# Distributed Estimation of Gaussian Correlations

Uri Hadar and Ofer Shayevitz[*]

June 26, 2018

**Abstract**

We study a distributed estimation problem in which two remotely located parties, Alice and Bob, observe an unlimited number of i.i.d. samples corresponding to two different parts of a random vector. Alice can send $k$ bits on average to Bob, who in turn wants to estimate the cross-correlation matrix between the two parts of the vector. In the case where the parties observe jointly Gaussian scalar random variables with an unknown correlation $\rho$, we obtain two constructive and simple unbiased estimators attaining a variance of $(1 - \rho^2)/(2k \ln 2)$, which coincides with a known but non-constructive random coding result of Zhang and Berger. We extend our approach to the vector Gaussian case, which has not been treated before, and construct an estimator that is uniformly better than the scalar estimator applied separately to each of the correlations. We then show that the Gaussian performance can essentially be attained even when the distribution is completely unknown. This in particular implies that in the general problem of distributed correlation estimation, the variance can decay at least as $O(1/k)$ with the number of transmitted bits. This behavior, however, is not tight: we give an example of a rich family of distributions for which local samples reveal essentially nothing about the correlations, and where a slightly modified estimator attains a variance of $2^{-\Omega(k)}$.

## 1 Introduction and Main Results

Estimating the parameters of an unknown distribution from its samples is a basic task in many scientific problems. The vast majority of research in this field has been dedicated to the centralized setup, where a number of independent samples are being observed by the estimating entity [1]. However, in many cases the data for the estimation task might be collected by remote terminals, who then need to communicate information regarding their observations in order to perform (or improve) estimation. When the budget for communication is limited, the parties

must judiciously encode their observations and send a compressed version that is as useful as possible, creating a tension between communication and estimation.

In this paper, we study the following distributed estimation setup. Let $\mathbf{X}$ and $\mathbf{Y}$ be a pair of jointly distributed random vectors taking values in Euclidean spaces of dimensions $d_X$ and $d_Y$ respectively. Assume the distribution of the pair is only known to belong to a given family of distributions, but is otherwise arbitrary. Two remotely located parties, Alice and Bob, draw i.i.d. samples $\{(\mathbf{X}_i, \mathbf{Y}_i)\}$ from this distribution, where the $\mathbf{X}$ component is observed only by Alice and the $\mathbf{Y}$ component is observed only by Bob. The parties are interested in estimating the set of correlations between the entries of $\mathbf{X}$ and $\mathbf{Y}$ using their local samples and limited communication. Specifically, we focus on the regime where the number of samples locally available to each party is essentially *unlimited*, but only a *fixed* number of $k$ bits can be transmitted on average from (say) Alice to Bob. In this extremal regime there is no coupling between data collection and communication (typically captured by the notion of *rate*, of communication bits per data sample), and the only constraint in the system stems from its distributive nature. Moreover, we restrict attention to cases where the correlations cannot be estimated locally (e.g. Gaussian marginals do not depend on the cross-correlation parameters), which further distills the distributive aspect of the problem.

In what follows we focus mainly on the Gaussian case, i.e., where $\mathbf{X}$ and $\mathbf{Y}$ are jointly Gaussian random vectors. We begin our discussion with the scalar $d_X = d_Y = 1$ case, where our goal is to estimate the correlation coefficient $\rho$. The only work we are aware of that deals with distributed estimation of the bivariate normal correlation under communication constraints is by Zhang and Berger [2], who studied the problem as an application of a more general result. Using random coding techniques, they proved the existence of an asymptotically unbiased estimator whose variance they provided as a function of the number of samples and the rate $R$ of communication bits per sample. Specializing to our setup by plugging in $k/R$ as the number of samples, the Zhang-Berger variance is given by

$$\mathsf{Var}\,\hat{\rho}_{ZB} = \frac{R}{k}\left(1 + \rho^2 + \frac{1-\rho^2}{2^{2R}-1} + o(1)\right). \tag{1}$$

Since we do not impose a rate constraint in our setup, we can minimize the variance over $R$ to obtain

$$\inf_{R>0} \mathsf{Var}\,\hat{\rho}_{ZB} = \frac{1}{k}\left(\frac{1-\rho^2}{2\ln 2} + o(1)\right), \tag{2}$$

which is attained (not surprisingly) in the zero-rate limit as $R \to 0$. It should be noted that this estimator was not claimed to be optimal in any sense. Furthermore, as the authors themselves indicate, the results in [2] apply only to the single scalar parameter case, and it is not clear how to extend this approach to the vector case.

In this Gaussian scalar setup, addressed in Section 2, we introduce the following constructive scheme: Alice sends to Bob the index $J$ of the largest sample

2

among her first $2^k$ samples, and Bob computes the unbiased estimator

$$\hat{\rho}_{\max} = \frac{Y_J}{\mathbb{E}\, X_J} \approx \frac{Y_J}{\sqrt{2k \ln 2}}. \tag{3}$$

In Theorem 1, we show that this simple estimator attains the same variance as the non-constructive Zhang-Berger estimator (2), i.e.,

$$\mathsf{Var}\,\hat{\rho}_{\max} = \frac{1}{k}\left(\frac{1-\rho^2}{2\ln 2} + o(1)\right). \tag{4}$$

Then, in preparations for the vector case, we describe a simple variation of this estimator: Alice scans her samples sequentially and finds the index $J$ of the first sample to pass a suitably chosen threshold. She then compresses this index using an optimal lossless variable-length code and sends the encoded version to Bob, who computes an estimator using his corresponding $Y$ sample, in a way similar to the maximum estimator above. This threshold estimator is unbiased, and also attains the Zhang-Berger variance. We note that the maximal/threshold-passing value of a scalar i.i.d. Gaussian sequence has been employed before in problems of writing on dirty paper [3], [4], and Gaussian lossy source coding [5].

We proceed to consider the vector Gaussian setup (Section 3). Without loss of generality, we assume that both parties know the distribution of Alice's vector, since she can estimate it arbitrarily well from her local samples and send a sufficiently accurate quantization to Bob with what can be shown to be a negligible cost in communication. In the case where $d_X = 1$ and $d_Y > 1$ we can trivially extend the scalar estimator by having Alice perform the same encoding (maximal or threshold) and have Bob apply the same type of estimation to each of the entries of $\mathbf{Y}$ using the single index obtained from Alice. The case of $d_X > 1$ and $d_Y = 1$ is more interesting. Of course, one could simply estimate each one of the correlations $\rho_\ell$ between $(\mathbf{X})_\ell$ and $Y$ separately by repeating the scalar method. A worthy goal is therefore to find an estimator that *dominates* the scalar approach, uniformly for all correlation values. In Proposition 3, we show that performing general linear operations (e.g., whitening the signal) before applying the scalar estimator, does not dominate the scalar approach. We then describe a multidimensional estimator that *does* dominate the scalar approach, by generalizing the scalar threshold to an appropriately chosen $d_X$-dimensional *stopping set*. We show that the resulting (constructive) estimator $\hat{\boldsymbol{\rho}}$ attains a total mean squared error that is a function of *the highest correlation only*, and is given by

$$\mathbb{E}\,\|\hat{\boldsymbol{\rho}} - \boldsymbol{\rho}\|^2 \leq \frac{1}{k}\left(\frac{d_X^2}{2\ln 2} \min_{\ell \in [d]}\{1 - \rho_\ell^2\} + o(1)\right). \tag{5}$$

This is proved Theorem 4. We note that the case of $d_X, d_Y > 1$ is again a trivial extension of the $d_X > 1, d_Y = 1$ case.

Returning to the general non-Gaussian setup (Section 4), we provide two additional results. In Section 4.1 we show how our estimators above can essentially

be used to obtain the *same variance guarantees* when $(\mathbf{X}, \mathbf{Y})$ are *arbitrarily distributed*, subject only to uniform integrability fourth moment conditions. This in particular means that one can always get a $O(1/k)$ variance in distributed correlation estimation with $k$ transmitted bits on average. Recall that in centralized estimation problems, when the family of distributions is sufficiently smooth in the parameter of interest, the Cramér–Rao lower bound implies that the optimal estimation variance is $\Theta(1/n)$, where $n$ is the number of samples. Thus, the centralized number of samples required to achieve the same variance as in the distributed case is at least linear in the number of communication bits, i.e., each communication bit is worth at least a constant number of samples. It is perhaps tempting to guess that this relation is fundamental, i.e., that a bit is equivalent to a constant number of samples, hence that the variance cannot decrease faster than $\Omega(1/k)$, assuming that the family of distributions is such that Bob cannot estimate the correlations from his local samples. While we conjecture this is true in the Gaussian case, it does not hold in general: In Subsection 4.2 we give an example of a rich family of distributions for which local samples reveal essentially nothing about the correlations, and where the variance of our (slightly modified) estimator is $2^{-\Omega(k)}$.

## 1.1 Related Work

The problem of distributed estimation under communication constraints has been studied in the last couple of decades by several authors. Zhang and Berger [2] used random coding techniques to establish the existence of an asymptotically unbiased estimator whose variance is *upper* bounded by a single-letter expression. Their results are limited to a certain family of joint distributions (that must satisfy an *additivity* condition) that depend on a one-dimensional parameter. Ahlswede and Burnashev [6] gave a multi-letter lower bound on the minimax estimation variance in the one-dimensional case. Han and Amari [7] (see also the survey paper [8]) suggested a rate constrained encoding scheme, and obtained the likelihood equation based on the decoded statistic. They also showed that the estimation variance asymptotically achieves the inverse of the Fisher information of that statistic. Their results only apply to finite alphabets. Amari [9] discussed optimal compression in the specific setting of estimating the correlation between two binary sources. He showed that under linear-threshold encoding, there does not exist a single scheme that is uniformly optimal for all correlation values. A similar setup was discussed by Haim and Kochman [10] in the context of hypothesis testing between two correlation values. Zhang *et al* [11] provided minimax lower bounds for a distributed estimation setting in which all terminals observe samples from the same distribution. El Gamal and Lai [12] showed that Slepian-Wolf rates are not necessary for distributed estimation over finite alphabets.

There is a rich literature addressing other aspects of the distributed estimation problem. Xiao *et al* [13] and Lou [14] considered distributed estimation of a location parameter under energy and bandwidth constraints. Gubner [15] considered a Bayesian distributed estimation setting and suggested a local quan-

tization algorithm. Xu and Raginsky [16] provided lower bounds on the risk in a distributed Bayesian estimation setting with noisy channels between the data collection terminals and the estimation entity. Braverman *et al* [17] provided lower bounds for some high dimensional distributed estimation problems, again when the samples of all terminals are from the same distribution, e.g. for distributed estimation of the multivariate Guassian mean when it is known to be sparse. The authors of [18], [19], [20] and [21] addressed various distributed estimation setups where the measurements across the sensors are assumed to be independent.

## 1.2 Notations and preliminaries

The standard normal density is denoted by $\phi(x) = e^{-x^2/2}/\sqrt{2\pi}$, and the tail probability by $Q(x) \triangleq \int_x^\infty \phi(t)dt$. For $Z \sim \mathcal{N}(0,1)$ the inverse Mills ratio is denoted by

$$s(x) \triangleq \mathbb{E}(Z \mid Z > x) = \frac{\phi(x)}{Q(x)}. \tag{6}$$

We write log and ln for the base 2 and natural logarithm, respectively. The *entropy* of the geometric distribution with parameter $p$ is given by $h_g(p) \triangleq h(p)/p$, where $h(p) = -p \log p - (1-p) \log(1-p)$ is the binary entropy function. Note that $h_g(p) = -\log(p)(1 + o(1))$ as $p \to 0$. Recall also that any discrete random variable (e.g. in our case, a geometric r.v.) can be losslessly encoded using a prefix-free code with expected length exceeding its entropy by at most one bit [22]. Since in the setups we consider the entropy grows large, this excess one bit has vanishing effect on our results, hence for the sake of readability we disregard it throughout.

For any natural $n$ we denote $[n] \triangleq \{1, \ldots, n\}$. For a vector $\mathbf{v}$, the $i$-th coordinate is denoted by $(\mathbf{v})_i$. Similarly, $(M)_{ij}$ denotes the $ij$-th entry of the matrix $M$. The $d \times d$ identity matrix is denoted by $\mathbf{I}_d$. We use the standard order notation; in the following, $f$ and $g$ are positive functions with discrete or continuous domain. We write $f = o(g)$ to indicate that $\lim f/g = 0$, and $f = O(g)$ to indicate that $\limsup f/g < \infty$, where the arguments and implied limits should be clear from the context. Writing $f = \Omega(g)$ means that $g = O(f)$, and $f = \Theta(g)$ means that both $f = O(g)$ and $f = \Omega(g)$.

Given a statistic $T$, and a scalar parameter $\theta$ we wish to estimate, The Fisher information of estimating $\theta$ from $T$ (see e.g. [1]) is given by

$$\mathrm{I}_T(\theta) \triangleq \mathbb{E}\left[\left(\frac{\partial \log f(T \mid \theta)}{\partial \theta}\right)^2\right], \tag{7}$$

where $f(t \mid \theta)$ is the p.d.f. of $T$ for the given value of $\theta$. The Cramér–Rao lower bound (CRLB) states that, under some regularity conditions (see e.g. [1]) that are trivially satisfied in our setups, any unbiased estimator $\hat{\theta} = \hat{\theta}(T)$ of $\theta$ satisfies

$$\mathsf{Var}\,\hat{\theta} \geq 1/\mathrm{I}_T(\theta). \tag{8}$$

An estimator $\hat{\theta}$ that satisfies (8) with equality is said to be *efficient*. We emphasize that the efficiency is with respect to the statistic $T$ by saying it is *efficient given $T$*. The estimators and statistics in this paper depend on the number of communicated bits, $k$. We call an estimator $\hat{\theta}$ *asymptotically efficient given $T$* if $\mathbb{E}\,\hat{\theta} \to \theta$, and $\mathrm{I}_T(\theta) \cdot \mathsf{Var}\,\hat{\theta} \to 1$ as $k \to \infty$. The estimated parameter may be vector valued, in which case $\mathrm{I}_T(\boldsymbol{\theta})$ is a matrix given by

$$\mathrm{I}_T(\boldsymbol{\theta}) \triangleq \mathbb{E}\left[\left(\frac{\partial \log f(T \mid \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)^T \cdot \left(\frac{\log f(T \mid \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right)\right], \tag{9}$$

and the CRLB reads $\mathsf{Cov}\,\hat{\boldsymbol{\theta}} \geq \mathrm{I}_T^{-1}$ where the inequality is in the positive semidefinite sense. In the vector case we say that an estimator $\hat{\boldsymbol{\theta}}(T)$ is asymptotically efficient if the estimator $\mathbf{v}^T\hat{\boldsymbol{\theta}}(T)$ of $\mathbf{v}^T\boldsymbol{\theta}$ is asymptotically efficient for any $\mathbf{v} \in \mathbb{R}^{\dim(\boldsymbol{\theta})}$. We note that since the aforementioned regularity conditions are satisfied in our Gaussian setups, then (asymptotic) efficiency of an estimator implies that it is *(asymptotically) minimum variance unbiased*.

