

Quantized Compressive Sensing

Wei Dai, Hoa Vinh Pham, and Olgica Milenkovic
 Department of Electrical and Computer Engineering
 University of Illinois at Urbana-Champaign

Abstract

We study the average distortion introduced by scalar, vector, and entropy coded quantization of compressive sensing (CS) measurements. The asymptotic behavior of the underlying quantization schemes is either quantified exactly or characterized via bounds. We adapt two benchmark CS reconstruction algorithms to accommodate quantization errors, and empirically demonstrate that these methods significantly reduce the reconstruction distortion when compared to standard CS techniques.

I. INTRODUCTION

Compressive sensing (CS) is a linear sampling method that converts unknown input signals, embedded in a high dimensional space, into signals that lie in a space of significantly smaller dimension. In general, it is not possible to uniquely recover an unknown signal using measurements of reduced-dimensionality. Nevertheless, if the input signal is sufficiently sparse, exact reconstruction is possible. In this context, assume that the unknown signal $\mathbf{x} \in \mathbb{R}^N$ is K -sparse, i.e., that there are at most K nonzero entries in \mathbf{x} . A naive reconstruction method is to search among all possible signals and find the sparsest one which is consistent with the linear measurements. This method requires only $m = 2K$ random linear measurements, but finding the sparsest signal representation is an NP-hard problem. On the other hand, Donoho and Candès et. al. demonstrated in [1]–[4] that sparse signal reconstruction is a polynomial time problem if more measurements are taken. This is achieved by casting the reconstruction problem as a linear programming problem and solving it using the *basis pursuit* (BP) method. More recently, the authors proposed the *subspace pursuit* (SP) algorithm in [5] (see also the independent work [6] for a closely related approach). The computational complexity of the SP algorithm is linear in the signal dimension, and the required number of linear measurements is of the same order as that for the BP method.

For most practical applications, it is reasonable to assume that the measurements are quantized and therefore do not have infinite precision. When the quantization error is bounded and known in advance, upper bounds on the reconstruction distortion were derived for the BP method in [7] and the SP algorithm in [5], [6], respectively. For bounded compressible signals, which have transform coefficients with magnitudes that decay according to a power law, an upper bound on the reconstruction distortion introduced by a uniform quantizer was derived in [8]. The same quantizer was studied in [9] for exactly K -sparse signals and it was shown that a large fraction of quantization regions is not used [9]. All of the above approaches focus on the worst case analysis, or simple one-bit quantization [10]. An exception includes the overview paper [11], which focuses on the average performance of uniform quantizers, assuming that the support set of the sparse signal is available at the quantizer.

As opposed to the worst case analysis, we consider the average distortion introduced by quantization. We study the asymptotic distortion rate functions for scalar quantization, entropy coded scalar quantization, and vector quantization of the measurement vectors. Exact asymptotic distortion rate functions are derived for scalar quantization when both the measurement matrix and the sparse signals obey a certain probabilistic model. Lower and upper bounds on the asymptotic distortion rate functions are also derived for other quantization scenarios, and the problem of compressive sensing matrix quantization is briefly discussed as well. In addition, two benchmark CS reconstruction algorithms are adapted to accommodate quantization errors. Simulations show that the new algorithms offer significant performance improvement over classical CS reconstruction techniques that do not take quantization errors into consideration.

This paper is organized as follows. Section II contains a brief overview of CS theory, the BP and SP reconstruction algorithms, and various quantization techniques. In Section III, we analyze the CS distortion rate function and examine the

*Part of the material in this paper was submitted to the IEEE Information Theory Workshop (ITW), 2009, and the IEEE International Symposium on Information Theory (ISIT), 2009.

influence of quantization errors on the BP and SP reconstruction algorithms. In Section IV, we describe two modifications of the aforementioned algorithms, suitable for quantized data, that offer significant performance improvements when compared to standard BP and SP techniques. Simulation results are presented in Section V.

