

# Appendices

## A Additional Details on Motivation in Introduction

### A.1 Preprocessing all client vectors by the same random matrix does not improve performance

Consider  $n$  clients. Suppose client  $i$  holds a vector  $\mathbf{x}_i \in \mathbb{R}^d$ . We want to apply Rand- $k$  or Rand- $k$ -Spatial, while also making the encoding process more flexible than just randomly choosing  $k$  out of  $d$  coordinates. One naïve way of doing this is for each client to pre-process its vector by applying an orthogonal matrix  $\mathbf{G} \in \mathbb{R}^{d \times d}$  that is the *same* across all clients. Such a technique might be helpful in improving the performance of quantization because the MSE due to quantization often depends on how uniform the coordinates of  $\mathbf{x}_i$ 's are, i.e. whether the coordinates of  $\mathbf{x}_i$  have values close to each other.  $\mathbf{G}$  is designed to be the random matrix (e.g. SRHT) that rotates  $\mathbf{x}_i$  and makes its coordinates uniform.

Each client sends the server  $\hat{\mathbf{x}}_i = \mathbf{E}_i \mathbf{G} \mathbf{x}_i$ , where  $\mathbf{E}_i \in \mathbb{R}^{k \times d}$  is the subsampling matrix. If we use Rand- $k$ , the server can decode each client vector by first applying the decoding procedure of Rand- $k$  and then rotating it back to the original space, i.e.,  $\hat{\mathbf{x}}_i^{(\text{Naïve})} = \frac{d}{k} \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i$ . Note that

$$\begin{aligned} \mathbb{E}[\hat{\mathbf{x}}_i^{(\text{Naïve})}] &= \frac{d}{k} \mathbb{E}[\mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i] \\ &= \frac{d}{k} \mathbf{G}^T \frac{k}{d} \mathbf{I}_d \mathbf{G} \mathbf{x}_i \\ &= \mathbf{x}_i. \end{aligned}$$

Hence,  $\hat{\mathbf{x}}_i^{(\text{Naïve})}$  is unbiased. The MSE of  $\hat{\mathbf{x}}^{(\text{Naïve})} = \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{x}}_i^{(\text{Naïve})}$  is given as

$$\begin{aligned} \mathbb{E} \left\| \bar{\mathbf{x}} - \hat{\mathbf{x}}^{(\text{Naïve})} \right\|_2^2 &= \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i - \frac{1}{n} \frac{d}{k} \sum_{i=1}^n \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i \right\|_2^2 \\ &= \frac{1}{n^2} \mathbb{E} \left\| \sum_{i=1}^n \mathbf{x}_i - \frac{d}{k} \sum_{i=1}^n \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i \right\|_2^2 \\ &= \frac{1}{n^2} \left\{ \frac{d^2}{k^2} \mathbb{E} \left\| \sum_{i=1}^n \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i \right\|_2^2 - \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right\} \\ &= \frac{1}{n^2} \left\{ \frac{d^2}{k^2} \left( \sum_{i=1}^n \mathbb{E} \|\mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i\|_2^2 + \sum_{i \neq j} \mathbb{E} \langle \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i, \mathbf{G} \mathbf{E}_j^T \mathbf{E}_j \mathbf{G} \mathbf{x}_j \rangle \right) - \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right\}. \end{aligned} \tag{12}$$

Next, we bound the first term in Eq. [12](#)

$$\begin{aligned} \mathbb{E} \|\mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i\|_2^2 &= \mathbb{E}[\mathbf{x}_i^T \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i] = \mathbb{E}[\mathbf{x}_i^T \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i] \\ &= \mathbf{x}_i^T \mathbf{G}^T \mathbb{E}[(\mathbf{E}_i^T \mathbf{E}_i)^2] \mathbf{G} \mathbf{x}_i \\ &= \mathbf{x}_i^T \frac{k}{d} \mathbf{I}_d \mathbf{x}_i \quad (\because (\mathbf{E}_i^T \mathbf{E}_i)^2 = \mathbf{E}_i^T \mathbf{E}_i) \\ &= \frac{k}{d} \|\mathbf{x}_i\|_2^2 \end{aligned} \tag{13}$$

The second term in Eq. [12](#) can also be simplified as follows.

$$\begin{aligned} &\mathbb{E}[\langle \mathbf{G}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{G} \mathbf{x}_i, \mathbf{G}^T \mathbf{E}_j^T \mathbf{E}_j \mathbf{G} \mathbf{x}_j \rangle] \\ &= \langle \mathbf{G}^T \mathbb{E}[\mathbf{E}_i^T \mathbf{E}_i] \mathbf{G} \mathbf{x}_i, \mathbf{G}^T \mathbb{E}[\mathbf{E}_j^T \mathbf{E}_j] \mathbf{G} \mathbf{x}_j \rangle \\ &= \langle \mathbf{G}^T \frac{k}{d} \mathbf{I}_d \mathbf{G} \mathbf{x}_i, \mathbf{G}^T \frac{k}{d} \mathbf{I}_d \mathbf{G} \mathbf{x}_j \rangle \end{aligned}$$

$$= \frac{k^2}{d^2} \langle \mathbf{x}_i, \mathbf{x}_l \rangle. \quad (14)$$

Plugging Eq. 13 and Eq. 14 into Eq. 12, we get the MSE is

$$\begin{aligned} & \mathbb{E} \|\bar{\mathbf{x}} - \hat{\mathbf{x}}^{(\text{Naive})}\|_2^2 \\ &= \frac{1}{n^2} \left\{ \frac{d^2}{k^2} \left( \sum_{i=1}^n \frac{k}{d} \|\mathbf{x}_i\|_2^2 + 2 \sum_{i=1}^n \sum_{l=i+1}^n \frac{k^2}{d^2} \langle \mathbf{x}_i, \mathbf{x}_l \rangle \right) - \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right\} \\ &= \frac{1}{n^2} \left( \frac{d}{k} - 1 \right) \sum_{i=1}^n \|\mathbf{x}_i\|_2^2, \end{aligned}$$

which has exactly the same MSE as that of Rand- $k$ . The problem is that if each client applies the same rotational matrix  $\mathbf{G}$ , simply rotating the vectors will not change the  $\ell_2$  norm of the decoded vector, and hence the MSE. Similarly, if one applies Rand- $k$ -Spatial, one ends up having exactly the same MSE as that of Rand- $k$ -Spatial as well. Hence, we need to design a new decoding procedure when the encoding procedure at the clients are more flexible.

## A.2 $nk \gg d$ is not interesting

One can rewrite  $\sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i$  in the Rand-Proj-Spatial estimator (Eq. 5) as  $\sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i = \sum_{j=1}^{nk} \mathbf{g}_j \mathbf{g}_j^T$ , where  $\mathbf{g}_j \in \mathbb{R}^d$  and  $\mathbf{g}_{ik}, \mathbf{g}_{ik+1}, \dots, \mathbf{g}_{(i+1)k}$  are the rows of  $\mathbf{G}_i$ . Since when  $nk \gg d$ ,  $\sum_{j=1}^{nk} \mathbf{g}_j \mathbf{g}_j^T \rightarrow \mathbb{E}[\sum_{j=1}^{nk} \mathbf{g}_j \mathbf{g}_j^T]$  due to Law of Large Numbers, one way to see the limiting MSE of Rand-Proj-Spatial when  $nk$  is large is to approximate  $\sum_{i=1}^n \sum_{j=1}^{nk} \mathbf{g}_j \mathbf{g}_j^T$  by its expectation.

By Lemma 4.1 when  $\mathbf{G}_i = \mathbf{E}_i$ , Rand-Proj-Spatial recovers Rand- $k$ -Spatial. We now discuss the limiting behavior of Rand- $k$ -Spatial when  $nk \gg d$  by leveraging our proposed Rand-Proj-Spatial. In this case, each  $\mathbf{g}_j$  can be viewed as a random based vector  $\mathbf{e}_w$  for  $w$  randomly chosen in  $[d]$ .  $\sum_{i=1}^{nk} \mathbf{g}_j \mathbf{g}_j^T \rightarrow \mathbb{E}[\sum_{i=1}^{nk} \mathbf{g}_j \mathbf{g}_j^T] = \sum_{i=1}^{nk} \frac{1}{d} \mathbf{I}_d = \frac{nk}{d} \mathbf{I}_d$ . And so the scalar  $\bar{\beta}$  in Eq. 5 to ensure an unbiased estimator is computed as

$$\begin{aligned} \bar{\beta} \mathbb{E} \left[ \left( \frac{nk}{d} \mathbf{I}_d \right)^\dagger \mathbf{G}_i^T \mathbf{G}_i \right] &= \mathbf{I}_d \\ \bar{\beta} \frac{d}{nk} \mathbf{I}_d \mathbb{E} [\mathbf{G}_i^T \mathbf{G}_i] &= \mathbf{I}_d \\ \bar{\beta} \frac{d}{nk} \frac{k}{d} &= \mathbf{I}_d \\ \bar{\beta} &= n \end{aligned}$$

And the MSE is now

$$\begin{aligned} \mathbb{E} \left[ \|\bar{\mathbf{x}} - \hat{\mathbf{x}}\| \right] &= \mathbb{E} \left[ \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i - \frac{1}{n} \bar{\beta} \frac{d}{nk} \mathbf{I}_d \sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i \right\|_2^2 \right] \\ &= \frac{1}{n^2} \left\{ \bar{\beta}^2 \frac{d^2}{n^2 k^2} \mathbb{E} \left[ \left\| \sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i \right\|_2^2 \right] - \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right\} \\ &= \frac{1}{n^2} \left\{ n^2 \frac{d^2}{n^2 k^2} \left( \sum_{i=1}^n \mathbb{E} \left[ \left\| \mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i \right\|_2^2 \right] + 2 \sum_{i=1}^n \sum_{l=i+1}^n \langle \mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i, \mathbf{E}_l^T \mathbf{E}_l \mathbf{x}_l \rangle \right) - \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right\} \\ &= \frac{1}{n^2} \left\{ \frac{d^2}{k^2} \left( \sum_{i=1}^n \mathbb{E} \left[ \mathbf{x}_i^T (\mathbf{E}_i^T \mathbf{E}_i)^2 \mathbf{x}_i \right] + 2 \sum_{i=1}^n \sum_{l=i+1}^n \frac{k^2}{d^2} \langle \mathbf{x}_i, \mathbf{x}_l \rangle \right) - \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right\} \\ &= \frac{1}{n^2} \left\{ \frac{d^2}{k^2} \left( \sum_{i=1}^n \frac{k}{d} \|\mathbf{x}_i\|_2^2 + 2 \sum_{i=1}^n \sum_{l=i+1}^n \frac{k^2}{d^2} \langle \mathbf{x}_i, \mathbf{x}_l \rangle \right) - \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 - 2 \sum_{i=1}^n \sum_{l=i+1}^n \langle \mathbf{x}_i, \mathbf{x}_l \rangle \right\} \\ &= \frac{1}{n^2} \left( \frac{d}{k} - 1 \right) \sum_{i=1}^n \|\mathbf{x}_i\|_2^2 \end{aligned}$$

which is exactly the same MSE as Rand- $k$ . This implies when  $nk$  is large, the MSE of Rand- $k$ -Spatial does not get improved compared to Rand- $k$  with correlation information. Intuitively, this implies when  $nk \gg d$ , the server gets enough amount of information from the client, and does not need correlation to improve its estimator. Hence, we focus on the more interesting case when  $nk < d$  — that is, when the server does not have enough information from the clients, and thus wants to use additional information, i.e. cross-client correlation, to improve its estimator.