The Fisher information matrix of a Gaussian vector with mean $\mu$ and covariance matrix $\Sigma$, where both are functions of a parameter vector $\boldsymbol{\theta}$, is given by (see e.g. [23])

$$(\mathrm{I})_{ij} = \frac{\partial \mu^T}{\partial (\boldsymbol{\theta})_i}\Sigma^{-1}\frac{\partial \mu}{\partial (\boldsymbol{\theta})_j} + \frac{1}{2}\,\mathrm{tr}\left(\Sigma^{-1}\frac{\partial \Sigma}{\partial (\boldsymbol{\theta})_i}\Sigma^{-1}\frac{\partial \Sigma}{\partial (\boldsymbol{\theta})_j}\right). \tag{10}$$

A common setup throughout is where a parameter is estimated from $(\mathbf{X}, \mathbf{Y})$ where $\mathbf{Y}|\mathbf{X}$ is Gaussian, and the distribution of $\mathbf{X}$ does not depend on the parameter. In this case we have

$$\mathrm{I}_{\mathbf{X},\mathbf{Y}} = \mathbb{E}_{\mathbf{X}}\,\mathrm{I}_{\mathbf{Y}|\mathbf{X}} \tag{11}$$

where $\mathrm{I}_{\mathbf{Y}|\mathbf{X}}$ is obtained via (10).

## 2 Estimating a single correlation

In this section, we consider the case where $X$ and $Y$ are both scalar, jointly Gaussian r.v.s, with unknown parameters satisfying only $\mathbb{E}\,X^2, \mathbb{E}\,Y^2 < u$ for some known $u$. Since the number of local samples available to Alice and Bob is unlimited, they can both estimate their own mean and variance arbitrarily well (taking $u$ into account) and normalize their samples accordingly. Hence, without loss of generality we can assume that $X, Y \sim \mathcal{N}(0, 1)$, and that the only unknown parameter is their correlation coefficient $\rho$. This model can be written as

$$Y = \rho X + \sqrt{1 - \rho^2}Z \tag{12}$$

where $Z \sim \mathcal{N}(0, 1)$ is statistically independent of $X$.

Alice, who observes the i.i.d. samples $\{X_i\}$, can transmit $k$ bits on average to Bob, who observes the corresponding $\{Y_i\}$ samples and would like to obtain a

good estimate of $\rho$ in the mean squared error sense. We note that the conditional Fisher information of $\rho$ from $Y$, given that $X = x$, is

$$\mathrm{I}_{Y|X=x}(\rho) = \frac{(1-\rho^2)x^2 + 2\rho^2}{(1-\rho^2)^2},$$
(13)

which is linear in $x^2$. This motivates using an estimator based on a measurement for which $|x|$ is as large as possible. The same can also be intuitively deduced from (12), since if one controls $X$, then picking it as large as possible would "maximize the SNR". For simplicity, we look at large positive values of $x$ rather than large value of $|x|$. Our derivations can be easily modified to hold in the latter case (with one extra bit describing the sign) without affecting the results.

## 2.1 Max estimator

Following the heuristic discussion above, consider the following scheme. Given the constraint $k$ on the expected number of communication bits, Alice looks at her first $2^k$ samples, finds the maximal one, and sends its index

$$J = \operatorname*{argmax}_{i \in [2^k]} X_i$$
(14)

to Bob, using exactly $k$ bits. Bob now looks at $Y_J$, his sample that corresponds to the same index, which we refer to as the *co-max*[1]. If Bob were in possession of $X_J$ as well, and observing the model (12) again, a natural estimator for $\rho$ he could have used is $Y_J/X_J$. In fact, it can be shown that this estimator is an approximated solution to the maximum likelihood equation, which is third a degree polynomial in this case (see appendix A.7). However, since $X_J$ is not available, Bob uses the estimator

$$\hat{\rho}_{\max} = \frac{Y_J}{\mathbb{E}\, X_J}$$
(15)

that depends only on $J$ (communicated by Alice) and on his own samples. The following Theorem shows that this simple estimator attains the same variance as the non-constructive Zhang-Berger estimator (2), and also that knowing the value of $X_J$ does not help.

**Theorem 1.** *The estimator $\hat{\rho}_{\max}$ is unbiased with*

$$\mathsf{Var}\, \hat{\rho}_{\max} = \frac{1}{k}\left(\frac{1-\rho^2}{2\ln 2} + o(1)\right)$$
(16)

*where $k$ is the number of transmitted bits. Furthermore, $\hat{\rho}_{\max}$ is asymptotically efficient given $(X_J, Y_J)$.*

---

[1] This is also known in the literature as the *max concomitant*, see e.g. [24]

7

*Proof.* It is easy to check that $\hat{\rho}_{\max}$ is unbiased. In order to compute its variance, we need to compute the mean and variance of $X_J$, which is the maximum of $2^k$ i.i.d. standard normal r.v.s. From extreme value theory (see e.g. [24]) applied to the normal distribution case, we obtain:

$$\mathbb{E}\, X_J = \sqrt{2\ln(2^k)}(1 + o(1)) \tag{17}$$

$$\mathbb{E}\, X_J^2 = 2\ln(2^k)(1 + o(1)) \tag{18}$$

$$\mathsf{Var}\, X_J = O\left(\frac{1}{\ln(2^k)}\right). \tag{19}$$

Therefore, we have that

$$\mathsf{Var}\, \hat{\rho}_{\max} = \frac{1}{(\mathbb{E}\, X_J)^2}\, \mathsf{Var}(\rho X_J + \sqrt{1-\rho^2}Z) \tag{20}$$

$$= \frac{1}{(\mathbb{E}\, X_J)^2}(\rho^2\, \mathsf{Var}\, X_J + 1 - \rho^2) \tag{21}$$

$$= \frac{1}{2k\ln 2}(1 - \rho^2 + o(1)). \tag{22}$$

Now, recalling (13), the Fisher Information of $\rho$ from $(X_J, Y_J)$ is given by

$$\mathrm{I}_{X_J Y_J}(\rho) = \frac{(1-\rho^2)\,\mathbb{E}\, X_J^2 + 2\rho^2}{(1-\rho^2)^2} \tag{23}$$

$$= 2k\ln 2\left(\frac{1}{1-\rho^2} + o(1)\right), \tag{24}$$

and hence $\hat{\rho}_{\max}$ is asymptotically efficient given $(X_J, Y_J)$. $\qquad\square$

Theorem 1 suggests that using a better estimator of $X_J$ in lieu of its expectation (by having Alice send some quantization of $\hat{X}_J$ and having Bob compute $\hat{\rho} = Y_J/\hat{X}_J$) would not improve the performance asymptotically, as $\hat{\rho}_{\max}$ is optimal among all unbiased estimators that use both the max and co-max. In Section 4.2, we will see that this observation does not extend to some other additive models.

We note that the random coding Zhang-Berger estimator only deals with the scalar case, and as the authors themselves indicate [2], it remains unclear whether it could be extended to the the case of multiple correlations. In contrast, our constructive approach can also be naturally extended to the multidimensional case. To that end, it is instructive to first describe a simple variation of our scalar estimator.

## 2.2    Threshold estimator

We now introduce a simple modification to max estimator that will be useful in the sequel. Instead of taking the maximum of a fixed number of measurements, Alice sequentially scans her samples until she finds a sample that exceeds some

fixed threshold, to be determined later. She then sends the index of this sample to Bob, who proceeds similarly to the max method. The main difference is that using the max method Alice sends a fixed number of bits, whereas using the threshold method she sends a random number of bits. In this subsection, we introduce and analyze the threshold estimator and demonstrate that it is asymptotically equivalent to the max estimator, in terms of how the estimation variance is related to the expected number of bits transmitted. As mentioned above, the main motivation for studying the threshold estimator is that in contrast to the max estimator, it can be naturally extended to the multidimensional case.

Precisely, let

$$J = \min\{i : X_i > t\}, \tag{25}$$

and consider the estimator

$$\hat{\rho}_{\text{th}} = \frac{Y_J}{\mathbb{E}\, X_J}. \tag{26}$$

Note that the index $J$ is distributed geometrically with parameter $p = \Pr(X > t) = Q(t)$. Alice can therefore represent $J$ using a prefix-free code (e.g., Huffman) with at most $h_g(p) + 1$ bits on average, where $h_g(p)$ is the entropy of this geometric distribution [22]. For brevity of exposition, we assume that the expected number of bits is exactly $k = h_g(p)$, as this does not affect the asymptotic behavior. Therefore, to satisfy the communication constraint the threshold must be set to

$$t = Q^{-1}(h_g^{-1}(k)). \tag{27}$$

We later show that $t = \sqrt{2k \ln 2}(1 + o(1))$ as $k$ grows large. The following Theorem shows that as the max estimator, the threshold estimator also attains the same variance as the non-constructive Zhang-Berger estimator (2), and also that knowing the value of $X_J$ again does not help.

**Theorem 2.** *The estimator $\hat{\rho}_{\text{th}}$ is unbiased with*

$$\mathsf{Var}\, \hat{\rho}_{\text{th}} = \frac{1}{k}\left(\frac{1 - \rho^2}{2 \ln 2} + o(1)\right) \tag{28}$$

*where $k$ is the expected number of transmitted bits. Furthermore, $\hat{\rho}_{\text{th}}$ is asymptotically efficient given $(X_J, Y_J)$.*

*Proof.* It is immediate to verify that $\hat{\rho}_{\text{th}}$ is unbiased. We have from (6) that $\mathbb{E}\, X_J = s(t)$, and straightforward calculations give $\mathbb{E}\, X_J^2 = 1 + ts(t)$. Also it is known that $t \le s(t) \le t + t^{-1}$, and that (see e.g. [25])

$$\frac{1}{s(t)} = \frac{1}{t} - \frac{1}{t^3} + \frac{3}{t^5} + O\left(\frac{1}{t^7}\right). \tag{29}$$

9

Combining the above yields

$$\mathbb{E}\, X_J^2 = t^2(1 + o(1)), \quad \mathsf{Var}\, X_J = 1/t^2 + O(1/t^4). \qquad (30)$$

Let us now express the threshold $t$ in terms of $k$. We have $h_g(p) = -\log(p)(1 + o(1))$ as $p \to 0$, and also that $-\ln Q(t) = \frac{t^2}{2}(1 + o(1))$. Therefore the expected number of bits sent by Alice is

$$\begin{align}
k &= h_g(Q(t)) &&(31)\\
&= -\log(Q(t))(1 + o(1)) &&(32)\\
&= t^2\left(\frac{1}{2\ln 2} + o(1)\right), &&(33)
\end{align}$$

which yield $t = \sqrt{2k\ln 2}/(1 + o(1))$. Combining this with (30) and recalling the model (12), we obtain

$$\begin{align}
\mathsf{Var}\, \hat{\rho}_{\text{th}} &= \frac{1}{s^2}\left(\rho^2\, \mathsf{Var}\, X_J + 1 - \rho^2\right) &&(34)\\
&= \frac{1 - \rho^2}{t^2}(1 + o(1)) &&(35)\\
&= \frac{1}{k}\left(\frac{1 - \rho^2}{2\ln 2} + o(1)\right). &&(36)
\end{align}$$

Recalling (13), the Fisher information is given by

$$\mathrm{I}_{X_J Y_J} = \frac{t^2}{1 - \rho^2}(1 + o(1)) = \frac{2k\ln 2}{1 - \rho^2}(1 + o(1)), \qquad (37)$$

concluding the proof. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad\square$

Note that unlike the maximum estimator, the threshold estimator's variance admits an exact non-asymptotic expression:

$$\mathsf{Var}\, \hat{\rho}_{\text{th}} = \frac{1}{s^2(t)}(1 - \rho^2(s(t) - t) \cdot s(t)), \qquad (38)$$

where $t = Q^{-1}(h_g^{-1}(k))$.

## 3 Estimating multiple correlations

We proceed to address the more challenging multidimensional case where $\mathbf{X}, \mathbf{Y}$ are jointly Gaussian random vectors with unknown parameters. As in the scalar case, we only assume that the variances of all the entries of both $\mathbf{X}$ and $\mathbf{Y}$ are bounded by some known constant, hence Alice and Bob can compute the means and variances of their samples, and normalize them accordingly. Thus, without loss of generality we can assume that all the entries of $\mathbf{X}$ and $\mathbf{Y}$ have zero mean

and unit variance. In fact, for the same reasons we can assume that Alice knows $\mathsf{Cov}\,\mathbf{X}$ and Bob knows $\mathsf{Cov}\,\mathbf{Y}$.

As before, Alice observes the i.i.d. samples $\{\mathbf{X}_i\}$ and can transmit $k$ bits on average to Bob, who observes the corresponding $\{\mathbf{Y}_i\}$ samples and would like to obtain a good estimate of $\mathbb{E}\,\mathbf{Y}\mathbf{X}^T$, the collection of all the correlations between the different entries of $\mathbf{X}$ and $\mathbf{Y}$. For simplicity, our performance measure will be the expected sum of squared estimation errors across all such correlations.

Below we discuss the two extremal setups: The case where $X$ is a scalar and $\mathbf{Y}$ is a vector, and the opposite case where $\mathbf{X}$ is a vector and $Y$ is a scalar. This is sufficient since estimators for the general setup where both $\mathbf{X}, \mathbf{Y}$ are vectors are straightforward to construct by combining the two extremal setups, hence discussing this more general setup adds no useful insight. Clearly, the scalar methods suggested in Section 2 can be directly applied to the multidimensional case, by allocating the bits between the tasks of estimating each correlation separately. It is therefore interesting to try and find a truly multidimensional scheme that *dominates* the scalar method, i.e., performs at least as good uniformly for all possible values of the correlations.

