# Collaborative Compressors in Distributed Mean Estimation with Limited Communication Budget

Harsh Vardhan hharshvardhan@ucsd.edu

Department of Computer Science and Engineering University of California, San Diego

Arya Mazumdar arya@ucsd.edu

Halicioğlu Data Science Institute University of California, San Diego

Reviewed on OpenReview: https://openreview.net/forum?id=AtCKHCoMA7

### **Abstract**

Distributed high dimensional mean estimation is a common aggregation routine used often in distributed optimization methods. Most of these applications call for a communication-constrained setting where vectors, whose mean is to be estimated, have to be compressed before sharing. One could independently encode and decode these to achieve compression, but that overlooks the fact that these vectors are often close to each other. To exploit these similarities, recently Suresh et al., 2022, Jhunjhunwala et al., 2021, Jiang et al, 2023, proposed multiple correlation-aware compression schemes. However, in most cases, the correlations have to be known for these schemes to work. Moreover, a theoretical analysis of graceful degradation of these correlation-aware compression schemes with increasing dissimilarity is limited to only the  $\ell_2$ -error in the literature. In this paper, we propose four different collaborative compression schemes that agnostically exploit the similarities among vectors in a distributed setting. Our schemes are all simple to implement and computationally efficient, while resulting in big savings in communication. The analysis of our proposed schemes show how the  $\ell_2$ ,  $\ell_\infty$  and cosine estimation error varies with the degree of similarity among vectors.

### 1 Introduction

We study the problem of estimating the empirical mean, or average, of a set of high-dimensional vectors in a communication-constrained setup. We assume a distributed problem setting, where m clients, each with a vector  $g_i \in \mathbb{R}^d$ , are connected to a single server (see, Fig. 1a). Our goal is to estimate their mean g on the server, where

$$g \triangleq \frac{1}{m} \sum_{i \in [m]} g_i. \tag{1}$$

We use [m] to denote the set  $\{1, 2, ..., m\}$ . The clients can communicate with the server via a communication channel which allows limited communication. The server does not have access to data but has relatively more computational power than individual clients.

This problem, referred to as distributed mean estimation (DME), is an important subroutine in several distributed learning applications. One common application is distributed training or federated learning (McMahan et al., 2016; McMahan & Ramage, 2017), when different clients correspond to different edge devices.

The typical learning task for DME is supervised learning via gradient-based methods (Bottou & Bousquet, 2007; Robbins & Monro, 1951). In this case, the vectors  $g_i$  correspond to the gradient updates for each

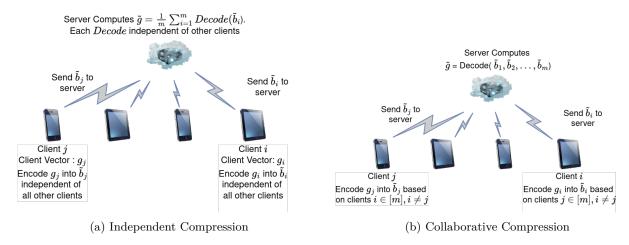


Figure 1: Compression for Distributed Mean Estimation

client i computed on its local training data, and g is the average gradient over all clients. On the other hand, distributed mean estimation is also used in unsupervised learning problems such as distributed KMeans (Liang et al., 2013), distributed PCA (Liang et al., 2014) or distributed power iteration (Li et al., 2021). In distributed KMeans and distributed power iteration,  $g_i$  corresponds to estimates of the cluster center and the top eigenvector, respectively, on the  $i^{th}$  client.

The naive strategy of clients sending their vectors  $g_i$  to the server for DME incurs no error; however, it has a high communication cost, rendering it burdensome in most of the real-world network applications. A principled way to tackle this is to use compression: each client  $i \in [m]$  compresses its vector  $g_i$  into an efficient encoding  $\tilde{b}_i \in \mathcal{B}_i$  which can then be sent to the server; The server forms an estimate  $\tilde{g}$  of the mean g using the encodings  $\{\tilde{b}_i\}_{i \in [m]}$ . We can then compute the error of the estimate  $\tilde{g}$  and the number of bits required to communicate  $\tilde{b}_i$  (i.e.,  $\log_2 |\mathcal{B}_i|$ ) to analyze the efficiency of the compression scheme. As opposed to distributed statistical inference (Braverman et al., 2016; Garg et al., 2014), we do not assume that  $g_i$  are sampled from a distribution, and instead the estimation error of these schemes is computed in terms of  $g_i$ .

One way to approach this compression paradigm is when each client compresses its vector oblivious to others, and the server separately decodes the vectors before aggregating (Figure 1a). We call this *independent compression* and several existing works (Konečný & Richtárik, 2018; Suresh et al., 2017; Safaryan et al., 2021; Gandikota et al., 2022; Vargaftik et al., 2021; Ghosh et al., 2021) use such a compression scheme. The simplest example of this scheme is RandK (Konečný & Richtárik, 2018), where each client sends only  $K \in \mathbb{N}$  coordinates as  $\tilde{b}_i$ , and the server estimates  $\tilde{g}$  as the average of K-sparse vectors from each client. As K < d, this scheme requires less communication than sending the full vector  $g_i$  from each client  $i \in [m]$ .

However, independent compressors suffer from a significant drawback, especially when the vectors to be aggregated are similar/not-too-far, which is often the case for gradient aggregation in distributed learning. Consider the case when two distinct clients  $i, j \in [m]$  have different vectors  $g_i \neq g_j$ , but they differ in only one coordinate. Then, independent compressors like RandK will end up sending  $\tilde{b}_i$  and  $\tilde{b}_j$  which are very similar (in fact, same with high probability) to each other, and therefore wasting communication. Collaborative compressors (Suresh et al., 2022; Szlendak et al., 2021; Jhunjhunwala et al., 2021; Jiang et al., 2023) can alleviate this problem. Figure 1b describes a collaborative compressor, where the encodings  $\{\tilde{g}_i\}_{i\in[m]}$  may not be independent of each other and a decoding function jointly decodes all encodings to obtain the mean estimate  $\tilde{g}$ . Clearly, this opens up more possibilities to reduce communication - but also the error of collaborative compressors can be made to scale as the variability of the vectors.

The amount of required communication also depends on the metric for estimation error. Among the existing schemes for collaborative compressors, most provide guarantees on the  $\ell_2$  error  $||\tilde{g} - g||_2^2$  (Suresh et al., 2022; Szlendak et al., 2021; Jhunjhunwala et al., 2021; Jiang et al., 2023). Also, in collaborative compressors, the error must ideally be dependent on *some measure of correlation/distance* among the vectors, which is indeed

Compressor	Error metric	Error	# Bits/client
NoisySign (Algorithm 1)	$  \tilde{g} - g  _{\infty}$	$\exp(\frac{  g  _{\infty}^2}{\sigma^2}) \cdot (\Delta_{\phi} + \sqrt{\frac{\log m}{m}})$	d
HadamardMultiDim (Algorithm 3)	$\mathbb{E}[  \tilde{g} - g  _{\infty}]$	$rac{B}{2^{m-1}} + \Delta_{ ext{Hadamard}}$	d
SparseReg (Algorithm 4)	$\mathbb{E}[  \tilde{g} - g  _2^2]$	$B^2 \exp\left(-\frac{2m\log L}{d}\right) + \Delta_{\text{reg}}$	$\log L$ $(L \ge 1 \text{ tunable})$
OneBit (Algorithm 5)	$\arccos{\langle \tilde{g}, g \rangle}$	$\pi(\Delta_{\mathrm{corr}} + \frac{d}{mt})$	$t $ ( $t \ge 1 \text{ tunable}$ )

Table 1: Theoretical results for our proposed collaborative compression schemes.  $\Delta_{\Phi}$ ,  $\Delta_{\text{Hadamard}}$ ,  $\Delta_{\text{reg}}$  and  $\Delta_{\text{corr}}$  are measures of average dissimilarity between vectors  $\{g_i\}_{i\in[m]}$  defined in Theorems 1, 2, 3 and Lemma 1 respectively. For NoisySign,  $\sigma>0$  is an algorithm parameter. For HadamardMultiDim, we assume  $||g_i||_{\infty}\leq B, \forall i\in[m]$ . For SparseReg, we assume  $||g_i||_2\leq B, \forall i\in[m]$  and L is an algorithm parameter. For OneBit, g is the unit vector along the average  $\frac{1}{m}\sum_{i=1}^m g_i$  and  $\tilde{g}$  is also a unit vector.

the case for all of these schemes. In this paper, the measure of such a distance is denoted with  $\Delta$ , with some subscript signifying the exact measure; the vectors in question have high similarity as  $\Delta \to 0$ . The estimation error naturally grows with the dimension d, and decays with the number of clients m (due to an averaging). One of our major contributions is to design a compression scheme that has significantly improved dependence on the number of clients m to counter the effect of growing dimension d.

If one were to estimate the unit vector in the direction of the average vector  $\frac{1}{m}\sum_{i=1}^{m}g_{i}$ , which is often important for gradient descent applications, using an estimate of the mean with low  $\ell_{2}$  error can be highly sub-optimal as the  $\ell_{2}$  error might be large even if all the vectors point in the same direction but have different norms. For this, the cosine distance  $\arccos(\frac{\langle \bar{g},g\rangle}{\|\bar{g}\|\|g\|})$  is a better measure, which has not been studied in the literature. We also give a compression scheme specifically tailored for this error metric. Another interesting metric is the  $\ell_{\infty}$ -error which has also not been studied except for in (Suresh et al., 2022). There as well, we give an improved dependence of the estimation error on m.

Further drawback of existing collaborative compressors such as, (Jhunjhunwala et al., 2021; Jiang et al., 2023) is that they require the knowledge of correlation between vectors before employing their compression. Without this knowledge, their error guarantees do not hold.

**Notation.** We use [a, b] to denote  $[a, a+1, a+2, \ldots, b]$  for  $a, b \in \mathbb{N}$ . We use  $g^{(j)}, j \in [d]$  to denote the  $j^{th}$  coordinate of a vector  $g \in \mathbb{R}^d$ . For a permutation  $\rho$  on [m],  $\rho^{(i)}$  denotes mapping of  $i \in [m]$  under  $\rho$ . We use  $\Pi_m$  to denote the uniform distribution over all permutations of [m]. We use  $\exp(a)$  to denote  $e^{\mathcal{O}(a)}$  for any  $a \in \mathbb{R}$ . For any matrix  $A \in \mathbb{R}^{m \times n}$ ,  $A_k$  for  $k \in [m]$  refers to the  $k^{th}$  row of A.

Our contributions. We provide four different collaborative compressors, which are communication-efficient, give error guarantees for different error metrics ( $\ell_2$  error,  $\ell_\infty$  error and cosine distance), and exhibit optimal dependence on the number of clients m and the diameter of ambient space B. To see the advantage of collaboration, we define few natural similarity metrics. All our schemes show graceful degradation of error with the similarity metric between different clients. Our schemes have three subroutines: Init which corresponds to initialization/setting up a protocol, Encode which is performed individually at each client to obtain their encoding  $\tilde{b}_i$ , and Decode which is performed at the server on all the encodings to obtain estimate of mean  $\tilde{g}$ .

Below our main contributions are summarized. The theoretical guarantees for our algorithms are provided in Table 1. For each of our proposed compressors, increasing m increases the total communication required (which is a product of m and the # Bits/client), but reduces the estimation error in the appropriate error metric up to a constant  $\Delta$  that depends on the deviation between client vectors.

1. We provide a simple collaborative scheme based on the popular signSGD (Bernstein et al., 2018a) scheme, NoisySign (Algorithm 1), where sign of each coordinate of a vector is sent after adding Gaussian noise. An advantage of this scheme, compared to others is that we can infer the vector g with an  $\ell_{\infty}$  error guarantee increasing with  $||g||_{\infty}$  and decreasing with m, without the knowledge of  $||g||_{\infty}$  itself. The dissimilarity is

Compressor	Error	# Bits/client	Notes
RandK (Konečný & Richtárik, 2018)	$\mathcal{O}(rac{d}{K} ilde{B}^2)$	$32K + K \log d$	Independent
SRQ (Suresh et al., 2017)	$\mathcal{O}(rac{\log d}{m(K-1)^2} ilde{B}^2)$	Kd	Independent
Kashin (Safaryan et al., 2021)	$\mathcal{O}(\left(\frac{10\sqrt{\lambda}}{\sqrt{\lambda}-1)}\right)^4 \tilde{B}^2)$	$31 + \lambda d$	Independent
Drive (Vargaftik et al., 2021)	$\mathcal{O}( ilde{B}^2)$	32+d	Independent
PermK (Szlendak et al., 2021)	$\mathcal{O}((1-\max\{0,\frac{m-d}{m-1}\})\Delta_2)$	$32K + K \log d$	Collaborative
RandKSpatial (Jhunjhunwala et al., 2021)	$\mathcal{O}(rac{d}{mK}\Delta_2)$	$32K + K \log d$	Needs Correlation
RandKSpatialProj (Jiang et al., 2023)	$\mathcal{O}(\frac{d}{mK}\Delta_2)$	$32K + K \log d$	Needs Correlation
Correlated SRQ (Suresh et al., 2022)	$\mathcal{O}\left(\frac{1}{m}\min\left\{\frac{\sqrt{d}\Delta_{\infty}^{d}B}{K},\frac{dB^{2}}{K^{2}}\right\}\right)$	$2d\log K + K\log d$	$  g_i  _2 \le B, \forall i \in [m]$

Table 2: Comparison of existing independent and collaborative compressors in terms of  $\ell_2$  error and bits communicated. K is the number of coordinates communicated for sparsification methods (RandK, PermK, RandKSpatial, RandKSpatialProj) and the number of quantization levels for quantization methods (SRQ, vqSGD, Correlated SRQ). The constant  $\lambda$  is a parameter of the Kashin scheme. Further,  $\tilde{B}^2 = \frac{1}{m} \sum_{i=1}^m ||g_i||_2^2$ ,  $\Delta_2 = \frac{1}{m} \sum_{i=1}^m ||g_i - g||_2^2$ , and  $\Delta_{\infty} = \max_{j \in [d]} \frac{1}{m} \sum_{i=1}^m |g_i^{(j)} - g^{(j)}|$ . It is also assumed that a real is equivalent to 32 bits, which is an informal norm in this literature.

 $\Delta_{\Phi} = \mathcal{O}(\frac{1}{m\sigma}\sum_{i=1}^{m}||g-g_i||_{\infty})$ , where  $\sigma$  is the variance of the noise added (Theorem 1). This scheme is described in Section 2.