## B Additional Details on the Rand-Proj-Spatial Family Estimator

### B.1 $\bar{\beta}$ is a scalar

From Eq. 20 in the proof of Theorem 4.3 and Eq. 25 in the proof of Theorem 4.4, it is evident that the unbiasedness of the mean estimator  $\hat{\mathbf{x}}^{\text{Rand-Proj-Spatial}}$  is ensured collectively by

- The random sampling matrices  $\{\mathbf{E}_i\}$ .
- The orthogonality of scaled Hadamard matrices  $\mathbf{H}^T \mathbf{H} = d\mathbf{I}_d = \mathbf{H} \mathbf{H}^T$ .
- The rademacher diagonal matrices, with the property  $(\mathbf{D}_i)^2 = \mathbf{I}_d$ .

### B.2 Alternative motivating regression problems

#### Alternative motivating regression problem 1.

Let  $\mathbf{G}_i \in \mathbb{R}^{k \times d}$  and  $\mathbf{W}_i \in \mathbb{R}^{d \times k}$  be the encoding and decoding matrix for client  $i$ . One possible alternative estimator that translates the intuition that the decoded vector should be close to the client's original vector, for all clients, is by solving the following regression problem,

$$\begin{aligned} \hat{\mathbf{x}} &= \arg \min_{\mathbf{W}} f(\mathbf{W}) = \mathbb{E}[\|\bar{\mathbf{x}} - \frac{1}{n} \sum_{i=1}^n \mathbf{W}_i \mathbf{G}_i \mathbf{x}_i\|_2^2] \\ \text{subject to } \bar{\mathbf{x}} &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i] \end{aligned} \quad (15)$$

where  $\mathbf{W} = (\mathbf{W}_1, \mathbf{W}_2, \dots, \mathbf{W}_n)$  and the constraint enforces unbiasedness of the estimator. The estimator is then the solution of the above problem. However, we note that optimizing a decoding matrix  $\mathbf{W}_i$  for each client leads to performing individual decoding of each client's compressed vector instead of a joint decoding process that considers all clients' compressed vectors. Only a joint decoding process can achieve the goal of leveraging cross-client information to reduce the estimation error. Indeed, we show as follows that solving the above optimization problem in Eq. 15 recovers the MSE of our baseline Rand- $k$ . Note

$$\begin{aligned} f(\mathbf{W}) &= \mathbb{E}[\|\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{W}_i \mathbf{G}_i \mathbf{x}_i)\|_2^2] = \mathbb{E}[\|\frac{1}{n} \sum_{i=1}^n (\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i\|_2^2] \\ &= \mathbb{E}\left[\frac{1}{n^2} \left( \sum_{i=1}^n \|(\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i\|_2^2 + \sum_{i \neq j} \langle (\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i, (\mathbf{I}_d - \mathbf{W}_j \mathbf{G}_j) \mathbf{x}_j \rangle \right)\right] \\ &= \frac{1}{n^2} \left( \sum_{i=1}^n \mathbb{E}[\|(\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i\|_2^2] + \sum_{i \neq j} \mathbb{E}[\langle (\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i, (\mathbf{I}_d - \mathbf{W}_j \mathbf{G}_j) \mathbf{x}_j \rangle] \right). \end{aligned} \quad (16)$$

By the constraint of unbiasedness, i.e.,  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i]$ , there is

$$\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i] = 0 \Leftrightarrow \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i] = 0.$$

We now show that a sufficient and necessary condition to satisfy the above unbiasedness constraint is that for all  $i \in [n]$ ,  $\mathbb{E}[\mathbf{W}_i \mathbf{G}_i] = \mathbf{I}_d$ .

*Sufficiency.* It is obvious that if for all  $i \in [n]$ ,  $\mathbb{E}[\mathbf{W}_i \mathbf{G}_i] = \mathbf{I}_d$ , then we have  $\frac{1}{n} \mathbb{E}[(\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i] = 0$ .

*Necessity.* Consider the special case that for some  $i \in [n]$  and  $\lambda \in [d]$ ,  $\mathbf{x}_i = n \mathbf{e}_\lambda$ , where  $\mathbf{e}_\lambda$  is the  $\lambda$ -th canonical basis vector, and  $\mathbf{x}_j = 0$ , and for all  $j \in [n] \setminus \{i\}$ . Then,

$$\mathbf{e}_\lambda = \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i] = \frac{1}{n} \mathbb{E}[\mathbf{W}_i \mathbf{G}_i] \mathbf{e}_\lambda = [\mathbb{E}[\mathbf{W}_i \mathbf{G}_i]]_\lambda,$$

where  $[\cdot]_\lambda$  denotes the  $\lambda$ -th column of matrix  $\mathbb{E}[\mathbf{W}_i \mathbf{G}_i]$ .

Since our approach is agnostic to the choice of vectors, we need this choice of decoder matrices, by varying  $\lambda$  over  $[d]$ , we see that we need  $\mathbb{E}[\mathbf{W}_i \mathbf{G}_i] = \mathbf{I}_d$ . And by varying  $i$  over  $[n]$ , we see that we need  $\mathbb{E}[\mathbf{W}_j \mathbf{G}_j] = \mathbf{I}_d$  for all  $j \in [n]$ .

Therefore,  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i] \Leftrightarrow \forall i \in [n], \mathbb{E}[\mathbf{W}_i \mathbf{G}_i] = \mathbf{I}_d$ .

This implies the second term of  $f(\mathbf{W})$  in Eq. [16](#) is 0, that is,

$$\sum_{i \neq j} \mathbb{E} \left[ \left\langle (\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i, (\mathbf{I}_d - \mathbf{W}_j \mathbf{G}_j) \mathbf{x}_j \right\rangle \right] = 0.$$

Hence, we only need to solve

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{W}} f_2(\mathbf{W}) = \sum_{i=1}^n \mathbb{E} \left[ \left\| (\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i \right\|_2^2 \right] \quad (17)$$

Since each  $\mathbf{W}_i$  appears in  $f_2(\mathbf{W})$  separately, each  $\mathbf{W}_i$  can be optimized separately, via solving

$$\min_{\mathbf{W}_i} \mathbb{E} \left[ \left\| (\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i \right\|_2^2 \right] \quad \text{subject to } \mathbb{E}[\mathbf{W}_i \mathbf{G}_i] = \mathbf{I}_d.$$

One natural solution is to take  $\mathbf{W}_i = \frac{d}{k} \mathbf{G}_i^\dagger$ ,  $\forall i \in [n]$ . For  $i \in [n]$ , let  $\mathbf{G}_i = \mathbf{V}_i \Lambda_i \mathbf{U}_i^T$  be its SVD, where  $\mathbf{V}_i \in \mathbb{R}^{k \times d}$  and  $\mathbf{U}_i \in \mathbb{R}^{d \times d}$  are orthogonal matrices. Then,

$$\mathbf{W}_i \mathbf{G}_i = \frac{d}{k} \mathbf{U}_i \Lambda_i^\dagger \mathbf{V}_i^T \mathbf{V}_i \Lambda_i \mathbf{U}_i^T = \frac{d}{k} \mathbf{U}_i \Lambda_i^\dagger \Lambda_i \mathbf{U}_i^T = \frac{d}{k} \mathbf{U}_i \Sigma \mathbf{U}_i^T,$$

where  $\Sigma$  is a diagonal matrix with 0s and 1s on the diagonal.

For simplicity, we assume the random matrix  $\mathbf{U}_i$  follows a continuous distribution.  $\mathbf{U}_i$  being discrete follows a similar analysis. Let  $\mu(\mathbf{U}_i)$  be the measure of  $\mathbf{U}_i$ .

$$\begin{aligned} \mathbb{E}[\mathbf{W}_i \mathbf{G}_i] &= \frac{d}{k} \mathbb{E}[\mathbf{U}_i \Sigma \mathbf{U}_i^T] = \frac{d}{k} \int_{\mathbf{U}_i} \mathbb{E}[\mathbf{U}_i \Sigma_i \mathbf{U}_i^T \mid \mathbf{U}_i] \cdot d\mu(\mathbf{U}_i) \\ &= \frac{d}{k} \int_{\mathbf{U}_i} \mathbf{U}_i \mathbb{E}[\Sigma_i \mid \mathbf{U}_i] \mathbf{U}_i^T \cdot \mu(\mathbf{U}_i) \\ &= \frac{d}{k} \int_{\mathbf{U}_i} \mathbf{U}_i \frac{k}{d} \mathbf{I}_d \mathbf{U}_i^T \cdot d\mu(\mathbf{U}_i) \\ &= \frac{d}{k} \frac{k}{d} \mathbf{I}_d = \mathbf{I}_d, \end{aligned}$$

which means the estimator  $\frac{1}{n} \sum_{i=1}^n \frac{k}{d} \mathbf{G}_i^\dagger \mathbf{G}_i$  satisfies unbiasedness. The MSE is now

$$\begin{aligned} MSE &= \mathbb{E} \left[ \left\| \bar{\mathbf{x}} - \frac{1}{n} \sum_{i=1}^n \mathbf{W}_i \mathbf{G}_i \mathbf{x}_i \right\|_2^2 \right] = \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \left[ \left\| (\mathbf{I}_d - \mathbf{W}_i \mathbf{G}_i) \mathbf{x}_i \right\|_2^2 \right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \left( \|\mathbf{x}_i\|_2^2 + \mathbb{E}[\|\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i\|_2^2] - 2 \langle \mathbf{x}_i, \mathbb{E}[\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i] \rangle \right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \left( \|\mathbf{x}_i\|_2^2 + \mathbb{E}[\|\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i\|_2^2] - 2 \langle \mathbf{x}_i, \mathbf{x}_i \rangle \right) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n^2} \sum_{i=1}^n \left( \mathbb{E}[\|\mathbf{W}_i \mathbf{G}_i \mathbf{x}_i\|_2^2 - \|\mathbf{x}_i\|_2^2] \right) \\
&= \frac{1}{n^2} \sum_{i=1}^n \left( \mathbf{x}_i \mathbb{E}[(\mathbf{W}_i \mathbf{G}_i)^T (\mathbf{W}_i \mathbf{G}_i)] \mathbf{x}_i - \|\mathbf{x}_i\|_2^2 \right).
\end{aligned}$$