## 3.1  $X$ is a scalar, $\mathbf{Y}$ is a vector

Suppose $(X, \mathbf{Y})$ are jointly Gaussian, where $X \sim \mathcal{N}(0, 1)$, $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{Y}})$ is a $d$-dimensional (column) vector, and $\Sigma_{\mathbf{Y}}$ has an all-ones diagonal and is known to Bob, who is interested in estimating the column correlation vector

$$\boldsymbol{\rho} = \mathbb{E}\,\mathbf{Y}X = [\rho_1, \ldots, \rho_d]^T. \tag{39}$$

The natural extension of the two scalar methods of Section 2 to this case is obvious. Here we analyze the threshold method, yet the max method is as simple and would yield the same results. Alice waits until $X_i$ passes a threshold $t > 0$ and transmits the resulting index

$$J = \min\{i : X_i > t\} \tag{40}$$

to Bob, where $t = Q^{-1}(h_g^{-1}(k))$. The estimator is then

$$\hat{\boldsymbol{\rho}} = \frac{1}{\mathbb{E}\,X_J}\mathbf{Y}_J = \frac{1}{s(t)}\mathbf{Y}_J, \tag{41}$$

which is an unbiased approximation of the maximum likelihood estimator (see Appendix A.7).

**Theorem 3.** *The estimator $\hat{\boldsymbol{\rho}}$ in* (41) *is unbiased with*

$$\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}} = \frac{1}{k}\left(\frac{1}{2\ln 2}\sum_{\ell=1}^{d}(1 - \rho_\ell^2) + o(1)\right) \tag{42}$$

*where $k$ is the expected number of transmitted bits. Furthermore, $\hat{\boldsymbol{\rho}}$ is asymptotically efficient given $(X_J, \mathbf{Y}_J)$.*

11

*Proof.* This is simple consequence of Theorem 2, except for asymptotic efficiency which we prove in Appendix A.1. $\square$

This method (trivially) dominates the scalar method applied separately to each of the correlations, as the latter would yield $\sum \mathsf{Var}\, \hat{\rho}_i = \frac{1}{k}\left(\frac{d}{2\ln 2}\sum_{\ell=1}^{d}(1-\rho_\ell^2) + o(1)\right)$.

## 3.2  X is a vector, $Y$ is a scalar

Consider the setup where $(\mathbf{X}, Y)$ are jointly Gaussian where $Y \sim \mathcal{N}(0,1)$ and $\mathbf{X}$ is a $d$-dimensional (column) vector $\sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathbf{X}})$ where $\Sigma_{\mathbf{X}}$ is known to Alice and has an all-ones diagonal. Alice observes $\{\mathbf{X}_i\}$ and transmits $k$ bits to Bob on average, who observes $\{Y_i\}$ and wishes to estimate the *row* vector

$$\boldsymbol{\rho} = \mathbb{E}\, Y\mathbf{X}^T = [\rho_1, \ldots, \rho_d]. \tag{43}$$

The model can be written as (see e.g. [26])

$$Y = \boldsymbol{\rho}\,\Sigma_{\mathbf{X}}^{-1}\mathbf{X} + \sigma Z \tag{44}$$

where $Z \sim \mathcal{N}(0,1)$ is independent of $\mathbf{X}$, and $\sigma^2 = 1 - \boldsymbol{\rho}\,\Sigma_{\mathbf{X}}^{-1}\boldsymbol{\rho}^T$.

A naive extension of the scalar method to this setup would be to allocate the bits between the correlations and apply the scalar (max or threshold) scheme $d$ times, using the fact that the model can also be written as

$$Y = \rho_\ell(\mathbf{X})_\ell + \sqrt{1-\rho_\ell^2}\,Z \tag{45}$$

for any $\ell \in [d]$. One could suggest to improve performance by having Alice locally perform some general linear operation on $\mathbf{X}$ before applying the scalar method, then having Bob perform the inverse operation. While this can indeed help for certain correlation values, it cannot improve the performance uniformly, even if the linear operation can depend on $\Sigma_{\mathbf{X}}$ (hence can e.g. whiten $\mathbf{X}$). See Appendix A.2 for details.

We now introduce an estimator that *does* dominate the scalar method. In fact, the mean squared error attained by this estimator is dictated by the single "best" entry of $\boldsymbol{\rho}$, namely by the highest correlation only. Our method is based on replacing the scalar one-dimensional threshold by $d$-dimensional *stopping sets* $A_1, \ldots, A_d \subset \mathbb{R}^d$. Similarly to the scalar case, Alice waits until $\mathbf{X}_i \in A_1$ for the first time, then again until $\mathbf{X}_i \in A_2$, and so on[2] until $\mathbf{X}_i \in A_d$. Alice then describes the resulting indices $J_1, \ldots J_d$ to Bob using an optimal variable-rate prefix-free code of expected length equal to the entropy of the associated geometric distribution (again, we neglect the excess one bit). Defining Alice's corresponding sample matrix $\mathbf{X_J} = [\mathbf{X}_{J_1}, \ldots, \mathbf{X}_{J_d}] \in \mathbb{R}^{d\times d}$, Alice creates some quantization $\hat{\mathbf{X}}_{\mathbf{J}}$ of $\mathbf{X_J}$, as further discussed below. Writing

---

[2]The communication cost can be slightly improved if Alice first seeks $\mathbf{X}_i$ that lies in the union of all sets, then $\mathbf{X}_i$ that lies in the union of the remaining sets, and so on. The difference is negligible for small $\Pr(A_1), \ldots, \Pr(A_d)$.

$Y_{\mathbf{J}} = [Y_{J_1}, \ldots, Y_{J_d}] \in \mathbb{R}^d$ for the corresponding sample vector on Bob's side, we consider the estimator

$$\hat{\boldsymbol{\rho}} = Y_{\mathbf{J}} \hat{\mathbf{X}}_{\mathbf{J}}^{-1} \Sigma_{\mathbf{X}} \qquad (46)$$

Note that in order to compute this estimator, Bob needs to know Alice's covariance matrix $\Sigma_{\mathbf{X}}$. Recall however that we have assumed without loss of generality that this is in fact a correlation matrix, hence all its entries have absolute value at most 1. Using a uniform quantizer of $[-1, 1]$ with (say) $\sqrt{k}$ bits, each entry of this matrix can be described to Bob with a resolution of roughly $2^{-\sqrt{k}}$, using only $d^2 \sqrt{k}$ bits overall. It is simple to check that this results in a negligible cost both in communication and in the mean squared error, and hence we disregard this issue below.

The general task is the following. Given a specified average number of bits $k$, find some quantization scheme $\mathbf{X}_{\mathbf{J}} \to \hat{\mathbf{X}}_{\mathbf{J}}$ using $k_q$ bits per entry, and sets $A_1, \ldots, A_d \in \mathbb{R}^d$, that

$$\text{minimize } \mathbb{E} \|Y_{\mathbf{J}} \hat{\mathbf{X}}_{\mathbf{J}}^{-1} \Sigma_{\mathbf{X}} - \boldsymbol{\rho}\|^2 \qquad (47)$$

$$\text{subject to } \sum_{\ell=1}^{d} h_g(\Pr(\mathbf{X} \in A_\ell)) + d^2 \cdot k_q = k \qquad (48)$$

Since the model (44) is linear with $d$ parameters, it is clear that we need at least $d$ different samples in order to obtain an estimator with a vanishing mean squared error. Furthermore, since Alice is given some control over the choice of $\mathbf{X}$ via her ability to pick samples from a large random set, it makes sense to try and make the problem as "well-posed" as possible, e.g., by striving to make the matrix $\mathbf{X}_{\mathbf{J}}$ have the smallest possible condition number while satisfying the communication constraints, which essentially dictate the number of samples we can choose from. A reasonable choice is therefore to try and make $\mathbf{X}_{\mathbf{J}}$ as diagonal as possible, by waiting each time for one coordinate to be strong and the others weak.

To make the problem tractable we apply the rationale above to a whitened version of $\mathbf{X}$, which allows us to directly compute the stopping probability. Let

$$\mathbf{W} = \Sigma_{\mathbf{X}}^{-\frac{1}{2}} \mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d) \qquad (49)$$

be the whitened version of $\mathbf{X}$, and $\{\mathbf{W}_i = \Sigma_{\mathbf{X}}^{-\frac{1}{2}} \mathbf{X}_i\}_i$ the associated whitened samples. We define the stopping sets

$$A_\ell^w = \left\{ \mathbf{w} \in \mathbb{R}^d \ : |w_\ell| > a, \ |w_j| < b \ \forall \ j \neq \ell \right\}, \qquad (50)$$

and the corresponding time indices

$$J_\ell^w = \min\{i > J_{\ell-1}^w : \mathbf{W}_i \in A_\ell^w\} \qquad (51)$$

for $\ell \in [d]$, with $J_0^w = 0$ by definition.

Note that by construction, $\Pr(\mathbf{W} \in A_\ell^w) = 2Q(a)(1 - 2Q(b))^{d-1}$ for any $\ell$. Alice then creates the matrix

$$\mathbf{W_J} = [\mathbf{W}_{J_1^w}, \ldots, \mathbf{W}_{J_d^w}] \in \mathbb{R}^{d \times d}, \tag{52}$$

and transmits to Bob the indices $J_1^w, \ldots, J_d^w$ using

$$k_l = h_g \left( 2Q(a)(1 - 2Q(b))^{d-1} \right) \tag{53}$$

bits per index on average, and $\hat{\mathbf{W}}_\mathbf{J}$, which is a quantized version of $\mathbf{W_J}$.

Note that in this method, in contrast to the ones considered thus far, Alice transmits to Bob some information regarding the actual values of her observations, rather than their locations alone. The reason is that the variance of the off-diagonal entries of $\mathbf{W_J}$ do not vanish as $k$ gets large. Nevertheless, we will show that a very simple quantizer using only a negligible number of bits is enough to represent $\mathbf{W_J}$ with sufficient accuracy for our purposes. Precisely, Alice quantizes $\mathbf{W_J}$ using exactly $k_q$ bits per entry, as follows: The diagonal entries are truncated to a maximal absolute value of $c = \sqrt{3}a$, and the double segment $[-c, -a] \cup [a, c]$ is uniformly quantized into $2^{k_q}$ levels. Off-diagonal entries, that all lie in the segment $[-b, b]$, are uniformly quantized into $2^{k_q}$ levels.

Given a communication constraint of $k$ bits on average, we need to choose $k_l, k_q$ that satisfy

$$d \cdot k_l + d^2 \cdot k_q = k, \tag{54}$$

and thresholds $a, b$ that satisfy (53). Furthermore, for reasons explained in the proof of Theorem 4, we need both $a^2$ and $(a - b)^2$ to increase with $k_l$, and $k_q, k_l$ to satisfy $k_l = k(1/d - o(1))$ and $k_l 2^{-k_q} \to 0$. One such choice is

$$k_l = \frac{1}{d} \left( \sqrt{k+1} - 1 \right)^2, \quad k_q = \sqrt{\frac{4k_l}{d^3}}, \tag{55}$$

and

$$a = Q^{-1} \left( \frac{h_g^{-1}(k_l)}{2(1 - 2Q(b_0))^{d-1}} \right), \quad b = b_0, \tag{56}$$

for some small fixed $b_0$.

After receiving $J_1^w, \ldots, J_d^w$ and $\hat{\mathbf{W}}_\mathbf{J}$, Bob creates the vector

$$Y_\mathbf{J} = [Y_{J_1^w}, \ldots, Y_{J_d^w}] \tag{57}$$

and performs estimation. The model (44) can be written as

$$Y = \boldsymbol{\rho} \Sigma_\mathbf{X}^{-\frac{1}{2}} \mathbf{W} + \sigma Z, \tag{58}$$

and thus the estimator is

$$\hat{\boldsymbol{\rho}} = Y_\mathbf{J} \hat{\mathbf{W}}_\mathbf{J}^{-1} \Sigma_\mathbf{X}^{\frac{1}{2}}. \tag{59}$$

14

**Theorem 4.** *The estimator $\hat{\boldsymbol{\rho}}$ in (59) satisfies*

$$\mathbb{E} \|\hat{\boldsymbol{\rho}} - \boldsymbol{\rho}\|^2 \leq \frac{1}{k} \left( \frac{d^2}{2\ln 2} \min_{\ell \in [d]} \{1 - \rho_\ell^2\} + o(1) \right), \tag{60}$$

*where $k$ is the expected number of transmitted bits. Furthermore, $\hat{\boldsymbol{\rho}}$ is asymptotically efficient given $(\mathbf{W_J}, Y_{\mathbf{J}})$.*

We prove this theorem in the next subsection.

**Corollary 1.** $\hat{\boldsymbol{\rho}}$ *in (59) dominates the scalar estimator.*

*Proof.* Allocating $k/d$ bits per correlation and using the scalar estimator (max or threshold), results in a sum of variances

$$\frac{1}{k} \left( \frac{d^2}{2\ln 2} \frac{1}{d} \sum_{\ell \in [d]} (1 - \rho_\ell^2) + o(1) \right) \tag{61}$$

which is greater than (60) for all values of $\rho_1, \dots, \rho_d$ (except when they are all equal). One could also use a nonuniform bit allocation for the scalar estimation, in which case the average in (61) would be replaced by a weighted average, which also is always greater than the minimum. $\qquad\square$

**Remark 1.** Theorem 4 implies in particular that when (say) $|\rho_1| = 1$, then the variance of our estimator decays faster than $\Omega(1/k)$. This is intuitively reasonable, since in this case $Y$ is equal to $\pm X_1$, hence $\Sigma_{\mathbf{X}}$ itself provides all the information about $\boldsymbol{\rho}$, which can be locally computed by Alice and communicated to Bob with variance of $2^{-\Omega(k)}$. Note however that Alice *cannot know* that $|\rho_1| = 1$, and neither can Bob (though he may have good reason to suspect so), hence it is still a bit surprising that our estimator allows this situation to nevertheless be exploited.

**Remark 2.** It is interesting to compare the performance of the estimator discussed in this subsection, to the performance of the estimator in the other extremal setup of Subsection 3.1, where $X$ is a scalar and $\mathbf{Y}$ is a vector. While both dominate the naive scheme of applying the scalar method $d$ times, neither dominates the other. The difference between them, essentially, is the difference between $\sum(1 - \rho_\ell^2)/d$ and $d \min\{1 - \rho_\ell^2\}$. For example, the former outperforms the latter if all correlations are equal, whereas the latter outperforms the former if any of the correlations is $\pm 1$.

## 3.3 Proof of Theorem 4

Consider the estimator

$$\hat{\boldsymbol{\rho}}_0 = Y_{\mathbf{J}} \mathbf{W_J}^{-1} \Sigma_{\mathbf{X}}^{\frac{1}{2}}. \tag{62}$$

Note that this estimator uses the non-quantized $\mathbf{W_J}$ which cannot be described to Bob with a finite number of bits, and hence is unrealizable. Nevertheless, as the following lemma shows, the loss incurred by employing $\hat{\boldsymbol{\rho}}$ instead, which uses the quantized $\mathbf{W_J}$, is small.