- 2. ( $\ell_{\infty}$ -guarantee) For vectors with  $\ell_{\infty}$  norm bounded by B, we propose a collaborative compression scheme, HadamardMultiDim (Algorithm 3) which performs coordinate-wise collaborative binary search. We obtain the best dependence on m and B for the  $\ell_{\infty}$  error ( $\mathcal{O}(B \cdot \exp(-m))$ ) while suffering from an extra error term  $\Delta_{\text{Hadamard}}$ , which is a measure of average dissimilarity between compressed vectors.  $\Delta_{\text{Hadamard}}$  lies in the range  $[\Delta_{\infty}, \Delta_{\infty, \max}]$  where  $\Delta_{\infty} = \max_{j \in [d]} \frac{1}{m} \sum_{i=1}^{m} |g_{i}^{(j)} g^{(j)}|$  and  $\Delta_{\infty, \max} = \max_{j \in [d], i \in [m]} |g_{i}^{(j)} g^{(j)}|$  (Theorem 2). In Section 3.3, we provide a practical example where value of  $\Delta_{\text{Hadamard}}$  can be approximated and use it compare theoretical guarantees of HadamardMultiDim with those of baselines in Table 2.
- 3. ( $\ell_2$ -guarantee) For vectors with  $\ell_2$  norm bounded by B, we provide a collaborative compression scheme SparseReg (Algorithm 4) based on Sparse Regression Codes (Venkataramanan et al., 2014b;a). We obtain the best dependence on B and m for the  $\ell_2$  error ( $\mathcal{O}(B\exp(-m/d))$ ) while compressing to much less than d bits (in fact, to a constant number of bits) per client. The error consists of a penalty for the dissimilarity,  $\Delta_{\text{reg}}$ , the average dissimilarity between compressed vectors which lies in the range  $[\Delta_2, \Delta_{2,\text{max}}]$  where  $\Delta_2 = \frac{1}{m} \sum_{i=1}^m ||g g_i||_2^2$  and  $\Delta_{2,\text{max}} = \max_{i \in [m]} ||g g_i||_2^2$  (see, Theorem 3).
- 4. (cosine-guarantee) For unit norm vectors  $\{g_i\}_{i\in[m]}$ , we estimate the unit vector g in the direction of the average  $\frac{1}{m}\sum_{i=1}^{m}g_i$ . For this, motivated by one-bit compressed sensing (Boufounos & Baraniuk, 2008), our collaborative compression scheme, OneBit (Algorithm 5), sends the sign of the inner product between the vector  $g_i$  and a random Gaussian vector. By establishing an equivalence to halfspace learning with malicious noise, we propose two decoding schemes: the first one is based on (Shen, 2023) which is optimal for halfspace learning but harder to implement and a second one, based on (Kalai et al., 2008) which is easy to implement. Both schemes are computationally efficient, and have an extra dissimilarity term in the error,  $\Delta_{\text{corr}} = \frac{1}{m\pi} \sum_{i=1}^{m} \cos^{-1}(\langle g, g_i \rangle)$ , which is the appropriate dissimilarity between unit vectors (see Theorem 4).
- 5. (Experiments) We perform a simulation for DME with our schemes as the dissimilarities vary and compare the three different error metrics from above with various existing baselines (Fig 2a-2c). We also used our DME subroutines in the downstream tasks of KMeans, power iteration, and linear regression on real (and federated) datasets (Fig 2d-2i). Our schemes have lowest error in all metrics for low dissimilarity regime.

**Organization.** In the next subsection, we present related works in distributed mean estimation. The NoisySign algorithm is given in Algorithm 1, and its analysis can be found in Section 2. In Section 3, we present the two schemes obtaining optimal dependence on m, HadamardMultiDim in Subsection 3.1 and

SparseReg in Subsection 3.2. In Section 4, we analyze the OneBit compression scheme. Finally, in Section 5, we provide experimental results for our schemes.

### 1.1 Related Works

Compressors in Distributed Learning. Starting from (Konečný et al., 2016) most compression schemes in distributed learning involve either quantization or sparsification. In quantization schemes, the real valued input space is quantized to specific levels, and each input is mapped to one of these quantization levels. A theoretical analysis for unbiased quantization was provided in (Alistarh et al., 2017). Subsequently, the distributed mean estimation problem with limited communication was formulated in (Suresh et al., 2017) where two schemes, stochastic rotated quantization (SRQ) and variable length coding, were proposed. These schemes matched the lower bound for communication and  $\ell_2$  error in terms of  $\tilde{B}^2 = \frac{1}{m} \sum_{i=1}^{m} ||g_i||_2^2$ . Performing a coordinate-wise sign is also a quantization operation, introduced in (Bernstein et al., 2018b). Further advances in quantization include multiple quantization levels (Wen et al., 2017), probabilistic quantization with noise (Chen et al., 2020; Jin et al., 2021; Safaryan & Richtarik, 2021), vector quantization (Gandikota et al., 2022; Zhu et al., 2023), and applying structured rotation before quantization (Vargaftik et al., 2021; Safaryan et al., 2021). Sparsification involves selecting only a subset of coordinates to communicate. Common examples include RandK (Konečný & Richtárik, 2018), TopK (Stich et al., 2018) and their combinations (Beznosikov et al., 2022). Note, for all independent compressors, the  $\ell_2$  error scales as  $\tilde{B}^2$ .

Collaborative Compressors. PermK (Szlendak et al., 2021) was the first collaborative compressor, where each client would send a different set of K coordinates. Their error scales with the empirical variance,  $\Delta_2 = \frac{1}{m} \sum_{i=1}^{m} ||g_i - g||_2^2$ . If  $\Delta_2$  is known, or one of the vectors  $g_i$  is known, the lattice-based quantizer in (Davies et al., 2021) and correlated noise based quantizer in (Mayekar et al., 2021) obtains  $\ell_2$  error in terms of  $\Delta_2$ . Further, RandKSpatial (Jhunjhunwala et al., 2021) and RandKSpatialProj (Jiang et al., 2023) utilize the correlation information to obtain the correct normalization coefficients for RandK with rotations, obtaining guarantees in terms of  $\Delta_2$ . In absence of correlation information, they propose a heuristic. A quantizer also based on correlated noise, was proposed in (Suresh et al., 2022) which achieves the lower bound for scalars. However, for d-dimensional vectors of  $\ell_2$ -norm at most B, their dependence on dimension d and number of clients m can be improved by our schemes.

We provide a summary of existing compressors in Table 2, along with their error guarantees.

# 2 NoisySign for unbounded $||g_i||_{\infty}$

In this section, we utilize collaborative compression to use sign compressor when  $||g_i||_{\infty}$  is unbounded. The sign-compressor (Bernstein et al., 2018a) applies the sign function coordinate-wise, where  $\mathrm{sign}(x) = +1$  if  $x \geq 0$  and -1 otherwise. For this section, we will focus on a single coordinate  $j \in [d]$ . Note that for any  $i \in [m]$ ,  $\mathrm{sign}(g_i^{(j)})$  does not have information about  $|g_i^{(j)}|$ . Existing compressors (Karimireddy et al., 2020) remedy this by sending  $|g_i^{(j)}|$  separately, or assuming that  $|g_i^{(j)}|$  is bounded by some constant B (Safaryan & Richtarik, 2021; Jin et al., 2023; Tang et al., 2023). In the second case, the maximum error that can be incurred is  $\frac{B}{2}$ . This can be improved by adding uniform symmetric noise before taking signs (Chen et al., 2020; Chzhen & Schechtman, 2023) depending on g. However, if no information is available about  $|g_i^{(j)}|$ , we cannot provide an estimate of  $g_i^{(j)}$ .

We utilize the concept of adding noise before taking signs, however, to accommodate possibly unbounded  $|g_i^{(j)}|$ , we add symmetric noise with unbounded support. One choice for such noise is the Gaussian distribution  $\mathcal{N}(0,\sigma^2)$  for some  $\sigma \in \mathbb{R}$ . For  $\xi_i^{(j)} \sim \mathcal{N}(0,\sigma^2)$ , we send  $\tilde{b}_i^{(j)} = \text{sign}(g_i^{(j)} + \xi_i^{(j)})$  as the encoding. Note that  $\mathbb{E}[\tilde{b}_i^{(j)}] = \Phi_{\sigma}(g_i^{(j)})$ , where  $\Phi_{\sigma}(t) = 2 \Pr_{x \sim \mathcal{N}(0,\sigma^2)}[x \geq -t] - 1 = \text{erf}(\frac{t}{\sqrt{2}\sigma})$ , and erf is the error function for the unit normal distribution. A single  $\tilde{b}_i^j$  gives us information about  $g_i^{(j)}$ , however, using it to decode  $g_i^{(j)}$  might incur a very large variance. However, assuming that all  $g_i^{(j)}$  are close to  $g_i^{(j)}$  for  $i \in [m]$ ,  $\frac{1}{m} \sum_{i=1}^m \tilde{b}_i^{(j)}$  is a good estimator for  $\Phi_{\sigma}(g^{(j)})$ . So, to estimate  $g^{(j)}$ , we can use  $\Phi_{\sigma}^{-1}(\frac{1}{m} \sum_{i=1}^m \tilde{b}_i^{(j)})$ . This scheme performed

coordinate-wise is the NoisySign algorithm described in Algorithm 1. We provide estimation error for recovering g using this scheme.

**Theorem 1** (Estimation error of noisy sign). With probability  $1 - 2dm^{-c}$ , for some constant c > 0, the estimation error of Algorithm 1 is

$$||\tilde{g} - g||_{\infty} \le \sqrt{\frac{\pi}{2}} \left( \left( 1 - \frac{\Delta_{\Phi} + \sqrt{\frac{2c \log m}{m}}}{\alpha(||g||_{\infty})} \right)^{-1} - 1 \right), \tag{2}$$

where  $\Delta_{\Phi} \triangleq \max_{j \in [d]} \left| \frac{1}{m} \sum_{i=1}^{m} \Phi_{\sigma}(g_i^{(j)}) - \Phi_{\sigma}(g^{(j)}) \right|$  and  $\alpha(u) \triangleq 1 - \Phi_{\sigma}(u)$ . Further, for small  $\Delta_{\phi}$  and  $\frac{||g||_{\infty}}{\sigma}$ , and large m,

$$||\tilde{g} - g||_{\infty} = \mathcal{O}\left(\exp(||g||_{\infty}^{2}/\sigma^{2}) \cdot (\Delta_{\phi} + \sqrt{\log m/m})\right).$$
(3)

Applying  $\Phi_{\sigma}^{-1}$  to estimate g makes our scheme collaborative, as an independent scheme would have individually decoded each client's  $\tilde{b}_i^j$  and thus would not take advantage of concentration over m machines. To gain insight into the error, note that  $(1-x)^{-1}-1\approx x$ , for small x. The error increases with the increase in  $||g||_{\infty}$  as we are compressing unbounded variables  $g_i$  into the bounded domain [-1,1] which is the range of the function  $\Phi_{\sigma}$ . The number of clients m determines the resolution with which we can measure on this domain, as the value  $\frac{1}{m}\sum_{i=1}^{m}\tilde{b_i}$  can only be in multiples of  $\frac{1}{m}$ . Therefore, increasing m decreases the error. As  $m\to\infty$ , the  $\ell_{\infty}$ -error approaches  $\Delta_{\Phi} \exp(\frac{||g||_{\infty}}{\sigma^2})$ . Note that  $\Delta_{\Phi}$  determines the average separation between vectors in terms of the  $\Phi_{\sigma}$  operator. If vectors  $g_i$  are similar to each other,  $\Delta_{\Phi}$  is small, and the error is small as a result. Further,  $\Delta_{\Phi}$  can be bounded by more interpretable quantities if the average separation between  $g_i$  and g is small:

$$\Delta_{\Phi} \le \sqrt{\frac{2}{\pi}} \frac{1}{m\sigma} \sum_{i \in [m]} ||g_i - g||_{\infty}. \tag{4}$$

Note that  $\Delta_{\Phi}$  is always  $\leq 2$ , so if the average error in terms of  $\ell_{\infty}$  norm is much smaller than  $\sigma$ , then the above bound makes sense. Combining Eq (4) with Eq (3), we can find an asymptotic expression for the estimation error. A large value of  $\sigma$  can help us reduce this upper bound on  $\Delta_{\phi}$  as well as the impact of  $\alpha(||g||_{\infty})$  on the estimation error. One can tune the value of  $\sigma$  if additional information is known, for instance a large upper bound on  $||g||_{\infty}$  and  $\frac{1}{m}\sum_{i=1}^{m}||g_{i}-g||_{\infty}$ . However, even without this information, we can run our algorithm. In contrast, the value of  $\max_{i\in[m]}||g||_{\infty}$  is required to even run vanilla sign compression algorithms. Additionally, it yields a constant error of  $\mathcal{O}(\max_{i\in[m]}||g_{i}||_{\infty})$ , as each sign would need to be accurate. If an upper bound is known for  $\max_{i\in[m]}||g_{i}||_{\infty}$ , then the error scales with this upper bound. In contrast, for large m and small  $\Delta_{\Phi}$ , knowledge of a large upper bound helps us to set a large  $\sigma$  and thus our collaborative compressor performs much better. Proofs for this section are provided in Appendix A.

### Algorithm 1 NoisySign

$$\begin{split} & \frac{\mathrm{Encode}(g_i)}{\mathrm{Sample}} \underbrace{\xi_i \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_d)}_{\tilde{b}_i = \mathrm{sign}(g_i + \xi_i)} \\ & \tilde{b}_i = \mathrm{sign}(g_i + \xi_i) \\ & \mathbf{return} \quad \tilde{b}_i. \\ & \underline{\mathrm{Decode}(\{\tilde{b}_i\}_{i \in [m]})}_{\tilde{g}^{(j)} \leftarrow \Phi_\sigma^{-1}(\frac{1}{m} \sum_{i=1}^m \tilde{b}_i^{(j)}), j = 1, \ldots, d} \\ & \mathbf{return} \quad \tilde{g} \end{split}$$

### Algorithm 2 Hadamard1DEnc

Input: Scalar s, Level K  $S_K^- = \bigcup_{k=0}^{2^{K-1}-1} \left[ -B + \frac{2kB}{2^{K-1}}, -B + \frac{(2k+1)B}{2^{K-1}} \right]$ return -1 if  $s \in S_K^-$  else +1

### Algorithm 3 HadamardMultiDim

 $\begin{array}{c} \underline{\mathbf{Init}()} \\ \mathbf{Sample} \ \mathbf{a} \ \mathbf{random} \ \mathbf{permutation} \ \rho \sim \Pi_m. \\ \mathbf{Clients} \ \mathbf{and} \ \mathbf{server} \ \mathbf{share} \ \rho. \\ \underline{\mathbf{Encode}(g_i)} \\ \overline{\mathbf{for} \ j \in [d]} \ \mathbf{do} \\ \tilde{b_i}^{(j)} \leftarrow \mathbf{Hadamard1DEnc}(g_i^{(j)}, \rho^{(i)}) \\ \mathbf{end} \ \mathbf{for} \\ \mathbf{return} \ \tilde{b_i} \\ \underline{\mathbf{Decode}(\{\tilde{b_i}\}_{i \in [m]})} \\ \overline{\mathbf{for} \ j \in [d]} \ \mathbf{do} \\ \tilde{g}^{(j)} = \sum_{i=1}^m \tilde{b_i}^{(j)} \cdot \frac{B}{2\rho^{(i)}-1} \\ \mathbf{end} \ \mathbf{for} \\ \mathbf{return} \ \tilde{g} \\ \end{array}$ 

For NoisySign, collaboration ensured concentration over m. In the next section, we utilize collaboration to obtain optimal dependence on m.

### 3 Optimal Dependence on m

If  $||g||_{\infty}$  or  $||g||_2$  is bounded, we can obtain an almost optimal exponential decay with m. We provide two schemes that obtain optimal  $\ell_{\infty}$  (by modifying the sign compressor) and  $\ell_2$  error dependence in terms of m and the diameter of the space B. The proofs for this section are provided in Appendix B.

### 3.1 HadamardMultiDim

When the vectors have bounded  $\ell_{\infty}$  norm, instead of obliviously using the sign compressor on every coordinate on every client, one may be able to divide their range and cleverly select bits to encode the most information. We call our algorithm Hadamard scheme, because the binary-search method involved is akin to the rows of a Hadamard-type matrix.

**Assumption 1** (Bounded domain).  $||g_i||_{\infty} \leq B, \forall i \in [m].$ 

This would imply that for any  $j \in [d]$ ,  $g_i^{(j)} \in [-B,B]$ ,  $\forall i \in [m]$ . Now, consider the  $i^{th}$  client and the scalar  $g_i^{(j)}$  and assume that we are allowed to encode this using m bits. The best error that we can achieve is  $\frac{B}{2^{m-1}}$ , by performing a binary search on the range [-B,B] for  $g_i^{(j)}$ , sending one bit per level of the binary search. However, this scheme is not collaborative. To obtain a collaborative scheme, for some permutation  $\rho$  on the set of clients [m], the  $i^{th}$  client can perform binary search until level  $\rho^{(i)}$  and sends its decision at level  $\rho^{(i)}$ . In this case, each client sends only 1 bit per coordinate. To decode  $\tilde{g}^{(j)}$ , we take a weighted sum of the signs obtained from different clients weighed by their coefficients  $\frac{B}{2\rho^{(i)}-1}$ . This is the core subroutine (Algorithm 2). The full compression scheme for d coordinates applies this coordinate-wise in Algorithm 3. Note that the clients and the server should share the permutation  $\rho$  before encoding and decoding, which need not change over different instantiations of the mean estimation problem. To understand the core idea of the scheme, consider the case when all vectors  $g_i = g$ . Then, sending a different level from a different client is equivalent to doing a full binary search to quantize g. As long as  $g_i$ s are close to g, we hope that this scheme should give us a good estimate of g. Let  $\tilde{b}_{i,k}^{(j)}$  denotes the encoding of  $g_i^{(j)}$  at level  $k \ \forall i, k \in [m], j \in [d]$ . Therefore, if we were to use binary search to encode  $g_i^{(j)}$ , then  $\sum_{k=1}^m \frac{B\tilde{b}_{i,k}^{(j)}}{2^{k-1}}$  would be it's decoded value.