Again, let  $\mathbf{G}_i = \mathbf{V}_i \Lambda_i \mathbf{U}_i^T$  be its SVD and consider  $\mathbf{W}_i \mathbf{G}_i = \frac{d}{k} \mathbf{U}_i \Sigma_i \mathbf{U}_i^T$ , where  $\Sigma_i$  is a diagonal matrix with 0s and 1s. Then,

$$\begin{aligned}
MSE &= \frac{1}{n^2} \sum_{i=1}^n \sum_{i=1}^n \left( \mathbf{x}_i^T \frac{d^2}{k^2} \mathbb{E}[\mathbf{U}_i \Sigma_i \mathbf{U}_i^T \mathbf{U}_i \Sigma_i \mathbf{U}_i^T] \mathbf{x}_i - \|\mathbf{x}_i\|_2^2 \right) \\
&= \frac{1}{n^2} \sum_{i=1}^n \left( \frac{d^2}{k^2} \mathbf{x}_i^T \mathbb{E}[\mathbf{U}_i \Sigma_i^2 \mathbf{U}_i^T] \mathbf{x}_i - \|\mathbf{x}_i\|_2^2 \right).
\end{aligned}$$

Since  $\mathbf{G}_i$  has rank  $k$ ,  $\Sigma_i$  is a diagonal matrix with  $k$  out of  $d$  entries being 1 and the rest being 0. Let  $\mu(\mathbf{U}_i)$  be the measure of  $\mathbf{U}_i$ . Hence, for  $i \in [n]$ ,

$$\begin{aligned}
\mathbb{E}[\mathbf{U}_i \Sigma_i^2 \mathbf{U}_i^T] &= \int_{\mathbf{U}_i} \mathbb{E}[\mathbf{U}_i \Sigma_i^2 \mathbf{U}_i^T \mid \mathbf{U}_i] d\mu(\mathbf{U}_i) \\
&= \int_{\mathbf{U}_i} \mathbf{U}_i \mathbb{E}[\Sigma_i^2 \mid \mathbf{U}_i] \mathbf{U}_i^T d\mu(\mathbf{U}_i) \\
&= \int_{\mathbf{U}_i} \frac{k}{d} \mathbf{U}_i \mathbf{I}_d \mathbf{U}_i^T d\mu(\mathbf{U}_i) \\
&= \frac{k}{d} \int_{\mathbf{U}_i} \mathbf{I}_d d\mu(\mathbf{U}_i) \\
&= \frac{k}{d} \mathbf{I}_d.
\end{aligned}$$

Therefore, the MSE of the estimator, which is the solution of the optimization problem in Eq. [15](#) is

$$MSE = \frac{1}{n^2} \sum_{i=1}^n \left( \frac{d^2}{k^2} \mathbf{x}_i^T \frac{k}{d} \mathbf{I}_d \mathbf{x}_i - \|\mathbf{x}_i\|_2^2 \right) = \frac{1}{n^2} \left( \frac{d}{k} - 1 \right) \sum_{i=1}^n \|\mathbf{x}_i\|_2^2,$$

which is the same MSE as that of Rand- $k$ .

### Alternative motivating regression problem 2.

Another motivating regression problem based on which we can design our estimator is

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \mathbf{x} - \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \mathbf{x}_i \right\|_2^2 \quad (18)$$

Note that  $\mathbf{G}_i \in \mathbb{R}^{k \times d}$ ,  $\forall i \in [n]$ , and so the solution to the above problem is

$$\hat{\mathbf{x}}^{(\text{solution})} = \left( \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \right)^\dagger \left( \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \mathbf{x}_i \right),$$

and to ensure unbiasedness of the estimator, we can set  $\bar{\beta} \in \mathbb{R}$  and have the estimator as

$$\hat{\mathbf{x}}^{(\text{estimator})} = \bar{\beta} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \right)^\dagger \left( \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \mathbf{x}_i \right).$$

It is not hard to see this estimator does not lead to an MSE as low as Rand-Proj-Spatial does. Consider the full correlation case, i.e.,  $\mathbf{x}_i = \mathbf{x}$ ,  $\forall i \in [n]$ , for example, the estimator is now

$$\hat{\mathbf{x}}^{(\text{estimator})} = \bar{\beta} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \right)^\dagger \left( \frac{1}{n} \sum_{i=1}^n \mathbf{G}_i \right) \mathbf{x}.$$

Note that  $\text{rank}(\frac{1}{n} \sum_{i=1}^n \mathbf{G}_i)$  is at most  $k$ , since  $\mathbf{G}_i \in \mathbb{R}^{k \times d}$ ,  $\forall i \in [k]$ . This limits the amount of information of  $\mathbf{x}$  the server can recover.

While recall that in this case, the Rand-Proj-Spatial estimator is

$$\hat{\mathbf{x}}^{(\text{Rand-Proj-Spatial})} = \bar{\beta} \left( \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i \right)^\dagger \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i \mathbf{x} = \bar{\beta} \mathbf{S}^\dagger \mathbf{S} \mathbf{x},$$

where  $\mathbf{S}$  can have rank at most  $nk$ .

### B.3 Why deriving the MSE of Rand-Proj-Spatial with SRHT is hard

To analyze Eq. [11] one needs to compute the distribution of eigendecomposition of  $\mathbf{S} = \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i$ , i.e. the sum of the covariance of SRHT. To the best of our knowledge, there is no non-trivial closed form expression of the distribution of eigen-decomposition of even a single  $\mathbf{G}_i^T \mathbf{G}_i$ , when  $\mathbf{G}_i$  is SRHT, or other commonly used random matrices, e.g. Gaussian. When  $\mathbf{G}_i$  is SRHT, since  $\mathbf{G}_i^T \mathbf{G}_i = \mathbf{D}_i \mathbf{H} \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i$  and the eigenvalues of  $\mathbf{E}_i^T \mathbf{E}_i$  are just diagonal entries, one might attempt to analyze  $\mathbf{H} \mathbf{D}_i$ . While the hardmard matrix  $\mathbf{H}$ 's eigenvalues and eigenvectors are known<sup>3</sup>, the result can hardly be applied to analyze the distribution of singular values or singular vectors of  $\mathbf{H} \mathbf{D}_i$ .

Even if one knows the eigen-decomposition of a single  $\mathbf{G}_i^T \mathbf{G}_i$ , it is still hard to get the eigen-decomposition of  $\mathbf{S}$ . The eigenvalues of a matrix  $\mathbf{A}$  can be viewed as a non-linear function in the  $\mathbf{A}$ , and hence it is in general hard to derive closed form expressions for the eigenvalues of  $\mathbf{A} + \mathbf{B}$ , given the eigenvalues of  $\mathbf{A}$  and that of  $\mathbf{B}$ . One exception is when  $\mathbf{A}$  and  $\mathbf{B}$  have the same eigenvector and the eigenvalues of  $\mathbf{A} + \mathbf{B}$  becomes a sum of the eigenvalues of  $\mathbf{A}$  and  $\mathbf{B}$ . Recall when  $\mathbf{G}_i = \mathbf{E}_i$ , Rand-Proj-Spatial recovers Rand- $k$ -Spatial. Since  $\mathbf{E}_i^T \mathbf{E}_i$ 's all have the same eigenvectors (i.e. same as  $\mathbf{I}_d$ ), the eigenvalues of  $\mathbf{S} = \sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i$  are just the sum of diagonal entries of  $\mathbf{E}_i^T \mathbf{E}_i$ 's. Hence, deriving the MSE for Rand- $k$ -Spatial is not hard compared to the more general case when  $\mathbf{G}_i^T \mathbf{G}_i$ 's can have different eigenvectors.

Since one can also view  $\mathbf{S} = \sum_{i=1}^{nk} \mathbf{g}_i \mathbf{g}_i^T$ , i.e. the sum of  $nk$  rank-one matrices, one might attempt to recursively analyze the eigen-decomposition of  $\sum_{i=1}^{n'} \mathbf{g}_i \mathbf{g}_i^T + \mathbf{g}_{n'+1} \mathbf{g}_{n'+1}^T$  for  $n' \leq n$ . One related problem is eigen-decomposition of a low-rank updated matrix in perturbation analysis: Given the eigen-decomposition of a matrix  $\mathbf{A}$ , what is the eigen-decomposition of  $\mathbf{A} + \mathbf{V} \mathbf{V}^T$ , where  $\mathbf{V}$  is low-rank matrix (or more commonly rank-one)? To compute the eigenvalues of  $\mathbf{A} + \mathbf{V} \mathbf{V}^T$  directly from that of  $\mathbf{A}$ , the most effective and widely applied solution is to solve the so-called secular equation, e.g. [59, 60, 61]. While this can be done computationally efficiently, it is hard to get a closed form expression for the eigenvalues of  $\mathbf{A} + \mathbf{V} \mathbf{V}^T$  from the secular equation.

The previous analysis of SRHT in e.g. [37, 38, 39, 45, 55] is based on asymptotic properties of SRHT, such as the limiting eigen-spectrum, or concentration bounds that bounds the singular values. To analyze the MSE of Rand-Proj-Spatial, however, we need an exact, non-asymptotic analysis of the distribution of SRHT. Concentration bounds does not apply, since computing the pseudo-inverse in Eq. [5] naturally bounds the eigenvalues, and applying concentration bounds will only lead to a loose upper bound on MSE.