**Lemma 1.** *For any $a, b$ such that $a > d(b + 1)$,*

$$\mathbb{E}\,\|\hat{\boldsymbol{\rho}} - \boldsymbol{\rho}\,\|^2 \le \mathbb{E}\,\|\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho}\,\|^2 + (2d)^6 \left(e^{-\frac{a^2}{2}} + 2^{-k_q}\right) \tag{63}$$

*where $k_q$ bits are used to represent each entry in $\hat{\mathbf{W}}_{\mathbf{J}}$.*

*Proof.* See Appendix A.3. $\qquad\square$

The estimator (62) can be written as

$$\hat{\boldsymbol{\rho}}_0 = \boldsymbol{\rho} + \sigma Z_{\mathbf{J}} \mathbf{W}_{\mathbf{J}}^{-1} \Sigma_{\mathbf{X}}^{\frac{1}{2}}. \tag{64}$$

where $Z_{\mathbf{J}} = [Z_{J_1^w}, \ldots, Z_{J_d^w}] \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ is independent of $\mathbf{W}_{\mathbf{J}}$. It follows that $\hat{\boldsymbol{\rho}}_0$ is *unbiased* with

$$\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_0 = \sigma^2 \Sigma_{\mathbf{X}}^{\frac{1}{2}}\, \mathbb{E}\left((\mathbf{W}_{\mathbf{J}} \mathbf{W}_{\mathbf{J}}^T)^{-1}\right) \Sigma_{\mathbf{X}}^{\frac{1}{2}}. \tag{65}$$

In view of Lemma 1, it is sufficient to analyze the performance of the unrealizable estimator $\hat{\boldsymbol{\rho}}_0$. For the purpose of analyzing $\hat{\boldsymbol{\rho}}_0$ only, we can assume that Bob is given the value of $\mathbf{W}_{\mathbf{J}}$ for free, and the only cost is in transmitting the indices $J_1^w, \ldots, J_d^w$. Given $k_l$ bits for the representation of each of the locations, our general goal is to find $a, b$ that

$$\text{minimize } \operatorname{tr} \mathsf{Cov}\,\hat{\boldsymbol{\rho}}_0 \tag{66}$$

$$\text{subject to } h_g(2Q(a)(1 - 2Q(b))^{d-1}) = k_l. \tag{67}$$

Before proceeding to the analysis of $\hat{\boldsymbol{\rho}}_0$, we need the following two technical lemmas.

**Lemma 2.** *Let $\mathbf{M}$ be a square random matrix with independent entries, where the diagonal entries are i.i.d. with one distribution, and the off-diagonal entries are i.i.d. with another, symmetric distribution. Then $\mathbb{E}\,\mathbf{M}$, $\mathbb{E}\,\mathbf{M}\mathbf{M}^T$ and $\mathbb{E}\left((\mathbf{M}\mathbf{M}^T)^{-1}\right)$ are scalar multiples of the identity matrix.*

*Proof.* The claim for $\mathbb{E}\,\mathbf{M}$ and $\mathbb{E}\,\mathbf{M}\mathbf{M}^T$ is trivial. For $\mathbb{E}\left((\mathbf{M}\mathbf{M}^T)^{-1}\right)$ see Appendix A.4. $\qquad\square$

**Lemma 3** (Johnson [27])**.** *For any $n$-by-$m$ matrix $B = (b_{ij})$, $n \le m$, the smallest singular value is bounded below by*

$$\min_{i \in [n]} \left\{ |b_{ii}| - \frac{1}{2} \left( \sum_{j \in [n] \setminus i} |b_{ij}| + \sum_{j \in [n] \setminus i} |b_{ji}| \right) \right\} \tag{68}$$

The following lemma provides a simplified expression and bounds for $\mathbb{E}\left((\mathbf{W}_{\mathbf{J}} \mathbf{W}_{\mathbf{J}}^T)^{-1}\right)$, that will aid in proving Proposition 1 below.

**Lemma 4.** *The following claims hold for*

$$\alpha = d^{-1} \operatorname{tr} \mathbb{E}\,\mathbf{W}_{\mathbf{J}} \mathbf{W}_{\mathbf{J}}^T, \quad \beta = d^{-1} \operatorname{tr} \mathbb{E}\left((\mathbf{W}_{\mathbf{J}} \mathbf{W}_{\mathbf{J}}^T)^{-1}\right). \tag{69}$$

16

*(i)* $\mathbb{E}\,\mathbf{W_J}\mathbf{W_J}^T = \alpha\mathbf{I}_d$

*(ii)* $\mathbb{E}\left((\mathbf{W_J}\mathbf{W_J}^T)^{-1}\right) = \beta\mathbf{I}_d$

*(iii)* For any $a, b$ such that $a > (d-1)b$,

$$(a^2 + d + 1)^{-1} \leq \alpha^{-1} \leq \beta \leq (a - (d-1)b)^{-2}. \tag{70}$$

*Proof.* Recall that the vectors $\mathbf{W}_i$ are i.i.d. across the time index $i$, and that the entries of each one are i.i.d. with a symmetric distribution. Taking into account the rectangular structure of the stopping sets $A_\ell^w$ we see that $\mathbf{W_J}$ has independent entries where diagonal elements have one distribution, and off-diagonal elements follow another, symmetric distribution. Thus, the matrix $\mathbf{W_J}$ satisfies the conditions of Lemma 2. This proves claim (i) (which also holds trivially by construction) and claim (ii).

We proceed to prove claim (iii). Denoting the singular values of $\mathbf{W_J}$ by $\sqrt{\lambda_1} \geq \ldots \geq \sqrt{\lambda_d}$, we have that

$$\operatorname{tr}(\mathbf{W_J}\mathbf{W_J}^T)^{-1} = \sum_{\ell \in [d]} \lambda_\ell^{-1} \leq d\lambda_d^{-1}. \tag{71}$$

By construction, the diagonal entries of $\mathbf{W_J}$ are larger than $a$ in absolute value, and the off-diagonal entries are smaller than $b$ in absolute value. Therefore Lemma 3 yields

$$\sqrt{\lambda_d} \geq a - (d-1)b. \tag{72}$$

We thus have

$$\beta d = \operatorname{tr} \beta \mathbf{I}_d = \operatorname{tr} \mathbb{E}(\mathbf{W_J}\mathbf{W_J}^T)^{-1} \leq d(a - (d-1)b)^{-2}, \tag{73}$$

which establishes the rightmost inequality in claim (iii). The middle inequality holds since

$$\beta d = \sum_{\ell \in [d]} \mathbb{E}\,\frac{1}{\lambda_\ell} \geq \sum_{\ell \in [d]} \frac{1}{\mathbb{E}\,\lambda_\ell} \geq \frac{d^2}{\sum_{\ell \in [d]} \mathbb{E}\,\lambda_\ell} \tag{74}$$

$$= \frac{d^2}{\operatorname{tr} \mathbb{E}\,\mathbf{W_J}\mathbf{W_J}^T} = \frac{d^2}{d\alpha}, \tag{75}$$

where the two inequalities follow from Jensen's inequality applied to the function $1/x$. Note that the rows and columns of $\mathbf{W_J}$ have the same distribution. Therefore

$$\alpha = \mathbb{E}\,\|\mathbf{W}_{I_1}\|^2 \tag{76}$$

$$= \mathbb{E}\left((\mathbf{W})_1^2 \big| |(\mathbf{W})_1| > a\right) + (d-1)\,\mathbb{E}\left((\mathbf{W})_2^2 \big| |(\mathbf{W})_2| < b\right) \tag{77}$$

$$= 1 + as(a) + (d-1)\,\mathbb{E}\left((\mathbf{W})_2^2 \big| |(\mathbf{W})_2| < b\right) \tag{78}$$

$$\leq 1 + a(a + a^{-1}) + (d-1)1 \tag{79}$$

$$= a^2 + 1 + d \tag{80}$$

which completes the proof. $\qquad\square$

Lemma 4 and (65) implies that the optimization problem (66)-(67) can be written as

$$\text{minimize } \beta \tag{81}$$

$$\text{subject to } h_g(2Q(a)(1 - 2Q(b))^{d-1}) = k_l. \tag{82}$$

Note that both $Q(\cdot)$ and $h_g(\cdot)$ are monotonically decreasing. Therefore from (82) it is clear that increasing $k_l$ means increasing $a$ and/or decreasing $b$. From (70) we get that $\beta$ decreases as $a$ increases *and* gets farther away from $b$. We conclude therefore that a reasonable approximation to the solution of the optimization problem above, for large $k_l$, ia as given in (56). Note that the proposed approximated solution satisfies the constraint exactly.

**Proposition 1.** *The estimator* (62) *is unbiased and, for the choice of $a, b$ given in* (56), *it satisfies*

$$\text{tr Cov } \hat{\boldsymbol{\rho}}_0 \leq \frac{1}{k_l} \left( \frac{d}{2 \ln 2} \min_{\ell \in [d]} \{1 - \rho_\ell^2\} + o(1) \right) \tag{83}$$

*where $k_l$ is the expected number of bits used to describe each of the locations $J_1^w, \ldots, J_d^w$. Furthermore, $\hat{\boldsymbol{\rho}}_0$ is asymptotically efficient given $(\mathbf{W_J}, Y_\mathbf{J})$.*

*Proof.* In light of Lemma 4, (65) can be written as

$$\text{Cov } \hat{\boldsymbol{\rho}}_0 = \beta \sigma^2 \Sigma_\mathbf{X}. \tag{84}$$

Using (10)-(11) with $\mu = \boldsymbol{\rho} \Sigma_\mathbf{X}^{-\frac{1}{2}} \mathbf{W_J}$ and $\Sigma = \sigma^2 \mathbf{I}_d$, we get that the Fisher information matrix of $(\mathbf{W_J}, Y_\mathbf{J})$ is

$$\mathrm{I}_{\mathbf{W_J} Y_\mathbf{J}} = \frac{1}{\sigma^2} \Sigma_\mathbf{X}^{-\frac{1}{2}} \mathbb{E} \, \mathbf{W_J} \mathbf{W_J}^T \Sigma_\mathbf{X}^{-\frac{1}{2}} + \frac{2d}{\sigma^4} \Sigma_\mathbf{X}^{-1} \boldsymbol{\rho}^T \boldsymbol{\rho} \Sigma_\mathbf{X}^{-1} \tag{85}$$

$$= \frac{\alpha}{\sigma^2} \Sigma_\mathbf{X}^{-1} + \frac{2d}{\sigma^4} \Sigma_\mathbf{X}^{-1} \boldsymbol{\rho}^T \boldsymbol{\rho} \Sigma_\mathbf{X}^{-1}, \tag{86}$$

and using the Sherman–Morrison formula (e.g. [28]) we get

$$\mathrm{I}_{\mathbf{W_J} Y_\mathbf{J}}^{-1} = \frac{\sigma^2}{\alpha} \left( \Sigma_\mathbf{X} - \frac{2d}{\alpha \sigma^2 + 2d(1 - \sigma^2)} \boldsymbol{\rho}^T \boldsymbol{\rho} \right). \tag{87}$$

We take $a, b$ of (56). Note that $b$ is fixed and that $a$ increases with $k_l$. From Lemma 4 we have that

$$\alpha^{-1} = a^{-2}(1 + o(1)), \quad \beta = a^{-2}(1 + o(1)), \tag{88}$$

which implies

$$\text{Cov } \hat{\boldsymbol{\rho}}_0 = \frac{\sigma^2}{a^2}(1 + o(1)) \Sigma_\mathbf{X} \tag{89}$$

$$\mathrm{I}_{\mathbf{W_J} Y_\mathbf{J}}^{-1} = \frac{\sigma^2}{a^2}(1 + o(1)) \Sigma_\mathbf{X} \tag{90}$$

18

and thus $\hat{\boldsymbol{\rho}}_0$ is asymptotically efficient. For large $k_l$ we have

$$k_l = h_g(Q(a)2(1 - 2Q(b))^{d-1}) \tag{91}$$
$$= -\log(Q(a)2(1 - 2Q(b))^{d-1})(1 + o(1)) \tag{92}$$
$$= -\log(Q(a))(1 + o(1)) \tag{93}$$
$$= \frac{a^2}{2\ln 2}(1 + o(1)) \tag{94}$$
$$\tag{95}$$

and thus

$$\text{tr}\,\text{Cov}\,\hat{\boldsymbol{\rho}}_0 = d\beta\sigma^2 = \frac{d\sigma^2}{a^2}(1 + o(1)) \tag{96}$$
$$= \frac{1}{k_l}\left(\frac{d\sigma^2}{2\ln 2} + o(1)\right). \tag{97}$$

It remains to show that

$$\sigma^2 \le \min\{1 - \rho_\ell^2\}. \tag{98}$$

Note that $\sigma^2 = \text{Var}(Y|\mathbf{X})$ is the MMSE of estimating $Y$ from $\mathbf{X}$ (see e.g. [23]). Therefore it is not greater than $1 - \rho_\ell^2 = \text{Var}(Y|(\mathbf{X})_\ell)$, which is the MMSE of estimating $Y$ from the $\ell$-th coordinate only. $\qquad\square$

Theorem 4 now follows from Lemma 1 and Proposition 1.

# 4 Non-Gaussian Families

In this section, we move beyond the Gaussian setup and consider the problem of distributed correlation estimation in more general families of distributions, based on our Gaussian constructions. For brevity of exposition, we limit our discussion to the scalar case; the results can be extended in an obvious way to the vector case. We note that in contrast to the Gaussian setting, the marginal distribution of $X$ or $Y$ in other families of distributions may depend on the correlation, in which case Alice or Bob could use their (unlimited) local measurements to improve their inference (and in some cases to even learn $\rho$ exactly without any communication). For example, if $X$ is uniformly distributed over the interval $[-\sqrt{3}, \sqrt{3}]$, and $Y = \rho X + \sqrt{1 - \rho^2}Z$ where $Z$ is uniformly distributed over the discrete set $\{-1, 1\}$, then it is clear that the distribution of $Y$, which can be determined with arbitrary accuracy by Bob, determines $\rho$ up to its sign, reducing our problem to a binary hypothesis testing one. Such scenarios render our method useless, or, at the very least, degenerate.

Our interest, therefore, is in families of distributions where the marginals reveal little or nothing about the correlation. Specifically, we say that a family $\mathscr{F}$ of distributions on $(X, Y)$ is *correlation-hiding* if each pair of marginals can be associated with an infinite number of possible correlations; namely, for any

two marginals $p_X$ and $p_Y$ that are possible for some member of $\mathscr{F}$, there exists a countably infinite set $\mathscr{F}' \subseteq \mathscr{F}$ of joint distributions with marginals $p_X$ and $p_Y$, and with correlation coefficients that are all distinct.

