(62):1–11, 2018.
- [21] Jakub Konečný, H. Brendan McMahan, Felix Yu, Peter Richtárik, Ananda Theertha Suresh, and Dave Bacon. Federated learning: strategies for improving communication efficiency. In *NIPS Private Multi-Party Machine Learning Workshop*, 2016.
- [22] Zhize Li, Dmitry Kovalev, Xun Qian, and Peter Richtárik. Acceleration for compressed gradient descent in distributed and federated optimization. In *Int. Conf. Machine Learning (ICML)*, 2020.
- [23] Yujun Lin, Song Han, Huizi Mao, Yu Wang, and William J. Dally. Deep gradient compression: Reducing the communication bandwidth for distributed training. In *Int. Conf. Learning Representations (ICLR)*, 2018.
- [24] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, 2017.
- [25] Konstantin Mishchenko, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. Distributed learning with compressed gradient differences. preprint arXiv:1901.09269, 2019.
- [26] Mher Safaryan, Egor Shulgin, and Peter Richtárik. Uncertainty principle for communication compression in distributed and federated learning and the search for an optimal compressor. preprint arXiv:2002.08958, 2020.
- [27] Jürgen Schmidhuber. Deep learning in neural networks: An overview. In *Neural networks*, volume 61, page 85–117, 2015.
- [28] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and application to data-parallel distributed training of speech DNNs. In *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.
- [29] C.E. Shannon. Coding theorems for a discrete source with a fidelity criterion. *IRE Nat. Conv. Rec.*, 27:379–423,623–656, 1948.
- [30] C.E. Shannon. A mathematical theory of communication. *Bell Syst. Tech. J.*, 4:142–163, 1959.
- [31] Sebastian U. Stich. Local SGD converges fast and communicates little. In *Int. Conf. Learning Representations (ICLR)*, 2019.
- [32] Sebastian U. Stich and Sai Praneeth Karimireddy. The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication. preprint arXiv:1909.05350, 2019.
- [33] Ananda Theertha Suresh, Felix X. Yu, Sanjiv Kumar, and H. Brendan McMahan. Distributed mean estimation with limited communication. In *Int. Conf. Machine Learning (ICML)*, 2017.
- [34] Hanlin Tang, Chen Yu, Xiangru Lian, Tong Zhang, and Ji Liu. DoubleSqueeze: Parallel stochastic gradient descent with double-pass error-compensated compression. In *Int. Conf. Machine Learning*, volume PMLR 97, pages 6155–6165, 2019.

- [35] Sharan Vaswani, Francis Bach, and Mark Schmidt. Fast and faster convergence of SGD for over-parameterized models and an accelerated perceptron. In *Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, 2019.
- [36] Thijs Vogels, Sai Praneeth Karimireddy, and Martin Jaggi. PowerSGD: Practical low-rank gradient compression for distributed optimization. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2019.
- [37] Hongyi Wang, Scott Sievert, Shengchao Liu, Zachary Charles, Dimitris Papailiopoulos, and Stephen Wright. Atomo: Communication-efficient learning via atomic sparsification. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2018.
- [38] Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for communication-efficient distributed optimization. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2018.
- [39] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Terngrad: Ternary gradients to reduce communication in distributed deep learning. In *Neural Information Processing Systems Conf. (NeurIPS)*, page 1509–1519, 2017.
- [40] Hantian Zhang, Jerry Li, Kaan Kara, Dan Alistarh, Ji Liu, and Ce Zhang. ZipML: Training linear models with end-to-end low precision, and a little bit of deep learning. In *Int. Conf. Machine Learning (ICML)*, volume 70, page 4035–4043, 2017.
- [41] Yuchen Zhang, John Duchi, Michael I Jordan, and Martin J Wainwright. Information-theoretic lower bounds for distributed statistical estimation with communication constraints. In *Advances in Neural Information Processing Systems 26 (NIPS)*, pages 2328–2336, 2013.
- [42] Shuai Zheng, Ziyue Huang, and James T. Kwok. Communication-efficient distributed blockwise momentum SGD with error-feedback. In *Neural Information Processing Systems Conf. (NeurIPS)*, 2019.

Appendix

A Discussion on Finite Precision Floats

In the paper, we formally consider compression of arbitrary vectors of \mathbb{R}^d , but in practice, in computers, real numbers are represented with finite-precision floats, typically using 32 bits. As a consequence, a nonzero real cannot be too small, too large, and its precision is limited. Compression, like every operation, amounts to a sequence of elementary arithmetic operations, each being exact only up to so-called ‘machine precision’, which is difficult to model. Thus, we could restrict ourselves to vectors in $\{0\} \cup \{x : r \leq \|x\| \leq R\}$ for some $0 < r < R$, instead of the whole space \mathbb{R}^d , but this would not account for the finite precision of floats, and since this set is not stable by arithmetic operations, this would not be enough to model the setting in all rigor. So, we prefer to stick with the general setting of \mathbb{R}^d throughout the paper, since there is no issue with limit cases of very large or very small nonzero numbers, that would deserve a particular discussion; the finite precision makes them automatically irrelevant in practice. In other words, the finite representation of reals is not more problematic with compression than for any learning or optimization task, and more generally for the numerical implementation of any mathematical algorithm.

In particular, considering floats with 32 bits, a non-compressed vector x of \mathbb{R}^d is actually represented using $32d$ bits. When we decompose x into its ‘gain’ $\|x\|$ and ‘shape’ $x/\|x\| \in \mathbb{S}^d$, there is no trickery in considering that $\|x\|$ is represented using 31 bits (the sign bit can be omitted) and that $x/\|x\|$ is actually compressed. The multiplication by $\|x\|$ at decompression has finite precision, just like any arithmetic operation.