<sup>3</sup>See this note <https://core.ac.uk/download/pdf/81967428.pdf>

## B.4 More simulation results on incorporating various degrees of correlation

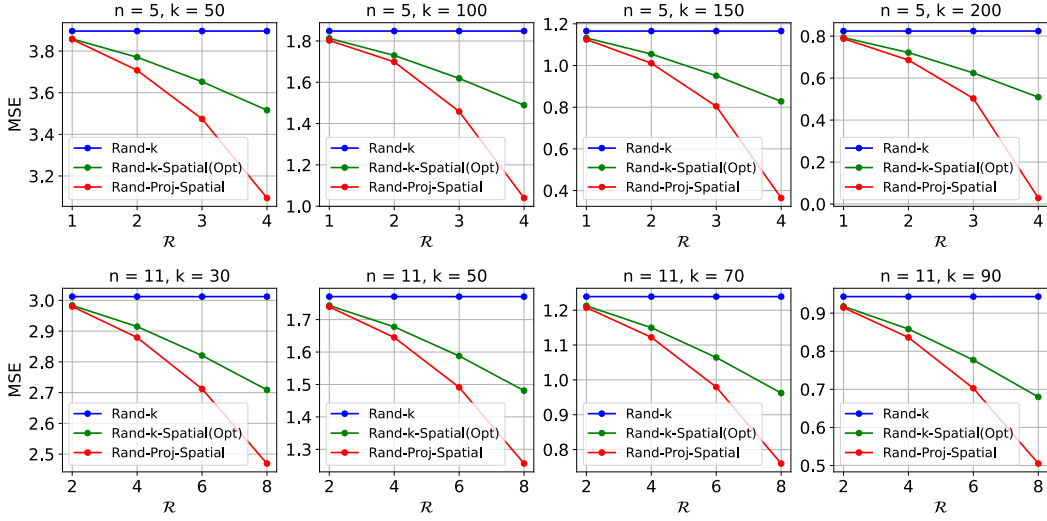


Figure 6: MSE comparison of estimators Rand- $k$ , Rand- $k$ -Spatial(Opt), Rand-Proj-Spatial, given the degree of correlation  $\mathcal{R}$ . Rand- $k$ -Spatial(Opt) denotes the estimator that gives the lowest possible MSE from the Rand- $k$ -Spatial family. We consider  $d = 1024$ , a smaller number of clients  $n \in \{5, 11\}$ , and  $k$  values such that  $nk < d$ . In each plot, we fix  $n, k, d$  and vary the degree of positive correlation  $\mathcal{R}$ . Note the range of  $\mathcal{R}$  is  $\mathcal{R} \in [0, n - 1]$ . We choose  $\mathcal{R}$  with equal space in this range.

## C All Proof Details

### C.1 Proof of Theorem 4.3

**Theorem 4.3** (MSE under Full Correlation). *Consider  $n$  clients, each holding the same vector  $\mathbf{x} \in \mathbb{R}^d$ . Suppose we set  $T(\lambda) = \lambda$ ,  $\bar{\beta} = \frac{d}{k}$  in Eq. 5 and the random linear map  $\mathbf{G}_i$  at each client to be an SRHT matrix. Let  $\delta$  be the probability that  $\mathbf{S} = \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i$  does not have full rank. Then, for  $nk \leq d$ ,*

$$\mathbb{E} \left[ \|\hat{\mathbf{x}}^{(\text{Rand-Proj-Spatial(Max)})} - \bar{\mathbf{x}}\|_2^2 \right] \leq \left[ \frac{d}{(1-\delta)nk + \delta k} - 1 \right] \|\mathbf{x}\|_2^2 \quad (19)$$

*Proof.* All clients have the same vector  $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_n = \mathbf{x} \in \mathbb{R}^d$ . Hence,  $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \mathbf{x}$ , and the decoding scheme is

$$\hat{\mathbf{x}}^{(\text{Rand-Proj-Spatial(Max)})} = \bar{\beta} \left( \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i \right)^\dagger \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i \mathbf{x} = \bar{\beta} \mathbf{S}^\dagger \mathbf{S} \mathbf{x},$$

where  $\mathbf{S} = \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i$ . Let  $\mathbf{S} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$  be its eigendecomposition. Since  $\mathbf{S}$  is a real symmetric matrix,  $\mathbf{U}$  is orthogonal, i.e.,  $\mathbf{U}^T \mathbf{U} = \mathbf{I}_d = \mathbf{U} \mathbf{U}^T$ . Also,  $\mathbf{S}^\dagger = \mathbf{U} \mathbf{\Lambda}^\dagger \mathbf{U}^T$ , where  $\mathbf{\Lambda}^\dagger$  is a diagonal matrix, such that

$$[\mathbf{\Lambda}^\dagger]_{ii} = \begin{cases} 1/[\mathbf{\Lambda}]_{ii} & \text{if } \Lambda_{ii} \neq 0, \\ 0 & \text{else.} \end{cases}$$

Let  $\delta_c$  be the probability that  $\mathbf{S}$  has rank  $c$ , for  $c \in \{k, k+1, \dots, nk-1\}$ . Note that  $\delta = \sum_{c=k}^{nk-1} \delta_c$ . For vector  $\mathbf{m} \in \mathbb{R}^d$ , we use  $\text{diag}(\mathbf{m}) \in \mathbb{R}^{d \times d}$  to denote the matrix whose diagonal entries correspond to the coordinates of  $\mathbf{m}$  and the rest of the entries are zeros.

**Computing  $\bar{\beta}$ .** First, we compute  $\bar{\beta}$ . To ensure that our estimator  $\hat{\mathbf{x}}^{(\text{Rand-Proj-Spatial(Max)})}$  is unbiased, we need  $\bar{\beta} \mathbb{E}[\mathbf{S}^\dagger \mathbf{S} \mathbf{x}] = \mathbf{x}$ . Consequently,

$$\mathbf{x} = \bar{\beta} \mathbb{E}[\mathbf{U} \mathbf{\Lambda}^\dagger \mathbf{U}^T \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T] \mathbf{x}$$

$$\begin{aligned}
&= \bar{\beta} \left[ \sum_{\mathbf{U}=\Phi} \Pr[\mathbf{U} = \Phi] \mathbb{E}[\mathbf{U} \Lambda^\dagger \Lambda \mathbf{U}^T \mid \mathbf{U} = \Phi] \right] \mathbf{x} \\
&= \bar{\beta} \left[ \sum_{\mathbf{U}=\Phi} \Pr[\mathbf{U} = \Phi] \mathbf{U} \mathbb{E}[\Lambda^\dagger \Lambda \mid \mathbf{U} = \Phi] \mathbf{U}^T \right] \mathbf{x} \\
&\stackrel{(a)}{=} \bar{\beta} \left[ \sum_{\mathbf{U}=\Phi} \Pr[\mathbf{U} = \Phi] \mathbf{U} \mathbb{E}[\text{diag}(\mathbf{m}) \mid \mathbf{U} = \Phi] \mathbf{U}^T \right] \mathbf{x} \\
&\stackrel{(b)}{=} \bar{\beta} \sum_{\mathbf{U}=\Phi} \Pr[\mathbf{U} = \Phi] \left[ \mathbf{U} \left( (1-\delta) \frac{nk}{d} \mathbf{I}_d + \sum_{c=k}^{nk-1} \delta_c \frac{c}{d} \mathbf{I}_d \right) \mathbf{U}^T \right] \mathbf{x} \\
&= \bar{\beta} \left[ (1-\delta) \frac{nk}{d} + \sum_{c=k}^{nk-1} \delta_c \frac{c}{d} \right] \mathbf{x} \\
\Rightarrow \bar{\beta} &= \frac{d}{(1-\delta)nk + \sum_{c=k}^{nk-1} \delta_c c} \tag{20}
\end{aligned}$$

where in (a),  $\mathbf{m} \in \mathbb{R}^d$  such that

$$\mathbf{m}_i = \begin{cases} 1 & \text{if } \Lambda_{jj} > 0 \\ 0 & \text{else.} \end{cases}$$

Also, by construction of  $\mathbf{S}$ ,  $\text{rank}(\text{diag}(\mathbf{m})) \leq nk$ . Further, (b) follows by symmetry across the  $d$  dimensions.

Since  $\delta k \leq \sum_{c=k}^{nk-1} \delta_c c \leq \delta(nk-1)$ , there is

$$\frac{d}{(1-\delta)nk + \delta(nk-1)} \leq \bar{\beta} \leq \frac{d}{(1-\delta)nk + \delta k} \tag{21}$$

**Computing the MSE.** Next, we use the value of  $\bar{\beta}$  in Eq. [20](#) to compute MSE.

$$\begin{aligned}
MSE(\text{Rand-Proj-Spatial(Max)}) &= \mathbb{E}[\|\hat{\mathbf{x}}^{(\text{Rand-Proj-Spatial(Max)})} - \bar{\mathbf{x}}\|_2^2] = \mathbb{E}[\|\bar{\beta} \mathbf{S}^\dagger \mathbf{S} \mathbf{x} - \mathbf{x}\|_2^2] \\
&= \bar{\beta}^2 \mathbb{E}[\|\mathbf{S}^\dagger \mathbf{S} \mathbf{x}\|_2^2] + \|\mathbf{x}\|_2^2 - 2 \langle \bar{\beta} \mathbb{E}[\mathbf{S}^\dagger \mathbf{S} \mathbf{x}], \mathbf{x} \rangle \\
&= \bar{\beta}^2 \mathbb{E}[\|\mathbf{S}^\dagger \mathbf{S} \mathbf{x}\|_2^2] - \|\mathbf{x}\|_2^2 \quad (\text{Using unbiasedness of } \hat{\mathbf{x}}^{(\text{Rand-Proj-Spatial(Max)})}) \\
&= \bar{\beta}^2 \mathbf{x}^T \mathbb{E}[\mathbf{S}^T (\mathbf{S}^\dagger)^T \mathbf{S}^\dagger \mathbf{S}] \mathbf{x} - \|\mathbf{x}\|_2^2. \tag{22}
\end{aligned}$$

Using  $\mathbf{S}^\dagger = \mathbf{U} \Lambda^\dagger \mathbf{U}^T$ ,

$$\begin{aligned}
\mathbb{E}[\mathbf{S}^T (\mathbf{S}^\dagger)^T \mathbf{S}^\dagger \mathbf{S}] &= \mathbb{E}[\mathbf{U} \Lambda \mathbf{U}^T \mathbf{U} \Lambda^\dagger \mathbf{U}^T \mathbf{U} \Lambda^\dagger \mathbf{U}^T \mathbf{U} \Lambda \mathbf{U}^T] \\
&= \mathbb{E}[\mathbf{U} \Lambda (\Lambda^\dagger)^2 \Lambda \mathbf{U}^T] \\
&= \sum_{\mathbf{U}=\Phi} \mathbf{U} \mathbb{E}[\Lambda (\Lambda^\dagger)^2 \Lambda] \mathbf{U}^T \cdot \Pr[\mathbf{U} = \Phi] \\
&= \sum_{\mathbf{U}=\Phi} \mathbf{U} \left[ (1-\delta) \frac{nk}{d} \mathbf{I}_d + \sum_{c=k}^{nk-1} \delta_c \frac{c}{d} \mathbf{I}_d \right] \mathbf{U}^T \cdot \Pr[\mathbf{U} = \Phi] \\
&= \left[ (1-\delta) \frac{nk}{d} + \sum_{c=k}^{nk-1} \delta_c \frac{c}{d} \right] \cdot \sum_{\mathbf{U}=\Phi} \mathbf{U} \mathbf{U}^T \cdot \Pr[\mathbf{U} = \Phi] \\
&= \left[ (1-\delta) \frac{nk}{d} + \sum_{c=k}^{nk-1} \delta_c \frac{c}{d} \right] \mathbf{I}_d \\
&= \frac{1}{\bar{\beta}} \mathbf{I}_d \tag{23}
\end{aligned}$$