**Theorem 2** (HadamardMultiDim Error). Under Assumptions 1, the estimation error for Algorithm 3 is

$$\mathbb{E}_{\rho \sim \Pi_m}[||\tilde{g} - g||_{\infty}] \le \frac{B}{2^{m-1}} + \min\{\Delta_{\text{Hadamard}}, \Delta_{\infty, \max}\},\tag{5}$$

$$where \ \Delta_{\text{Hadamard}} \equiv \max_{r \in [d]} \sqrt{\frac{1}{m^2} \sum_{1 \leq i \neq j \leq m} \sum_{k=1}^{m} \left( \frac{B(\tilde{b}_{i,k}^{(r)} - \tilde{b}_{j,k}^{(r)})}{2^{k-1}} \right)^2} \ , \ and \ \Delta_{\infty, \max} \equiv \max_{r \in [d], i \in [m]} |g_i^{(r)} - g^{(r)}|.$$

As we allow full collaboration between clients, in the worst case, we might have to incur a cost  $\Delta_{\infty,\text{max}}$  which is the worst case dissimilarity among clients. However, if client vectors are close, we might end up paying a much lower cost of  $\Delta_{\text{Hadamard}}$ .

### 3.2 Sparse regression coding

In this part, we extend the coordinate-wise guarantee of the HadamardMultiDim to  $\ell_2$  error between d-dimensional vectors of bounded  $\ell_2$ -norm.

**Assumption 2** (Norm Ball).  $||g_i||_2 \leq B, \forall i \in [m]$ .

To extend the idea of binary search and full collaboration from HadmardMultiDim, we first need a compression scheme which performs binary search on d dimensional vectors with  $\ell_2$  error guarantees. Sparse Regression codes (Venkataramanan et al., 2014b;a), which are known to achieve the optimal rate-distortion tradeoff

for a Gaussian source, fit our requirements. Let  $A \in \mathbb{R}^{mL \times d}$  be a matrix for some parameter L > 0, where  $A_k$  denotes the kth row of A. To compress a single vector g by Sparse Regression codes with the matrix A, we use  $m \log L$  bits to find the closest vector to g in the first L rows of A; say the index of this vector is  $\tilde{b}_1$ . Similar to binary search, we subtract  $c_1 A_{\tilde{b}_1}$  from g, where  $c_1$  is given in (9) to obtain an updated g. We repeat the process using the next set of L rows. Here, each set of L rows corresponds to a single level of binary search, with the coefficients  $c_i$  obtained from Eq (9) having a decaying exponent. Although (Venkataramanan et al., 2014b) requires A to be a random matrix with iid Gaussian entries, if A satisfies the following definition, then we can apply the theoretical guarantees of (Venkataramanan et al., 2014b).

**Definition 1** ( $(\delta_1, \delta_2, \Gamma)$ -Cover). A matrix  $A \in \mathbb{R}^{mL \times d}$  is a  $(\delta_1, \delta_2, \Gamma)$ -Cover for  $\delta_1, \delta_2 \in (0, 1)$  and  $\Gamma > 0$ , if for each  $k \in [m]$ ,  $\exists \epsilon_k, \gamma_k, \Gamma_k \in \mathbb{R}$ , such that

$$\max_{j \in [(k-1)L+1,kL]} ||A_j||_2 \le \sqrt{d} (1+\gamma_k); \quad \max_{j \ne j' \in [(k-1)L+1,kL]} ||A_j - A_{j'}||_2 \le \Gamma_k$$

$$\forall v \in \mathbb{S}^{d-1}, \quad \max_{j \in [(k-1)L+1,kL]} \langle v, A_j \rangle \ge \sqrt{2 \log L} (1+\epsilon_k),$$

$$\sum_{k=1}^m \frac{|\gamma_k|}{m} \le \delta_1, \quad \sum_{k=1}^m \frac{|\epsilon_k|}{m} \le \delta_2, \ and \ \max_{k \in [m]} \Gamma_k \le \Gamma.$$

By the above definition, for each  $k \in [m]$ , the L consecutive rows of A, which correspond to  $\{A_j, j \in [(k-1)L+1, kL]\}$ , form a  $\varepsilon$ -net (Vershynin, 2018, Chapter 4.2) of  $\mathbb{S}^{d-1}$  with  $\varepsilon = \sqrt{2(1-\sqrt{\frac{2\log L}{d}})}$ . By carefully selecting the parameters in the proof of (Venkataramanan et al., 2014b, Theorem 1), we can show that this scheme obtains  $\ell_2$  error  $B \exp(-m)$  as long as A is a  $(\delta_1, \delta_2, \Gamma)$ -cover. Note that the matrix A is deterministically constructed with m sets of  $\varepsilon$ -nets. The choice of Gaussian A in (Venkataramanan et al., 2014b) can now be understood from the fact that such random Gaussian matrices are in fact  $\varepsilon$ -nets. The following remark summarizes this.

Remark 1 (Gaussian A). If each entry of the matrix A is sampled iid from  $\mathcal{N}(0,1)$  then, for any  $\delta_1, \delta_2, \delta_3 \in [0,1)$  and  $\Gamma^2 = 2d + 4\sqrt{d\log(\frac{mL^2}{\delta_3})} + 4\log(\frac{mL^2}{\delta_3})$ , the matrix A is a  $(\delta_1, \delta_2, \Gamma)$ -Cover with probability  $1 - \delta_3 - 2m^2Le^{-\frac{d\delta_1^2}{8}} - m\left(\frac{L^{2\delta_2}}{\log L}\right)^{-m}$ .

Note that the expressions for  $\delta_1, \delta_2$  can be obtained from (Venkataramanan et al., 2014b) and the proof for  $\Gamma$  is provided in Appendix B.3. For  $d = \Omega(\log m)$ , the probability of error in the above Remark can be made arbitrarily close to 1 for large m.

We now extend the scheme based on sparse regression codes in (Venkataramanan et al., 2014b) to our setting of distributed mean estimation, similar to HadamardMultiDim. Each client  $i \in [m]$  encodes at level  $\rho^{(i)}$  where  $\rho$  is a permutation on [m] and the server computes the weighted sum of the encodings from each client with corresponding coefficients  $c_{\rho^{(i)}}$ . Here,  $\tilde{b}_{i,k}$  denotes the encoding of  $g_i$  at level  $k \ \forall i,k \in [m]$ . If we were to use Sparse Regression codes to encode only  $g_i$  then  $\sum_{k=1}^m c_k A_{(k-1)L+\tilde{b}_{i,k}}$  would be it's decoded value. The full scheme in described in Algorithm 4 and we compute it's  $\ell_2$  estimation error in the next theorem. In this scheme, each client only needs to send a single index  $\tilde{b}_{i,\rho^{(i)}}$ . As  $\rho$  is known to the server, this index is in the range [L] and each client uses only  $\ell$ 0 bits for communication.

**Theorem 3** (SparseReg Error). Under Assumption 2, the following statements hold for  $\tilde{g}$ , the output of Algorithm 4 for any matrix  $A \in \mathbb{R}^{mL \times d}$ .

$$\mathbb{E}_{\rho \sim \Pi_m}[||g - \tilde{g}||_2^2] = ||\bar{g} - g||_2^2 + \Delta_{\text{reg}}, \tag{6}$$

$$where, \ \Delta_{\text{reg}} \equiv \frac{1}{m^2} \sum_{i,j \in [m], i \neq j} \sum_{k=1}^m c_k^2 ||A_{(k-1)L + \tilde{b}_{i,k}} - A_{(k-1)L + \tilde{b}_{j,k}}||_2^2 \ and \ \bar{g} \equiv \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^m c_k A_{(k-1)L + \tilde{b}_{i,k}}.$$

In addition, if the matrix A is a  $(\delta_1, \delta_2, \Gamma)$ -Cover, then,

$$\mathbb{E}_{\rho \sim \Pi_m}[||g - \tilde{g}||_2^2] \le \Lambda + \Delta_{2,\max}, \quad and, \quad ||\bar{g} - g||_2^2 \le \Lambda,$$

$$where, \ \Lambda \equiv B^2 \left( 1 + \frac{10 \log L}{d} e^{\frac{m \log L}{d}} (\delta_1 + \delta_2) \right)^2 \left( 1 - \frac{2 \log L}{d} \right)^m.$$

$$(7)$$

$$\Delta_{\text{reg}} \le \frac{2B^2 \Gamma^2 \log L}{d^2 m^2} \sum_{i,j \in [m], i \ne j} \sum_{k=1}^m \left( 1 - \frac{2 \log L}{d} \right)^{k-1} \mathbf{1}(\tilde{b}_{i,k} \ne \tilde{b}_{j,k}). \tag{8}$$

**Remark 2** (Growth of  $\Lambda$  with m). By expanding the squares, we can see that  $\Lambda$  can be bounded by 3 terms.

$$\Lambda \le B^2 e^{-\frac{2m \log L}{d}} + \frac{20B^2 \log L}{d} (\delta_1 + \delta_2) e^{-\frac{m \log L}{d}} + \frac{100B^2 \log^2 L}{d^2} (\delta_1 + \delta_2)^2.$$

The first and the second terms decrease exponentially with m. The third term is a small constant independent of m, depending on  $\delta_1, \delta_2 \ll 1$  that can be made arbitrarily small.

The proof is provided in Appendix B.2. Similar to HadmardMultiDim, the first term in the estimation error,  $\Lambda$ , has an exponential dependence on m and is obtained from the existing results of Sparse Regression Codes from (Venkataramanan et al., 2014b). In terms of  $\ell_2$  error this dependence on m is better than all the prior methods. Note that if our matrix A is a  $(0,0,\Gamma)$ -Cover, i.e.,  $\delta_1 = \delta_2 = 0$ , then  $\Lambda$  decreases exponentially with m. As mentioned in Remark 1, we can can achieve very small (but nonzero) values for  $\delta_1, \delta_2$  with Gaussian matrices.

By combining Eqs (6) and (7), we can bound the squared  $\ell_2$  estimation error by  $\Lambda + \min\{\Delta_{\rm reg}, \Delta_{2,\rm max}\}$ . Note that this error resembles Theorem 2, as it contains an exponentially decreasing term of m along with a minimum of two terms, one of which is the maximum deviation between any client and the other ( $\Delta_{2,\rm max}$ ) and the other is the average difference of encoded values across all pairs of clients ( $(\Delta_{\rm reg})$ ). Further, even the dissimilarity term  $\Delta_{\rm reg}$  has a similar structure to  $\Delta_{\rm Hadamard}$  as it is the pairwise difference between encodings of two different vectors at all levels, where the difference at level k is proportional to  $e^{-k}$ . If the client vectors are not close to each other, we might incur the worst possible error  $\Delta_{2,\rm max}$ , but if they are close, we will pay an average price in terms of  $\Delta_{\rm reg}$ .

The upper bound on  $\Delta_{\text{reg}}$  in Eq (8) depends on A only via the encoded values  $\{\tilde{b}_{i,k}\}_{i,k\in[m]}$ . For this upper bound, the corresponding terms for each level, here denoted by k, decrease geometrically with k. The term  $\frac{2B^2\Gamma^2\log L}{d^2m^2}$ , ensures that  $\Delta_{\text{reg}}$  is of the correct scale as in the worst case when all encodings are different for all pairs of clients, the summation gives us  $\Theta(\frac{d}{2\log L})$  which results in an overall  $\Delta_{\text{reg}}$  of the scale of  $B^2$ . From (Venkataramanan et al., 2014a), for each  $k \in [m]$ , all encodings upto and including the  $k^{th}$  encoding, i.e.,  $\{\tilde{b}_{i,k'}\}_{k'\leq k}$ , form an estimate of  $g_i$  that has error scaling as  $Be^{-\frac{\log L(k-1)}{d}}$ , which is similar to the error scaling of HadamardMultiDim for a single dimension.

While both the HadmardMultiDim and SparseReg schemes achieve very low communication rate, that comes at the price of O(m) computing in the Encode step. This higher cost in computing is to be expected when one wants to exploit the full potential of collaborative compression (e.g., (Jiang et al., 2023), where the Decode step takes  $O(m^2)$  time).

### 3.3 Benefits of HadmardMultiDim and SparseReg

We now provide a example to show that for practical scenarios, the error terms  $\Delta_{reg}$  and  $\Delta_{\text{Hadamard}}$  are much smaller than their worst case values  $\Delta_{\infty,\text{max}}$  and  $\Delta_{2,\text{max}}$  and have similar behavior as their average values,  $\Delta_2 \equiv \frac{1}{m} \sum_{i=1}^m ||g_i - g||_2^2$  and  $\Delta_\infty \Delta_\infty \equiv \max_{r \in [d]} \frac{1}{m} \sum_{i=1}^m |g_i^{(r)} - g^{(r)}|$  respectively. Consider the scenario of Theorem 2 ( $\ell_\infty$  error) and set d=1. Assume that the first c vectors are  $g_1'$  and the remaining m-c vectors are  $g_2'$ , for some constant  $c \ll m$ . In this case,  $\Delta_{\infty,\text{max}} = (1 - \frac{c}{m})|g_1' - g_2'| \approx |g_1' - g_2'|$ , while  $\Delta_\infty \approx \frac{2c}{m}|g_1' - g_2'|$ . In this scenario, if the compressed values  $\tilde{b}$  for  $g_1'$  and  $g_2'$  according to the HadamardMultiDim differ at  $k \in \mathcal{K} \subseteq [m]$  levels, then,  $\Delta_{\text{Hadamard}} \approx \sqrt{\frac{c}{m} \sum_{k \in \mathcal{K}} (B/2^{k-1})^2} \leq \sqrt{\frac{cB^2}{m} \sum_{k=k^*}^\infty 2^{-2(k-1)}} \leq \sqrt{\frac{4c}{3m}} \frac{B}{2^{k^*-1}}$ , where

 $k^* = \min_{k \in \mathcal{K}} k$  and we upper bound the sum over the set  $\mathcal{K}$  by the sum of a geometric series. As  $\Delta_{\text{Hadamard}}$  averages over all machines, it decreases with m similar to  $\Delta_2$  and should be much smaller than  $\Delta_{\infty,\text{max}}$ . The only case when it is not smaller than  $\Delta_{\infty,\text{max}}$  is when  $g_1'$  and  $g_2'$  are very close, so that  $\Delta_{\infty,\text{max}} = \mathcal{O}(\sqrt{m^{-1}})$ , but the first level where they differ  $(\min_{k \in \mathcal{K}} k)$  is very small. One such example is when the quantized values of  $g_1'$  in the set  $\mathcal{K}$  sorted by the levels in increasing order are (+1, -1, -1, -1) and that of  $g_2'$  are (-1, +1, +1, +1). As the vectors are extremely close in this case, the estimation error with  $\Delta_{\infty,\text{max}}$  is not very large. Further, if we assume a distributional assumption on the vectors  $g_i$ , similar to how we generate Figure 2b, obtaining vectors where  $\Delta_{\text{Hadamard}} > \Delta_{\infty,\text{max}}$ , happens with low probability.