II. PRELIMINARIES

A. Compressive Sensing (CS)

In CS, one encodes a signal \mathbf{x} of dimension N by computing a measurement vector \mathbf{y} of dimension of $m \ll N$ via linear projections, i.e.,

$$\mathbf{y} = \Phi \mathbf{x},$$

where $\Phi \in \mathbb{R}^{m \times N}$ is referred to as the *measurement matrix*. In this paper, we assume that $\mathbf{x} \in \mathbb{R}^N$ is exactly K -sparse, i.e., that there are exactly K entries of \mathbf{x} that are nonzero. The reconstruction problem is to recover \mathbf{x} given \mathbf{y} and Φ .

The BP method is a technique that casts the reconstruction problem as a l_1 -regularized optimization problem, i.e.,

$$\min \|\mathbf{x}\|_1 \text{ subject to } \mathbf{y} = \Phi \mathbf{x}, \quad (1)$$

where $\|\mathbf{x}\|_1 = \sum_{i=1}^N |x_i|$ denotes the l_1 -norm of the vector \mathbf{x} . It is a convex optimization problem and can be solved efficiently by linear programming techniques. The reconstruction complexity equals $O(m^2 N^{3/2})$ if the convex optimization problem is solved using interior point methods [12].

The computational complexity of CS reconstruction can be further reduced by the SP algorithm, recently proposed by two research groups [5], [6]. It is an iterative algorithm drawing on the theory of list decoding. The computational complexity of this algorithm is upper bounded by $O(Km(N + K^2))$, which is significantly smaller than the complexity of the BP method whenever $K \ll N$. See [5] for a detailed performance and complexity analysis of this greedy algorithm.

A sufficient condition for both the BP and SP algorithms to perform exact reconstruction is based on the so called restricted isometry property (RIP) [2], formally defined as follows.

Definition 1: (RIP). A matrix $\Phi \in \mathbb{R}^{m \times N}$ is said to satisfy the Restricted Isometry Property (RIP) with coefficients (K, δ) for $K \leq m$, $0 \leq \delta \leq 1$, if for all index sets $I \subset \{1, \dots, N\}$ such that $|I| \leq K$ and for all $\mathbf{q} \in \mathbb{R}^{|I|}$, one has

$$(1 - \delta) \|\mathbf{q}\|_2^2 \leq \|\Phi_I \mathbf{q}\|_2^2 \leq (1 + \delta) \|\mathbf{q}\|_2^2.$$

The RIP parameter δ_K is defined as the infimum of all parameters δ for which the RIP holds, i.e.,

$$\delta_K := \inf \left\{ \delta : (1 - \delta) \|\mathbf{q}\|_2^2 \leq \|\Phi_I \mathbf{q}\|_2^2 \leq (1 + \delta) \|\mathbf{q}\|_2^2, \right. \\ \left. \forall |I| \leq K, \forall \mathbf{q} \in \mathbb{R}^{|I|} \right\}. \quad (2)$$

It was shown in [5], [7] that both BP and SP algorithms lead to exact reconstructions of K -sparse signals if the matrix Φ satisfies the RIP with a constant parameter, i.e., $\delta_{c_1 K} \leq c_0$ where both $c_1 \in \mathbb{R}^+$ and $c_0 \in (0, 1)$ are constants independent of K (although different algorithms may have different parameters c_0 s and c_1 s). Most known families of matrices satisfying the RIP property with optimal or near-optimal performance guarantees are random, including Gaussian random matrices with i.i.d. $\mathcal{N}(0, 1/m)$ entries, where $m \geq O(K \log N)$.

For completeness, we briefly describe the SP algorithm. For an index set $T \subset \{1, 2, \dots, N\}$, let Φ_T be the ‘‘truncated matrix’’ consisting of the columns of Φ indexed by T , and let $\text{span}(\Phi_T)$ denote the subspace in \mathbb{R}^m spanned by the columns of Φ_T . Suppose that $\Phi_T^* \Phi_T$ is invertible. For any given $\mathbf{y} \in \mathbb{R}^m$, the projection of \mathbf{y} onto $\text{span}(\Phi_T)$ is defined as

$$\mathbf{y}_p = \text{proj}(\mathbf{y}, \Phi_T) := \Phi_T (\Phi_T^* \Phi_T)^{-1} \Phi_T^* \mathbf{y}, \quad (3)$$

where Φ^* denotes the conjugate transpose of Φ .