Substituting Eq. 23 in Eq. 22, we get

$$\begin{aligned} MSE(\text{Rand-Proj-Spatial(Max)}) &= \bar{\beta}^2 \mathbf{x}^T \frac{1}{\bar{\beta}} \mathbf{I}_d \mathbf{x} - \|\mathbf{x}\|_2^2 = (\bar{\beta} - 1) \|\mathbf{x}\|_2^2 \\ &\leq \left[ \frac{d}{(1-\delta)nk + \delta k} - 1 \right] \|\mathbf{x}\|_2^2, \end{aligned}$$

where the inequality is by Eq. 21. □

### C.2 Comparing against Rand- $k$

Next, we compare the MSE of Rand-Proj-Spatial(Max) with the MSE of the baseline Rand- $k$  analytically in the full-correlation case. Recall that in this case,

$$MSE(\text{Rand-}k) = \frac{1}{n} \left( \frac{d}{k} - 1 \right) \|\mathbf{x}\|_2^2.$$

We have

$$\begin{aligned} MSE(\text{Rand-Proj-Spatial(Max)}) &\leq MSE(\text{Rand-}k) \\ \Leftrightarrow \frac{d}{(1-\delta)nk + \delta k} - 1 &\leq \frac{1}{n} \left( \frac{d}{k} - 1 \right) \\ \Leftrightarrow \frac{d}{k} \frac{n - (1-\delta)n - \delta}{n((1-\delta)n + \delta)} &\leq 1 - \frac{1}{n} \\ \Leftrightarrow \frac{d}{k} \cdot \frac{\delta - \delta/n}{(1-\delta)n + \delta} &\leq \frac{n-1}{n} \\ \Leftrightarrow d\delta \left(1 - \frac{1}{n}\right)n &\leq k(n-1) \cdot ((1-\delta)n + \delta) \\ \Leftrightarrow d\delta &\leq k \cdot ((1-\delta)n + \delta) \\ \Leftrightarrow d\delta + kn\delta - k\delta &\leq kn \\ \Leftrightarrow \delta &\leq \frac{kn}{d + kn - k} \\ \Leftrightarrow \delta &\leq \frac{1}{\frac{d}{kn} + 1 - \frac{1}{n}} \end{aligned}$$

Since  $nk \leq d$ , for  $n \geq 2$ , the above implies when

$$\delta \leq \frac{1}{1 + \frac{1}{2}} = \frac{2}{3},$$

the MSE of Rand-Proj-Spatial(Max) is always less than that of Rand- $k$ .

### C.3 $\mathcal{S}$ has full rank with high probability

We empirically verify that  $\delta \approx 0$ . With  $d \in \{32, 64, 128, \dots, 1024\}$  and 4 different  $nk$  value such that  $nk \leq d$  for each  $d$ , we compute  $\text{rank}(\mathcal{S})$  for  $10^5$  trials for each pair of  $(nk, d)$  values, and plot the results for all trials. All results are presented in Figure 7. As one can observe from the plots,  $\text{rank}(\mathcal{S}) = nk$  with high probability, suggesting  $\delta \approx 0$ .

This implies the MSE of Rand-Proj-Spatial(Max) is

$$MSE(\text{Rand-Proj-Spatial(Max)}) \approx \left( \frac{d}{nk} - 1 \right) \|\mathbf{x}\|_2^2,$$

in the full correlation case.

### C.4 Proof of Theorem 4.4

**Theorem 4.4** (MSE under No Correlation). *Consider  $n$  clients, each holding a vector  $\mathbf{x}_i \in \mathbb{R}^d$ ,  $\forall i \in [n]$ . Suppose we set  $T \equiv 1$ ,  $\bar{\beta} = \frac{d^2}{k}$  in Eq. 5 and the random linear map  $\mathbf{G}_i$  at each client to be*

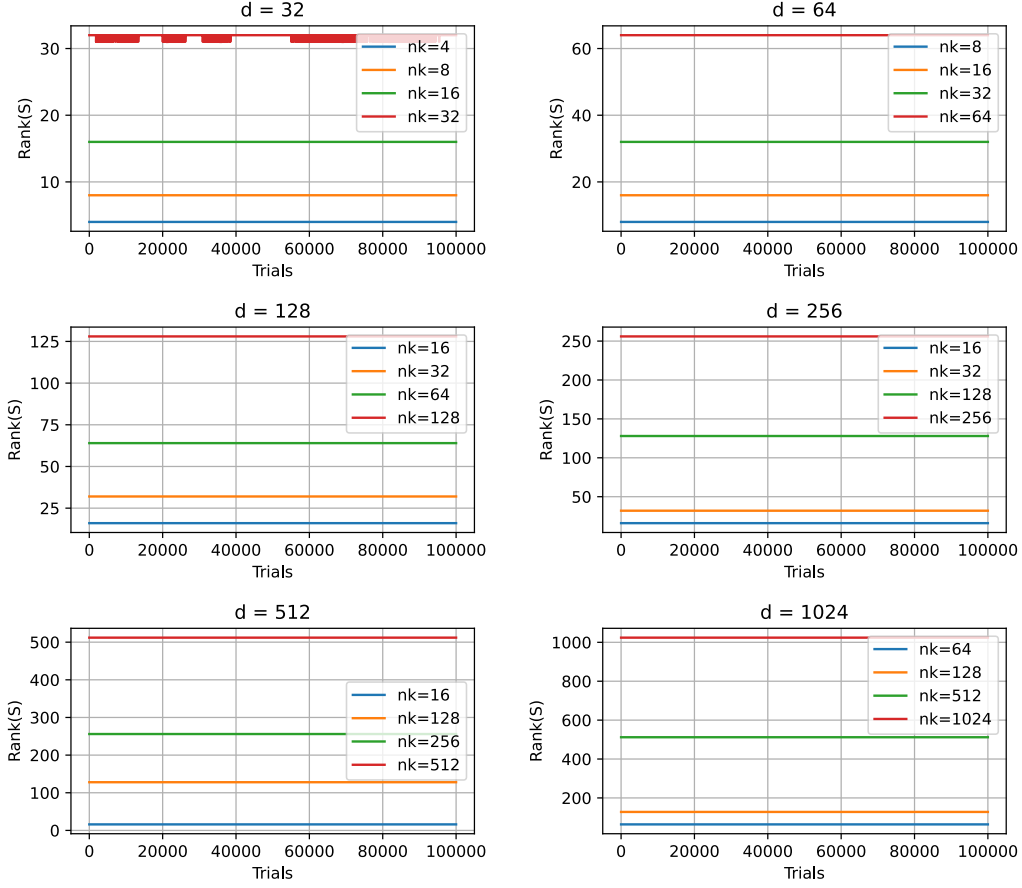


Figure 7: Simulation results of  $\text{rank}(\mathbf{S})$ , where  $\mathbf{S} = \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i$ , with  $\mathbf{G}_i$  being SRHT. With  $d \in \{32, 64, 128, \dots, 1024\}$  and 4 different  $nk$  values such that  $nk \leq d$  for each  $d$ , we compute  $\text{rank}(\mathbf{S})$  for  $10^5$  trials for each pairs of  $(nk, d)$  values and plot the results for all trials. When  $d = 32$  and  $nk = 32$  in the first plot,  $\text{rank}(\mathbf{S}) = 31$  in 2100 trials, and  $\text{rank}(\mathbf{S}) = nk = 32$  in all the rest of the trials. For all other  $(nk, d)$  pairs,  $\mathbf{S}$  always has rank  $nk$  in the  $10^5$  trials. This verifies that  $\delta = \Pr[\text{rank}(\mathbf{S}) < nk] \approx 0$ .

an SRHT matrix. Then, for  $nk \leq d$ ,

$$\mathbb{E} \left[ \|\hat{\mathbf{x}}^{(\text{Rand-Proj-Spatial})} - \bar{\mathbf{x}}\|_2^2 \right] = \frac{1}{n^2} \left( \frac{d}{k} - 1 \right) \sum_{i=1}^n \|\mathbf{x}_i\|_2^2.$$

*Proof.* When the client vectors are all orthogonal to each other, we define the transformation function on the eigenvalue to be  $T(\lambda) = 1, \forall \lambda \geq 0$ . We show that by considering the above constant  $T$ , SRHT becomes the same as rand  $k$ . Recall  $\mathbf{S} = \sum_{i=1}^n \mathbf{G}_i^T \mathbf{G}_i$  and let  $\mathbf{G}^T \mathbf{G} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$  be its eigendecomposition. Then,

$$T(\mathbf{S}) = \mathbf{U} T(\mathbf{\Lambda}) \mathbf{U}^T = \mathbf{U} \mathbf{I}_d \mathbf{U}^T = \mathbf{I}_d.$$

Hence,  $(T(\mathbf{S}))^\dagger = \mathbf{I}_d$ . And the decoded vector for client  $i$  becomes

$$\begin{aligned} \hat{\mathbf{x}}_i &= \bar{\beta} \left( T(\mathbf{G}^T \mathbf{G}) \right)^\dagger \mathbf{G}_i^T \mathbf{G}_i \mathbf{x}_i = \bar{\beta} \mathbf{G}_i^T \mathbf{G}_i \mathbf{x}_i = \bar{\beta} \frac{1}{d} \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i, \\ \hat{\mathbf{x}} &= \frac{1}{n} \sum_{i=1}^n \hat{\mathbf{x}}_i = \frac{1}{n} \bar{\beta} \sum_{i=1}^n \frac{1}{d} \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i \end{aligned} \quad (24)$$

$D_i$  is a diagonal matrix. Also,  $\mathbf{E}_i^T \mathbf{E}_i \in \mathbb{R}^{d \times d}$  is a diagonal matrix, where the  $i$ -th entry is 0 or 1.