# $$\begin{split} & \frac{\text{Algorithm 4 SparseReg}}{\text{Init ()}} \\ & \frac{\text{Init ()}}{\text{Clients and server share } A \in \mathbb{R}^{mL \times d}, \text{ and a permutation } \rho \sim \Pi_m. \\ & \frac{\text{Encode}(g_i)}{g_i' \leftarrow g_i} \\ & \text{for } j \in [\rho^{(i)}] \text{ do} \\ & \tilde{b}_{i,j} \leftarrow \arg\max_{r \in [L]} \langle A_{(j-1)L+r}, g_i' \rangle \\ & g_i' \leftarrow g_i' - c_j A_{(j-1)L+\tilde{b}_{i,j}} \\ & \text{end for} \\ & \tilde{b}_i \leftarrow \tilde{b}_{i,\rho^{(i)}} \\ & \text{return } \tilde{b}_i \\ & \frac{\text{Decode}(\{\tilde{b}_i\}_{i \in [m]})}{\tilde{g}} \leftarrow \sum_{i \in [m]} c_{\rho^{(i)}} A_{(\rho^{(i)}-1)L+\tilde{b}_i} \end{split}$$

$$c_i = B\sqrt{\frac{2\log L}{d^2} \left(1 - \frac{2\log L}{d}\right)^{i-1}} \tag{9}$$

### Algorithm 5 OneBit

$$\begin{split} &\frac{\text{Init}()}{\text{Clients and server share unit vectors }} \{z_i\}_{i \in [m]}. \\ &\frac{\text{Encode}(g_i)}{\tilde{b}_i \leftarrow \text{sign}(\langle g_i, z_i \rangle)} \\ &\text{return } \tilde{b}_{\underline{i}} \\ &\text{Decode}(\{\tilde{b}_i\}_{i \in [m]}) \\ &g' \leftarrow \left\{ \begin{array}{c} (\text{Shen, 2023, Algorithm 1) (Tech. I)} \\ &\frac{1}{m} \sum_{i=1}^m z_i \tilde{b}_i \text{ (Tech. II)} \\ \tilde{g} \leftarrow g'/||g'||_2 \end{array} \right. \end{split}$$

Using this example, we can compare the error of our proposed schemes to baselines mentioned in Table 2. Assume a modification of the above example where each vector is d-dimensional, and all of it's coordinates are either  $g_1' \mathbb{M}_d$  or  $g_2'$ . Consider any  $\ell_2$  compressor whose error is either proportional to  $\Psi \tilde{B}^2$  or  $\Psi \Delta_2$  and it sends  $\psi$  bits/client for some constants  $\psi, \Psi > 0$ . The  $\ell_2$  error is defined as  $\mathbb{E}[||\tilde{g} - g||_2^2]$  and the  $\ell_\infty$  error is defined as  $\mathbb{E}[||\tilde{g} - g||_\infty]$ , therefore the corresponding  $\ell_\infty$  error of these compressors is  $\sqrt{\Psi} \tilde{B}$  or  $\sqrt{\Psi} \Delta_2$ . Now, consider the example that we just presented with d > 1 and all coordinates being equal for each vector. Therefore,  $\Delta_2 \approx \frac{cd}{m} |g_2' - g_1'|^2$ , and plugging this in, the  $\ell_2$  error of the schemes is  $\sqrt{\Psi} \tilde{B}$  or  $\sqrt{\Psi} \frac{cd}{m} |g_2' - g_1'|$ . HadamardMultiDim sends d bits/client; therefore, to compare with any of these schemes, we set  $\psi = d$ . For RandK, this would mean setting  $K = \frac{d}{32 + \log d}$ . Now, if  $|g_1'|, |g_2'| \approx B$  but  $|g_2' - g_1'| \ll B$ , then  $\tilde{B} \approx \sqrt{d}B$ . Using these approximations, the error of RandK is  $\sqrt{(32 + \log d)d}B$  as  $\Psi = \frac{d}{K} = 32 + \log d$ . This is much larger than the  $\ell_\infty$  error of HadamardMultiDim, as the first term is  $B \cdot 2^{m-1}$  and the second term  $\Delta_{\text{Hadamard}} \approx \sqrt{\frac{c}{m}} |g_2' - g_1'|$ . A similar argument holds for all independent compression schemes, as their  $\ell_\infty$  error scales as  $\tilde{B}$ , which in the worst case is  $\sqrt{d}B$ .

For compressors whose error scales as  $\Psi\Delta_2$  (PermK, RandKSpatial, RandKSpatialProj), we obtain the same number of bits/client as HadamardMultiDim scheme by setting  $K = \frac{d}{32 + \log d}$ . Consider RandKSpatialProj, where  $\Psi = \frac{32 + \log d}{m}$ , and the error for our example is  $\sqrt{c\frac{(32 + \log d)d}{m^2}}|g_2' - g_1'|$ . As long as d > m, this error is larger than  $\Delta_{\text{Hadamard}}$  by constant terms. A similar argument holds for RandKSpatial and PermK. Additionally, note that the theoretical guarantees for RandKSpatial and RandKSpatialProj do not hold if the correlation is unknown. Without this information, the heuristics they use do not result in theoretical guarantees, and their error might be as bad as RandK. CorrelatedSRQ achieves the lower bound for collaborative compressors for d=1, and is based on a coordinate-wise scheme, hence the  $\Delta_{\infty}$  in its error guarantees. However, for  $d\gg 1$ , its error scales poorly. For the example described above,  $||g_i||_2 \leq \sqrt{d}B$ , therefore, the  $\ell_{\infty}$  error for CorrelatedSRQ is  $\sqrt{\frac{1}{m}} \min\{\frac{d\Delta_{\infty}^d B}{K}, \frac{d^3 B^2}{K^2}\}$ . Even for K=2, correlated SRQ requires double the number of bits/client as HadamardMultiDim. The first term of HadamardMultiDim is  $B\cdot 2^{m-1}$ , which is much smaller than any of these terms, while  $\Delta_{\text{Hadamard}} \approx \sqrt{\frac{m}{c}}\Delta_{\infty}$  for our example. Therefore, as long

as  $\left(\frac{m^2K}{cdB}\right)^{1/(2d-1)} < \Delta_{\infty} < \frac{\sqrt{c}dB}{mK}$ ,  $\Delta_{\text{Hadamard}}$  is smaller than  $\ell_{\infty}$  error of CorrelatedSRQ. The size of this interval for  $\Delta_{\infty}$  increases as d increases.

With the above example and analysis, we have specified the exact scenarios when HadamardMultiDim outperforms baselines, and this can be easily extended to SparseReg.

### 4 One-bit Schemes

In this section, our vectors are assumed to belong on the unit sphere  $\mathbb{S}^{d-1}$ . Further, our goal is to recover the unit vector in the direction of the average vector  $g = (\frac{1}{m} \sum_{i \in [m]} g_i) / ||\frac{1}{m} \sum_{i \in [m]} g_i||_2$ .

**Assumption 3** (Unit vectors).  $g_i \in \mathbb{S}^{d-1}, \forall i \in [m].$ 

Consider the collaborative compressor where each client has sample  $z_i \sim \text{Unif}(\mathbb{S}^{d-1})$  (which are also available to the server apriori). Client i sends the single bit  $\tilde{b}_i = \text{sign}(\langle g_i, z_i \rangle)$  to the server. To recover g, consider the trivial case when all vectors  $g_i$ s were equal. Then, each  $\tilde{b}_i = \text{sign}(\langle g, z_i \rangle)$ , and to recover g, the server needs to learn the halfspace corresponding to g from a set of m labeled datapoints. Applying the same method to when  $g_i$ s are not all the same, we can estimate g by solving the following optimization problem.

$$\min_{\tilde{q} \in \mathbb{S}^{d-1}} \frac{1}{m} \mathbf{1}(\tilde{b_i} \neq \operatorname{sign}(\langle z_i, \tilde{g} \rangle)). \tag{10}$$

Here,  $\mathbf{1}(\cdot)$  denotes the indicator function. We can intuitively view (10) as a halfspace learning problem with a groundtruth g, but in the presence of noise, as  $g_i \neq g$ . Learning halfspaces in the presence of noise is hard in general (Guruswami & Raghavendra, 2006). In our setting, if we sample  $z_i$  from the intersection of the halfspaces with normal vectors g and  $g_i$ , then the label is  $\operatorname{sign}(\langle g, z_i \rangle)$ , otherwise, it is  $-\operatorname{sign}(\langle g, z_i \rangle)$ . We can consider this to be under the malicious noise model, wherein a fraction of datapoints are corrupted.

**Lemma 1** (Malicious Noise). If  $z_i \sim \text{Unif}(\mathbb{S}^{d-1})$  and  $\tilde{b_i} = \text{sign}(\langle z_i, g_i \rangle)$ ,  $\forall i \in [m]$ , then, with probability  $1 - \mathcal{O}(\exp(-m\Delta_{\text{corr}}))$ ,  $\zeta$ , the fraction of the set of datapoints  $\{(z_i, \tilde{b_i})\}_{i \in [m]}$  satisfying  $\text{sign}(\langle z_i, g_i \rangle) \neq \text{sign}(\langle g, z_i \rangle)$  is equal to  $\Theta(\Delta_{\text{corr}})$ , where  $\Delta_{\text{corr}} \triangleq \frac{1}{m\pi} \sum_{i=1}^m \arccos(\langle g_i, g_i \rangle)$ .

The proof of the lemma is provided in Appendix C.1. Our methods will use  $\Delta_{\text{corr}}$  to measure the deviation between clients. For small  $\Delta_{\text{corr}}$ , we obtain better performance. If  $\langle g, g_i \rangle \geq 0, \forall i \in [m]$ , then

$$\cos(\pi \Delta_{\text{corr}}) \ge \sqrt{\frac{1}{m} + \frac{2}{m^2} \sum_{1 \le i < j \le m} \langle g_i, g_j \rangle}.$$
 (11)

As long as the corruption level,  $\zeta < \frac{1}{2}$ , we can hope to recover the halfspace g. We provide two techniques – Techniques I and II, to recover g, thus yielding two corresponding Decode procedures.

The first decoding procedure (Technique I) is a linear time algorithm for halfspace learning in the presence of malicious noise (Shen, 2023, Theorem 3) that provides optimal sample complexity and noise tolerance.

**Theorem 4** (Error of Technique I). If  $\zeta$  defined in Lemma 1 is less than  $\frac{1}{2}$ , after running Algorithm 5 with Technique I, with probability  $1 - \delta - \mathcal{O}(\exp(-m\Delta_{\text{corr}}))$ , we obtain a hyperplane  $\tilde{g}$  such that,  $\langle \tilde{g}, g \rangle \geq \cos(\pi(\Delta_{\text{corr}} + \frac{d}{m}\operatorname{polylog}(d, \frac{1}{m}, \frac{1}{\delta}))$ .

The algorithm itself is fairly complicated. It assigns weights to different points based on how likely they are to be corrupted. The algorithm proceeds in stages, wherein each stage decreases the weights of the corrupted points and solves the weighted version of (10). The key technique is to use matrix multiplicative weights update (MMWU) (Arora et al., 2012) to yield linear time implementation of both these steps, instead of (Awasthi et al., 2017) which used polynomial time linear programs for this purpose.

Technique II is the simple average algorithm of (Servedio, 2002), which obtains suboptimal error guarantees for halfspace learning.

**Theorem 5** (Error of Technique II). If  $\zeta$  defined in in Lemma 1 is less than  $\frac{1}{2}$ , after running Algorithm 5 with Technique II, with probability  $1 - \delta - \mathcal{O}(\exp(-m\Delta_{\text{corr}}))$ , we obtain a hyperplane  $\tilde{g}$  such that,  $\langle \tilde{g}, g \rangle \ge \cos(\pi(\sqrt{d}\Delta_{\text{corr}} + \frac{d}{\sqrt{m}}\sqrt{\log(\frac{d}{\delta})}))$ .

The performance of both techniques improves with decrease in  $\Delta_{\text{corr}}$ . Since we have only m bits to infer a d-dimensional vector, we require m > d, with Technique II requiring  $m > d^2$ . If we send t bits per client in OneBit, then the number of samples for the halfspace learning is mt, thus obtaining the guarantee in Table 1. The main benefit of OneBit schemes is their extreme communication efficiency. Existing quantization and sparsification schemes require sending at least  $\log K$  or  $\log d$ , where K is the number of quantization levels. The proofs for this section are provided in Appendix C.

Comparison to  $\ell_2$  compressor. Note that, we can use compressor for  $\ell_2$  error to first decode the mean and then normalize it to obtain its unit vector. If such a scheme uses t bits and has  $\ell_2$  error either  $\Lambda\Delta_2$  or  $\Lambda\tilde{B}^2$  then its cosine similarity  $\frac{\langle g,\tilde{g}\rangle}{||g'||_2||\tilde{g}||_2} \geq 1 - \frac{\Lambda}{2||g'||_2^2}$  for  $||g'||_2 \approx ||\tilde{g}||_2$ , where  $g' = \frac{1}{m} \sum_{i=1}^m g_i$  and  $\tilde{g}$  is the estimate of g'. To compare this with OneBit Technique I, we send  $\lambda$  bits per client to obtain the same communication budget. The cosine similarity of this scheme is  $\cos(\pi(\Delta_{\text{corr}} + \frac{d}{tm}))$ . We can lower bound this similarity by  $1 - 2\pi^2\Delta_{\text{corr}}^2 + 2\pi^2\frac{d^2}{m^2t^2}$  as  $\cos(x) \geq 1 - \frac{x^2}{2}$ . Comparing this cosine similarity with that obtained for  $\ell_2$ -compressor, as long as  $2\pi^2\Delta_{\text{corr}}^2 + 2\pi^2\frac{d^2}{m^2\beta^2} < \Lambda$ , OneBit Technique I performs better. For any sparsification scheme sending K coordinates,  $\Lambda$  is at least  $\frac{d}{mK}$ . If we set  $t = 32K + K \log d$ , OneBit Technique I outerperforms the sparsification scheme as long as  $\Delta_{\text{corr}}$  is small.

### 5 Experiments

Setup. To compare the performance of our proposed algorithms, we perform DME for three different distributions which correspond to the three error metrics covered by our schemes –  $\ell_2$ ,  $\ell_\infty$  and cosine distance. Then, we run our algorithms as the DME subroutine for four different downstream distributed learning tasks – KMeans, power iteration, linear regression and logistic regression. KMeans and power iteration are run on MNIST (LeCun & Cortes, 2010) and FEMNIST (Caldas et al., 2018) datasets, and we report the KMeans cost and top eigenvalue as the metrics. For linear regression, we run gradient descent on UJIndoorLoc (Torres-Sospedra et al., 2014) and a Synthetic mixture of regressions dataset, with low dissimilarity between the mixture components, and report the test MSE. For logistic regression, we run gradient descent on the HAR (Reyes-Ortiz et al., 2012) dataset and report the training loss and test accuracy for binary classification. We compare against all baselines in Table 2 for 5 random seeds and report the methods which perform the best. For DME experiemnts, the shaded corresponds to 2 standard deviations around the mean, while for all downstream tasks, it corresponds to 1 standard deviation. The results for KMeans, power iteration and linear regression are reported in Fig 2, while the results for logistic regression are reported in Fig 3. Additional details for our experimental setup are deferred to Appendix D.

Results. Distributed Mean Estimation. From Fig 2a and 2b, HadamardMultiDim and SparseReg, whose error is optimal in m, obtain the best performance in terms of  $\ell_{\infty}$  and  $\ell_2$  error for low dissimilarity. Especially, for HadamardMultiDim in Fig 2b, the gap in  $\ell_{\infty}$  error to next best scheme is very large. NoisySign obtains competitive performance to other baselines as we use a large  $\sigma$ . The performance of OneBit for cosine distance metric (Fig 2c) shows that compressors with  $\ell_2$  error guarantees perform poorly in terms of cosine distance. For all collaborative compression schemes, including our proposed schemes, performance degrades as dissmilarity increases. From Fig 2a and 2b, the rate of this decrease is more severe for SparseReg than HadamardMultiDim. For large dissimilarity, HadamardMultiDim and SparseReg can perform worse than certain baselines.

KMeans and Power iteration. For MNIST dataset, where dissimilarity is low, HadamardMultiDim performs best for KMeans and close to the best baseline for power iteration (Fig 2d and 2e). Most of our collaborative compression schemes do not perform as well as RandK on FEMNIST, due to higher client dissimilarity. OneBit is very communication-efficient, so running it for the same communication budget as our baselines ensures that it still remains competitive for KMeans(Fig 2g). By increase in confidence intervals for the same baselines from MNIST to FEMNIST, we can see that FEMNIST dataset is much more heterogeneous.