The corresponding projection residue vector \mathbf{y}_r and projection coefficient vector \mathbf{x}_p are defined as

$$\mathbf{y}_r = \text{resid}(\mathbf{y}, \Phi_T) := \mathbf{y} - \mathbf{y}_p, \quad (4)$$

and

$$\mathbf{x}_p = \text{pcoeff}(\mathbf{y}, \Phi_T) := (\Phi_T^* \Phi_T)^{-1} \Phi_T^* \mathbf{y}. \quad (5)$$

The steps of the SP algorithm are summarized below.

Algorithm 1 The Subspace Pursuit (SP) Algorithm

Input: K, Φ, \mathbf{y}

Initialization: Let $T^0 = \{K \text{ indices corresponding to entries of largest magnitude in } \Phi^* \mathbf{y}\}$ and $\mathbf{y}_r^0 = \text{resid}(\mathbf{y}, \Phi_{T^0})$.

Iteration: At the ℓ^{th} iteration, go through the following steps.

- 1) $\tilde{T}^\ell = T^{\ell-1} \cup \{K \text{ indices corresponding to entries of largest magnitude in } \Phi^* \mathbf{y}_r^{\ell-1}\}$.
- 2) Let $\mathbf{x}_p = \text{pcoeff}(\mathbf{y}, \Phi_{\tilde{T}^\ell})$ and $T^\ell = \{K \text{ indices corresponding to entries of largest magnitude in } \mathbf{x}_p\}$.
- 3) $\mathbf{y}_r^\ell = \text{resid}(\mathbf{y}, \Phi_{T^\ell})$.
- 4) If $\|\mathbf{y}_r^\ell\|_2 > \|\mathbf{y}_r^{\ell-1}\|_2$, let $T^\ell = T^{\ell-1}$ and quit the iteration.

Output: The vector $\hat{\mathbf{x}}$ satisfying $\hat{\mathbf{x}}_{\{1, \dots, N\} - T^\ell} = \mathbf{0}$ and $\hat{\mathbf{x}}_{T^\ell} = \text{pcoeff}(\mathbf{y}, \Phi_{T^\ell})$.

In what follows, we study the performance of the SP and BP reconstruction algorithms when the measurements are subjected to three different quantization schemes. We also discuss the issue of quantizing the measurement matrix values.

B. Scalar and Vector Quantization

Let $\mathcal{C} \subset \mathbb{R}^m$ be a finite discrete set, referred to as a codebook. A quantizer is a mapping from \mathbb{R}^m to the codebook \mathcal{C} with the property that

$$\begin{aligned} \mathfrak{q} : \mathbb{R}^m &\rightarrow \mathcal{C} \\ \mathbf{y} &\mapsto \boldsymbol{\omega} \in \mathcal{C} \text{ if } \mathbf{y} \in \mathcal{R}_\omega, \end{aligned} \quad (6)$$

where ω is referred to as a *level* and \mathcal{R}_ω is the *quantization region* corresponding to the level ω . The performance of a quantizer is often described by its distortion-rate function, defined as follows. Let the distortion measure be the squared Euclidean distance (i.e., mean squared error (MSE)). For a random source $\mathbf{Y} \in \mathbb{R}^m$, the distortion associated with a quantizer \mathfrak{q} is $D_{\mathfrak{q}} := \mathbb{E} \left[\|\mathbf{Y} - \mathfrak{q}(\mathbf{Y})\|_2^2 \right]$. For a given codebook \mathcal{C} , the optimal quantization function that minimizes the Euclidean distortion measure is given by

$$\mathfrak{q}^*(\mathbf{Y}) = \arg \min_{\boldsymbol{\omega} \in \mathcal{C}} \|\mathbf{Y} - \boldsymbol{\omega}\|_2^2.$$

As a result, the corresponding quantization region is given by

$$\mathcal{R}_\omega := \left\{ \mathbf{y} \in \mathbb{R}^m : \|\mathbf{y} - \boldsymbol{\omega}\|_2^2 \leq \|\mathbf{y} - \boldsymbol{\omega}'\|_2^2, \forall \boldsymbol{\omega}' \in \mathcal{C} \right\}, \quad (7)$$

and the distortion associated with this codebook \mathcal{C} equals

$$D(\mathcal{C}) := \mathbb{E} \left[\|\mathbf{Y} - \mathfrak{q}^*(\mathbf{Y})\|_2^2 \right].$$