**Computing  $\bar{\beta}$ .** To ensure that  $\hat{\mathbf{x}}$  is an unbiased estimator, from Eq. [24](#)

$$\begin{aligned}
\mathbf{x}_i &= \bar{\beta} \mathbb{E}[\mathbf{G}_i^T \mathbf{G}_i] \mathbf{x}_i \\
&= \frac{\bar{\beta}}{d} \mathbb{E}[\mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i] \mathbf{x}_i \\
&= \frac{\bar{\beta}}{d} \mathbb{E}_{D_i} \left[ \mathbf{D}_i \mathbf{H}^T \underbrace{\mathbb{E}[\mathbf{E}_i^T \mathbf{E}_i]}_{=(k/d)\mathbf{I}_d} \mathbf{H} \mathbf{D}_i \right] \mathbf{x}_i && (\because \mathbf{E}_i \text{ is independent of } D_i) \\
&= \frac{\bar{\beta}}{d} k \mathbb{E}_{D_i} [\mathbf{D}_i^2] \mathbf{x}_i && (\because \mathbf{H}^T \mathbf{H} = d\mathbf{I}_d) \\
&= \frac{\bar{\beta} k}{d} \mathbf{x}_i && (\because D_i^2 = \mathbf{I} \text{ is now deterministic.}) \\
\Rightarrow \bar{\beta} &= \frac{d}{k}. && (25)
\end{aligned}$$

**Computing the MSE.**

$$\begin{aligned}
MSE &= \mathbb{E} \left\| \hat{\mathbf{x}} - \bar{\mathbf{x}} \right\|_2^2 \\
&= \mathbb{E} \left\| \frac{1}{n} \bar{\beta} \sum_{i=1}^n \frac{1}{d} \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \\
&= \frac{1}{n^2} \left\{ \mathbb{E} \left\| \bar{\beta} \sum_{i=1}^n \frac{1}{d} \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i \right\|_2^2 + \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right. \\
&\quad \left. - 2 \left\langle \bar{\beta} \mathbb{E} \left[ \sum_{i=1}^n \frac{1}{d} \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i \right], \sum_{i=1}^n \mathbf{x}_i \right\rangle \right\} \\
&= \frac{1}{n^2} \left\{ \bar{\beta}^2 \mathbb{E} \left\| \sum_{i=1}^n \frac{1}{d} \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i \right\|_2^2 - \left\| \sum_{i=1}^n \mathbf{x}_i \right\|_2^2 \right\} && (\because \mathbb{E}[\hat{\mathbf{x}}] = \bar{\mathbf{x}}) \\
&= \frac{1}{n^2} \left\{ \sum_{i=1}^n \frac{\bar{\beta}^2}{d^2} \mathbb{E} \left\| \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i \right\|_2^2 - \sum_{i=1}^n \left\| \mathbf{x}_i \right\|_2^2 \right. && (26) \\
&\quad \left. + 2 \sum_{i=1}^n \sum_{l=i+1}^n \frac{\bar{\beta}^2}{d^2} \left\langle \mathbb{E}[\mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i], \mathbb{E}[\mathbf{D}_l \mathbf{H}^T \mathbf{E}_l^T \mathbf{E}_l \mathbf{H} \mathbf{D}_l \mathbf{x}_l] \right\rangle - 2 \sum_{i=1}^n \sum_{l=i+1}^n \left\langle \mathbf{x}_i, \mathbf{x}_l \right\rangle \right\}.
\end{aligned}$$

Note that in Eq. [26](#)

$$\begin{aligned}
\mathbb{E} \left\| \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i \right\|_2^2 &= \mathbb{E}[\mathbf{x}_i^T \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i] \\
&= d \mathbb{E}[\mathbf{x}_i^T \mathbf{D}_i \mathbf{H}^T (\mathbf{E}_i^T \mathbf{E}_i)^2 \mathbf{H} \mathbf{D}_i \mathbf{x}_i] && (\because D_i^2 = \mathbf{I}_d; \mathbf{H}^T \mathbf{H} = \mathbf{H} \mathbf{H}^T = d\mathbf{I}_d) \\
&= d \mathbf{x}_i^T \mathbb{E}_{D_i} [\mathbf{D}_i \mathbf{H}^T \mathbb{E}[\mathbf{E}_i^T \mathbf{E}_i] \mathbf{H} \mathbf{D}_i] \mathbf{x}_i && (\mathbf{E}_i, \mathbf{D}_i \text{ are independent; } (\mathbf{E}_i^T \mathbf{E}_i)^2 = \mathbf{E}_i^T \mathbf{E}_i) \\
&= kd \|\mathbf{x}_i\|_2^2, && (27)
\end{aligned}$$

since  $\mathbb{E}[\mathbf{E}_i^T \mathbf{E}_i] = (k/d)\mathbf{I}_d$ ,  $\mathbf{H}^T \mathbf{H} = d\mathbf{I}_d$  and for  $i \neq l$

$$\left\langle \mathbb{E}[\mathbf{D}_i \mathbf{H}^T \mathbf{E}_i^T \mathbf{E}_i \mathbf{H} \mathbf{D}_i \mathbf{x}_i], \mathbb{E}[\mathbf{D}_l \mathbf{H}^T \mathbf{E}_l^T \mathbf{E}_l \mathbf{H} \mathbf{D}_l \mathbf{x}_l] \right\rangle = \left\langle k \mathbf{x}_i, k \mathbf{x}_l \right\rangle = k^2 \left\langle \mathbf{x}_i, \mathbf{x}_l \right\rangle. \quad (28)$$

Substituting Eq. [27](#), [28](#) in Eq. [26](#), we get

$$\begin{aligned}
MSE &= \frac{1}{n^2} \left\{ \left( \frac{\bar{\beta}^2}{d^2} \sum_{i=1}^n kd \|\mathbf{x}_i\|_2^2 + 2 \sum_{i=1}^n \sum_{l=i+1}^n \frac{\bar{\beta}^2 k^2}{d^2} \left\langle \mathbf{x}_i, \mathbf{x}_l \right\rangle \right) - \sum_{i=1}^n \left\| \mathbf{x}_i \right\|_2^2 - 2 \sum_{i=1}^n \sum_{l=i+1}^n \left\langle \mathbf{x}_i, \mathbf{x}_l \right\rangle \right\} \\
&= \frac{1}{n^2} \left( \frac{d}{k} - 1 \right) \sum_{i=1}^n \|\mathbf{x}_i\|_2^2,
\end{aligned}$$

which is exactly the same as the MSE of rand  $k$ .  $\square$

### C.5 Rand-Proj-Spatial recovers Rand- $k$ -Spatial (Proof of Lemma 4.1)

**Lemma 4.1** (Recovering Rand- $k$ -Spatial). *Suppose client  $i$  generates a subsampling matrix  $\mathbf{E}_i = [\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_k}]^\top$ , where  $\{\mathbf{e}_j\}_{j=1}^d$  are the canonical basis vectors, and  $\{i_1, \dots, i_k\}$  are sampled from  $\{1, \dots, d\}$  without replacement. The encoded vectors are given as  $\hat{\mathbf{x}}_i = \mathbf{E}_i \mathbf{x}_i$ . Given a function  $T$ ,  $\hat{\mathbf{x}}$  computed as in Eq. [5](#) recovers the Rand- $k$ -Spatial estimator.*

*Proof.* If client  $i$  applies  $\mathbf{E}_i \in \mathbb{R}^{k \times d}$  as the random matrix to encode  $\mathbf{x}_i$  in Rand-Proj-Spatial, by Eq. [5](#) client  $i$ 's encoded vector is now

$$\hat{\mathbf{x}}_i^{(\text{Rand-Proj-Spatial})} = \bar{\beta} \left( T \left( \sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i \right) \right)^\dagger \mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i \quad (29)$$

Notice  $\mathbf{E}_i^T \mathbf{E}_i$  is a diagonal matrix, where the  $j$ -th diagonal entry is 1 if coordinate  $j$  of  $\mathbf{x}_i$  is chosen. Hence,  $\mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i$  can be viewed as choosing  $k$  coordinates of  $\mathbf{x}_i$  without replacement, which is exactly the same as Rand- $k$ -Spatial's (and Rand- $k$ 's) encoding procedure.

Notice  $\sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i$  is also a diagonal matrix, where the  $j$ -th diagonal entry is exactly  $M_j$ , i.e. the number of clients who selects the  $j$ -th coordinate as in Rand- $k$ -Spatial [\[12\]](#). Furthermore, notice  $\left( T \left( \sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i \right) \right)^\dagger$  is also a diagonal matrix, where the  $j$ -th diagonal entry is  $\frac{1}{T(M_j)}$ , which recovers the scaling factor used in Rand- $k$ -Spatial's decoding procedure.

Rand-Proj-Spatial computes  $\bar{\beta}$  as  $\bar{\beta} \mathbb{E} \left[ \left( T \left( \sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i \right) \right)^\dagger \mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i \right] = \mathbf{x}_i$ . Since  $\left( T \left( \sum_{i=1}^n \mathbf{E}_i^T \mathbf{E}_i \right) \right)^\dagger$  and  $\mathbf{E}_i^T \mathbf{E}_i \mathbf{x}_i$  recover the scaling factor and the encoding procedure of Rand- $k$ -Spatial, and  $\bar{\beta}$  is computed in exactly the same way as Rand- $k$ -Spatial does,  $\bar{\beta}$  will be exactly the same as in Rand- $k$ -Spatial.

Therefore,  $\hat{\mathbf{x}}_i^{(\text{Rand-Proj-Spatial})}$  in Eq. [29](#) with  $\mathbf{E}_i$  as the random matrix at client  $i$  recovers  $\hat{\mathbf{x}}_i^{(\text{Rand-}k\text{-Spatial})}$ . This implies Rand-Proj-Spatial recovers Rand- $k$ -Spatial in this case.  $\square$

## D Additional Experiment Details and Results

**Implementation.** All experiments are conducted in a cluster of 20 machines, each of which has 40 cores. The implementation is in Python, mainly based on `numpy` and `scipy`. All code used for the experiments can be found at <https://github.com/11hifish/Rand-Proj-Spatial>.

**Data Split.** For the non-IID dataset split across the clients, we follow [62] to split Fashion-MNIST, which is used in distributed power iteration and distributed  $k$ -means. Specifically, the data is first sorted by labels and then divided into  $2n$  shards with each shard corresponding to the data of a particular label. Each client is then assigned 2 shards (i.e., data from 2 classes). However, this approach only works for datasets with discrete labels (i.e. datasets used in classification tasks). For the other dataset UJIIndoor, which is used in distributed linear regression, we first sort the dataset by the ground truth prediction and then divides the sorted dataset across the clients.

### D.1 Additional experimental results

For each one of the three tasks, distributed power iteration, distributed  $k$ -means, and distributed linear regression, we provide additional results when the data split is IID across the clients for smaller  $n, k$  values in Section D.1.1 and when the data split is Non-IID across the clients in Section D.1.2 For the Non-IID case, we use the same settings (i.e.  $n, k, d$  values) as in the IID case.