B Proofs for Section 2

B.1 Relaxed classes of compression operators

As mentioned in the paper and in Appendix A, any operator from $\mathbb{U}(\omega)$, $\mathbb{B}(\alpha)$, $\mathbb{C}(\alpha)$ cannot be encoded with a finite number of bits. For example, in the case of $\mathbb{B}(\alpha)$, the inequality (4) breaks near $x = 0$ and $\|x\| \rightarrow \infty$. However, in practice, machine floats have finite precision and we do not deal with nonzero values that are too small and too large. To reflect this practical aspect into the theory, we adjust the definition of α -contractive compressors $\mathbb{B}(\alpha)$ and consider the following class instead. For the sake of concreteness, we carry out the discussion for the class $\mathbb{B}(\alpha)$ only and note that analogous observations can be adopted for other two classes.

Definition 5 (Practical α -contractive compressions). *Let $\alpha \in (0, 1)$ and $R \geq 1$ be fixed. We denote by $\mathbb{B}^1(\alpha, R)$ the class of (possibly randomized) operators $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that*

$$\begin{aligned} \mathbb{E} [\|\mathcal{C}(x) - x\|^2] &\leq \alpha \|x\|^2, & \text{if } 1/R \leq \|x\| \leq R \\ \mathcal{C}(x) &= 0, & \text{if } \|x\| < 1/R \text{ or } \|x\| > R. \end{aligned}$$

Note that, for simplicity, we take $r = 1/R$.

The class of all α -contractive operators $\mathbb{B}(\alpha)$ can be seen as the limit of the class $\mathbb{B}^1(\alpha, R) \rightarrow \mathbb{B}(\alpha)$ as $R \rightarrow \infty$. The advantage of the class $\mathbb{B}^1(\alpha, R)$ compared to $\mathbb{B}(\alpha)$ is that it allows an encoding with finite number of bits. Next, we relax the definition of $\mathbb{B}^1(\alpha, R)$ as follows:

Definition 6 (Weak α -contractive compressions). *Let $\alpha \in (0, 1)$ and $R > 0$ be fixed. We denote by $\mathbb{B}^2(\alpha, R)$ the class of (possibly randomized) operators $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that*

$$\begin{aligned} \mathbb{E} [\|\mathcal{C}(x) - x\|^2] &\leq \alpha R^2, & \text{if } \|x\| \leq R \\ \mathcal{C}(x) &= 0, & \text{if } \|x\| > R. \end{aligned}$$

The following simple lemma shows that the latter class is much more general and contains the first class.

Lemma 1. *If $R \geq \alpha^{-1/4}$ then $\mathbb{B}^1(\alpha, R) \subset \mathbb{B}^2(\alpha, R)$.*

Proof. Let $\mathcal{C} \in \mathbb{B}^1(\alpha, R)$ with $\alpha R^4 \geq 1$. If $\|x\| < 1/R$ then

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] = \|x\|^2 < \frac{1}{R^2} \leq \alpha R^2.$$

If $1/R \leq \|x\| \leq R$ then

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq \alpha \|x\|^2 \leq \alpha R^2.$$

□

Now, the lower bound $\alpha 4^{b/d} \geq 1$ was proved for any $\mathcal{C} \in \mathbb{B}^2(\alpha, R)$ and hence for any $\mathcal{C} \in \mathbb{B}^1(\alpha, R)$ with sufficiently large R . Since this lower bound is independent of R , it can be associated with the limit class $\mathbb{B}(\alpha)$ under the described practical caveat.

Lastly, we define another class of contractive compression operators which will be used to provide some examples related to the optimality.

Definition 7 (Spherical α -contractive compressions). *Let $\alpha \in (0, 1)$ and $R \geq 1$ be fixed. We denote by $\mathbb{B}^3(\alpha, R)$ the class of (possibly randomized) operators $\mathcal{C}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ such that*

$$\begin{aligned} \mathbb{E} [\|\mathcal{C}(x) - x\|^2] &\leq \alpha, & \text{if } \|x\| = 1 \\ \mathcal{C}(x) &= \|x\| \mathcal{C}(x/\|x\|), & \text{if } 1/R \leq \|x\| \leq R \\ \mathcal{C}(x) &= 0, & \text{if } \|x\| < 1/R \text{ or } \|x\| > R. \end{aligned}$$

The advantage of this class is that any operator $\mathcal{C} \in \mathbb{B}^3(\alpha, R)$ can be uniquely identified by its restriction $\mathcal{C}: \mathbb{S}^d \rightarrow \mathbb{R}^d$ to the unit sphere. To compress a given vector $x \in \mathbb{R}^d$, we compress its projection $x/\|x\| \in \mathbb{S}^d$ by applying \mathcal{C} and we send $\mathcal{C}(x/\|x\|)$ together with the norm $\|x\| \in \mathbb{R}$.

Subsequently, we will concentrate on the compression of unit vectors with as few bits as possible.

Lemma 2. $\mathbb{B}^3(\alpha, R) \subset \mathbb{B}^1(\alpha, R)$.