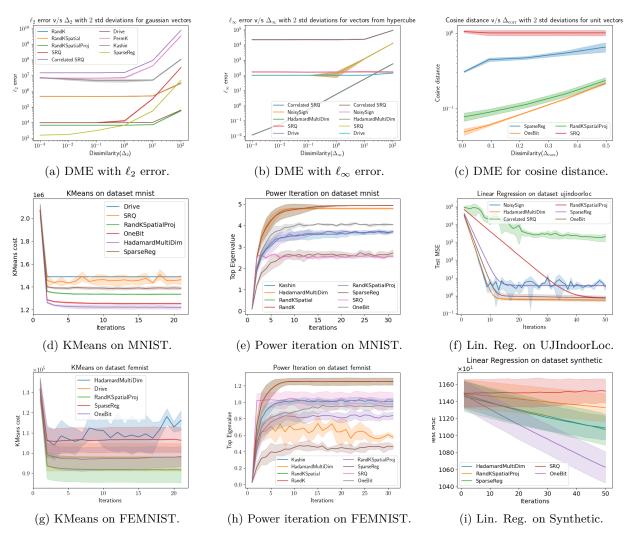
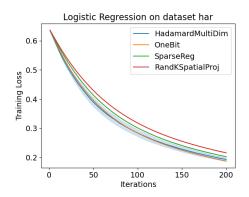


Figure 2: Performance of DME(Distributed Mean Estimation), KMeans, Power iteration and linear regression for the same communication budget. For each experiment, we report the best compressors. Lin. Reg. refer to Linear Regression. For power iteration, higher top eigenvalue is better. For all other experiments, we report the error, so lower is better.

Linear Regression. From Fig 2f and2i, all collaborative compressors perform better than independent compressors as UJIndoorLoc and synthetic datasets have low dissimilarity among clients as compared to FEMNIST. Our schemes can take full advantage of this low dissimilarity, so HadamardMultiDim and SparseReg outperform baselines on both datasets. As the Synthetic dataset has lower dissimilarity than UJIndoorLoc, even the NoisySign performs better than other baselines. Further, OneBit obtains best performance.

Logistic Regression. From Fig 3, our compressors, Onebit, Sparsereg, and HadamardMultDim, are the best, second, and fourth best compressors, respectively, in terms of both training loss and test accuracy. Further, among the top 4 best-performing schemes, only one baseline, RandKSpatialProj, comes in third. This shows the benefit of using collaborative compressors.

Stability of Our Algorithms. Among all our proposed algorithms, OneBit and SparseReg are the most stable with respect to different random seeds followed by HadamardMultiDim and NoisySign. The most stable algorithms among all baselines seems to be RandKSpatialProj. OneBit and Sparsereg are slightly less stable



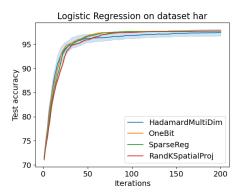


Figure 3: Performance of compressors for Logistic regression on HAR (Reyes-Ortiz et al., 2012) dataset. Left: Training Logistic loss, Right: Test Accuracy.

than it. However, all our algorithms except NoisySign are much more stable than most other baselines like SRQ, Correlated SRQ and RandK.

### 6 Conclusion

We proposed four communication-efficient collaborative compression schemes to obtain error guarantees in  $\ell_2$ -error (SparseReg),  $\ell_{\infty}$ -error (NoisySign, HadamardMultiDim) and cosine distance (OneBitAvg). The estimation error of our schemes improves with number of clients, and degrades with increasing dissimilarity between clients. Our schemes are biased and our dissimilarity metrics ( $\Delta_{\text{reg}}$ ,  $\Delta_{\text{Hadamard}}$ ) depend on the quantization levels. However, these drawbacks can be removed by using existing techniques for converting biased compressors to unbiased ones (Beznosikov et al., 2022), and adding noise before quantization (Tang et al., 2023; Chzhen & Schechtman, 2023). Schemes such as error feedback (Karimireddy et al., 2019) reduces the error of independent compressors, and it will be interesting to check if it works for our collaborative compressors.

**Acknowledgment.** This work is supported in part by NSF awards 2217058 (EnCORE Inst.), 2112665 (TILOS Inst.), and a grant from Adobe Research.

### References

Ahmad Ajalloeian and Sebastian U. Stich. Analysis of SGD with biased gradient estimators. CoRR, abs/2008.00051, 2020.

Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. QSGD: Communication-Efficient SGD via Gradient Quantization and Encoding. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: a meta-algorithm and applications. *Theory Comput.*, 8:121–164, 2012.

Pranjal Awasthi, Maria Florina Balcan, and Philip M. Long. The power of localization for efficiently learning linear separators with noise. *J. ACM*, 63(6), jan 2017. ISSN 0004-5411. doi: 10.1145/3006384.

Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd: Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*, pp. 560–569. PMLR, 2018a.

Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signSGD: Compressed optimisation for non-convex problems. In Jennifer Dy and Andreas Krause (eds.), *Proceedings* 

- of the 35th International Conference on Machine Learning, volume 80 of Proceedings of Machine Learning Research, pp. 560–569. PMLR, 10–15 Jul 2018b.
- Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, and Mher Safaryan. On Biased Compression for Distributed Learning, December 2022. arXiv:2002.12410 [cs, math, stat].
- Léon Bottou and Olivier Bousquet. The Tradeoffs of Large Scale Learning. In J. Platt, D. Koller, Y. Singer, and S. Roweis (eds.), *Advances in Neural Information Processing Systems*, volume 20. Curran Associates, Inc., 2007.
- Petros T Boufounos and Richard G Baraniuk. 1-bit compressive sensing. In 2008 42nd Annual Conference on Information Sciences and Systems, pp. 16–21. IEEE, 2008.
- Mark Braverman, Ankit Garg, Tengyu Ma, Huy L. Nguyen, and David P. Woodruff. Communication lower bounds for statistical estimation problems via a distributed data processing inequality. In *Proceedings of the Forty-Eighth Annual ACM Symposium on Theory of Computing*, STOC '16, pp. 1011–1020, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450341325. doi: 10.1145/2897518.2897582.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. Found. Trends Mach. Learn., 8(3–4): 231–357, November 2015. ISSN 1935-8237. doi: 10.1561/2200000050.
- Sebastian Caldas, Peter Wu, Tian Li, Jakub Konečný, H. Brendan McMahan, Virginia Smith, and Ameet Talwalkar. LEAF: A benchmark for federated settings. *CoRR*, abs/1812.01097, 2018.
- Seok-Ho Chang, Pamela C. Cosman, and Laurence B. Milstein. Chernoff-type bounds for the gaussian error function. *IEEE Transactions on Communications*, 59(11):2939–2944, 2011. doi: 10.1109/TCOMM.2011. 072011.100049.
- Xiangyi Chen, Tiancong Chen, Haoran Sun, Steven Z. Wu, and Mingyi Hong. Distributed Training with Heterogeneous Data: Bridging Median- and Mean-Based Algorithms. In *Advances in Neural Information Processing Systems*, volume 33, pp. 21616–21626. Curran Associates, Inc., 2020.
- Evgenii Chzhen and Sholom Schechtman. SignSVRG: fixing SignSGD via variance reduction, May 2023. arXiv:2305.13187 [math, stat].
- Peter Davies, Vijaykrishna Gurunanthan, Niusha Moshrefi, Saleh Ashkboos, and Dan Alistarh. New bounds for distributed mean estimation and variance reduction. In *International Conference on Learning Representations*, 2021.
- Venkata Gandikota, Daniel Kane, Raj Kumar Maity, and Arya Mazumdar. vqsgd: Vector quantized stochastic gradient descent. *IEEE Transactions on Information Theory*, 68(7):4573–4587, 2022. doi: 10.1109/TIT.2022.3161620.
- Ankit Garg, Tengyu Ma, and Huy Nguyen. On Communication Cost of Distributed Statistical Estimation and Dimensionality. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger (eds.), *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- Avishek Ghosh, Raj Kumar Maity, Swanand Kadhe, Arya Mazumdar, and Kannan Ramchandran. Communication-efficient and byzantine-robust distributed learning with error feedback. *IEEE Journal on Selected Areas in Information Theory*, 2(3):942–953, 2021.
- Venkatesan Guruswami and Prasad Raghavendra. Hardness of learning halfspaces with noise. In 2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06), pp. 543–552, 2006. doi: 10.1109/FOCS.2006.33.
- Divyansh Jhunjhunwala, Ankur Mallick, Advait Gadhikar, Swanand Kadhe, and Gauri Joshi. Leveraging spatial and temporal correlations in sparsified mean estimation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 14280–14292. Curran Associates, Inc., 2021.

- Shuli Jiang, Pranay Sharma, and Gauri Joshi. Correlation aware sparsified mean estimation using random projection. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Richeng Jin, Yufan Huang, Xiaofan He, Huaiyu Dai, and Tianfu Wu. Stochastic-Sign SGD for Federated Learning with Theoretical Guarantees, September 2021. arXiv:2002.10940 [cs, stat].
- Richeng Jin, Xiaofan He, Caijun Zhong, Zhaoyang Zhang, Tony Quek, and Huaiyu Dai. Magnitude Matters: Fixing SIGNSGD Through Magnitude-Aware Sparsification in the Presence of Data Heterogeneity, February 2023. arXiv:2302.09634 [cs].
- Adam Tauman Kalai, Adam R. Klivans, Yishay Mansour, and Rocco A. Servedio. Agnostically Learning Halfspaces. SIAM Journal on Computing, 37(6):1777–1805, January 2008. ISSN 0097-5397. doi: 10.1137/060649057. Publisher: Society for Industrial and Applied Mathematics.
- Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. Error feedback fixes SignSGD and other gradient compression schemes. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), Proceedings of the 36th International Conference on Machine Learning, volume 97 of Proceedings of Machine Learning Research, pp. 3252–3261. PMLR, 09–15 Jun 2019.
- Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. SCAFFOLD: Stochastic controlled averaging for federated learning. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 5132–5143. PMLR, 13–18 Jul 2020.
- Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: distributed machine learning for on-device intelligence. arXiv preprint arXiv:1610.02527, 2016.
- Jakub Konečný and Peter Richtárik. Randomized Distributed Mean Estimation: Accuracy vs. Communication. Frontiers in Applied Mathematics and Statistics, 4, 2018. ISSN 2297-4687.
- B. Laurent and P. Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5), October 2000. ISSN 0090-5364. doi: 10.1214/aos/1015957395. URL http://dx.doi.org/10.1214/aos/1015957395.
- Yann LeCun and Corinna Cortes. MNIST handwritten digit database. 2010.
- Xiang Li, Kaixuan Huang, Wenhao Yang, Shusen Wang, and Zhihua Zhang. On the convergence of fedavg on non-iid data. In *International Conference on Learning Representations*, 2020.
- Xiang Li, Shusen Wang, Kun Chen, and Zhihua Zhang. Communication-efficient distributed svd via local power iterations. In *International Conference on Machine Learning*, pp. 6504–6514. PMLR, 2021.
- Yingyu Liang, Maria-Florina Balcan, and Vandana Kanchanapally. Distributed pca and k-means clustering. 2013.
- Yingyu Liang, Maria-Florina F Balcan, Vandana Kanchanapally, and David Woodruff. Improved distributed principal component analysis. *Advances in neural information processing systems*, 27, 2014.
- S. Lloyd. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982. doi: 10.1109/TIT.1982.1056489.
- Prathamesh Mayekar, Ananda Theertha Suresh, and Himanshu Tyagi. Wyner-Ziv Estimators: Efficient Distributed Mean Estimation with Side-Information. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, pp. 3502–3510. PMLR, March 2021. ISSN: 2640-3498.
- Brendan McMahan and Daniel Ramage. Federated learning: Collaborative machine learning without centralized training data. https://research.googleblog.com/2017/04/federated-learning-collaborative.html, 2017.

- H Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. arXiv preprint arXiv:1602.05629, 2016.
- Jorge Reyes-Ortiz, Davide Anguita, Alessandro Ghio, Luca Oneto, and Xavier Parra. Human Activity Recognition Using Smartphones. UCI Machine Learning Repository, 2012. DOI: https://doi.org/10.24432/C54S4K.
- Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. The Annals of Mathematical Statistics, 22(3):400 407, 1951. doi: 10.1214/aoms/1177729586.
- Mher Safaryan and Peter Richtarik. Stochastic Sign Descent Methods: New Algorithms and Better Theory. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 9224–9234. PMLR, July 2021. ISSN: 2640-3498.
- Mher Safaryan, Egor Shulgin, and Peter Richtárik. Uncertainty principle for communication compression in distributed and federated learning and the search for an optimal compressor. *Information and Inference: A Journal of the IMA*, 11(2):557–580, 04 2021. ISSN 2049-8772. doi: 10.1093/imaiai/iaab006.
- Rocco A. Servedio. Perceptron, winnow, and pac learning. *SIAM Journal on Computing*, 31(5):1358–1369, 2002. doi: 10.1137/S0097539798340928.
- Jie Shen. Pac learning of halfspaces with malicious noise in nearly linear time. In Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent (eds.), Proceedings of The 26th International Conference on Artificial Intelligence and Statistics, volume 206 of Proceedings of Machine Learning Research, pp. 30–46. PMLR, 25–27 Apr 2023.
- Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified SGD with Memory. In Advances in Neural Information Processing Systems, volume 31. Curran Associates, Inc., 2018.
- Ananda Theertha Suresh, Felix X. Yu, Sanjiv Kumar, and H. Brendan McMahan. Distributed Mean Estimation with Limited Communication. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 3329–3337. PMLR, July 2017. ISSN: 2640-3498.
- Ananda Theertha Suresh, Ziteng Sun, Jae Ro, and Felix Yu. Correlated Quantization for Distributed Mean Estimation and Optimization. In *Proceedings of the 39th International Conference on Machine Learning*, pp. 20856–20876. PMLR, June 2022. ISSN: 2640-3498.
- Rafał Szlendak, Alexander Tyurin, and Peter Richtárik. Permutation compressors for provably faster distributed nonconvex optimization, 2021.
- Zhiwei Tang, Yanmeng Wang, and Tsung-Hui Chang. \$z\$-SignFedAvg: A Unified Stochastic Sign-based Compression for Federated Learning. February 2023.
- Joaqun Torres-Sospedra, Raul Montoliu, Adolfo Martnez-Us, Tomar Arnau, and Joan Avariento. UJIIndoorLoc. UCI Machine Learning Repository, 2014. DOI: https://doi.org/10.24432/C5MS59.
- Shay Vargaftik, Ran Ben-Basat, Amit Portnoy, Gal Mendelson, Yaniv Ben-Itzhak, and Michael Mitzenmacher. DRIVE: One-bit Distributed Mean Estimation. In *Advances in Neural Information Processing Systems*, volume 34, pp. 362–377. Curran Associates, Inc., 2021.
- Ramji Venkataramanan, Antony Joseph, and Sekhar Tatikonda. Lossy Compression via Sparse Linear Regression: Performance Under Minimum-Distance Encoding. *IEEE Transactions on Information Theory*, 60(6):3254–3264, June 2014a. ISSN 1557-9654. doi: 10.1109/TIT.2014.2313085.
- Ramji Venkataramanan, Tuhin Sarkar, and Sekhar Tatikonda. Lossy Compression via Sparse Linear Regression: Computationally Efficient Encoding and Decoding. *IEEE Transactions on Information Theory*, 60(6):3265–3278, June 2014b. ISSN 1557-9654. doi: 10.1109/TIT.2014.2314676.
- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018.

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 *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17, pp. 1508–1518, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.

Heng Zhu, Avishek Ghosh, and Arya Mazumdar. Optimal compression of unit norm vectors in the high distortion regime. In 2023 IEEE International Symposium on Information Theory (ISIT), pp. 719–724. IEEE, 2023.