**Discussion.** For smaller  $n, k$  values compared to the data dimension  $d$ , there is less information or less correlation from the client vectors. Hence, both Rand- $k$ -Spatial and Rand-Proj-Spatial perform better as  $nk$  increases. When  $n, k$  is small, one might notice Rand-Proj-Spatial performs worse than Rand- $k$ -Wangni in some settings. However, Rand- $k$ -Wangni is an *adaptive* estimator, which optimizes the sampling weights for choosing the client vector coordinates through an iterative process. That means Rand- $k$ -Wangni requires more computation from the clients, while in practice, the clients often have limited computational power. In contrast, our Rand-Proj-Spatial estimator is *non-adaptive* and the server does more computation instead of the clients. This is more practical since the central server usually has more computational power than the clients in applications like FL. See the introduction for more discussion.

In most settings, we observe the proposed Rand-Proj-Spatial has a better performance compared to Rand- $k$ -Spatial. Furthermore, as one would expect, both Rand- $k$ -Spatial and Rand-Proj-Spatial perform better when the data split is IID across the clients since there is more correlation among the client vectors in the IID case than in the Non-IID case.

#### D.1.1 More results in the IID case

**Distributed Power Iteration and Distributed  $K$ -Means.** We use the Fashion-MNIST dataset for both distributed power iteration and distributed  $k$ -means, which has a dimension of  $d = 1024$ . We consider more settings for distributed power iteration and distributed  $k$ -means here:  $n = 10, k \in \{5, 25, 51\}$ , and  $n = 50, k \in \{5, 10\}$ .

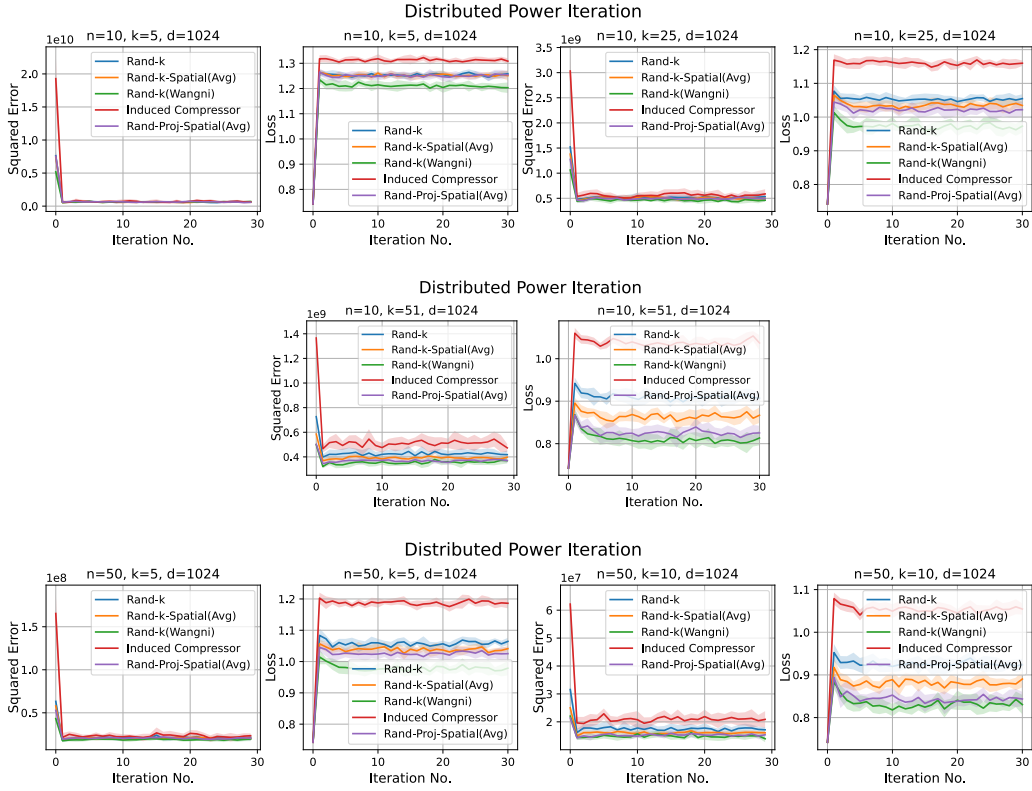


Figure 8: More results of distributed power iteration on Fashion-MNIST (IID data split) with  $d = 1024$  when  $n = 10, k \in \{5, 25, 51\}$  and when  $n = 50, k \in \{5, 10\}$ .

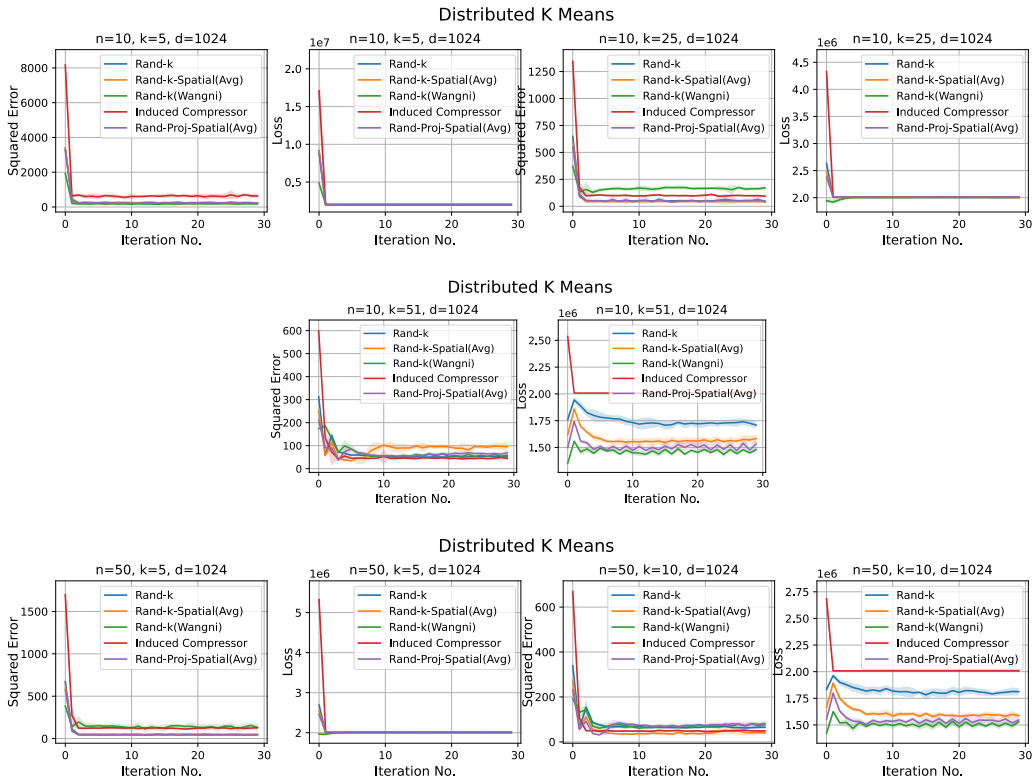


Figure 9: More results on distributed  $k$ -means on Fashion-MNIST (IID data split) with  $d = 1024$  when  $n = 10, k \in \{5, 25, 51\}$  and when  $n = 50, k \in \{10, 51\}$ .

**Distributed Linear Regression.** We use the UJIndoor dataset distributed linear regression, which has a dimension of  $d = 512$ . We consider more settings here:  $n = 10, k \in \{5, 25\}$  and  $n = 50, k \in \{1, 5\}$ .

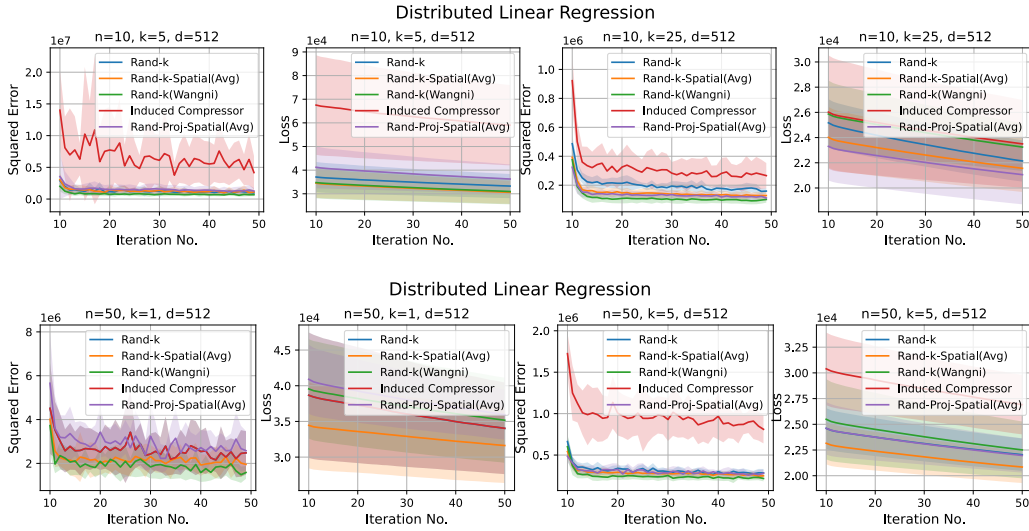


Figure 10: More results of distributed linear regression on UJIndoor (IID data split) with  $d = 512$ , when  $n = 10, k \in \{5, 25\}$  and when  $n = 50, k \in \{1, 5\}$ . Note when  $k = 1$ , the Induced estimator is the same as Rand- $k$ .

### D.1.2 Additional results in the Non-IID case

In this section, we report results when the dataset split across the clients are Non-IID, using the same datasets as in the IID case. We choose exactly the same set of  $n, k$  values as in the IID case.

**Distributed Power Iteration and Distributed  $K$ -Means.** Again, both distributed power iteration and distributed  $k$ -means use the Fashion-MNIST dataset, with a dimension  $d = 1024$ . We consider the following settings for both tasks:  $n = 10, k \in \{5, 25, 51, 102\}$  and  $n = 50, k \in \{5, 10, 20\}$ .