Below, we discuss two types of correlation-hiding families. The first is the family of all possible distributions (subject only to mild moment constraints), which is obviously correlation-hiding. We show that for this family, the Gaussian performance can be uniformly attained. The idea is very simple: we perform "Gaussianization" of the samples using the Central Limit Theorem (CLT), and then apply the Gaussian estimators; showing that this indeed works, however, is somewhat technically involved. The second type of families that we consider are ones where $p_X$ is known, and where $Y = \alpha X + Z$ for some unknown coefficient $\alpha$ and unknown independent noise $Z$. We show that such families are correlation-hiding, and that we can sometimes (depending on $p_X$) obtain a variance that decays much faster with $k$ than the Gaussian one.

## 4.1 Unknown Distributions

In this subsection, we consider the case where the joint distribution of $X$ and $Y$ is completely unknown, subject only to mild moment conditions. We show how the threshold method of Subsection 2.2 can be extended to this setup, using the CLT, to yield the same performance guarantees. The basic idea is to use the unlimited number of samples in order to create Gaussian r.v.s with the same correlation, by averaging over blocks of samples. Due to the CLT, it is intuitively clear that this approach works if Alice and Bob use infinite sized blocks. This is however impractical, and the main technical challenge is to show that using finite large enough blocks, i.e., changing the order of limits, still works.

Let $(X, Y)$ be drawn from the family

$$\mathscr{F} = \{p_{XY} : \mathbb{E}\, X^2, \mathbb{E}\, Y^2 < u, \mathbb{E}\, Y^4 < \infty\}. \tag{99}$$

where $u$ is some known constant. Again, since we assume that local measurements are essentially unlimited, and the second moments have known upper bounds, we can assume without loss of generality that $\mathbb{E}\, X = \mathbb{E}\, Y = 0$, and $\mathbb{E}\, X^2 = \mathbb{E}\, Y^2 = 1$. The following claim is immediate from the fact that $\mathscr{F}$ contains in particular the Gaussian distributions.

**Corollary 2.** *The family $\mathscr{F}$ in* (99) *is correlation-hiding.*

Let us now proceed to describe our estimator. Alice and Bob first locally sum over their measurements to create the new i.i.d. sequences $\{\bar{X}_i\}_i, \{\bar{Y}_i\}_i$, given by

$$\bar{X}_i = \frac{1}{\sqrt{m}} \sum_{j \in S_i} X_j, \quad \bar{Y}_i = \frac{1}{\sqrt{m}} \sum_{j \in S_i} Y_j \tag{100}$$

where the $S_i$'s are disjoint index sets of size $m$. For brevity, we suppress the dependence of these new r.v.s on $m$. The sequence of pairs $\{(\bar{X}_i, \bar{Y}_i)\}_i$ is clearly

i.i.d. Denoting by $(\bar{X}, \bar{Y})$ a generic pair in this sequence, the correlation between $\bar{X}$ and $\bar{Y}$ is clearly the same as the correlation between $X$ and $Y$. Alice and Bob can therefore apply the threshold method to the sequence $\{(\bar{X}_i, \bar{Y}_i)\}_i$ in order to estimate the original $\rho$. We now show that the performance of this estimator approaches the Gaussian performance as $m \to \infty$. Given a communication constraint of $k$ bits, the threshold $t$ is chosen (as in the Gaussian case) such that $h_g(Q(t)) = k$. We denote

$$\bar{J} = \min\{i : \bar{X}_i > t\} \tag{101}$$

and the estimator

$$\hat{\rho}_{\text{th}}^{(m)} = \frac{\bar{Y}_{\bar{J}}}{s(t)}, \tag{102}$$

where $s(t)$ is given in (6). Note we cannot normalize by $\mathbb{E}\,\bar{X}_{\bar{J}}$ to get a strictly unbiased estimator since we assume unknown distributions and thus $\mathbb{E}\,\bar{X}_{\bar{J}}$ is not known for finite $m$. The expected number of bits needed to describe $\bar{J}$ is

$$k^{(m)} = h_g(\Pr(\bar{X} > t)). \tag{103}$$

**Remark 3.** Note that practical scenarios would require the choice of some fixed $m$. Therefore, in cases where the support of $X$ is finite, we might get that $\Pr(\bar{X} > t) = 0$ which means Alice waits forever and the estimator is undefined. Therefore, while the distribution of $(X, Y)$ need not be known in general, such a practical scenario requires some knowledge regarding the support of $X$ in the form of a number $x$ such that $\Pr(X_i > x) > 0$ (which must exist since $\mathbb{E}\,X = 0$). Then we can take $m > t^2/x^2$ to assure $\Pr(\bar{X} > t) > 0$.

**Theorem 5.** *Let* $t = Q^{-1}(h_g^{-1}(k))$. *Then for the family* $\mathscr{F}$ *in* (99) *it holds that* $\lim_{m \to \infty} k^{(m)} = k$ *and*

$$\lim_{m \to \infty} \mathbb{E}(\hat{\rho}_{\text{th}}^{(m)} - \rho)^2 = \frac{1}{k}\left(\frac{1 - \rho^2}{2\ln 2} + o(1)\right) \tag{104}$$

*Proof.* Due to the CLT we have for any fixed $t > 0$ that

$$\lim_{m \to \infty} \Pr(\bar{X} > t) = Q(t) \tag{105}$$

and thus, since $h_g$ is smooth, the communication constraint is asymptotically satisfied. We have

$$\mathbb{E}(\hat{\rho}_{\text{th}}^{(m)} - \rho)^2 = \frac{\mathbb{E}\,\bar{Y}_{\bar{J}}^2}{s^2(t)} - 2\rho\frac{\mathbb{E}\,\bar{Y}_{\bar{J}}}{s(t)} + \rho^2 \tag{106}$$

and thus it suffices to show that the first two moments of $\bar{Y}_{\bar{J}}$ converge to their values under the Gaussian distribution. Denoting by $(X^{\mathcal{N}}, Y^{\mathcal{N}})$ and $Y_{\bar{J}}^{\mathcal{N}}$ the associated r.v.s under a Gaussian distribution, it is enough to show that $\bar{Y}_{\bar{J}}$

21

converges in distribution to $Y_{\bar{J}}^{\mathcal{N}}$ as $m \to \infty$, and that $\bar{Y}_{\bar{J}}^2$ is uniformly integrable [29]. To show convergence in distribution, observe that

$$\lim_{m \to \infty} \Pr(\bar{Y}_{\bar{J}} > y) = \lim_{m \to \infty} \Pr(\bar{Y} > y | \bar{X} > t) \tag{107}$$

$$= \lim_{m \to \infty} \frac{\Pr(\bar{Y} > y, \bar{X} > t)}{\Pr(\bar{X} > t)} \tag{108}$$

$$= \frac{\lim_{m \to \infty} \Pr(\bar{Y} > y, \bar{X} > t)}{\lim_{m \to \infty} \Pr(\bar{X} > t)} \tag{109}$$

$$= \frac{\Pr(Y^{\mathcal{N}} > y, X^{\mathcal{N}} > t)}{\Pr(X^{\mathcal{N}} > t)} \tag{110}$$

$$= \Pr(Y^{\mathcal{N}} > y | X^{\mathcal{N}} > t) \tag{111}$$

$$= \Pr(Y_{\bar{J}}^{\mathcal{N}} > y) \tag{112}$$

where (109) holds since the denominator is not zero, and (110) holds by virtue of the CLT.

It follows from (105) that there exist some $m_0$ and $c > 0$ (e.g., $c = Q(t)/2$) such that

$$\Pr(\bar{X} > t) > c \quad \forall\, m \geq m_0, \tag{113}$$

and therefore we assume without loss of generality that $m \geq m_0$. To prove uniform integrability of $\bar{Y}_{\bar{J}}^2$ it suffices to show that $\sup_m \mathbb{E}\,|\bar{Y}_{\bar{J}}|^\gamma < \infty$ for some $\gamma > 2$ [29]. For simplicity, we set $\gamma = 4$:

$$\mathbb{E}\,|\bar{Y}_{\bar{J}}|^4 = \mathbb{E}(|\bar{Y}|^4 \mid \bar{X} > t) \tag{114}$$

$$\leq \frac{\mathbb{E}\,|\bar{Y}|^4}{\Pr(\bar{X} > t)} \tag{115}$$

$$= \frac{\mathbb{E}(\frac{1}{\sqrt{m}} \sum_j Y_j)^4}{\Pr(\bar{X} > t)} \tag{116}$$

$$= \frac{\frac{1}{m}\mathbb{E}\,Y^4 + 3\frac{m-1}{m}(\mathbb{E}\,Y^2)^2}{\Pr(\bar{X} > t)} \tag{117}$$

$$< \frac{1}{c}\left(\frac{\mathbb{E}\,Y^4}{m} + 3\right), \tag{118}$$

which is finite since $\mathbb{E}\,Y^4 < \infty$. $\qquad\square$

**Example 1** (Doubly symmetric binary r.v.s)**.** Consider the family of distributions where $X \sim \text{Bernoulli}(1/2)$ and $Y = X \oplus Z$ where $Z \sim \text{Bernoulli}(p)$ is independent of $X$, $p \in [0,1]$ is unknown, and $\oplus$ is the binary XOR operation. The associated Gaussian version of these r.v.s (after removing the mean) are the jointly normal, zero mean unit norm r.v.s $\bar{X}$ and $\bar{Y}$, with correlation $\rho = 1 - 2p$. Our unbiased estimator can therefore obtain a variance of $\frac{1-(1-2p)^2}{2k\ln 2} = \frac{2p(1-p)}{k\ln 2}$ for the estimation of $\rho$, which corresponds to a variance of

$\frac{p(1-p)}{2k\ln 2}$ for the estimation of $p$. This can be juxtaposed with the straightforward approach of simply sending $X_1, \ldots, X_k$ to Bob and applying the (efficient) estimator $\hat{p} = \frac{1}{k}\sum_{j=1}^{k} X_j \oplus Y_j$. This unbiased estimator has a variance of $\frac{p(1-p)}{k}$, which is interestingly slightly worse than what we got using the Gaussian approach. It may be possible to improve the former by using lossy compression, but we do not explore this direction here.

**Remark 4.** Estimating the joint probability mass function of general discrete distributions on $X, Y$ can be similarly cast as a correlation estimation problem. However, the gain observed in the binary case above does not carry over to the general case. This is however not unexpected, since our estimator does not assume any bound on the cardinality of $X$ and $Y$.

## 4.2 Additive Noise Families

In this subsection, we consider a more restricted model where the distribution $p_X$ of $X$ is fixed (but not necessarily Gaussian) and has bounded variance, and where

$$Y = \alpha X + Z \tag{119}$$

for some unknown bounded constant $\alpha$, where $Z$ is an arbitrary r.v. with bounded variance that is independent[3] of $X$. Let us denote this family of distributions by $\mathscr{F}(p_X)$. First, we note:

**Corollary 3.** $\mathscr{F}(p_X)$ *is correlation-hiding for any* $p_X$.

*Proof.* See Appendix A.6. $\qquad\square$

We now show that the threshold estimator proposed for the Gaussian case applies to $\mathscr{F}(p_X)$ as well, and that its performance can be better or worse, depending on $p_X$. Specifically, we show that the $O(1/k)$ decay of the variance with the number of bits is not fundamental, as for some (heavier tailed) choices of $p_X$ we obtain a behavior of $O(1/k^2)$ using the same threshold estimator, and $2^{-\Omega(k)}$ using a slightly modified estimator. The latter is essentially the best possible using our approach, since we utilize $O(2^k)$ samples (with high probability), which corresponds to a variance of $\Omega(2^{-k})$ even in the centralized case.

As in the previous sections, Alice and Bob can normalize their measurements locally. Therefore, we can assume without loss of generality that (119) can be written as

$$Y = \rho X + \sqrt{1-\rho^2}Z \tag{120}$$

where $X$ and $Z$ are independent, zero mean unit variance r.v.s, and the correlation is $\rho = \mathbb{E}\,XY$. We assume that $p_Z$ is arbitrary and unknown, and that

---

[3]It is in fact sufficient for our purposes to assume only that $\mathbb{E}(Z|X)$ and $\mathsf{Var}(Z|X)$ do not depend on $X$

$p_X$ is arbitrary but known. Applying the threshold method of Subsection 2.2 to this non-Gaussian setup, we denote as usual $J = \min\{i : X_i > t\}$ the first index to pass the threshold $t$, where $t$ is chosen such that $h_g(\Pr(X > t)) = k$. Our estimator is

$$\hat{\rho}_{\text{th}} = \frac{Y_J}{\mathbb{E}\,X_J}. \tag{121}$$

The following claim is immediate.

**Corollary 4.** $\hat{\rho}_{\text{th}}$ is unbiased, and

$$\text{Var}\,\hat{\rho}_{\text{th}} = \frac{\rho^2\,\text{Var}(X \mid X > t) + 1 - \rho^2}{(\mathbb{E}(X \mid X > t))^2}. \tag{122}$$

Let us compute (122) for some specific choices of $p_X$.

**Example 2** (Laplace Distribution). Let $p_X$ be a zero-mean, unit-variance Laplace distribution, hence $\Pr(X > x) = \frac{1}{2}e^{-\sqrt{2}x}$ for $x > 0$. Thus,

$$\mathbb{E}(X \mid X > t) = t + \frac{1}{\sqrt{2}}, \quad \text{Var}(X \mid X > t) = \frac{1}{2}, \tag{123}$$

and

$$k = h_g\left(\frac{1}{2}e^{-\sqrt{2}t}\right) \tag{124}$$

$$= -\log\left(\frac{1}{2}e^{-\sqrt{2}t}\right)(1 + o(1)) \tag{125}$$

$$= \sqrt{2}t\log(e)(1 + o(1)). \tag{126}$$

Therefore (122) becomes

$$\text{Var}\,\hat{\rho}_{\text{th}} = \frac{1}{k^2}\left(\frac{2 - \rho^2}{(\ln 2)^2} + o(1)\right), \tag{127}$$

which yields a variance of $O(1/k^2)$, in contrast to the slower $O(1/k)$ attained in the Gaussian case.