Proof. Let $\mathcal{C} \in \mathbb{B}^3(\alpha, R)$. If $1/R \leq \|x\| \leq R$ then

$$\begin{aligned} \mathbb{E} [\|\mathcal{C}(x) - x\|^2] &= \mathbb{E} [\|\mathcal{C}(x/\|x\|) - x/\|x\|\|^2] \|x\|^2 \\ &= \mathbb{E} [\|\mathcal{C}(x/\|x\|) - x/\|x\|\|^2] \|x\|^2 \\ &\leq \alpha \|x\|^2. \end{aligned}$$

The other cases are trivial. □

Lemma 3 (Lemma 1 in [26]). *If $\mathcal{C} \in \mathbb{U}(\omega)$, then $\frac{1}{\omega+1}\mathcal{C} \in \mathbb{B}(\frac{\omega}{\omega+1})$.*

B.2 Two senses of optimality for compression: Proof of Proposition 1

If $\mathcal{C} \in \mathbb{B}(\alpha)$ with $\alpha \in [0, 1)$, then to minimize L -smooth and μ -strongly convex function f , CGD needs $\mathcal{O}\left(\frac{1}{1-\alpha}\kappa \log \frac{1}{\epsilon}\right)$ steps for ϵ -accuracy, where $\kappa = \frac{L}{\mu}$ is the condition number of f (see e.g. Theorem 13 of [6]). If we choose to not use compression operator and send uncompressed gradients ($\alpha = 0$) then we get iteration complexity of GD $\mathcal{O}\left(\kappa \log \frac{1}{\epsilon}\right)$, which is $\frac{1}{1-\alpha}$ times smaller than for CGD. If compression operator is unbiased with variance $\alpha \geq 0$, then the iteration complexity becomes $\mathcal{O}\left((1+\alpha)\kappa \log \frac{1}{\epsilon}\right)$ (see e.g. [18]). Alternatively, for an unbiased compression operator $\mathcal{C} \in \mathbb{U}(\alpha)$ one has $\frac{1}{1+\alpha}\mathcal{C} \in \mathbb{B}\left(\frac{\alpha}{1+\alpha}\right)$, which implies the iteration complexity $\mathcal{O}\left(\frac{1}{1-\frac{\alpha}{1+\alpha}}\kappa \log \frac{1}{\epsilon}\right) = \mathcal{O}\left((1+\alpha)\kappa \log \frac{1}{\epsilon}\right)$.

B.3 Dimension-tolerant compression schemes: Proof of Theorem 1

Statement (i) directly follows from (1) and Lemma 3, since $b \geq d \log_4 \frac{1}{\alpha}$ in the biased case and $b \geq d \log_4 \frac{\omega}{\omega+1}$ in the unbiased case.

Statement (ii): first we construct an unbiased compression operator on the unit sphere, which together with $\|x\|$ factor will prove the unbiased case. It follows from [19] (see also [26], Section 3) that one can construct an unbiased compression operator $\mathcal{C}: \mathbb{S}^d \rightarrow \mathbb{R}^d$ with $\omega = \mathcal{O}\left(\frac{d}{\log m/d}\right)$ variance and $\log m$ bits where the dependence of m from d can be up to exponential. Choosing $m = 2^{cd-31}$, we obtain a number of cd bits to encode $\mathcal{C}(x/\|x\|)$ together with $\|x\|$ and variance

$$\omega = \mathcal{O}\left(\frac{d}{\log m - \log d}\right) = \mathcal{O}\left(\frac{d}{cd - \log d - 31}\right) = \mathcal{O}(1/c).$$

For the biased case, Lemma 3 implies that the operator $\frac{1}{\omega+1}\mathcal{C}$ has variance

$$\alpha = 1 - \frac{1}{\omega + 1} = 1 - \frac{1}{\mathcal{O}(1/c) + 1} = \frac{1}{1 + \Omega(c)}$$

and uses the same number cd of bits as \mathcal{C} .

C Proofs for Section 3

C.1 Asymptotic tightness of the lower bound (1): Proof of Theorem 2

First of all, note that to construct a α -contractive compression operator $\mathcal{C}: \mathbb{S}^d \rightarrow \mathbb{R}^d$ on the unit sphere, it is sufficient to cover the unit sphere \mathbb{S}^d by spherical caps generated from balls of radius $\sqrt{\alpha}$. To see this, let $B^d(x^0, \sqrt{\alpha})$ be the ball of radius $\sqrt{\alpha}$ and center $x^0 \in \mathbb{R}^d$ and $C^d(x^0, \sqrt{\alpha}) := B^d(x^0, \sqrt{\alpha}) \cap \mathbb{S}^d$ be the corresponding spherical cap. Then compressing all points $x \in C^d(x^0, \sqrt{\alpha})$ to the center x^0 (i.e. $\mathcal{C}(x) = x^0$) we preserve α -contractive property $\|\mathcal{C}(x) - x\| \leq \alpha$ since $\|\mathcal{C}(x) - x\| = \|x^0 - x\| \leq \sqrt{\alpha}$.

It can be shown that in order to maximize the surface area of $C^d(x^0, \sqrt{\alpha})$, the center x^0 should be on the sphere of radius $\sqrt{1-\alpha}$, namely $\|x^0\| = \sqrt{1-\alpha}$. Based on the formula³ for the surface area of spherical caps, we compute the normalized surface area of $C^d(x^0, \sqrt{\alpha})$ to be $\frac{1}{2}I_\alpha\left(\frac{d-1}{2}, \frac{1}{2}\right)$. Thus, $C^d(x^0, \sqrt{\alpha})$ covers $\frac{1}{2}I_\alpha\left(\frac{d-1}{2}, \frac{1}{2}\right)$ portion of the unit sphere \mathbb{S}^d , where I is the regularized incomplete beta function

$$I_p(a, b) = \frac{B(p; a, b)}{B(a, b)} = \frac{\int_0^p t^{a-1}(1-t)^{b-1} dt}{\int_0^1 t^{a-1}(1-t)^{b-1} dt}, \quad a, b > 0, p \in [0, 1]. \quad (8)$$

Next, we use the following result on covering the sphere with balls:

Theorem 9 (see Theorem 1 in [7]). *For any $d \geq 3$ and $r \in (0, 1)$, the unit sphere \mathbb{S}^d can be covered with balls of radius r in a way that no point of \mathbb{S}^d is covered more than $400 d \ln d$ times.*

Let m be the number of balls of radius $\sqrt{\alpha}$ that cover the whole unit sphere with density at most $400 d \ln d$. This implies that

$$m \frac{1}{2} I_\alpha\left(\frac{d-1}{2}, \frac{1}{2}\right) \leq 400 d \ln d.$$

Now these m balls can be encoded using $b = \lceil \log m \rceil$ bits, so that $m \geq 2^{b-1}$. Therefore

$$2^b I_\alpha\left(\frac{d-1}{2}, \frac{1}{2}\right) \leq 1600 d \ln d. \quad (9)$$

³see https://en.wikipedia.org/wiki/Spherical_cap#Hyperspherical_cap

It remains to lower bound the function I , which we do as follows

$$\begin{aligned} I_\alpha \left(\frac{d-1}{2}, \frac{1}{2} \right) &= \frac{1}{B \left(\frac{d-1}{2}, \frac{1}{2} \right)} \int_0^\alpha t^{\frac{d-3}{2}} (1-t)^{-\frac{1}{2}} dt = \frac{\Gamma \left(\frac{d}{2} \right)}{\Gamma \left(\frac{d-1}{2} \right) \Gamma \left(\frac{1}{2} \right)} \int_0^\alpha t^{\frac{d-3}{2}} (1-t)^{-\frac{1}{2}} dt \\ &\geq \frac{1}{\sqrt{\pi}} \int_0^\alpha t^{\frac{d-3}{2}} dt = \frac{2}{\sqrt{\pi}(d-1)} \alpha^{\frac{d-1}{2}} \geq \frac{\alpha^{d/2}}{d}. \end{aligned}$$

Applying this lower bound to (9) and making some simplifications we get

$$2^b \alpha^{d/2} \leq 1600 d^2 \ln d, \quad (10)$$

which is equivalent to (6). Finally, since $1 \leq \sqrt[d]{d} \rightarrow 1$ as $d \rightarrow \infty$, we can make the right hand side of (6) smaller than $1 + \epsilon$ for any fixed ϵ just by choosing a large d .

C.2 Deterministic-biased version of SD: Proof of Theorem 3

Compression operator and variance bound. To compress a given nonzero vector $x \in \mathbb{R}^d$, we first compress the normalized vector $u = x/\|x\| \in \mathbb{S}^d$ and then rescale it. To quantize the coordinates of unit vector u , we apply dithering with levels $2k_i h$, $k_i \geq 0$, where $h = \sqrt{\nu/d}$ is the half-step and $\nu \geq 0$ is a free parameter of the compression operator. For each $u_i, i \in [d]$ we choose the nearest level so that $\|u_i - 2k_i h\| \leq h$. Letting $\hat{u}_i = \text{sign}(u_i) 2k_i h$ we have $|u_i - \hat{u}_i| \leq h$ for all $i \in [d]$. Therefore

$$\|u - \hat{u}\|^2 = \sum_{i=1}^d (u_i - \hat{u}_i)^2 \leq dh^2 = \nu.$$

Note that, after rescaling with $\|x\|$, this gives a compression with variance at most ν . However, $\|x\|$ is not always the best option. Specifically, we can choose the scaling factor $\gamma > 0$ so as to minimize the variance $\|x - \gamma \hat{u}\|^2$, which yields the optimal factor $\gamma^* = \frac{\langle x, \hat{u} \rangle}{\|\hat{u}\|^2}$ with optimal variance $\|x - \gamma^* \hat{u}\|^2 = \sin^2 \varphi \|x\|^2$, where $\varphi \in [0, \pi/2]$ is the angle between x and \hat{u} . Hence, defining the compression operator as $\mathcal{C}(x) = \gamma^* \hat{u}$, we have the following bound on the variance:

$$\|\mathcal{C}(x) - x\|^2 \leq \min(\nu, \sin^2 \varphi) \|x\|^2,$$

where $\varphi \in [0, \pi/2]$ is the angle between x and $\mathcal{C}(x)$, and in the case $\mathcal{C}(x) = 0$, we let $\varphi = \pi/2$.

Encoding. Now, we describe the encoding scheme itself; that is, how many and which bits we need to communicate for $\gamma^* \hat{u}$. We introduce the following notations:

$$\gamma := 2h\gamma^* \in \mathbb{R}_+, \quad s := (\text{sign}(u_i k_i))_{i=1}^d \in \{-1, 0, 1\}^d, \quad k := (k_i)_{i=1}^d \in \mathbb{N}_+^d.$$

Note that $\mathcal{C}(x) = \gamma^* \hat{u} = 2h\gamma^* \text{sign}(u) k = \gamma s k$. So, we need to encode the triple (γ, s, k) . Since $\gamma \in \mathbb{R}_+$, we need only $\lceil 31 \rceil$ bits for the scaling factor. Next we encode s . Let

$$n_0 := \#\{i \in [d] : s_i = 0\} = \#\{i \in [d] : k_i = 0\}$$

be the number of coordinates u_i that are compressed to 0. To communicate s , we first send the locations of those n_0 coordinates and then $\lceil d - n_0 \rceil$ bits for the values ± 1 . Sending n_0 positions can be done by sending $\lceil \log d \rceil$ bits⁴ representing the number n_0 , afterwards sending $\lceil \log \binom{d}{n_0} \rceil$ bits for the positions. Finally, it remains to encode k for which we only need to send nonzero entries since the positions of $k_i = 0$ are already encoded. We encode $k_i \geq 1$ with k_i bits: $k_i - 1$ ones followed by 0. Hence, encoding k required $\lceil \sum k_i \rceil$ additional bits.