Let $R := \frac{1}{m} \log_2 |\mathcal{C}|$ be the rate of the codebook \mathcal{C} . For a given code rate R , the distortion rate function is given by

$$D^*(R) := \inf_{\mathcal{C}: \frac{1}{m} \log_2 |\mathcal{C}| \leq R} D(\mathcal{C}). \quad (8)$$

For simplicity, assume that the random source \mathbf{Y} does not have mass points, and that the levels in the quantization codebook are all distinct. With these assumptions, though different quantization regions (7) may overlap, the ties can be broken arbitrarily as they happen with probability zero.

We study both vector quantization and scalar quantization. Scalar quantization has lower computational complexity than vector quantization. It is a special case of vector quantization when $m = 1$. To distinguish the two schemes, we use the subscripts SQ and VQ to refer to scalar and vector quantization, respectively. For quantized compressive sensing, we assume

that the quantization functions for all the coordinate of \mathbf{Y} are the same. The corresponding distortion rate function is therefore of the form

$$D_{SQ}^*(R) := \inf_{\mathcal{C}_{SQ}: \log_2 |\mathcal{C}_{SQ}| \leq R} \mathbb{E}_{\mathbf{Y}} \left[\sum_{i=1}^m |Y_i - q_{SQ}(Y_i)|^2 \right]. \quad (9)$$

Necessary conditions for optimal scalar quantizer design can be found in [13]. The quantization region for the level $\omega_i \in \mathcal{C}$, $i = 1, 2, \dots, 2^R$, can be written in the form $\mathcal{R}_{\omega_i} = \overline{(t_{i-1}, t_i)}$, where $t_{i-1}, t_i \in \mathbb{R} \cup \{-\infty\} \cup \{+\infty\}$ and $\overline{(t_{i-1}, t_i)}$ is the closure of the open interval (t_{i-1}, t_i) . An optimal quantizer satisfies the following conditions:

- 1) If the optimal quantizer has levels ω_{i-1} and ω_i , then the threshold that minimizes the mean square error (MSE) is

$$t_i = \frac{1}{2} (\omega_i + \omega_{i+1}). \quad (10)$$

- 2) If the optimal quantizer has thresholds t_{i-1} and t_i , then the level that minimizes the MSE is

$$\omega_i = \mathbb{E} \left[Y | Y \in \overline{(t_{i-1}, t_i)} \right]. \quad (11)$$

Lloyd's algorithm [13] for quantizer codebook design is based on the above necessary conditions. Lloyd's algorithm starts with an initial codebook, and then in each iteration, computes the thresholds t_i s according to (10) and updates the codebook via (11). Although Lloyd's algorithm is not guaranteed to find a global optimum for the quantization regions, it produces locally optimal codebooks.

As a low-complexity alternative to non-uniform quantizers, uniform scalar quantizers are widely used in practice. A uniform scalar quantizer is associated with a "uniform codebook" $\mathcal{C}_{u,SQ} = \{\omega_1 < \omega_2 < \dots < \omega_M\}$, for which $\omega_i - \omega_{i-1} = \omega_j - \omega_{j-1}$ for all $1 < i \neq j \leq M$. The difference between adjacent levels is often referred to as the step size, and denoted by $\Delta_{u,SQ}$. The corresponding distortion rate function is given by

$$D_{u,SQ}^*(R) := \inf_{\mathcal{C}_{u,SQ}: \log_2 |\mathcal{C}_{u,SQ}| \leq R} \mathbb{E}_{\mathbf{Y}} \left[\sum_{i=1}^m |Y_i - q_{SQ}(Y_i)|^2 \right]. \quad (12)$$

where \mathcal{C}_{SQ} in (9) is replaced by $\mathcal{C}_{u,SQ}$.

Definitions (9) and (12) are consistent with (8) as a Cartesian product of scalar quantizers can be viewed as a special form of a vector quantizer.