### A Proofs for Section 2

### A.1 Proof of Theorem 1

As all operations are coordinate-wise, we restrict our focus to only a single dimension  $j \in [d]$ .

$$\mathbb{E}_{\xi_i}[\tilde{b}_i^{(j)}] = \Phi_{\sigma}(g_i^{(j)}), \, \forall i \in [m]$$

Note that  $\Phi_{\sigma}(t) = \text{erf}(\frac{t}{\sqrt{2}\sigma})$  and  $\Phi_{\sigma}^{-1}(t) = \sqrt{2}\sigma \,\text{erf}^{-1}(t)$ . By Hoeffding's inequality for bounded random variables we have,

$$\Pr[|\frac{1}{m}\sum_{i=1}^{m}(\tilde{b_i}^{(j)} - \Phi_{\sigma}(g_i^{(j)}))| \ge t] \le 2e^{-\frac{mt^2}{2}}$$

If we set  $t = \sqrt{\frac{2c \log(m)}{m}}$ , for some c > 0 in the above inequality, then with probability  $1 - 2m^{-c}$ , we have,

$$\left|\frac{1}{m}\sum_{i=1}^{m}(\tilde{b_i}^{(j)} - \Phi_{\sigma}(g_i^{(j)}))\right| \le t$$

We can represent  $\frac{1}{m}\sum_{i=1}^{m} \tilde{b}_i = \Phi_{\sigma}(\tilde{g})$ , as  $\Phi_{\sigma}$  is an invertible function. To find the difference between  $\tilde{g}$  and g, we find the difference  $\Phi_{\sigma}(\tilde{g}) - \Phi_{\sigma}(g)$ . With probability  $1 - 2m^{-c}$ , we have,

$$|\Phi_{\sigma}(\tilde{g}^{(j)}) - \Phi_{\sigma}(g^{(j)})| \le \frac{1}{m} \sum_{i=1}^{m} |\Phi_{\sigma}(g_i^{(j)}) - \Phi_{\sigma}(g^{(j)})| + t$$

To remove the terms of  $\Phi_{\sigma}$ , we can apply the function  $\Phi_{\sigma}^{-1}$  on  $\tilde{g}^{(j)}$ . As  $\Phi_{\sigma}^{-1}$  is not Lipschitz, we need to perform its Taylor's expansion around  $\Phi_{\sigma}(g^{(j)})$  to account for the linear terms in the error. If  $\Delta_{\Phi} = \frac{1}{m} \sum_{i=1}^{m} |\Phi_{\sigma}(g_{i}^{(j)}) - \Phi_{\sigma}(g^{(j)})|$ , then we obtain,

$$|\tilde{g}^{(j)} - g^{(j)}| \le \max_{u \in [\Phi_{\sigma}(g^{(j)}) - \Delta_{\Phi} - t, \Phi_{\sigma}(g^{(j)}) + \Delta_{\Phi} + t]} |(\Phi_{\sigma}^{-1})'(u)|(\Delta_{\Phi} + t)$$
(12)

We now obtain an appropriate upper bound on  $(\Phi_{\sigma}^{-1})'(u)$  as we do not have a closed-form expression for it. We will use the properties of erf to obtain a suitable bound. First, note that  $\Phi_{\sigma}$  and  $\Phi_{\sigma}^{-1}$  are both odd functions, therefore,  $|\Phi^{-1}(u)| = |\Phi^{-1}(|u|)|$ , so we consider the bound for u > 0. Note that  $(\Phi^{-1})'(u) = \frac{1}{\Phi'(\Phi^{-1}(u))}$ . For u > 0, we have,

$$1 - \operatorname{erf}(u) \le e^{-u^2}$$

$$\operatorname{erf}(u) \ge 1 - e^{-u^2}$$

$$\operatorname{erf}^{-1}(u) \le \sqrt{-\log(1 - u)}$$

$$\Phi_{\sigma}^{-1}(u) = \sqrt{2}\sigma \operatorname{erf}^{-1}(u) \le \sigma \sqrt{-2\log(1 - u)}$$

$$(\Phi_{\sigma}^{-1})'(u) = \sqrt{\frac{\pi}{2}} e^{(\Phi_{\sigma}^{-1}(u))^2/(2\sigma^2)} \le \sqrt{\frac{\pi}{2}} e^{-2\log(1 - u)/2} = \sqrt{\frac{\pi}{2}} \frac{1}{1 - u}$$

For the first step, we use an upper bound on the complementary error function. For the third step, we use the fact that if  $f(x) \le g(x)$ , then  $f^{-1}(y) \ge g^{-1}(y)$ .

Using the following upper bound in Eq (12), we obtain,

$$|\tilde{g}^{(j)} - g^{(j)}| \leq \max_{u \in [\Phi_{\sigma}(g^{(j)}) - \Delta_{\Phi} - t, \Phi_{\sigma}(g^{(j)}) + \Delta_{\Phi} + t]} \sqrt{\frac{\pi}{2}} \frac{\Delta_{\Phi} + t}{1 - |u|}$$

$$\leq \sqrt{\frac{\pi}{2}} \frac{\Delta_{\Phi} + t}{1 - \max\{|\Phi_{\sigma}(g^{(j)}) - \Delta_{\Phi} - t|, |\Phi_{\sigma}(g^{(j)}) + \Delta_{\Phi} + t|\}}$$

We use  $\max\{|\Phi_{\sigma}(g^{(j)}) - \Delta_{\Phi} - t|, |\Phi_{\sigma}(g^{(j)}) + \Delta_{\Phi} + t|\} \leq \Phi_{\sigma}(|g^{(j)}|) + \Delta_{\Phi} + t$ , as  $\Phi_{\sigma}$  is an increasing odd function.

$$|\tilde{g}^{(j)} - g^{(j)}| \le \sqrt{\frac{\pi}{2}} \left( \left( 1 - \frac{\Delta_{\Phi} + t}{1 - \Phi_{\sigma}(|g^{(j)}|)} \right)^{-1} - 1 \right)$$

We extend the bound to d dimensions by taking a union bound, yielding a probability of error  $2dm^{-c}$ .

### A.2 Proof of Equation (3)

First, note that the term of  $(1+x)^{-1}-1=\frac{x}{1-x}\leq \mathcal{O}(x)$  when  $x\leq c$  for some constant  $c\in[0,1]$ . The conditions required for this are  $\Delta_{\Phi}$  to be small and both m and  $\alpha(||g||_{\infty})$  to be large. Note that  $\alpha(||g||_{\infty})=1-\Phi_{\sigma}(||g||_{\infty})$ , so it is small when the ratio  $\frac{||g||_{\infty}}{\sigma}$  is small. The bound on estimation error after this step is the following.

$$||\tilde{g} - g||_{\infty} = \mathcal{O}\left(\frac{\Delta_{\Phi} + \sqrt{\frac{\log m}{m}}}{\alpha(||g||_{\infty})}\right)$$

Now, we lower bound the term  $\alpha(||g||_{\infty})$ . Note that  $\alpha(x) = 1 - \Phi_{\sigma}(x) = \frac{1}{2}(1 - \text{erf}(\frac{x}{\sqrt{2}\sigma}))$  for any  $x \in \mathbb{R}_+$  where erf is the error function. Here, 1 - erf(x) is the complementary error function for any  $x \ge 0$  and it is lower bounded by  $\Omega(e^{-\beta'x^2})$  for  $\beta' > 1$  (Chang et al., 2011). Plugging in this lower bound, we complete the proof for a  $\beta = \beta'/2 > \frac{1}{2}$ .

### A.3 Proof of Equation (4)

The proof follows from using the triangle inequality and a Taylor's expansion for each  $\Phi_{\sigma}(g_i^{(j)})$  around  $g^{(j)}$ . Note that, for some  $u_i^{(j)}$  between  $g^{(j)}$  and  $g_i^{(j)}$ , we have,

$$\begin{split} &\Phi_{\sigma}(g_{i}^{(j)}) = \Phi_{\sigma}(g^{(j)}) + \sqrt{\frac{2}{\pi}} \frac{(g^{(j)} - g_{i}^{(j)})e^{-\frac{(u_{i}^{(j)})^{2}}{2\sigma^{2}})}}{\sigma} \\ &|\Phi_{\sigma}(g_{i}^{(j)}) - \Phi_{\sigma}(g^{(j)})| \leq \sqrt{\frac{2}{\pi}} \frac{|g^{(j)} - g_{i}^{(j)}|}{\sigma} \end{split}$$

We use the fact that  $e^{-\frac{(u_i^{(j)})^2}{2\sigma^2}} \leq 1$ . By using triangle inequality for any coordinate  $j \in [m]$ , we obtain,

$$\Delta_{\Phi} \leq \max_{j \in [d]} \frac{1}{m} \sum_{i \in [m]} |\Phi_{\sigma}(g_i^{(j)}) - \Phi_{\sigma}(g^{(j)})| \leq \frac{1}{m} \sum_{i \in [m]} \max_{j \in [d]} |\Phi_{\sigma}(g_i^{(j)}) - \Phi_{\sigma}(g^{(j)})|$$
$$\leq \sqrt{\frac{2}{\pi}} \frac{1}{m} \sum_{i \in [m]} \max_{j \in [d]} \frac{|g^{(j)} - g_i^{(j)}|}{\sigma} \leq \sqrt{\frac{2}{\pi}} \frac{1}{m} \sum_{i \in [m]} \frac{||g - g_i||_{\infty}}{\sigma}$$

## B Proofs for Section 3

### B.1 Proof of Theorem 2

To prove this Lemma, we first consider a single dimension  $j \in [d]$ . Note that the collaborative compression scheme has precision  $\frac{B}{2^{m-1}}$ . Since, for each client  $g_i^{(j)}$  lies in the range  $[g^{(j)} - \Delta_{\max}, g^{(j)} + \Delta_{\max}]$ , therefore, we must decode one of the elements which has error at most  $\Delta_{\max} + \frac{B}{2^{m-1}}$ . Since this is true for each coordinate individually, therefore, it holds for the  $\ell_{\infty}$  norm.

Consider a single dimension  $j \in [d]$ . Let  $g_i^{(j)}$  be the  $j^{th}$  coordinate of  $g_i$  and  $\rho_j$  be the permutation selected for the coordinate j. We omit j from  $g_i^{(j)}$  and  $\rho_j$  to simplify the notation. Let  $\tilde{b}_{i,p}$  be the estimate of  $g_i$  after decoding it for p levels where  $p \in [m]$ . Therefore, the estimator  $\tilde{g} = \sum_{i=1}^m \frac{\tilde{b}_{i,\rho_i}B}{2^{\rho_i-1}}$ . Let  $\tilde{g}_i = \sum_{k=1}^m \frac{\tilde{b}_{i,k}B}{2^{k-1}}$  be the decoded value of  $g_i$  till level m and  $\bar{g} = \frac{1}{m} \sum_{i=1}^m \tilde{g}_i = \sum_{k=1}^m \frac{\bar{b}_{k}B}{2^{k-1}}$ , where  $\bar{b}_k = \frac{1}{m} \sum_{i=1}^m \tilde{b}_{i,k}$ .

We compute the expected error for coordinate j, where the expectation is wrt the permutation  $\rho_j$ . Note that  $\mathbb{E}_{\rho \sim \Pi_m}[\tilde{g}_i] = \bar{g}$ .

$$\mathbb{E}_{\rho \sim \Pi_m}[|g - \tilde{g}|] = \sqrt{(\mathbb{E}_{\rho \sim \Pi_m}[|g - \tilde{g}|])^2} \leq \sqrt{\mathbb{E}_{\rho \sim \Pi_m}|g - \tilde{g}|^2} \leq \sqrt{\mathbb{E}_{\rho \sim \Pi_m}|\tilde{g} - \bar{g}|^2 + |g - \bar{g}|^2} \\
\leq \sqrt{\mathbb{E}_{\rho \sim \Pi_m}|\tilde{g} - \bar{g}|^2} + |g - \bar{g}| \leq \frac{1}{m} \sum_{i=1}^m |g_i - \tilde{g}_i| + \sqrt{\mathbb{E}_{\rho \sim \Pi_m}|\tilde{g} - \bar{g}|^2} \\
\leq \frac{B}{2^{m-1}} + \sqrt{\mathbb{E}_{\rho \sim \Pi_m}|\tilde{g} - \bar{g}|^2}$$

We use Jensen's inequality for the first inequality. For the second inequality, we use bias-variance decomposition for the random variable  $\tilde{g}$ , where the first term is its variance, and the second term is its bias wrt the term g. We then use  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for any  $a,b \geq 0$ . To handle the term  $|g-\bar{g}|$ , we expand both terms as a summation over m clients, followed by a triangle inequality. As each estimator  $\tilde{g}_i$  is at least  $\frac{B}{2^{m-1}}$  away from  $g_i$ , each term in the difference  $|g_i - \tilde{g}_i|$  has the upperbound  $\frac{B}{2^{m-1}}$ .

We now bound the variance term separately. Note that

$$\mathbb{E}_{\rho \sim \Pi_m} |\tilde{g} - \bar{g}|^2 = \mathbb{E}_{\rho \sim \Pi_m} |\tilde{g}|^2 - \bar{g}^2$$

We first evaluate the second moment  $\mathbb{E}_{\rho \sim \Pi_m} |\tilde{g}|^2$ .

$$\begin{split} \mathbb{E}_{\rho \sim \Pi_{m}} |\tilde{g}|^{2} &= \mathbb{E}_{\rho \sim \Pi_{m}} \left| \sum_{i=1}^{m} \frac{\tilde{b}_{i,\rho_{i}}}{2^{\rho_{i}-1}} \right|^{2} = \sum_{i=1}^{m} \mathbb{E}_{\rho \sim \Pi_{m}} \left[ \frac{\tilde{b}_{i,\rho_{i}}^{2}}{2^{2\rho_{i}-2}} \right] + B^{2} \sum_{1 \leq i \neq j \leq m} \mathbb{E}_{\rho \sim \Pi_{m}} \left[ \frac{\tilde{b}_{i,\rho_{i}}}{2^{\rho_{i}-1}} \frac{\tilde{b}_{i,\rho_{j}}}{2^{\rho_{j}-1}} |\rho_{i} \right] \\ &= \sum_{k=1}^{m} \frac{B^{2}}{2^{2k-2}} + B^{2} \sum_{1 \leq i \neq j \leq m} \mathbb{E}_{\rho_{i}} \left[ \frac{\tilde{b}_{i,\rho_{i}}}{2^{\rho_{i}-1}} \frac{\tilde{b}_{i,\rho_{j}}}{2^{\rho_{i}-1}} \frac{\tilde{b}_{i,\rho_{j}}}{2^{\rho_{j}-1}} |\rho_{i} \right] \\ &= \sum_{k=1}^{m} \frac{B^{2}}{2^{2k-2}} + B^{2} \sum_{1 \leq i \neq j \leq m} \mathbb{E}_{\rho_{i}} \left[ \frac{\tilde{b}_{i,\rho_{i}}}{2^{\rho_{i}-1}} \frac{1}{m-1} \sum_{l=1,l \neq \rho_{i}} \frac{\tilde{b}_{j,l}}{2^{l-1}} \right] \\ &= \sum_{k=1}^{m} \frac{B^{2}}{2^{2k-2}} + \frac{B^{2}}{m(m-1)} \sum_{1 \leq i \neq j \leq m} \sum_{k=1}^{m} \left[ \frac{\tilde{b}_{i,k}}{2^{k-1}} \sum_{l=1,l \neq \rho_{i}} \frac{\tilde{b}_{j,l}}{2^{l-1}} \right] \\ &= \sum_{k=1}^{m} \frac{B^{2}}{2^{2k-2}} + \frac{1}{m(m-1)} \sum_{1 \leq i \neq j \leq m} \sum_{k=1}^{m} \frac{\tilde{b}_{i,k}\tilde{b}_{j,k}}{2^{2k-2}} \\ &= \sum_{k=1}^{m} \frac{B^{2}}{2^{2k-2}} + \frac{1}{m(m-1)} \sum_{1 \leq i \neq j \leq m} \sum_{k=1}^{m} \frac{B^{2}\tilde{b}_{i,k}\tilde{b}_{j,k}}{2^{2k-2}} \\ &= \sum_{k=1}^{m} \frac{B^{2}}{2^{2k-2}} + \frac{1}{m(m-1)} \sum_{1 \leq i \neq j \leq m} \sum_{k=1}^{m} \frac{B^{2}(|\tilde{b}_{i,k}|^{2} + |\tilde{b}_{j,k}|^{2} - 2\tilde{b}_{i,k}\tilde{b}_{j,k})}{2^{2k-1}} \\ &= \frac{m^{2}|\tilde{g}|^{2} - \sum_{i=1}^{m}|\tilde{g}_{i}|^{2}}{m(m-1)} + \frac{1}{2m(m-1)} \sum_{1 \leq i \neq j \leq m} \sum_{k=1}^{m} \left( \frac{B(\tilde{b}_{i,k} - \tilde{b}_{j,k})}{2^{2k-1}} \right)^{2} \end{aligned}$$

Note that we expand the square of the sum of terms where  $\tilde{b}_{i,j}^2 = 1$ . For the second term, we use the law of total expectation by conditioning on the value of  $\rho_i$ . To evaluate the inner expectation, we note that  $\rho_j$  can take any value other than that of  $\rho_i$  with equal probability. To evaluate the outer expectation, note that  $\rho_i$  can take any value in [m] with equal probability. In the fourth equation, we subtract the term where l = k. Then, we can factorize the remaining terms to obtain  $\tilde{g}_i$  and  $\tilde{g}_j$ . Note that the sum of the product terms  $\tilde{g}_i\tilde{g}_j$  can be expressed as  $|\sum_{i=1}^m \tilde{g}_i|^2$ , with the square terms subtracted. Further, we express the term  $\frac{B^2}{2^{2k-2}} = \sum_{1 \le i \ne j \le m} \frac{B^2(|\tilde{b}_{i,k}|^2 + |\tilde{b}_{j,k}|^2)}{2^{2k-1}}$  as  $|\tilde{b}_{i,k}|^2 = 1$ . Finally, we complete the squares for each term k.