Thus, our encoding scheme for $\mathcal{C}(x) = \gamma s k$ is as follows

⁴We can further optimize this with Elias- ω encoding by sending $\approx \log n_0$ bits instead of $\log d$. However, both are negligible in the overall encoding and we will not complicate the analysis for this small improvement.

- scaling factor γ : 31 bits,
- signs s : $\log d + \log \binom{d}{n_0} + d - n_0$ bits,
- dithering levels k : $\sum_{i=1}^d k_i$ bits,
- total number of bits $b = 31 + \log d + \log \binom{d}{n_0} + d - n_0 + \sum_{i=1}^d k_i$.

Upper bound on b . We continue by giving a theoretical upper bound for the bits b needed to communicate $\mathcal{C}(x)$. Below, we derive an upper bound for $\sum k_i$. Since each $|u_i|$ is quantized to the nearest $2k_i h$, then we have this double bound $\max(0, (2k_i - 1)h) \leq |u_i| \leq (2k_i + 1)h$. Using this with the Cauchy–Schwarz inequality we get

$$1 = \left(\sum_{i=1}^d u_i^2 \right)^{1/2} \geq h \left(\sum_{k_i \geq 1} (2k_i - 1)^2 \right)^{1/2} \geq \frac{h}{\sqrt{d - n_0}} \sum_{k_i \geq 1} (2k_i - 1),$$

which implies the following bound on $\sum k_i$:

$$\sum_{i=1}^d k_i = \sum_{k_i \geq 1} k_i \leq \frac{1}{2} \left(\frac{\sqrt{d - n_0}}{h} + d - n_0 \right) = \frac{d}{2} \left(\sqrt{\frac{1 - n_0/d}{\nu}} + 1 - n_0/d \right).$$

Setting $\tau = n_0/d \in [0, 1]$, we further upper bound it using the AM-GM inequality

$$\sum_{i=1}^d k_i \leq \frac{d}{2} \left(\frac{1 - \tau/2}{\sqrt{\nu}} + 1 - \tau \right).$$

Let us consider the extreme cases $n_0 = 0$ and $n_0 = d$ separately. If $n_0 = d$, then $b = 31 + \log d$. If $n_0 = 0$, then $b \leq 31 + \log d + \left(\frac{3}{2} + \frac{1}{2\sqrt{\nu}} \right) d$. Note that these extreme cases are the best cases in terms of the number of bits. In the sequel, we assume that $1 \leq n_0 \leq d - 1$ and hence $\tau \in [1/d, 1 - 1/d]$. Next, we upper bound the term $\log \binom{d}{n_0}$, for which it is known the following tight estimate⁵

$$\frac{2^{dH_2(\tau)}}{\sqrt{8d\tau(1-\tau)}} \leq \binom{d}{\tau d} \leq \frac{2^{dH_2(\tau)}}{\sqrt{2\pi d\tau(1-\tau)}}, \quad 0 < \tau < 1,$$

where $H_2(\tau) = -\tau \log \tau - (1 - \tau) \log(1 - \tau)$ is the binary entropy function in bits. Hence

$$\log \binom{d}{n_0} = \log \binom{d}{\tau d} \leq -\frac{1}{2} \log(2\pi d\tau(1-\tau)) + dH_2(\tau).$$

The first term with negative sign saves at least $\frac{1}{2} \log 2\pi \approx 1.32$ bits and up to $\frac{1}{2} \log \frac{\pi d}{2}$ bits. In further estimations we upper bound it by -1 . So far, the following upper bound is obtained

$$\begin{aligned} b &\leq 30 + \log d + dH_2(\tau) + d(1 - \tau) + \frac{d}{2} \left(\frac{1 - \tau/2}{\sqrt{\nu}} + 1 - \tau \right) \\ &= 30 + \log d + \left(H_2(\tau) + \frac{3}{2}(1 - \tau) + \frac{1 - \tau/2}{2\sqrt{\nu}} \right) d \\ &:= 30 + \log d + \beta(\tau, \nu)d. \end{aligned}$$

⁵“The Theory of Error-Correcting Codes” by MacWilliams and Sloane (Chapter 10, Lemma 7, p. 309)

It remains to find an upper bound for $\beta(\tau, \nu)$ with respect to τ . As the entropy function H_2 and any linear function are concave, we can find the maximum by solving first order optimality condition. The equation $\frac{d}{d\tau}\beta(\tau, \nu) = 0$ gives the solution

$$\tau^* = \frac{1}{1 + 2^{\frac{1}{4}\left(6 + \frac{1}{\sqrt{\nu}}\right)}}.$$

Setting $\beta(\nu) := \beta(\tau^*, \nu)$, we upper bound the number of bits b as

$$b \leq 30 + \log d + \beta(\nu)d.$$

It can be shown that, with $\nu = 1/10$, one has $\beta(\nu) \approx 3.3495 < 3.35$. This completes the proof of Theorem 3.

C.3 Tighter bounds on minimal communication: Proofs of Theorems 4 and 5

The first motivation for this is that even though the uncertainty principle (1) is strong for constant α , it is not tight when α goes to 1 as d goes to infinity. In particular, for $\alpha = \frac{d-1}{d}$, it says that the number of bits is at least $d \log(1 - 1/d)/2$, which is constant. However, we can show that when $\alpha < 1$, one needs at least $\log(d)$ bits. This explains why there is no way to only communicate a fixed number of bits per round while still having $\alpha < 1$. Moreover, we will compute an explicit estimate of $b^*(\alpha, d)$, as a function of d and α only, with a very low error of $\frac{1}{2} \log \log d + C$ for some absolute constant C .