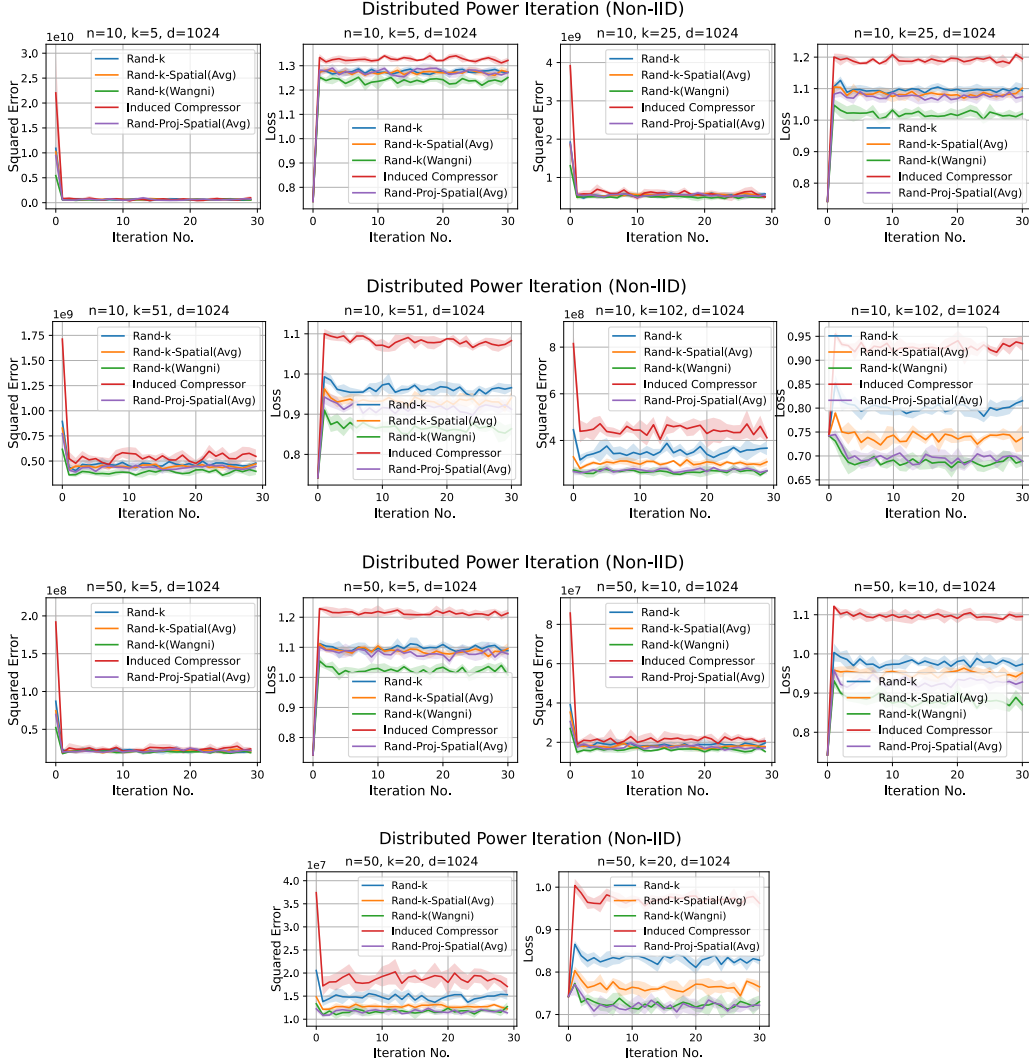


Figure 11: Results of distributed power iteration when the data split is Non-IID.  $n = 10, k \in \{5, 25, 51, 102\}$  and  $n = 50, k \in \{5, 10, 20\}$ .

**Distributed Linear Regression.** Again, we use the UJIndoor dataset for distributed linear regression, which has a dimension  $d = 512$ . We consider the following settings:  $n = 10, k \in \{5, 25, 50\}$  and  $n = 50, k \in \{1, 5, 50\}$ .



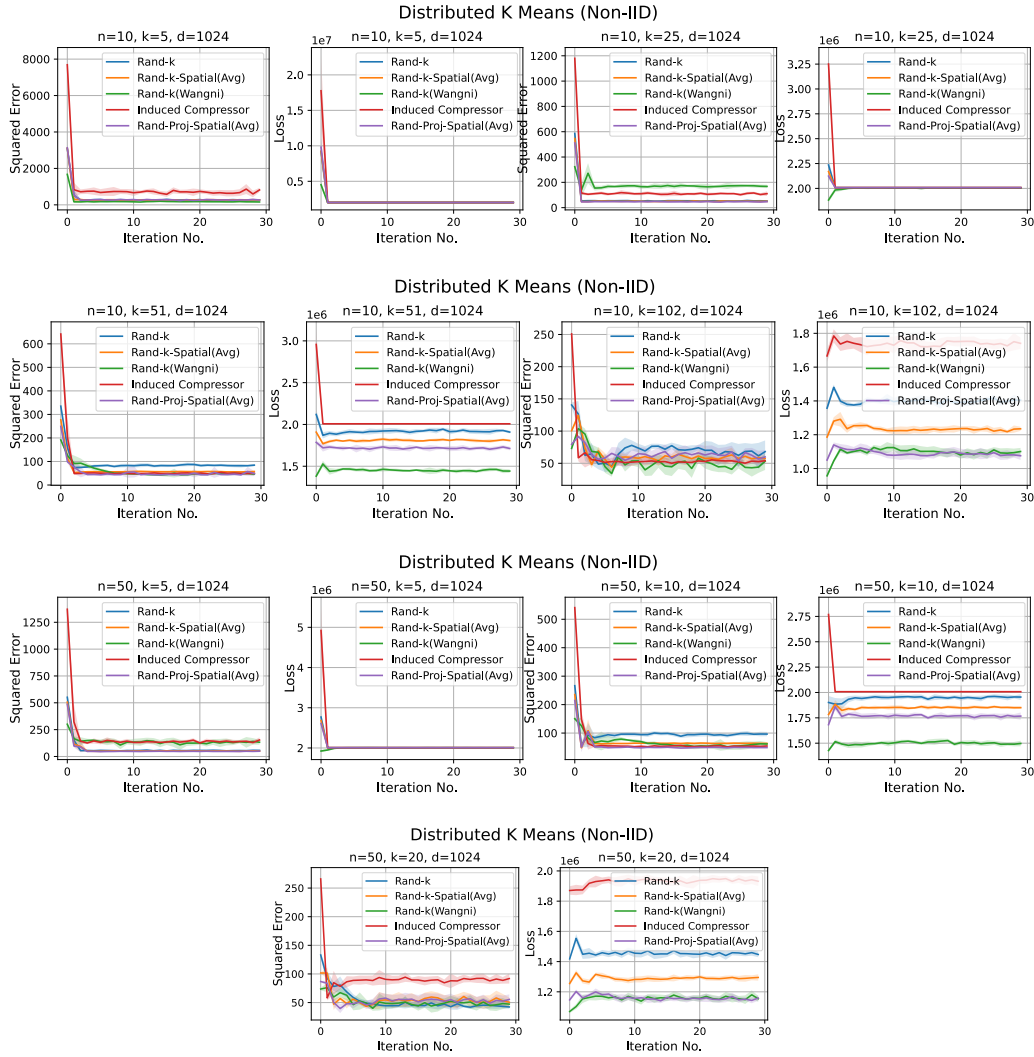


Figure 12: Results of distributed  $k$ -means when the data split is Non-IID.  $n = 10, k \in \{5, 25, 51, 102\}$  and  $n = 50, k \in \{5, 10, 20\}$ .

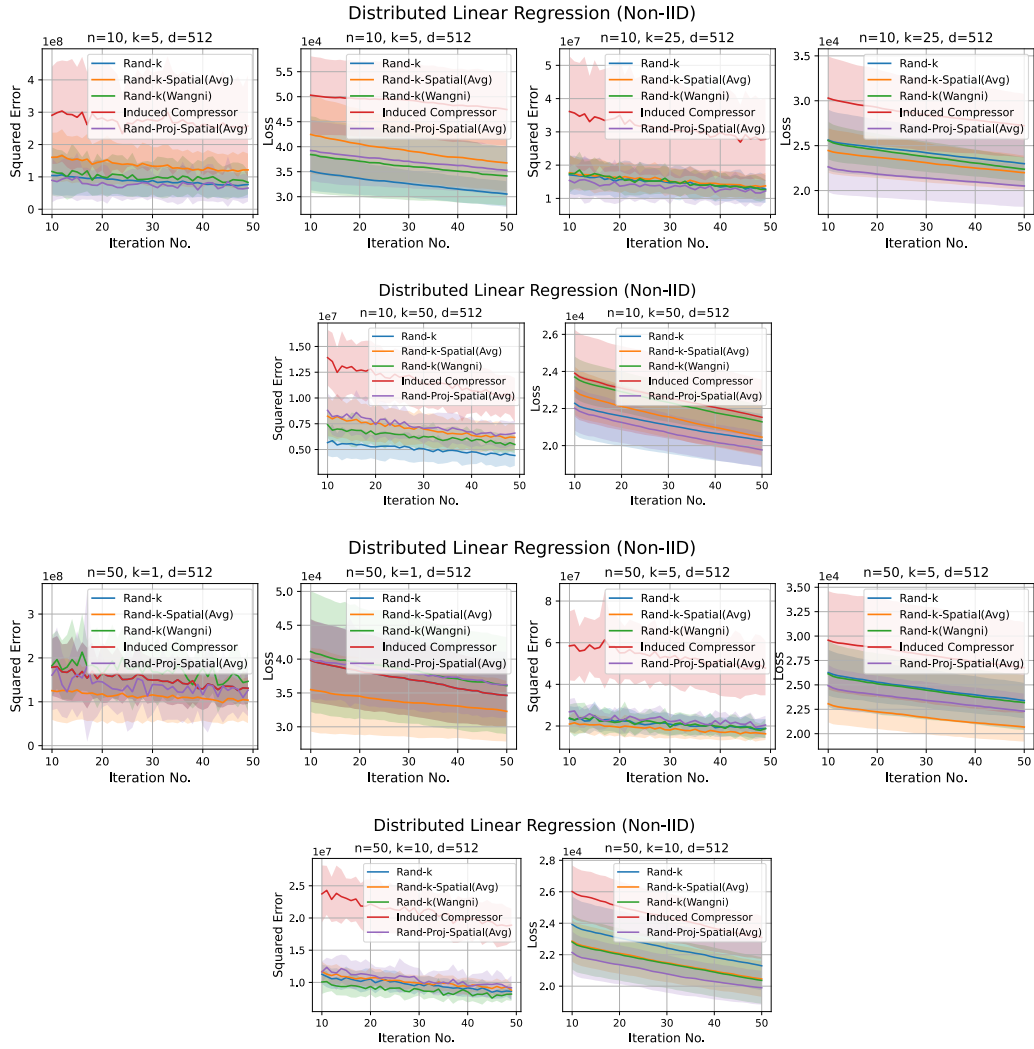


Figure 13: Results of distributed linear regression when the data split is Non-IID.  $n = 10, k \in \{5, 25, 50\}$  and  $n = 50, k \in \{1, 5, 50\}$ .