**Example 3** (Pareto Distribution). Motivated by the Laplace example which indicates that a heavier tail of $X$ may yield a faster decay of $\text{Var}\,\hat{\rho}_{\text{th}}$, we investigate the heaviest tail possible with finite variance. Suppose that $p_X$ is the (double-sided, zero mean) Pareto distribution, i.e.,

$$\Pr(X > x) = \Pr(X < -x) = \frac{1}{2}\left(\frac{x_0}{x}\right)^\alpha \tag{128}$$

for any $x > x_0$, where $\alpha > 2$ and $x_0 > 0$ is set such that $\text{Var}\,X = 1$. Then for any $t > x_0$

$$\mathbb{E}(X \mid X > t) = \frac{\alpha t}{\alpha - 1}, \quad \text{Var}(X \mid X > t) = \frac{\alpha t^2}{(\alpha - 1)^2(\alpha - 2)} \tag{129}$$

24

and (122) becomes

$$\text{Var}\,\hat{\rho}_{\text{th}} = \frac{\rho^2}{\alpha(\alpha - 2)} + O(1/t^2). \tag{130}$$

Thus, the variance of our threshold estimator does not vanish with the number of bits. This flaw can nevertheless be fixed in a very strong way, as we show next. Before we proceed, we note that for $p_X$ with a tail of the form $\Pr(X > x) \propto e^{-x^{\frac{1}{m}}}$, i.e. in between Pareto and Laplace, the threshold estimator yields $\text{Var}\,\hat{\rho}_{\text{th}} = O(1/k^2)$ for any natural $m$. Also, tails that decay faster than Gaussian may yield worse performance, e.g. the tail $e^{-x^4}$ yields $\text{Var}\,\hat{\rho}_{\text{th}} = O(1/\sqrt{k})$.

Getting back to the double-sided Pareto distribution, recall that in the Gaussian case it was shown that describing the value of $X_J$ does not improve estimation performance. This was due to the fact that for the Gaussian family, $\text{Var}\,X_J = \text{Var}(X \mid X > t) \to 0$. This is however not true in general; in fact, the Pareto distribution is an extreme case in which $\text{Var}(X \mid X > t) \to \infty$. Therefore, providing some information regarding the value of $X_J$ at the expense of the number of bits used to describe the index $J$, might improve performance. With that in mind, we consider the estimator

$$\hat{\rho}_{\text{th-q}} = \frac{Y_J}{\hat{X}_J} \tag{131}$$

that allocates $k_l$ bits to describe $J$, and $k_q$ bits to describe the value of $\hat{X}_J$, where $k_l + k_q = k$. We apply the following simple quantizer. For some $u > t$ we divide the region $[t, u]$ to $2^{k_q}$ equal segments of length $\Delta = 2^{-k_q}(u - t)$. For $x > u$ we set $\hat{x} = u$. In the following, we show that this estimator attains a variance that decays exponentially fast with $k$.

**Proposition 2.** *Consider the family $\mathscr{F}(p_X)$ where $p_X$ be the double-sided Pareto distribution. Then the estimator $\hat{\rho}_{\text{th-q}}$ in (131) satisfies*

$$\mathbb{E}(\hat{\rho}_{\text{th-q}} - \rho)^2 \leq (1 + \rho^2) \cdot 2^{-\frac{2}{\alpha}\frac{\alpha-2}{\alpha-1}k(1-o(1))}, \tag{132}$$

*where $k$ is the average number of transmitted bits.*

*Proof.* See Appendix A.5. $\qquad\square$

## 5 Conclusions

We have discussed the problem of estimating the correlations between remotely observed random vectors with unlimited local samples, under one-way communication constraints. For the case where the vectors are jointly Gaussian, we provided simple constructive unbiased estimators for the correlations; our estimators attain the best known non-constructive Zhang-Berger upper bound on

the variance in the scalar case, and use the local correlations to uniformly improve performance in the vector case, where the Zhang-Berger approach seems inapplicable. Loosely speaking, our approach is based on Alice scanning her the local observations and sending the index of suitably "large" samples that induce good signal-to-noise ratio for the estimation for Bob, who uses the corresponding samples on his end. We then showed that using the CLT, this approach can be applied to the case of estimating correlations for completely unknown distributions, with the exact same variance guarantees. While the Gaussian approach yields a variance that is inversely proportional to the expected number of transmitted bits, we show that for joint distributions generated via unknown fading channels with unknown additive noise, whose correlations cannot be estimated locally, a slightly modified estimator attains a variance decaying *exponentially fast* with the expected number of transmitted bits. It remains interesting to try and obtain lower bounds on the variance as a function of the number of bits and the richness of the family of distributions under consideration. We conjecture that the inversely proportional behavior of our Gaussian estimator is order-wise optimal in the Gaussian case, hence also for the case of unknown distributions.

# A  Appendix

## A.1  Proof of Theorem 3

The model can be written as (see e.g. [26])

$$\mathbf{Y} = \boldsymbol{\rho}\, X + \Sigma^{\frac{1}{2}} \mathbf{Z} \tag{133}$$

where $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ is independent of $X$, and $\Sigma = \Sigma_{\mathbf{Y}} - \boldsymbol{\rho}\,\boldsymbol{\rho}^T$. We have $\hat{\boldsymbol{\rho}} = (\boldsymbol{\rho}\, X_J + \Sigma^{\frac{1}{2}} \mathbf{Z}_J)/\mathbb{E}\, X_J$. Therefore $\mathbb{E}\,\hat{\boldsymbol{\rho}} = \boldsymbol{\rho}$ and

$$\mathsf{Cov}\,\hat{\boldsymbol{\rho}} = \frac{1}{(\mathbb{E}\, X_J)^2}\left(\Sigma + \mathsf{Var}(X_J)\,\boldsymbol{\rho}\,\boldsymbol{\rho}^T\right). \tag{134}$$

Using (10)-(11) with $\mu = \boldsymbol{\rho}\, X_J, \Sigma = \Sigma_{\mathbf{Y}} - \boldsymbol{\rho}\,\boldsymbol{\rho}^T$ we get that the Fisher information matrix pertaining to $(X_J, \mathbf{Y}_J)$ is

$$\mathrm{I}_{X_J\mathbf{Y}_J} = \Sigma^{-1}(\mathbb{E}\, X_J^2 + \boldsymbol{\rho}^T \Sigma^{-1} \boldsymbol{\rho}) + \Sigma^{-1}\,\boldsymbol{\rho}\,\boldsymbol{\rho}^T \Sigma^{-1}, \tag{135}$$

and applying the Sherman–Morrison formula (e.g. [28]) yields

$$\mathrm{I}_{X_J\mathbf{Y}_J}^{-1} = \frac{1}{EX_J^2 + \boldsymbol{\rho}^T \Sigma^{-1} \boldsymbol{\rho}}\left(\Sigma - \frac{\boldsymbol{\rho}\,\boldsymbol{\rho}^T}{EX_J^2 + 2\,\boldsymbol{\rho}^T \Sigma^{-1} \boldsymbol{\rho}}\right). \tag{136}$$

Using the arguments of Theorem 2 yields that both $\mathsf{Cov}\,\hat{\boldsymbol{\rho}}$ and $\mathrm{I}_{X_J\mathbf{Y}_J}^{-1}$ are $\frac{1}{t^2}(1 + o(1))\Sigma$ and thus $\hat{\boldsymbol{\rho}}$ is asymptotically efficient. Theorem 2 also implies that $k = t^2(2\ln 2 + o(1))$, and noting that $\mathrm{tr}\,\Sigma = \mathrm{tr}(\Sigma_{\mathbf{Y}} - \boldsymbol{\rho}\,\boldsymbol{\rho}^T) = d - \|\boldsymbol{\rho}\|^2$ concludes the proof.

## A.2 The scalar method with linear transformations

In this subsection we show that any method based on $d$ scalar transmissions cannot uniformly beat the scheme of applying the basic scalar method $d$ times. In this sense, the joint method proposed in Theorem 4 *is superior* because it *does* uniformly beat the simple scalar scheme. Specifically, let $M$ be some invertible $d \times d$ matrix known to both Alice and Bob, and let $\widetilde{\mathbf{X}} = M\mathbf{X}$. Suppose Alice and Bob apply the scalar method separately to obtain an estimator $\hat{\boldsymbol{\rho}}_M$ for the correlation vector $\boldsymbol{\rho}_M = \mathbb{E}\, Y\widetilde{\mathbf{X}}^T$, and then use the estimator $M^{-1}\hat{\boldsymbol{\rho}}_M$ to estimate $\boldsymbol{\rho}$. As it turns out, this family of estimators does not dominate the naive approach of estimating each correlation separately (i.e., $M = \mathbf{I}_d$).

**Proposition 3.** *For any two invertible $d \times d$ matrices $M_1, M_2$ (that can arbitrarily depend on $\Sigma_{\mathbf{X}}$, and are known to both Alice and Bob), $M_1^{-1}\hat{\boldsymbol{\rho}}_{M_1}$ does not dominate $M_2^{-1}\hat{\boldsymbol{\rho}}_{M_2}$.*

*Proof.* We need to show that any linear transformation applied to $\mathbf{X}$, followed by the scalar method, cannot be uniformly better than the scalar method itself. It suffices to show that for the two-dimensional case.

Alice creates the following two scalar sequences.

$$U_i = [a_1, b_1]\mathbf{X}_i, \text{ for } i = 1, \ldots, n_1 \text{ and} \tag{137}$$

$$V_i = [b_2, a_2]\mathbf{X}_i, \text{ for } i = n_1 + 1, \ldots, n_1 + n_2 \tag{138}$$

and allocates $k_1$ bits for $U$, and $k_2$ bits for $V$ (Note we can use either max or threshold method, and that $n_1, n_2$ can be arbitrarily large). One special case of the above is the "successive refinement" approach described in the introduction (for $b_1 = 0$), and another special case is the naive scalar method (for $a_1 = a_2 = 1$, $b_1 = b_2 = 0$ and $k_1 = k_2 = k/2$). Without loss of generality we assume $a_1, b_1$ are such that $\mathbb{E}\, U^2 = 1$, and $a_2, b_2$ are such that $\mathbb{E}\, V^2 = 1$. We denote

$$\alpha_1 = \mathbb{E}\, Y_iU_i = a_1\rho_1 + b_1\rho_2 \tag{139}$$

$$\alpha_2 = \mathbb{E}\, Y_iV_i = b_2\rho_1 + a_2\rho_2, \tag{140}$$

and $\boldsymbol{\alpha} = [\alpha_1, \alpha_2]^T$. We also denote $\boldsymbol{\rho} = [\rho_1, \rho_2]^T$ and

$$M = \begin{bmatrix} a_1 & b_1 \\ b_2 & a_2 \end{bmatrix}, \tag{141}$$

and therefore we have $\boldsymbol{\alpha} = M\,\boldsymbol{\rho}$. The best Bob can do (recall $U, V$ are independent) is to estimate $\alpha_1$ using $U$ and $\alpha_2$ using $V$ to obtain

$$\mathsf{Var}\,\hat{\alpha}_1 = \frac{1}{k_1}\left(\frac{1 - \alpha_1^2}{2\ln 2} + o(1)\right) \tag{142}$$

$$\mathsf{Var}\,\hat{\alpha}_2 = \frac{1}{k_2}\left(\frac{1 - \alpha_2^2}{2\ln 2} + o(1)\right) \tag{143}$$

and then take $\hat{\boldsymbol{\rho}}_{\mathrm{trn}} = M^{-1}\hat{\boldsymbol{\alpha}}$. The resulting sum of variances (note $\mathsf{Cov}(\hat{\alpha}_1, \hat{\alpha}_2) = 0$) is

$$\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{trn}} = \mathrm{tr}\,M^{-1}\,\mathsf{Cov}(\hat{\alpha})M^{-T} \tag{144}$$

$$= \mathrm{tr}\,M^{-1}\begin{bmatrix}\mathsf{Var}\,\hat{\alpha}_1 & 0 \\ 0 & \mathsf{Var}\,\hat{\alpha}_2\end{bmatrix}M^{-T} \tag{145}$$

$$= \frac{(a_2^2 + b_2^2)\,\mathsf{Var}\,\hat{\alpha}_1 + (a_1^2 + b_1^2)\,\mathsf{Var}\,\hat{\alpha}_2}{(a_1 a_2 - b_1 b_2)^2} \tag{146}$$

Applying the simple scalar method twice yields

$$\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{scl}} = \frac{1}{k}\left(\frac{k}{k_1'}\frac{1 - \rho_1^2}{2\ln 2} + \frac{k}{k_2'}\frac{1 - \rho_2^2}{2\ln 2} + o(1)\right). \tag{147}$$

with $k_1' + k_2' = k_1 + k_2 = k$. We want to show that $\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{trn}}$ cannot be uniformly better than $\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{scl}}$, namely, show that for any choice of $a_1, b_1, a_2, b_2, k_1, k_2$ (that do not depend on $\rho_1, \rho_2$) we can find $\rho_1, \rho_2$ such that $\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{scl}} < \mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{trn}}$. This is easy because we can always take $\rho_1, \rho_2 \in \{-1, 1\}$ (or arbitrarily close to $\pm 1$) which makes $\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{scl}} \approx 0$ and $\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{trn}} \neq 0$. If $\mathrm{tr}\,\mathsf{Cov}\,\hat{\boldsymbol{\rho}}_{\mathrm{trn}} = 0$ (i.e. $\alpha_1^2 = \alpha_2^2 = 1$), we can flip the sign of $\rho_2$ to obtain either $\alpha_1^2 \neq 1$ or $\alpha_2^2 \neq 1$. $\qquad\square$

## A.3   Proof of Lemma 1

Writing $\mathbf{W} = \mathbf{W_J}$ and $\hat{\mathbf{W}} = \hat{\mathbf{W}}_\mathbf{J}$, we have

$$\hat{\boldsymbol{\rho}}_0 = Y_\mathbf{J}\mathbf{W}^{-1}\Sigma_\mathbf{X}^{\frac{1}{2}} = \boldsymbol{\rho} + \sigma Z_\mathbf{J}\mathbf{W}^{-1}\Sigma_\mathbf{X}^{\frac{1}{2}} \tag{148}$$

$$\hat{\boldsymbol{\rho}} = Y_\mathbf{J}\hat{\mathbf{W}}^{-1}\Sigma_\mathbf{X}^{\frac{1}{2}} = \boldsymbol{\rho}\,\Sigma_\mathbf{X}^{-\frac{1}{2}}\mathbf{W}\hat{\mathbf{W}}^{-1}\Sigma_\mathbf{X}^{\frac{1}{2}} + \sigma Z_\mathbf{J}\hat{\mathbf{W}}^{-1}\Sigma_\mathbf{X}^{\frac{1}{2}} \tag{149}$$

$$\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho} = \sigma Z_\mathbf{J}\mathbf{W}^{-1}\Sigma_\mathbf{X}^{\frac{1}{2}} \tag{150}$$

$$\hat{\boldsymbol{\rho}} - \hat{\boldsymbol{\rho}}_0 = \boldsymbol{\rho}\,\Sigma_\mathbf{X}^{-\frac{1}{2}}(\mathbf{W}\hat{\mathbf{W}}^{-1} - I)\Sigma_\mathbf{X}^{\frac{1}{2}} + \sigma Z_\mathbf{J}(\hat{\mathbf{W}}^{-1} - \mathbf{W}^{-1})\Sigma_\mathbf{X}^{\frac{1}{2}}. \tag{151}$$