Proof of Theorem 4. Proving the result is equivalent to proving that the surface of the unit sphere cannot be covered by less than d smaller, identical balls. We can prove this easily by induction. To make the induction step, let us assume, without loss of generality, that one of the smaller balls is centered on the positive x_1 axis. Since the radius of this smaller ball is less than 1, the unit $(d-1)$ -dimensional sphere with $x_1 = 0$ is disjoint from the first smaller ball, which means, by induction, that it will itself require at least $d-1$ additional smaller balls, leading to the desired result. \square

In fact, the previous result can be used to obtain the following result.

Definition 8. For a covering of the surface of the unit sphere using identical spherical caps, we define the density of the cover to be the average number of caps covering a point on the surface of the unit sphere. Identically, this is equal to the number of spherical caps used multiplied by the fraction of the unit sphere covered by a single spherical cap.

Theorem 10. There exists an absolute constant B such that if the surface of the unit sphere is covered with identical smaller spherical caps, then the density of the covering is at least Bd .

Proof. We split this into two cases. The first case is when the radius of the spherical cap is larger than $\sqrt{1 - \frac{1}{d}}$. In this case, each spherical cover will cover at least a fraction of $A(d) = \text{Prob}\left(x_1 \geq \frac{1}{\sqrt{d}}\right)$ where x is chosen uniformly from the surface of the unit sphere. One can easily show that there exists $C_1 > 0$ such that $A(d) > C_1$ for all d . Indeed, it is enough to see that $A(d) > 0$ for all d and that $A(d)$ approaches $1 - \phi(1)$ where ϕ is the CDF of a standard normal. Combining this bound with the previous result of requiring at least d caps to cover the surface of the unit sphere, the density is at least C_1d when the radius of the spherical cap is at least $\sqrt{1 - \frac{1}{d}}$.

In the second case, when the radius of the spherical cap is less than $\sqrt{1 - \frac{1}{d}} < \sqrt{1 - \frac{1}{d+1}}$, The Coxeter-Few-Rogers ‘‘simplex’’ bound shows [11] that the density is at least C_2d for some absolute constant C_2 . Choosing $B = \min(C_1, C_2)$ gives us the desired result. \square

Theorem 11 (see [11]). There exists an absolute constant A such that for any d and any spherical radius $r < 1$, there exists a cover for the surface of the unit sphere with smaller, identical spherical caps of radius r such that the density of the covering is at most $Ad \log d$.

Lemma 4. *If $\|x\| = 1$ and $\|v - x\|^2 \leq \alpha$ for some t , then for $v' = \frac{\sqrt{1-\alpha}v}{\|v\|}$, one has $\|x - v'\|^2 \leq \alpha$. In other words, if some balls of radius $\sqrt{\alpha}$ cover the surface of the unit sphere, then projecting them onto the sphere of radius $\sqrt{1-\alpha}$ will still cover entirely the surface of the unit sphere.*

Proof. The initial condition is equivalent to $1 - \alpha + \|v\|^2 \leq 2\langle v, x \rangle$ which, using AM-GM, implies that $2\sqrt{1-\alpha}\|v\| \leq 2\langle v, x \rangle$, which can be rearranged to look like $1 - \alpha + \|v'\|^2 \leq 2\langle v', x \rangle$ or $\|v' - x\|^2 \leq \alpha$, as desired. \square

The above discussion leads us to our next result on $b^*(\alpha, d)$, which is an important quantity to study. First, it tells us that operators $\mathcal{C} \in \mathbb{B}(\alpha)$ cannot be compressed with less than b^* bits. Moreover, it tells us that this bound is tight, because there is at least one operator in $\mathbb{B}(\alpha)$ that can be compressed to no more than $b^*(\alpha, d)$ bits. Thus, we proceed with estimating $b^*(\alpha, d)$ explicitly with a very small estimation error of $\frac{1}{2} \log \log d + \mathcal{O}(1)$.

Proof of Theorem 5. Recall that $P(\alpha, d) = \frac{1}{2}I_\alpha(\frac{d-1}{2}, \frac{1}{2})$ is also equal to the fraction of the surface area of the surface of the unit sphere with $\sqrt{1-\alpha} \leq x_1$. This can be viewed as the probability that a point x chosen uniformly on the unit sphere satisfies $\sqrt{1-\alpha} \leq x_1$. In order to prove the theorem, we will prove the upper and lower bounds on $b^*(\alpha, d)$ separately.

We first prove the lower bound. Let \mathcal{C} be an arbitrary operator in $\mathbb{B}(\alpha)$ that can be encoded with no more than b bits. This means that at most 2^b possible values can be communicated. Let c_1, \dots, c_V be all the possible decodings, with $V \leq 2^b$. Now, if we consider the balls $B(c_i, \sqrt{\alpha})$, the surface of the unit sphere must be covered. Indeed, if $\|x\| = 1$ and x is not covered, then all the possible encodings of $\mathcal{C}(x)$ will produce a point whose distance from x is more than $\sqrt{\alpha}$, which contradicts the fact that the operator is in $\mathbb{B}(\alpha)$.

Now, since these small balls cover the surface of the unit sphere, one can use Lemma 4 to show that the balls centered at $B(C_i, \sqrt{\alpha})$, where $C_i = \frac{\sqrt{1-\alpha}c_i}{\|c_i\|}$, should also be a covering. Using Theorem 10, we know that the density of this new coverage is at least Bd , while it is at most $F2^b$, where F is the fraction of the surface area of the unit sphere that each one of these balls cover. In fact, one can compute F explicitly as $P(\alpha, d) = \text{Prob}(x_1 \geq \sqrt{1-\alpha})$ where x is chosen uniformly on the surface of the unit sphere. This gives us the lower bound $b \geq -\log P(\alpha, d) + \log d + \log B$.