Using the above value of second moment  $\mathbb{E}_{\rho \sim \Pi_m} |\tilde{g}|^2$ , we can compute the variance,

$$\mathbb{E}_{\rho \sim \Pi_m} |\tilde{g} - \bar{g}|^2 = \mathbb{E}_{\rho \sim \Pi_m} |\tilde{g}|^2 - |\bar{g}|^2 = \frac{|\bar{g}|^2 - \frac{1}{m} \sum_{i=1}^m |\tilde{g}_i|^2}{m-1} + \frac{1}{2m(m-1)} \sum_{1 \le i \ne j \le m} \sum_{k=1}^m \left( \frac{B(\tilde{b}_{i,k} - \tilde{b}_{j,k})}{2^{k-1}} \right)^2$$

$$= \frac{1}{2m^2} \sum_{1 \le i \ne j \le m} \sum_{k=1}^m \left( \frac{B(\tilde{b}_{i,k} - \tilde{b}_{j,k})}{2^{k-1}} \right)^2$$

We use 
$$\bar{g}^2 \le \frac{1}{m} \sum_{i=1}^m |\tilde{g}_i|^2 = \frac{1}{2m^2} \sum_{1 \le i \ne j \le m} (\tilde{g}_i - \tilde{g}_j)^2 \ge \frac{1}{2m^2} \sum_{1 \le i \ne j \le m} \sum_{k=1}^m \left( \frac{B(\tilde{b}_{i,k} - \tilde{b}_{j,k})}{2^{k-1}} \right)^2$$
.

To simplify this bound, we need to incorporate the difference in the actual gradient vectors. For this purpose, we try to bound the differences  $|\tilde{b}_{i,k} - \tilde{b}_{j,k}|$  in terms of  $\Delta_{ij} \triangleq |g_i - g_i|$ . If  $\Delta_{ij} = |g_i - g_j|$ , then  $\tilde{b}_{i,k} = \tilde{b}_{j,k}, \forall k \geq \log\left(\frac{B}{\Delta_{ij}}\right)$ .

### B.2 Proof for Theorem 3

We first show how to obtain the coefficients  $c_i$  and the bound on the term  $\Lambda$  in Eq (7). To obtain the coefficients  $c_i$ , we replace set  $L=m, n=d, R=\frac{m\log L}{d}$  and  $\sigma^2=\frac{B^2}{d}$  in (Venkataramanan et al., 2014b, Eq 2). Note that the terms  $\delta_1$  and  $\delta_2$  can be directly applied in (Venkataramanan et al., 2014b, Section V) with the same values of  $\epsilon_k$  and  $\gamma_k$ . Note that the bound on  $\Lambda$  in (Venkataramanan et al., 2014b, Theorem 1) is for gaussian sources with a fixed variance, however, their proof in (Venkataramanan et al., 2014b, Section V) only uses the fact that the squared  $\ell_2$  norm of the source is bounded, which is exactly our case from Assumption 2.

The proof of Eq (7) is same as Theorem 2 for a single dimension, with the coefficients  $\frac{B}{2^{j-1}}$  replaced by  $c_j$  and  $\tilde{b}_{i,k}^{(r)}$  replaced by  $A_{(k-1)L+\tilde{b}_{i,k}}$ . Following Appendix B.1, we can write down the  $\ell_2$  error.

$$\mathbb{E}_{\rho \sim \Pi_m}[||\tilde{g} - g||_2^2] = \mathbb{E}_{\rho \sim \Pi_m}[||g - \mathbb{E}_{\rho \sim \Pi_m}[\tilde{g}]||_2^2] + \mathbb{E}_{\rho \sim \Pi_m}[||\tilde{g} - \mathbb{E}_{\rho \sim \Pi_m}[\tilde{g}]||_2^2]$$

 $\mathbb{E}[\tilde{g}] = \bar{g} = \frac{1}{m} \sum_{i=1}^{m} \bar{g}_i, \text{ where } \bar{g}_i = \sum_{j=1}^{m} c_j A_{(j-1)L + \tilde{b}_{i,j}}. \text{ By triangle inequality, the first term is } \frac{1}{m} \sum_{i=1}^{m} ||g_i - \bar{g}_i||_2^2, \text{ which is bounded individually by } \Lambda = B^2 (1 + \frac{10 \log L}{d} e^{\frac{m \log L}{d}} (\delta_1 + \delta_2))^2 \left(1 - \frac{2 \log L}{d}\right)^m \text{ by setting } L = m, n = d, R = \frac{m \log L}{d}, \sigma^2 = \frac{B^2}{d} \text{ and } \delta_0 = 0 \text{ in (Venkataramanan et al., 2014b, Theorem 1)}.$ 

For the second term, we need to bound  $\mathbb{E}[||\tilde{g}||_2^2]$ .

$$\begin{split} \mathbb{E}[||\tilde{g}||_{2}^{2}] &= \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{m} c_{i}^{2} ||A_{(j-1)L+\tilde{b}_{i,j}}||_{2}^{2} \\ &+ \sum_{1 \leq i \neq j \leq m} \mathbb{E}_{\rho \sim \Pi_{m}} \left[ c_{\pi(i)} c_{\pi(j)} \langle A_{(\pi(i)-1)L+\tilde{b}_{i,\pi(i)}}, A_{(\pi(j)-1)L+\tilde{b}_{j,\pi(j)}} \rangle \right] \\ &= \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{m} c_{i}^{2} ||A_{(j-1)L+\tilde{b}_{i,j}}||_{2}^{2} \\ &+ \frac{1}{m(m-1)} \sum_{1 \leq i \neq j \leq m} \mathbb{E}_{\rho \sim \Pi_{m}} \left[ c_{\pi(i)} c_{\pi(j)} \langle A_{(\pi(i)-1)L+\tilde{b}_{i,\pi(i)}}, A_{(\pi(j)-1)L+\tilde{b}_{j,\pi(j)}} \rangle \right] \\ &= \frac{m^{2} ||\bar{g}||_{2}^{2} - \sum_{i=1}^{m} ||\tilde{g}_{i}||_{2}^{2}}{m(m-1)} + \frac{1}{m(m-1)} \sum_{1 \leq i \neq j \leq m} \sum_{k=1}^{m} c_{k}^{2} ||A_{(k-1)L+\tilde{b}_{j,k}} - A_{(k-1)L+\tilde{b}_{i,k}}||_{2}^{2} \end{split}$$

The remainder of the proof follows proof of Theorem 2 with  $|\cdot|^2$  replaced by  $||\cdot||_{2}^2$ .

For the first equation in Eq (7), since the collaborative compression scheme has precision  $\Lambda$ , for each client, the vector  $g_i$  lies in a ball around g of radius  $\Delta_{2,\max} = \max_{i \in [m]} ||g_i - g||_2$ . Therefore, if we use sparse regression codes to individually encode each vector  $g_i$ , it's decoded value  $\tilde{g}_i$  must lie either in the same ball or be at most  $\Lambda$   $\ell_2$  distance away from some vector in this ball. Since this set is convex, and the vector  $\bar{g}$  is also obtained as a convex combination of vectors inside this ball or at most  $\Lambda$  away from it, the vector  $\bar{g}$  also lies inside this ball or is at most  $\Lambda$  away from it.

This proves Eq (7).

### B.3 Proof of Equation (8) and Remark 1

To bound  $\Delta_{reg}$ , we need to bound the terms of the form  $||A_{(k-1)L+\tilde{b}_{i,k}}-A_{(k-1)L+\tilde{b}_{j,k}}||_2^2$  for clients  $i \neq j \in [m]$  and at level  $k \in [m]$ . Note that if  $\tilde{b}_{j,k} = \tilde{b}_{i,k}$ , then this difference is 0. If  $\tilde{b}_{j,k} \neq \tilde{b}_{i,k}$ , then this is  $\Gamma^2$  from the

definition of  $(\delta_1, \delta_2, \Gamma)$ -Cover. This results in the following upper bound for  $\Delta_{reg}$ .

$$\Delta_{reg} \leq \frac{\Gamma^{2}}{m} \sum_{i,j \in [m], i \neq j} \sum_{k=1}^{m} c_{k}^{2} \cdot \mathbf{1}(\tilde{b}_{i,k} \neq \tilde{b}_{j,k}) = \frac{2\Gamma^{2}B^{2} \log L}{d^{2}m} \sum_{i,j \in [m], i \neq j} \sum_{k=1}^{m} \left(1 - \frac{2 \log L}{d}\right)^{k-1} \mathbf{1}(\tilde{b}_{i,k} \neq \tilde{b}_{j,k})$$

$$= \frac{2\Gamma^{2}B^{2} \log L}{dm^{2}} \sum_{i,j \in [m], i \neq j} \sum_{k=1}^{m} \left(1 - \frac{2 \log L}{d}\right)^{k-1} \mathbf{1}(\tilde{b}_{i,k} \neq \tilde{b}_{j,k})$$

We use the value of  $c_k$  in the second inequality and then simplify the terms.

To prove the bound of  $\Gamma$  for the case of Gaussian A in Remark 1, note that  $\Gamma$  is an upper bound on the the squared  $\ell_2$  norm of the difference of two gaussian random vectors. For any  $r_1, r_2 \in [L]$  with  $r_1 \neq r_2$ ,  $A_{(k-1)L+r_1} - A_{(k-1)L+r_2} \sim \mathcal{N}(0, 2\mathbb{I}_d)$  as each vector is d-dimensional, they are sampled independently, and each element in each vector is sampled from a unit normal. Therefore,  $\frac{1}{2}||A_{(k-1)L+r_1} - A_{(k-1)L+r_2}||_2^2$  is a  $\chi_d^2$  random variable. Note that the scaling of  $\frac{1}{2}$  is to normalize the covariance to  $\mathbb{I}_d$ . By concentration of  $\chi_d^2$  random variables from (Laurent & Massart, 2000), we know that, with probability  $1 - \delta$ ,

$$||A_{(k-1)L+r_1} - A_{(k-1)L+r_2}||_2^2 \le 2d + 4\sqrt{d\log(1/\delta)} + 4\log(1/\delta)$$

Since  $r_1, r_2 \in [L]$ , we can take a union bound over all  $\frac{L(L-1)}{2}$  unique pairs of  $(r_1, r_2)$ . Further, we can take a union bound over all levels  $k \in [m]$ . These would set the probability of error to  $\frac{mL(L-1)\delta}{2}$ . Setting  $\delta_3$  to this value, we obtain  $\delta = \frac{2\delta_3}{mL(L-1)}$  and the following bound  $\forall k \in [m], r_1 \neq r_2 \in [L]$ ,

$$||A_{(k-1)L+r_1} - A_{(k-1)L+r_2}||_2^2 \le 2d + 4\sqrt{d\log(\frac{mL^2}{\delta_3})} + 4\log(\frac{mL^2}{\delta_3}) \triangleq \Gamma^2$$

We upper bound  $\frac{L(L-1)}{2}$  by  $L^2$ .

### B.4 Improving HadamardMultiDim and SparseReg with Repetitions

Let  $\tilde{g}_R = \frac{1}{R} \sum_{r=1}^R \tilde{g}(r)$  be estimator after R repetitions where  $\tilde{g}(r)$  is the estimator after 1 repetition of either HadamardMultiDim or SparseReg. To compute the error for both Theorem 2 and Theorem 3, we use a bias-variance decomposition of the squared estimation error. Consider SparseReg where we bound  $\mathbb{E}_{\rho \sim \Pi_m}[||\tilde{g} - g||_2^2]$ . Note that even with the new definition of  $\tilde{g}$ , it's mean wrt  $\rho$  remains the same as  $\mathbb{E}_{\rho \sim \Pi_m}[\tilde{g}] = \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{\rho \sim \Pi_m}[\tilde{g}(r)] = \frac{1}{R} \sum_{r=1}^R \bar{g} = \bar{g}$ . This contributes to the term with exponential dependence on m and it remains the same. Note that the variance term, however, decreases with R.

$$\mathbb{E}_{\rho \sim \Pi_m}[||\tilde{g} - \bar{g}||_2^2] = \frac{1}{R^2} \sum_{r=1}^R \mathbb{E}_{\rho \sim \Pi_m}[||\tilde{g}(r) - \bar{g}||_2^2] = \frac{1}{R} \mathbb{E}_{\rho \sim \Pi_m}[||\tilde{g}(1) - \bar{g}||_2^2] = \frac{\Delta_{\text{reg}}}{R}$$

Similarly for HadamardMultiDim, the term  $\Delta_{\text{Hadamard}}$  is replaced by  $\frac{\Delta_{\text{Hadamard}}}{\sqrt{R}}$ , as we take square root of the variance in it's analysis. Note that the  $\Delta_{2,\text{max}}$  and  $\Delta_{\infty,\text{max}}$  terms are not affected by R as it's analysis does not use  $\rho$ .

### C Proofs for Section 4

### C.1 Proof of Lemma 1

To prove this Lemma, note that  $\tilde{b}_i = sign(\langle g_i, z_i \rangle) \neq sign(\langle g, z_i \rangle)$  only if  $z_i$  is sampled from the symmetric difference of  $g_i$  and g. The probability that a  $z_i$  sampled uniformly from  $\mathbb{S}^{d-1}$  lies in this symmetric difference is given by  $\arccos(\langle g, g_i \rangle)/\pi$ . If we set  $\Delta_{\text{corr}} = \frac{1}{m\pi} \sum_{i \in [m]} \arccos(\langle g, g_i \rangle)$ 

Let  $\zeta$  be the fraction of  $z_i$  such that  $\tilde{b}_i \neq sign(\langle g, z_i \rangle)$ . Then, by Chernoff bound, we have,

$$\Pr[\zeta \ge (1+\gamma)\Delta_{\text{corr}}] \le e^{-\frac{\gamma^2 m \Delta_{\text{corr}}}{2+\gamma}}$$

By setting  $\gamma$  to be any small constant, we obtain, with probability  $1 - \mathcal{O}(\exp(-m\Delta_{\text{corr}}))$ , at most  $\zeta = \Theta(\Delta_{\text{corr}})$  fraction of datapoints are not generated from the halfspace with normal g and are thus corrupted.

### C.2 Proofs of Theorem 4 and 5

To prove Theorem 4, we utilize the guarantees of (Shen, 2023, Theorem 3), where the sample complexity requirement ensures that the error is  $\tilde{O}(\frac{d}{m})$ . Further, (Shen, 2023, Theorem 3) obtains error guarantee linear in the noise rate of the samples which is obtained from Lemma 1. The error guarantee is in terms of the symmetric difference between  $\tilde{g}$  and g wrt the uniform distribution on the unit sphere. Since this is equal to the angle between these two vectors divided by  $\pi$ , this gives us a bound on the inner product of these two unit vectors.

To prove Theorem 5, from (Kalai et al., 2008, Theorem 12), the sample complexity provides the term  $\frac{d}{\sqrt{m}}$  while the noise tolerance provides the term  $\sqrt{d}\Delta_{\text{corr}}$ .