Recall $Z_\mathbf{J}$ is a row vector $\sim \mathcal{N}(0, \mathbf{I}_d)$ independent of $\mathbf{W}$. It follows that

$$\mathbb{E}\,\|\hat{\boldsymbol{\rho}} - \hat{\boldsymbol{\rho}}_0\|^2 = \mathbb{E}\,\|\,\boldsymbol{\rho}\,\Sigma_\mathbf{X}^{-\frac{1}{2}}(\mathbf{W}\hat{\mathbf{W}}^{-1} - I)\Sigma_\mathbf{X}^{\frac{1}{2}}\|^2 \tag{152}$$

$$+ \sigma^2\,\mathbb{E}\,\mathrm{tr}(\hat{\mathbf{W}}^{-1} - \mathbf{W}^{-1})\Sigma_\mathbf{X}(\hat{\mathbf{W}}^{-T} - \mathbf{W}^{-T}), \tag{153}$$

and

$$\mathbb{E}(\hat{\boldsymbol{\rho}} - \hat{\boldsymbol{\rho}}_0)(\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho})^T = \sigma^2\,\mathbb{E}\,\mathrm{tr}(\hat{\mathbf{W}}^{-1} - \mathbf{W}^{-1})\Sigma_\mathbf{X}\mathbf{W}^{-T}. \tag{154}$$

Therefore

$$\mathbb{E}\,\|\hat{\boldsymbol{\rho}} - \boldsymbol{\rho}\|^2 = \mathbb{E}\,\|(\hat{\boldsymbol{\rho}} - \hat{\boldsymbol{\rho}}_0) + (\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho})\|^2 \tag{155}$$

$$= \mathbb{E}\,\|\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho}\|^2 + \mathbb{E}\,\|\hat{\boldsymbol{\rho}} - \hat{\boldsymbol{\rho}}_0\|^2 + 2\,\mathbb{E}(\hat{\boldsymbol{\rho}} - \hat{\boldsymbol{\rho}}_0)(\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho})^T \tag{156}$$

$$= \mathbb{E}\,\|\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho}\|^2 + \mathbb{E}\,\|\,\boldsymbol{\rho}\,\Sigma_\mathbf{X}^{-\frac{1}{2}}(\mathbf{W}\hat{\mathbf{W}}^{-1} - I)\Sigma_\mathbf{X}^{\frac{1}{2}}\|^2 \tag{157}$$

28

$$+ \sigma^2 \, \mathbb{E} \, \text{tr} (\hat{\mathbf{W}}^{-1} - \mathbf{W}^{-1}) \Sigma_{\mathbf{X}} (\hat{\mathbf{W}}^{-T} + \mathbf{W}^{-T}), \tag{158}$$

and thus

$$\mathbb{E} \, \| \hat{\boldsymbol{\rho}} - \boldsymbol{\rho} \|^2 - \mathbb{E} \, \| \hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho} \|^2 \tag{159}$$

$$= \mathbb{E} \, \| \boldsymbol{\rho} \, \Sigma_{\mathbf{X}}^{-\frac{1}{2}} (\mathbf{W} - \hat{\mathbf{W}}) \hat{\mathbf{W}}^{-1} \Sigma_{\mathbf{X}}^{\frac{1}{2}} \|^2 \tag{160}$$

$$+ \sigma^2 \, \mathbb{E} \, \text{tr} (\mathbf{W} - \hat{\mathbf{W}}) \hat{\mathbf{W}}^{-1} \Sigma_{\mathbf{X}} (\hat{\mathbf{W}}^{-T} + \mathbf{W}^{-T}) \mathbf{W}^{-1}. \tag{161}$$

Let us upper bound the two terms separately. Recall that by (73) we have $\| \mathbf{W}^{-1} \|_F^2 \leq d / (a - (d-1)b)^2$, which also holds for $\hat{\mathbf{W}}^{-1}$. Furthermore, the assumption that $a > d(b+1)$ implies that $a - (d-1)b$ is lower bounded by either $a/d$ or $d$.

First term (160): The Frobenius norm is sub-multiplicative (see e.g. [30]), and therefore

$$\mathbb{E} \, \| \boldsymbol{\rho} \, \Sigma_{\mathbf{X}}^{-\frac{1}{2}} (\mathbf{W} - \hat{\mathbf{W}}) \hat{\mathbf{W}}^{-1} \Sigma_{\mathbf{X}}^{\frac{1}{2}} \|_F^2 \tag{162}$$

$$\leq \| \boldsymbol{\rho} \, \Sigma_{\mathbf{X}}^{-\frac{1}{2}} \|_F^2 \| \Sigma_{\mathbf{X}}^{\frac{1}{2}} \|_F^2 \, \mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2 \| \hat{\mathbf{W}}^{-1} \|_F^2 \tag{163}$$

$$\leq \frac{d \| \boldsymbol{\rho} \, \Sigma_{\mathbf{X}}^{-\frac{1}{2}} \|_F^2 \| \Sigma_{\mathbf{X}}^{\frac{1}{2}} \|_F^2}{(a - (d-1)b)^2} \, \mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2 \tag{164}$$

$$\leq \frac{d \| \boldsymbol{\rho} \, \Sigma_{\mathbf{X}}^{-\frac{1}{2}} \|_F^2 \| \Sigma_{\mathbf{X}}^{\frac{1}{2}} \|_F^2}{(a/d)^2} \, \mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2 \tag{165}$$

$$= \frac{d^4 (1 - \sigma^2)}{a^2} \, \mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2 \tag{166}$$

where for (166) we used the fact that $\| \boldsymbol{\rho} \, \Sigma_{\mathbf{X}}^{-\frac{1}{2}} \|^2 = \boldsymbol{\rho} \, \Sigma_{\mathbf{X}}^{-1} \boldsymbol{\rho}^T = 1 - \sigma^2$, and $\| \Sigma_{\mathbf{X}}^{\frac{1}{2}} \|_F^2 = \text{tr} \, \Sigma_{\mathbf{X}} = d$.

Second term (161): For any two $d \times d$ matrices $A, B$, it can be easily shown that $\text{tr} \, AB^T \leq d^2 \| A \|_F \| B \|_F$. Therefore,

$$\sigma^2 \, \mathbb{E} \, \text{tr} (\mathbf{W} - \hat{\mathbf{W}}) \hat{\mathbf{W}}^{-1} \Sigma_{\mathbf{X}} (\hat{\mathbf{W}}^{-T} + \mathbf{W}^{-T}) \mathbf{W}^{-1} \tag{167}$$

$$\leq \sigma^2 d^2 \, \mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F \| \hat{\mathbf{W}}^{-1} \Sigma_{\mathbf{X}} (\hat{\mathbf{W}}^{-T} + \mathbf{W}^{-T}) \mathbf{W}^{-1} \|_F \tag{168}$$

$$\leq \sigma^2 d^2 \sqrt{\mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2} \sqrt{\mathbb{E} \, \| \hat{\mathbf{W}}^{-1} \Sigma_{\mathbf{X}} (\hat{\mathbf{W}}^{-T} + \mathbf{W}^{-T}) \mathbf{W}^{-1} \|_F^2} \tag{169}$$

$$\leq \sigma^2 d^2 \frac{\sqrt{2} d^{\frac{3}{2}} \| \Sigma_{\mathbf{X}} \|_F}{(a - (d-1)b)^3} \sqrt{\mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2} \tag{170}$$

$$\leq \frac{\sqrt{2} d^{4.5} \sigma^2}{(a/d) d^2} \sqrt{\mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2} \tag{171}$$

where (169) is due to the Cauchy–Schwarz inequality. For (171) note that $\| \Sigma_{\mathbf{X}} \|_F^2 \leq d^2$ because all the entries of $\Sigma_{\mathbf{X}}$ are less than or equal to one.

We now proceed to upper bound $\mathbb{E} \, \| \mathbf{W} - \hat{\mathbf{W}} \|_F^2$. Consider the following uniform quantizer: The diagonal entries are truncated at some $c > a$. The

double segment $\pm[a, c]$ is divided into $l_1$ regions of width $\epsilon_1 = 2(c-a)/l_1$ each. For $|w| > c$ we take $\hat{w} = \text{sign}(w)c$. Therefore

$$\mathbb{E}(W_{11} - \hat{W}_{11})^2 \tag{172}$$

$$= \frac{Q(c)}{Q(a)} \mathbb{E}\left((W_{11} - \hat{W}_{11})^2 \,\middle|\, |W_{11}| > c\right) \tag{173}$$

$$+ \left(1 - \frac{Q(c)}{Q(a)}\right) \mathbb{E}\left((W_{11} - \hat{W}_{11})^2 \,\middle|\, |W_{11}| < c\right) \tag{174}$$

$$\leq \frac{Q(c)}{Q(a)} \mathbb{E}\left((W_{11} - c)^2 \,\middle|\, |W_{11}| > c\right) + \epsilon_1^2 \tag{175}$$

$$= \frac{Q(c)}{Q(a)}(1 + cs(c) + c^2) + \epsilon_1^2 \tag{176}$$

$$\leq \frac{a + a^{-1}}{c} e^{-\frac{c^2 - a^2}{2}} 2(1 + c^2) + \epsilon_1^2 \tag{177}$$

$$\leq 8c^2 e^{-\frac{c^2 - a^2}{2}} + \epsilon_1^2 \tag{178}$$

where (177) is obtained with some manipulations on $t \leq s(t) \leq t + t^{-1}$. For the off-diagonal entries, the segment $[-b, b]$ is divided into $l_2$ regions of width $\epsilon_2 = 2b/l_2$ each. Therefore

$$\mathbb{E}(W_{12} - \hat{W}_{12})^2 \leq \epsilon_2^2. \tag{179}$$

It follows that

$$\mathbb{E}\|\mathbf{W} - \hat{\mathbf{W}}\|_F^2 \tag{180}$$

$$= d\,\mathbb{E}(W_{11} - \hat{W}_{11})^2 + (d^2 - d)\,\mathbb{E}(W_{12} - \hat{W}_{12})^2 \tag{181}$$

$$\leq 8dc^2 e^{-\frac{c^2 - a^2}{2}} + d\epsilon_1^2 + d(d-1)\epsilon_2^2 \tag{182}$$

$$\leq 8dc^2 e^{-\frac{c^2 - a^2}{2}} + d^2(\epsilon_1 + \epsilon_2)^2. \tag{183}$$

We take $c = \sqrt{3}a$ and $l_1 = l_2$ and thus $\epsilon_1 + \epsilon_2 = 2(c - a + b)/l_1 \leq 4a/l_1$. The number of bits used for quantization is $k_q = \log l_1$ and therefore $\epsilon_1 + \epsilon_2 \leq 4a2^{-k_q}$. Now,

$$\sqrt{\mathbb{E}\|\mathbf{W} - \hat{\mathbf{W}}\|_F^2} \tag{184}$$

$$\leq \sqrt{24da^2 e^{-a^2} + 16d^2 a^2 2^{-2k_q}} \tag{185}$$

$$\leq \sqrt{(5ad)^2 (e^{-a^2} + 2^{-2k_q})} \tag{186}$$

$$\leq 5ad(e^{-\frac{a^2}{2}} + 2^{-k_q}), \tag{187}$$

and finally, combining (187) with (166) and (171) yields

$$\mathbb{E}\|\hat{\boldsymbol{\rho}} - \boldsymbol{\rho}\|^2 - \mathbb{E}\|\hat{\boldsymbol{\rho}}_0 - \boldsymbol{\rho}\|^2 \tag{188}$$

30

$$\leq 25d^6(1-\sigma^2)(e^{-\frac{a^2}{2}}+2^{-k_q})^2+5\sqrt{2}d^{4.5}\sigma^2(e^{-\frac{a^2}{2}}+2^{-k_q}) \qquad (189)$$

$$\leq 25d^6(e^{-\frac{a^2}{2}}+2^{-k_q}) \qquad (190)$$

which completes the proof.

## A.4 Proof of Lemma 2

Denote by $\mathcal{P}$ the set of all $d \times d$ signed permutation matrices, i.e. matrices with exactly one nonzero entry in every row and every column, that takes values in $\{-1,1\}$. For any $d \times d$ matrix $B$ and any $P \in \mathcal{P}$, the matrix $PBP^T$ is obtained from $B$ by performing the same permutation on the rows and columns of $B$, with possible sign changes. Specifically, the diagonal of $PBP^T$ is a permutation of the diagonal of $B$, and the off-diagonal of $PBP^T$ is a permutation of the off-diagonal entries of $B$ with possible sign changes.

Suppose that the random matrix $\mathbf{N}$ has the same distribution as $PNP^T$ for any $P \in \mathcal{P}$. It follows that $\mathbb{E}\,\mathbf{N}$ must be a scalar multiple of the identity matrix since for any $i \neq j$ there exist two matrices $P_1, P_2 \in \mathcal{P}$ such that

1. for some $i' \neq j'$,

$$(P_1\mathbf{N}P_1^T)_{i'j'} = (\mathbf{N})_{ij} \qquad (191)$$

$$(P_2\mathbf{N}P_2^T)_{i'j'} = -(\mathbf{N})_{ij} \qquad (192)$$

and thus $(\mathbb{E}\,\mathbf{N})_{ij} = -(\mathbb{E}\,\mathbf{N})_{ij}$.

2. for some $i'$,

$$(P_1\mathbf{N}P_1^T)_{i'i'} = (\mathbf{N})_{ii} \qquad (193)$$

$$(P_2\mathbf{N}P_2^T)_{i'i'} = (\mathbf{N})_{jj} \qquad (194)$$

and thus $(\mathbb{E}\,\mathbf{N})_{ii} = (\mathbb{E}\,\mathbf{N})_{jj}$.