For the upper bound, one can use a constant number of bits to communicate $\|x\|$, then one can use the covering from Theorem 11 with radius equal to $\sqrt{\alpha}$ and quantize x to the nearest spherical cap center, which is guaranteed to be within a distance of $\sqrt{\alpha}$, ensuring that this quantization is in $\mathbb{B}(\alpha)$. Now, since this covering has density no more than $Ad \log d$, and since its density is equal to $VP(\alpha, d)$, where v is the number of spherical caps used and P is, as defined above, the fraction of the surface area covered by one a spherical cap of radius $\sqrt{\alpha}$, one can conclude that $v \leq \frac{Ad \log d}{P(\alpha, d)}$, which means that the centers can be encoded using no more than $-\log P(\alpha, d) \log d + \log \log d + \log A$ bits, yielding the desired upper bound. \square

D Proofs for Section 4

D.1 Lower bound on average communication: Proof of Theorem 6

Let X be a random vector with uniform distribution over the unit sphere \mathbb{S}^d and $\hat{X} = \mathcal{C}(X)$ be the compressed (random) vector. Note that, \hat{X} has two source of randomness, one from the random vector X and the other coming from the compression operator \mathcal{C} . Based on the assumption of finiteness of B (otherwise the lower bound is trivial), we conclude that \mathcal{C} , and hence the random vector \hat{X} , are discrete; that is, the set of possible values they can take is finite or countably infinite. Note that \hat{X} can be encoded with B bits in expectation with respect to its own source of randomness, as

$$B = \sup_{\|x\|=1} \mathbb{E}_{\mathcal{C}} [|\mathcal{C}(x)|] \geq \mathbb{E}_{\mathcal{C}, X} [|\mathcal{C}(X)|] = \mathbb{E}_{\hat{X}} [|\mathcal{C}(X)|].$$

Thus, the discrete random source \hat{X} admits an encoding with expected binary description length of B . Applying Shannon's source coding theorem on lossless compression⁶, we get $B \geq H(\hat{X})$, where H indicates the entropy of the source⁷ in bits.

Next, using the mutual information and relative entropy of \hat{X} and X , we further lower bound it as follows:

$$B \geq H(\hat{X}) \geq H(\hat{X}) - H(\hat{X}|X) = I(\hat{X}, X) = H(X) - H(X|\hat{X}).$$

Now, we denote by A the surface area of the unit sphere \mathbb{S}^d . For a given point $v \in \mathbb{R}^d$, let $A'(v)$ be the surface area of the cap $C^d(v, \sqrt{\alpha}) = B^d(v, \sqrt{\alpha}) \cap \mathbb{S}^d$. Using Lemma 4, it can be shown that in order to maximize the surface area of $C^d(v, \sqrt{\alpha})$, the center v should be on the sphere of radius $\sqrt{1-\alpha}$, namely $\|v\|_2 = \sqrt{1-\alpha}$. Using the formula⁸ for the surface area of spherical caps, we compute the normalized surface area of $C^d(v, \sqrt{\alpha})$ to be $P(\alpha, d) = \frac{1}{2}I_\alpha(\frac{d-1}{2}, \frac{1}{2})$. Thus, at best $C^d(v, \sqrt{\alpha})$ covers the portion $P(\alpha, d)$ of the unit sphere \mathbb{S}^d , where I_α is the regularized incomplete beta function. Therefore, for an arbitrary $v \in \mathbb{R}^d$, one has the upper bound $A'(v) \leq P(\alpha, d)A$. Note that as I_α is upper bounded by 1 (which directly follows from the definition) we get $P(\alpha, d) < 1/2$.

Since X is uniform on the unit sphere, its probability density function is $1/A$ and so the entropy $H(X) = \log(A)$. Similarly, since the random vector X conditioned with $\hat{X} = v$ is uniform over $C^d(v, \sqrt{\alpha})$, we have $H(X|\hat{X} = v) = \log A'(v)$. Hence

$$H(X|\hat{X}) = \mathbb{E}_{\hat{X}} \left[H(X|Y = \hat{X}) \right] = \mathbb{E}_{\hat{X}} \left[\log A'(\hat{X}) \right] \leq \log(P(\alpha, d)A),$$

resulting in the desired lower bound

$$B \geq H(X) - H(X|\hat{X}) \geq \log(A) - \log(P(\alpha, d)A) = -\log P(\alpha, d).$$

D.2 Randomized-unbiased version of Sparse Dithering: Proof of Theorem 7

In this section, we randomize Sparse Dithering to make it unbiased.

Compression operator and variance bound. Again, to compress a given nonzero vector $x \in \mathbb{R}^d$, we decompose x into the scalar $\|x\|$ and the unit vector $u = x/\|x\|$. To quantize the coordinates of u , we round to one of the two nearest neighbors, so as to preserve unbiasedness; that is, if $2k_i h \leq |u_i| \leq 2(k_i + 1)h$ for some $k_i \geq 0$, then

$$\hat{u}_i = \text{sign}(u_i) 2\hat{k}_i h = \begin{cases} \text{sign}(u_i) 2k_i h & \text{with probability } \frac{2(k_i+1)h - |u_i|}{2h} \\ \text{sign}(u_i) 2(k_i + 1)h & \text{with probability } \frac{|u_i| - 2k_i h}{2h} \end{cases}$$