# C.3 Proof of Equation (11)

To prove this remark, note that arccos(x) is concave for  $x \ge 0$ . Therefore, by applying Jensen's inequality, we obtain,

$$\Delta_{\text{corr}} = \frac{1}{m\pi} \sum_{i \in [m]} \arccos(\langle g_i, g \rangle) \le \frac{1}{\pi} \arccos\left(\langle \frac{1}{m} \sum_{i=1}^m g_i, g \rangle\right) = \frac{1}{\pi} \arccos\left(||\frac{1}{m} \sum_{i=1}^m g_i||_2 \langle g, g \rangle\right)$$

$$\le \frac{1}{\pi} \arccos\left(\sqrt{||\frac{1}{m} \sum_{i \in [m]} g_i||_2^2}\right) = \frac{1}{\pi} \arccos\left(\sqrt{||\frac{\sum_{i \in [m]} \langle g_i, g_i \rangle}{m^2} + \frac{2}{m^2} \sum_{1 \le i < j \le m} \langle g_i, g_j \rangle||}\right)$$

$$= \frac{1}{\pi} \arccos\left(\sqrt{\frac{1}{m} + \frac{2}{m^2} \sum_{1 \le i < j \le m} \langle g_i, g_j \rangle}\right)$$

### D Additional Experiment Details

Baselines We implement all the baselines mentioned in Table 2. As all these baselines are suited to  $\ell_2$  error, for the DME experiment on gaussians, where  $\ell_2$  error is the correct metric, compare SparseReg (Algorithm 4) to all these baselines. For  $\ell_\infty$  error uniform distribution, we implement NoisySign (Algorithm 1) and HadamardMultiDim (Algorithm 3) and compare it to Correlated SRQ (Suresh et al., 2022), as it's guarantees hold in single dimensions. We also add comparisons to its independent variant, SRQ (Suresh et al., 2017), and Drive (Vargaftik et al., 2021), which performs coordinate-wise signs. For the unit vector case, we implement OneBit (Algorithm 5 Technique II) and SparseReg(Algorithm 4) and compare it with one independent compressor (SRQ (Suresh et al., 2017)) and one collaborative compressor (RandKSpatialProj (Jiang et al., 2023)). Note that we set d=512 throughout our experiments and tune the parameters (number of coordinates sent (Konečný & Richtárik, 2018; Jhunjhunwala et al., 2021) or the quantization levels in (Suresh et al., 2017; 2022)) so that all compressors have the same number of bits communicated. For compressors without tunable parameters, we repeat them to match the communication budget. The communication budget is  $2375 \pm 25$  bits/client for each compressor every round. It is often not possible for two different forms of compressors, like sparsification methods and quantization methods to achieve the exact same communication budget even after tuning hence we allow a small range of communication budgets.

**Datasets** For the distributed mean estimation task, we generate d dimensional vectors on m = 100 clients. To compare  $\ell_2$  error, we generate g with  $||g||_2 = 100$ . Then, each client generates  $g_i$  from a  $\mathcal{N}(0, \Delta_2^2)$ , where  $\Delta_2 \in [0.001, 100]$ . To compare  $\ell_{\infty}$  error, we generate g uniformly from a hypercube  $[-B, B]^d$  where B = 100. Each client generates  $g_i$  from a smaller hypercube  $[-\Delta_{\infty}, \Delta_{\infty}]^d$  centered at g where  $\Delta_{\infty} \in [10^{-3}, 10^2]$ . To compare cosine distance, we generate g uniformly from the unit sphere, and each client generates  $g_i$  uniformly from the set of unit vectors at a cosine distance  $\Delta_{corr}$  from the g, Here,  $\Delta_{corr} \in [0.01, 0.4]$ .

For KMeans and power iteration, we set m=50. FEMNIST is a real federated dataset where each client has handwritten digits from a different person. We apply dimensionality reduction to set d=512. We run 20 iterations of Lloyd's algorithm (Lloyd, 1982) for KMeans and 30 power iterations. For distributed linear regression, the Synthetic dataset is a mixture of linear regressions, with one mixture component per client. The true model  $w_i \in \mathbb{R}^d$  for each component is obtained from DME setup for gaussians with  $\Delta_2=4$ . Then, we generate n=1000 datapoints on each client, where the features x are sampled from standard normal, while the labels y are generated as  $y=\langle w_i,x\rangle+\xi$ , where  $\xi$  is the zero-mean gaussian noise with variance  $10^{-2}$ . For UJIndoorLoc, we use the first d=512 of the 520 features following (Jiang et al., 2023). The task for UJIndoorLoc dataset is to predict the longitude of a phone call. For both the linear regression datasets, we run 50 iterations of GD. For MNIST and UJIndoorLoc, we split the dataset uniformly into m chunks one per client.

For logistic regression for binary classification, we select the last 2 classes of the HAR and label them with  $\pm 1$ . We split the dataset into m=20 clients iid. HAR dataset has 561 features which we reduce by PCA to d=512. We run distributed gradient descent with learning rate 0.001 for T=200 iterations on the logistic loss, where the logistic loss for any data point  $(x,y) \in \mathbb{R}^d \times \{\pm 1\}$  is defined as  $\ell(w,(x,y)) = \log(1 + \exp(-\langle w, x \rangle \cdot y))$  for any weight  $w \in \mathbb{R}^d$ .

Metrics With the same number of bits, we can directly compare the error of baselines. For mean estimation, we measure  $\ell_2$  error,  $\ell_{\infty}$  error and cosine distance for gaussian, uniform and unit vectors respectively. For KMeans, we report the KMeans objective. For power iteration, we report the top eigenvalue. For linear regression, we provide the mean squared error on a test dataset. For logistic regression, we report the training logistic loss and test accuracy for binary classification. For all experiments except power iteration and test accuracy logistic regression, lower value of the reported metric implies better performance. For power iteration, higher implies better performance, as we need to find the eigenvector corresponding to the top eigenvalue. For logistic regression, higher test accuracy implies better performance.

We provide the code in the supplementary material and all the experiments took 6 days to run on a single 20 core machine with 25 GB RAM.

### E Distributed Gradient Descent with SparseReg Compressor

This section uses our  $\ell_2$  compressor, SparseReg, for running FedAvg. Each client  $i \in [m]$  contains a local objective function  $f_i : \mathcal{W} \to \mathbb{R}$ . We define the global objective function  $f(w) = \frac{1}{m} \sum_{i=1}^m f_i(w)$ ,  $\forall w \in \mathcal{W} \subset \mathbb{R}^d$ . The goal is to find  $w^* \in \arg\min_{w \in \mathcal{W}} f(w)$ . Note that  $\nabla f(w) = \frac{1}{m} \sum_{i=1}^m \nabla f_i(w)$ , therefore, in our case, the vector  $g_i$  correspond to  $\nabla f_i(w)$ . We describe the algorithm in Algorithm 6

### Algorithm 6 Distributed Projected Gradient Descent with SparseReg compressor

```
Require: Initial iterate w^0 \in \mathcal{W}, Step size \gamma > 0
\frac{\text{Server}}{\text{SparseReg-Init()}}
\text{for } t = 0 \text{ to } T - 1 \text{ do}
\text{Send } w^t \text{ to all clients } i \in [m].
\text{Receive } \tilde{b_i}^t \text{ from clients } i \in [m].
\tilde{g}^t \leftarrow \text{SparseReg-Decode}(\{\tilde{b}_i^t\}_{i \in [m]})
w^{t+1} \leftarrow \text{proj}_{\mathcal{W}}(w^t - \eta_t \tilde{g}^t)
end for
\frac{\text{Client(i) at iteration } t}{\text{Receive } w^t \text{ from server.}}
\tilde{b}_i \leftarrow \text{SparseReg-Encode}(\nabla f_i(w^t))
Send \tilde{b}_i^t to server.
```

We first state the assumptions required for applying the SparseReg compressor.

**Assumption 4** (Bounded Gradient). For all  $w \in \mathcal{W}, i \in [m]$ , we assume that  $||\nabla f_i(w)||_2 \leq B$ .

By this assumption, we ensure that for each iteration t in Algorithm 6,  $||g_i||_2 = ||\nabla f_i(w^t)||_2$  is bounded. Further, bounded gradients imply that each  $f_i$  is Lipschitz. By triangle inequality, we can also establish the following corollary.

Corollary 1. The objective function f(w) is B-Lipschitz,  $\forall w \in \mathcal{W}$ .

From the above assumptions, it is clear that local objective functions need to be Lipschitz. From (Bubeck, 2015, Theorem 3.2), if the domain of iterates, W is bounded and f(w) is also convex, then gradient descent can converge at a rate  $\mathcal{O}(1/\sqrt{T})$ . We use these two assumptions, and establish a  $\mathcal{O}(1/\sqrt{T})$  rate along with a error obtained from Theorem 3. We define  $\Delta_{\text{reg}}(t)$  and  $\Delta_{2,\text{max}}(t)$  from Theorem 3 to be the corresponding errors for  $g_i = \nabla f_i(w^t), \forall i \in [m]$  for any t > 0.

**Assumption 5** (Bounded domain). The set W is closed and convex with diameter  $R^2$ .

**Assumption 6** (Convexity). The objective function f(w) is convex  $\forall w \in \mathcal{W}$ .

We now state our convergence result.

**Theorem 6.** Under Assumptions 4, 5, 6, running Algorithm 6 for T iterations with step size  $\eta_t = \frac{R}{B\sqrt{T}}$  with probability  $1 - 2m^2LT \exp(-d\delta_1^2/8) - mT\left(\frac{L^{2\delta_2}}{\log L}\right)^{-m}$  we have,

$$\mathbb{E}[f(\bar{w}^T)] - f(w^*) \le \frac{R(2B^2 + \Gamma_1)}{2B\sqrt{T}} + \sqrt{\Gamma_1}R, \quad \text{where,} \quad \bar{w}^T = \frac{1}{T} \sum_{t=0}^{T-1} w^t$$

$$\Gamma_1 = B^2 \left( 1 + \frac{10 \log L}{d} \exp\left(\frac{m \log L}{d}\right) (\delta_1 + \delta_2) \right)^2 \left( 1 - \frac{2 \log L}{d} \right)^m,$$

$$\Gamma_2 = \max_{t \in \{0, 1, \dots, T-1\}} \quad \min\{\Delta_{\text{reg}}(t), \Delta_{2, \text{max}}(t)\}$$

$$(13)$$

From the above theorem, we can see that the high probability terms and  $\Gamma_1$  and  $\Gamma_2$  are obtained from Theorem 3. Note that  $\Gamma = \mathcal{O}(B^2 \exp(-m/d))$ , therefore, for large m, the additional bias term of  $R\sqrt{\Gamma_1}$  is very small. Further, the term  $\Gamma_2 \leq B^2$ , therefore,  $\Gamma_2$  only affects constant terms in the convergence rate due to  $\sqrt{T}$  in the denominator. If  $\exp(-m/d) = \mathcal{O}(1/\sqrt{T})$  or  $m = \Omega(d \log T)$ , the final convergence rate of Algorithm 6 is  $\mathcal{O}(RB/\sqrt{T})$  which is the rate for distributed GD without compression.

We provide the proof for the above theorem, which modifies the proof of (Bubeck, 2015, Theorem 3.2) to handle a biased gradient oracle. We can also extend our analysis to other function classes, for instance strongly convex functions, by using existing works on biased gradient oracles (Ajalloeian & Stich, 2020). Extending the proof to FedAvg from distributed GD would require using biased gradient oracles in (Li et al., 2020). Further, these proofs can also be extended to HadamardMultiDim compressor, with an additional  $\sqrt{d}$  factor in the corresponding error terms from Theorem 2 to account for conversion from  $\ell_{\infty}$  to  $\ell_{2}$  norm.

### E.1 Proof of Theorem 6

At any iteration t > 0, we use  $\tilde{g}^t$  to denote the estimate of  $\nabla f(w^t)$ . From the proof of Theorem 3,  $||\mathbb{E}_t[\tilde{g}^t] - \nabla f(w^t)||_2 \leq \sqrt{\Gamma_1}$ , and  $\mathbb{V}ar_t(\tilde{g}^t|w^t) \leq \Gamma_2, \forall t > 0$ , where  $\mathbb{E}_t$  and  $\mathbb{V}ar_t$  are the expectation and variance wrt the randomness in the SparseReg compressor at iteration t. We take a union bound over the high probability terms in Theorem 3 over all iterations t = 0 to t = 0.

We can write the following equation by convexity of  $f(w^t)$ .

$$\begin{split} f(w^t) - f(w^\star) \leq & \langle \nabla f(w^t), w^t - w^\star \rangle = \langle \tilde{g}^t, w^t - w^\star \rangle + \langle \nabla f(w^t) - \tilde{g}^t, w^t - w^\star \rangle \\ \leq & \frac{1}{2n} (||w^t - w^\star||_2^2 - ||w^t - \eta \tilde{g}^t - w^\star||_2^2) + \eta ||\tilde{g}^t||_2^2 / 2 + \langle \nabla f(w^t) - \tilde{g}^t, w^t - w^\star \rangle \end{split}$$

In the second line, we use  $2\langle a,b\rangle = ||a||_2^2 + ||b||_2^2 - ||a-b||_2^2$ . Now, taking expectation wrt the randomness in SparseReg at iteration t, we obtain,

$$\begin{split} \mathbb{E}_{t}[f(w^{t})] - f(w^{\star}) \leq & \frac{1}{2\eta} (||w^{t} - w^{\star}||_{2}^{2} - \mathbb{E}_{t}[||w^{t} - \eta \tilde{g}^{t} - w^{\star}||_{2}^{2}]) + \eta \mathbb{E}_{t}[||\tilde{g}^{t}||_{2}^{2}]/2 \\ & + \langle \nabla f(w^{t}) - \mathbb{E}_{t}[\tilde{g}^{t}], w^{t} - w^{\star} \rangle \\ \leq & \frac{1}{2\eta} (||w^{t} - w^{\star}||_{2}^{2} - \mathbb{E}_{t}[||w^{t+1} - w^{\star}||_{2}^{2}]) + \eta (||\mathbb{E}_{t}[\tilde{g}^{t}]||_{2}^{2} + \mathbb{V}ar_{t}(\tilde{g}^{t}))/2 \\ & + ||\nabla f(w^{t}) - \mathbb{E}_{t}[\tilde{g}^{t}]||_{2} \cdot ||w^{t} - w^{\star}||_{2} \\ \leq & \frac{1}{2\eta} (||w^{t} - w^{\star}||_{2}^{2} - \mathbb{E}_{t}[||w^{t+1} - w^{\star}||_{2}^{2}]) + \eta (B^{2} + \Gamma_{2})/2 + \sqrt{\Gamma_{1}}R \end{split}$$

In the second line, we use the non-expansiveness of projections on a convex set,  $||w^t - \eta \tilde{g}^t - w^*||_2 \ge ||\operatorname{proj}_{\mathcal{W}}(w^t - \eta \tilde{g}^t - w^*)||_2$ , the decomposition of  $2^{nd}$  moment into square of mean and variance, and cauchy-schwartz inequality. In the third line, we plug in bounds of  $\Gamma_1, \Gamma_2$ , diameter of the set and by triangle inequality, argue that  $\mathbb{E}[\tilde{g}^t]$  also lies in an  $\ell_2$  ball of radius B.

Finally, we take expectations wrt all random variables, unroll the recursion from t = 0 to T, and divide both sides by T.

$$\frac{1}{T} \sum_{t=0}^{T} \mathbb{E}[f(w^{t})] - f(w^{\star}) \le \frac{R^{2}}{2\eta T} + \frac{\eta(B^{2} + \Gamma_{2})}{2} + \sqrt{\Gamma_{1}}R \le \frac{R(2B^{2} + \Gamma_{1})}{2B\sqrt{T}} + \sqrt{\Gamma_{1}}R$$

We obtain the final inequality by plugging in the step size  $\eta = \frac{R}{B\sqrt{T}}$ . By convexity of f, for  $\bar{w}^T = \sum_{t=0}^{T-1} w^t$ , we obtain,

$$\mathbb{E}[f(\bar{w}^T)] - f(w^*) \le \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[f(w^t)] - f(w^*) \le \frac{R(2B^2 + \Gamma_1)}{2B\sqrt{T}} + \sqrt{\Gamma_1}R$$