The assumptions in the lemma imply that $\mathbf{M}$ and $P\mathbf{M}P^T$ have the same distribution for any $P \in \mathcal{P}$, thus $\mathbf{M}\mathbf{M}^T$ and $(P\mathbf{M}P^T)(P\mathbf{M}P^T)^T$ have the same distribution. Hence

$$(P\mathbf{M}P^T)(P\mathbf{M}P^T)^T = P\mathbf{M}\mathbf{M}^T P^T \qquad (195)$$

and thus $P\mathbf{M}\mathbf{M}^T P^T$ has the same distribution as $\mathbf{M}\mathbf{M}^T$. This implies that $(P\mathbf{M}\mathbf{M}^T P^T)^{-1}$ has the same distribution as $(\mathbf{M}\mathbf{M}^T)^{-1}$, and since

$$(P\mathbf{M}\mathbf{M}^T P^T)^{-1} = P(\mathbf{M}\mathbf{M}^T)^{-1}P^T \qquad (196)$$

we have that $P(\mathbf{M}\mathbf{M}^T)^{-1}P^T$ has the same distribution as $(\mathbf{M}\mathbf{M}^T)^{-1}$. Therefore $\mathbb{E}(\mathbf{M}\mathbf{M}^T)^{-1}$ is a scalar multiple of the identity matrix.

## A.5 Proof of Proposition 2

For any distribution on $X$, and for any $u, \Delta$,

$$\mathbb{E}(\hat{\rho}_{\text{th-q}} - \rho)^2 = \rho^2 \, \mathbb{E}\left(\frac{X_J - \hat{X}_J}{\hat{X}_J}\right)^2 + (1 - \rho^2) \, \mathbb{E}\frac{1}{\hat{X}_J^2} \tag{197}$$

$$\leq \frac{\rho^2}{t^2} \, \mathbb{E}\left(X_J - \hat{X}_J\right)^2 + \frac{1 - \rho^2}{t^2} \tag{198}$$

$$= \frac{\rho^2}{t^2} \Pr(X_J < u) \, \mathbb{E}\left(\left(X_J - \hat{X}_J\right)^2 \middle| X_J < u\right) \tag{199}$$

$$+ \frac{\rho^2}{t^2} \Pr(X_J > u) \, \mathbb{E}\left((X_J - u)^2 \middle| X_J > u\right) + \frac{1 - \rho^2}{t^2}$$

$$\leq \frac{\rho^2}{t^2}\Delta^2 + \frac{\rho^2}{t^2}\frac{\Pr(X > u)}{\Pr(X > t)} \, \mathbb{E}\left((X - u)^2 \middle| X > u\right) + \frac{1 - \rho^2}{t^2}. \tag{200}$$

In this Pareto example we have

$$\mathbb{E}\left((X - u)^2 \middle| X > u\right) = cu^2 \tag{201}$$

where $c = 2/((\alpha - 1)^2(\alpha - 2))$, and thus

$$\mathbb{E}(\hat{\rho}_{\text{th-q}} - \rho)^2 \leq \frac{1 - \rho^2 + \rho^2\Delta^2}{t^2} + c\rho^2\left(\frac{t}{u}\right)^{\alpha - 2}. \tag{202}$$

We take $u = t^{\frac{\alpha}{\alpha - 2}}$, and thus (202) becomes

$$\mathbb{E}(\hat{\rho}_{\text{th-q}} - \rho)^2 \leq \frac{1 - \rho^2 + \rho^2\Delta^2 + c\rho^2}{t^2}. \tag{203}$$

The bits are allocated by

$$k_q = \frac{1}{\alpha - 1}k, \quad k_l = \frac{\alpha - 2}{\alpha - 1}k, \tag{204}$$

and the threshold $t$ is determined by $k_l$ from the solution of $k_l = h_g(\Pr(X > t))$, which yields $t = 2^{\frac{k_l}{\alpha}(1 - o(1))}$. We have $\Delta \leq 1$ since

$$\Delta = 2^{-k_q}(u - t) = 2^{-k_q}(t^{\frac{\alpha}{\alpha - 2}} - t) = 2^{-k_q}t^{\frac{\alpha}{\alpha - 2}}(1 - t^{\frac{-2}{\alpha - 2}})$$

$$= 2^{-\frac{k}{\alpha - 1}}2^{\frac{k}{\alpha - 1}(1 - o(1))}(1 - t^{\frac{-2}{\alpha - 2}}) = 2^{-o(k)}(1 - t^{\frac{-2}{\alpha - 2}}).$$

Note that for $\alpha > 3$ we have $c < 1$ and thus

$$\mathbb{E}(\hat{\rho}_{\text{th-q}} - \rho)^2 \leq \frac{1 + \rho^2}{t^2} \leq \frac{1 + \rho^2}{2^{\frac{2}{\alpha}\frac{\alpha - 2}{\alpha - 1}k(1 - o(1))}}. \tag{205}$$

## A.6 Proof of Corollary 3

Set any real-valued sequence $\{\alpha_k\}_{k=1}^{\infty}$ such that all the elements are distinct, and $\sum \alpha_k^2 = 1$. Pick $Z = \sum_k \alpha_k Z_k$, where $Z_k \sim p_X$ are i.i.d and mutually independent of $X$. Then $Y = \alpha X + \sum_k \alpha_k Z_k$ is a weighted sum of i.i.d. r.v.s, hence knowing $p_X$ and $p_Y$, or even knowing all the weights $\alpha, \{\alpha_k\}$, there is no way to distinguish between the case where $X$ has coefficient $\alpha$ and where $X$ has coefficient $\alpha_k$ for some $k$. Thus, there is an infinite number of possible correlations.

## A.7 Maximum likelihood approximation

In this section we provide further justification for the estimator $\hat{\boldsymbol{\rho}} = \mathbf{Y}_J / \mathbb{E}\, X_J$ (or $Y_J / \mathbb{E}\, X_J$ in the scalar setup which is a special case), by showing that $\mathbf{Y}_J / X_J$ is an approximation of the maximum likelihood estimator. The model is

$$\mathbf{Y}_J = \boldsymbol{\rho}\, X_J + \Sigma^{\frac{1}{2}} \mathbf{Z}_J \qquad (206)$$

where either $J = \mathrm{argmax}_i \{X_i\}$ if we use the max method, or, if we use the threshold method, $J = \min\{i : X_i > t\}$. We wish to maximize $f_{X_J \mathbf{Y}_J}$ which is equivalent to maximizing $f_{\mathbf{Y}_J | X_J}$, since $f_{X_J}$ does not depend on $\boldsymbol{\rho}$. It follows that the actual distribution of $X_J$ is irrelevant. We have $\mathbf{Y}_J | X_J \sim \mathcal{N}(\boldsymbol{\rho}\, X_J, \Sigma)$ with $\Sigma = \Sigma_{\mathbf{Y}} - \boldsymbol{\rho}\,\boldsymbol{\rho}^T$ and thus

$$\frac{\partial}{\partial \boldsymbol{\rho}} \ln f_{X_J \mathbf{Y}_J} \qquad (207)$$

$$= -\frac{1}{2} \frac{\partial}{\partial \boldsymbol{\rho}} \left( \ln \det \Sigma + (\mathbf{Y}_J - \boldsymbol{\rho}\, X_J)^T \Sigma^{-1} (\mathbf{Y}_J - \boldsymbol{\rho}\, X_J) \right)$$

$$= \Sigma^{-1} \left( \boldsymbol{\rho} - (\mathbf{Y}_J - \boldsymbol{\rho}\, X_J) \left( (\mathbf{Y}_J - \boldsymbol{\rho}\, X_J)^T \Sigma^{-1} \boldsymbol{\rho} - X_J \right) \right)$$

meaning we want to solve

$$\boldsymbol{\rho} = (\mathbf{Y}_J - \boldsymbol{\rho}\, X_J) \left( (\mathbf{Y}_J - \boldsymbol{\rho}\, X_J)^T \Sigma^{-1} \boldsymbol{\rho} - X_J \right). \qquad (208)$$

Note that the rightmost term is a scalar and thus the solution must be of the form $\hat{\boldsymbol{\rho}}_{\mathrm{ML}} = C \mathbf{Y}_J$ where $C$ is a scalar that depends on $X_J$ and $\mathbf{Y}_J$. Plugging it yields that $C$ is obtained as the solution of the third degree polynomial

$$\mathbf{Y}_J^T \Sigma_{\mathbf{Y}}^{-1} \mathbf{Y}_J C^2 (X_J - C) - (X_J^2 - 1 + \mathbf{Y}_J^T \Sigma_{\mathbf{Y}}^{-1} \mathbf{Y}_J) C + X_J = 0.$$

In our setups $X_J$ takes large values. This implies in general that the entries of $\mathbf{Y}_J$ are also large, and thus $C$ should be small as we expect $\hat{\boldsymbol{\rho}}_{\mathrm{ML}} = C \mathbf{Y}_J$ to produce moderate values. Therefore we can assume that $X_J - C \approx X_J$ and that $X_J^2 - 1 \approx X_J^2$, which results in a quadratic equation in $C$ whose solutions are $X_J / (\mathbf{Y}_J^T \Sigma_{\mathbf{Y}}^{-1} \mathbf{Y}_J)$ and $1/X_J$. Note that (with either max or threshold) $\mathsf{Var}\, X_J$ approaches zero as the number of bits increases, and therefore the loss in replacing $X_J$ with $\mathbb{E}\, X_J$ is negligible (it is also evident in the optimality claims throughout where it is shown that the estimators, which do not use the actual value of $X_J$, achieve the CRLB that assumes $X_J$ is known).

# References

[1] E. L. Lehmann and G. Casella, *Theory of point estimation.* Springer Science & Business Media, 2006.

[2] Z. Zhang and T. Berger, "Estimation via compressed information," *IEEE transactions on Information theory*, vol. 34, no. 2, pp. 198–211, 1988.

[3] T. Liu and P. Viswanath, "Opportunistic orthogonal writing on dirty paper," *IEEE transactions on information theory*, vol. 52, no. 5, pp. 1828–1846, 2006.

[4] S. Borade and L. Zheng, "Writing on fading paper and causal transmitter csi," in *Information Theory, 2006 IEEE International Symposium on*, pp. 744–748, IEEE, 2006.

[5] A. No and T. Weissman, "Rateless lossy compression via the extremes," *IEEE Transactions on Information Theory*, vol. 62, no. 10, pp. 5484–5495, 2016.

[6] R. Ahlswede and M. Burnashev, "On minimax estimation in the presence of side information about remote data," *The Annals of Statistics*, pp. 141–171, 1990.

[7] T. S. Han and S.-i. Amari, "Parameter estimation with multiterminal data compression," *IEEE transactions on Information Theory*, vol. 41, no. 6, pp. 1802–1833, 1995.

[8] S. Amari *et al.*, "Statistical inference under multiterminal data compression," *IEEE Transactions on Information Theory*, vol. 44, no. 6, pp. 2300–2324, 1998.

[9] S.-i. Amari, "On optimal data compression in multiterminal statistical inference," *IEEE Transactions on Information Theory*, vol. 57, no. 9, pp. 5577–5587, 2011.

[10] E. Haim and Y. Kochman, "Binary distributed hypothesis testing via Körner-Marton coding," in *Information Theory Workshop (ITW), 2016 IEEE*, pp. 146–150, IEEE, 2016.

[11] Y. Zhang, J. Duchi, M. I. Jordan, and M. J. Wainwright, "Information-theoretic lower bounds for distributed statistical estimation with communication constraints," in *Advances in Neural Information Processing Systems*, pp. 2328–2336, 2013.

[12] M. El Gamal and L. Lai, "Are Slepian-Wolf rates necessary for distributed parameter estimation?," in *Communication, Control, and Computing (Allerton), 2015 53rd Annual Allerton Conference on*, pp. 1249–1255, IEEE, 2015.

[13] J.-J. Xiao, S. Cui, Z.-Q. Luo, and A. J. Goldsmith, "Joint estimation in sensor networks under energy constraints," in *Sensor and Ad Hoc Communications and Networks, 2004. IEEE SECON 2004. 2004 First Annual IEEE Communications Society Conference on*, pp. 264–271, IEEE, 2004.

[14] Z.-Q. Luo, "Universal decentralized estimation in a bandwidth constrained sensor network," *IEEE Transactions on information theory*, vol. 51, no. 6, pp. 2210–2219, 2005.

[15] J. A. Gubner, "Distributed estimation and quantization," *IEEE Transactions on Information Theory*, vol. 39, no. 4, pp. 1456–1459, 1993.

[16] A. Xu and M. Raginsky, "Information-theoretic lower bounds on Bayes risk in decentralized estimation," *IEEE Transactions on Information Theory*, vol. 63, no. 3, pp. 1580–1600, 2017.

[17] M. Braverman, A. Garg, T. Ma, H. L. Nguyen, and D. 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*, pp. 1011–1020, ACM, 2016.

[18] I. D. Schizas, A. Ribeiro, and G. B. Giannakis, "Consensus in ad hoc wsns with noisy links—part i: Distributed estimation of deterministic signals," *IEEE Transactions on Signal Processing*, vol. 56, no. 1, pp. 350–364, 2008.

[19] J.-J. Xiao, A. Ribeiro, Z.-Q. Luo, and G. B. Giannakis, "Distributed compression-estimation using wireless sensor networks," *IEEE Signal Processing Magazine*, vol. 23, no. 4, pp. 27–41, 2006.

[20] A. Ribeiro and G. B. Giannakis, "Bandwidth-constrained distributed estimation for wireless sensor networks-part i: Gaussian case," *IEEE transactions on signal processing*, vol. 54, no. 3, pp. 1131–1143, 2006.

[21] P. Venkitasubramaniam, L. Tong, and A. Swami, "Quantization for maximin are in distributed estimation," *IEEE Transactions on Signal Processing*, vol. 55, no. 7, pp. 3596–3605, 2007.

[22] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2006.

[23] S. M. Kay, "Fundamentals of statistical signal processing, volume i: estimation theory," 1993.

[24] H. David and H. Nagaraja, *Order Statistics*. Wiley Series in Probability and Statistics, Wiley, 2004.

[25] C. G. Small, *Expansions and asymptotics for statistics*. CRC Press, 2010.

[26] A. C. Rencher, *Methods of multivariate analysis*, vol. 492. John Wiley & Sons, 2003.

[27] C. R. Johnson, "A Gersgorin-type lower bound for the smallest singular value," *Linear Algebra and its Applications*, vol. 112, pp. 1–7, 1989.

[28] W. W. Hager, "Updating the inverse of a matrix," *SIAM review*, vol. 31, no. 2, pp. 221–239, 1989.

[29] P. Billingsley, *Probability and Measure*. Wiley Series in Probability and Statistics, Wiley, 1995.

[30] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge university press, 1990.

[31] U. Hadar and O. Shayevitz, "Distributed estimation of Gaussian correlations," in *Information Theory (ISIT), 2018 IEEE International Symposium on*, pp. 511–515, IEEE, 2018.