Clearly, $\mathbb{E}[\hat{u}] = u$ and defining $\mathcal{C}(x) = \|x\|\hat{u}$ we maintain unbiasedness $\mathbb{E}[\mathcal{C}(x)] = x$. Bounding the second moment

$$\begin{aligned} \mathbb{E}[\hat{u}_i^2] &= (2k_i h)^2 \frac{2(k_i + 1)h - |u_i|}{2h} + (2(k_i + 1)h)^2 \frac{|u_i| - 2k_i h}{2h} \\ &= u_i^2 + (|u_i| - 2k_i h)(2(k_i + 1)h - |u_i|) \\ &\leq u_i^2 + \left(\frac{|u_i| - 2k_i h + 2(k_i + 1)h - |u_i|}{2} \right)^2 = u_i^2 + h^2, \end{aligned}$$

we conclude that

$$\frac{\mathbb{E}[\|\mathcal{C}(x)\|^2]}{\|x\|^2} = \mathbb{E}[\|\hat{u}\|^2] \leq \sum_{i=1}^d (u_i^2 + h^2) \leq 1 + dh^2 = 1 + \nu.$$

⁶see e.g. Theorem 5.5.1+Corollary or Theorem 5.11.1 of [9]

⁷A discrete random vector can be mapped to a discrete random variable preserving the same probability distribution (and so we can extend the source-coding inequality), as entropy is defined through probability mass/density function, not the actual values of the random source.

⁸see https://en.wikipedia.org/wiki/Spherical_cap#Hyperspherical_cap

Table 2: Total communication savings due to unbiased compression method.

COMPRESSION METHOD	Bits $\mathbb{E}[b]$	$\times(1 + \omega)$	$\beta := \mathbb{E}[b]/32d$	SAVINGS $\times(1 + \omega)\beta$
NO COMPRESSION (BASE)	$32d$	1	1	1
RANDOM SPARSIFICATION	$32k + \log_2 \binom{d}{k}$	d/k	$> k/d$	> 1
TERNARY QUANTIZATION	$\approx d \log_2 3$	\sqrt{d}	$1/20.2$ (0.05)	$\sqrt{d}/20.2$
STANDARD DITHERING	$\approx 2.8d$	2	$1/11.4$ (0.087)	$1/5.7$ (0.175)
NATURAL COMPRESSION	$9d$	$9/8$ (1.125)	$1/3.5$ (0.28)	$1/3.1$ (0.31)
Randomized SD (new)	$\approx 2.6d$	$5/4$ (1.25)	$1/12.3$ (0.081)	$1/9.9$ (0.10)

Hence, the variance of compression operator \mathcal{C} is $\omega \leq \nu$.

Encoding. Next, we proceed to the encoding scheme, exactly like in the deterministic case. We introduce the following notations:

$$\gamma := 2h\|x\| \in \mathbb{R}_+, \quad s := \left(\text{sign}(u_i \hat{k}_i) \right)_{i=1}^d \in \{-1, 0, 1\}^d, \quad \hat{k} := (\hat{k}_i)_{i=1}^d \in \mathbb{N}_+^d.$$

Note that $\mathcal{C}(x) = \|x\|\hat{u} = 2h\|x\| \text{sign}(u) \hat{k} = \gamma s \hat{k}$. So, we need to encode the triple (γ, s, \hat{k}) . The encoding scheme and the formula for the number of bits are the same, so we need to upper bound

$$\hat{b} = 31 + \log d + \log \binom{d}{\hat{n}_0} + d - \hat{n}_0 + \sum_{i=1}^d \hat{k}_i$$

in expectation, where $\hat{n}_0 := \#\{i \in [d]: \hat{k}_i = 0\}$.

Upper bound on $\mathbb{E}[\hat{b}]$. First, notice that

$$\mathbb{E} \left[\sum_{i=1}^d \hat{k}_i \right] = \frac{1}{2h} \sum_{i=1}^d \mathbb{E} [\hat{u}_i] = \frac{\|u\|_1}{2h} \leq \frac{\sqrt{d}}{2h} = \frac{d}{2\sqrt{\nu}}.$$

Considering $\hat{n}_0 = 0$ and $\hat{n}_0 = d$ cases separately, we get $\mathbb{E}[\hat{b}] \leq 31 + \log d + (1 + 1/2\sqrt{\nu})d$ and $\hat{b} = 31 + \log d$ respectively. Next, we use the same upper bound for the log-term $\log \binom{d}{\hat{n}_0} \leq dH_2(\hat{\tau}) - 1$ with $\hat{\tau} = \hat{n}_0/d \in [1/d, 1 - 1/d]$. It remains to upper bound $H_2(\hat{\tau}) + (1 - \hat{\tau})$, which is maximized when $\hat{\tau} = 1/3$ with value $\log 3$, i.e. $H_2(\hat{\tau}) + (1 - \hat{\tau}) \leq \log 3$. Thus, we have proved the formula for the number bits in expectation:

$$\mathbb{E} [\hat{b}] \leq 31 + \log d + (dH_2(\hat{\tau}) - 1) + d(1 - \hat{\tau}) + \frac{d}{2\sqrt{\nu}} \leq 30 + \log d + \left(\log 3 + \frac{1}{2\sqrt{\nu}} \right) d.$$

The parameter $\nu = 1/4$ is approximately the maximizer for

$$\frac{32d}{(1 + \omega)\mathbb{E}[\hat{b}]} = \frac{32}{(1 + \nu) \left(\log 3 + \frac{1}{2\sqrt{\nu}} \right)} \approx 9.9,$$

which shows the gain in total communication complexity. In other words, the scheme communicates $(1 + \log 3) d \approx 2.6d$ bits in each iteration (about 12 times less than without compression), but needs $1 + \omega = 5/4$ times more iterations.