III. DISTORTION ANALYSIS

We analyze the asymptotic behavior of the distortion rate functions introduced in the previous section. We assume that the quantization codebook \mathcal{C} , for both scalar and vector quantization, is designed offline and fixed when the measurements are taken.

A. Distortion of Scalar Quantization

For scalar quantization, we consider the following two CS scenarios.

Assumptions I:

- 1) Let $\Phi = \frac{1}{\sqrt{m}} \mathbf{A} \in \mathbb{R}^{m \times N}$, where the entries of \mathbf{A} are i.i.d. Subgaussian random variables¹ with zero mean and unit variance.

¹A random variable X is said to be *Subgaussian* if there exist positive constants c_1 and c_2 such that

$$\Pr(|X| > x) \leq c_1 e^{-c_2 x^2} \quad \forall x > 0.$$

One property of Subgaussian distributions is that they have a well defined moment generating function. Note that the Gaussian and Bernoulli distributions are special cases of the Subgaussian distribution.

2) Let $\mathbf{X} \in \mathbb{R}^N$ be an exactly K -sparse vector, that is, a signal that has exactly K nonzero entries. We assume that the nonzero entries of \mathbf{X} are i.i.d. Subgaussian random variables with zero mean and unit variance, although more general models can be analyzed in a similar manner.

Assumptions II: Assume that $\mathbf{X} \in \mathbb{R}^n$ is exactly K -sparse, and that the nonzero entries of \mathbf{X} are i.i.d. standard Gaussian random variables.

The asymptotic distortion-rate function of the measurement vector under the first CS scenario is characterized in Theorem 1.

Theorem 1: Suppose that Assumptions I hold. Then

$$\lim_{R \rightarrow \infty} \lim_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{K} D_{SQ}^*(R) = \frac{\pi\sqrt{3}}{2}, \quad (13)$$

and

$$\lim_{R \rightarrow \infty} \lim_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{KR} D_{u, SQ}^*(R) = \frac{4}{3} \ln 2. \quad (14)$$

The proof is based on the fact that the distributions of $\sqrt{\frac{m}{K}} Y_i$, $1 \leq i \leq m$, weakly converge to standard Gaussian distributions. The detailed description is given in Appendix A.

To study the scenario described by Assumptions II, we need the following definitions. For a given matrix Φ , let

$$\mu_1 := \frac{1}{N} \sum_{i \in [m], j \in [N]} \varphi_{i,j}^2, \quad (15)$$

and

$$\mu_2 := \max_{i \in [m], T \in \binom{[N]}{K}} \frac{m}{K} \sum_{j \in T} \varphi_{i,j}^2, \quad (16)$$

where $[m] = \{1, 2, \dots, m\}$ and $\binom{[N]}{K}$ denotes the set of all subsets of $[N]$ with cardinality K . Note that if the matrix Φ is generated from the random ensemble described in Assumption I.1), then $\mu_1 \in (1 - \epsilon, 1 + \epsilon)$ with high probability, for all $\epsilon > 0$, and whenever m and N are sufficiently large. It is straightforward to verify that $\mu_1 \leq \mu_2$.

With these definitions at hand, bounds on the distortion rate function can be described as below.

Theorem 2: Suppose that Assumption II holds. Then

$$\begin{aligned} \frac{\pi\sqrt{3}}{2} \mu_1 &\leq \liminf_{R \rightarrow \infty} \frac{2^{2R}}{K} D_{SQ}^*(R) \\ &\leq \limsup_{R \rightarrow \infty} \frac{2^{2R}}{K} D_{SQ}^*(R) \leq \frac{\pi\sqrt{3}}{2} \mu_2, \end{aligned} \quad (17)$$

and

$$\frac{4 \ln 2}{3} \mu_1 \leq \liminf_{R \rightarrow \infty} \frac{2^{2R}}{KR} D_{u, SQ}^*(R). \quad (18)$$

The detailed proof is postponed to Appendix B. Here, we sketch the basic ideas behind the proof. In order to construct a lower bound, suppose that one has prior information about the support set T before taking the measurements. For a given value of i and for a given T , we calculate the corresponding asymptotic distortion-rate function. The lower bound is obtained by taking the average of these distortion-rate functions over all possible values of i and T . For the upper bound, we design a sequence of sub-optimal scalar quantizers, then apply them to all measurement components, and finally construct a uniform upper bound on their asymptotic distortion-rate functions, valid for all i and T . The uniform upper bound is given in (17).

Remark 1: Our results are based on the fundamental assumption that the sparsity level K is known in advance and that the statistics of the sparse vector \mathbf{x} is specified. Very frequently, however, this is not the case in practice. If we relax Assumptions I and II further by assuming that K is sufficiently large, it will often be the case that the statistics of the measurement Y_i is well approximated by a Gaussian distribution. Here, note that different Y_i variables may have different variances and these variances are generally unknown in advance. The problem of statistical mismatch has been analyzed in the proof of the upper

bound (17) (see Proposition 1 of Appendix B for details). In particular, non-uniform quantization with slightly over-estimated variance performs better than that with under-estimated variance [14, Chapter 8.6].

According to Theorem 1, if the quantization rate R is sufficiently large, the distortion of the optimal non-uniform quantizer is approximately only $1/R$ of that of the optimal uniform quantizer. This gap can be closed by using entropy coding techniques in conjunction with uniform quantizers.

B. Uniform Scalar Quantization with Entropy Encoding

Let $\mathcal{B}_{enc} = \{v_1, v_2, \dots, v_M\}$ be a binary codebook, where the codewords v_i , $1 \leq i \leq M$, are finite-length strings over the binary field with elements $\{0, 1\}$. The codebook \mathcal{B}_{enc} can, in general, contain codewords of variable length - i.e., the lengths of different codewords are allowed to be different. Let ℓ_i be the length of codeword v_i , $i = 1, 2, \dots, M$. Then $v_i \in \{0, 1\}^{\ell_i \times 1}$. For a given quantization codebook $\mathcal{C} = \{\omega_1, \omega_2, \dots, \omega_M\}$, the encoding function f_{enc} is a mapping from the quantization codebook \mathcal{C} to the binary codebook \mathcal{B}_{enc} , i.e., $f_{enc}(\omega) = v \in \mathcal{B}_{enc}$. The extension f_{enc}^* is a mapping from finite length strings of \mathcal{C} to finite length strings of \mathcal{B}_{enc} (a concatenation of the corresponding binary codewords):

$$f_{enc}^*(\omega_{i_1} \omega_{i_2} \dots \omega_{i_s}) = f_{enc}(\omega_{i_1}) f_{enc}(\omega_{i_2}) \dots f_{enc}(\omega_{i_s}).$$

The code \mathcal{B}_{enc} is called *uniquely decodable* if any concatenation of binary codewords $v_{i_1} v_{i_2} \dots v_{i_s}$ has only one possible preimage string $\omega_{j_1} \omega_{j_2} \dots \omega_{j_s}$ producing it. In practice, the code \mathcal{B}_{enc} is often chosen to be a *prefix* code, that is, no codeword is a prefix of any other codeword. A prefix code can be uniquely decoded as the end of a codeword is immediately recognizable without checking future encoded bits.

We consider the case in which scalar quantization is followed by variable-length encoding. The corresponding expected encoding length \bar{L} is defined by

$$\bar{L} = \mathbb{E}_Y [\mathcal{L} \circ f_{enc} \circ q_{SQ}(Y)],$$

where $\mathcal{L}(v)$ outputs the length of the encoding codeword $v \in \mathcal{B}_{enc}$. The goal is to *jointly* design q_{SQ} and f_{enc} to minimize the expected encoding length \bar{L} . We are interested in the distortion rate function defined by

$$D_{enc}^*(R) := \inf_{\bar{L} \leq R} \mathbb{E}_Y \left[\sum_{i=1}^m |Y_i - q_{SQ}(Y_i)|^2 \right]. \quad (19)$$

Theorem 3: Suppose that Assumptions I hold. Then

$$\begin{aligned} \frac{\pi e}{6} &\leq \liminf_{R \rightarrow \infty} \liminf_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{K} D_{enc}^*(R) \\ &\leq \limsup_{R \rightarrow \infty} \limsup_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{K} D_{enc}^*(R) \leq \frac{\pi e}{3}, \end{aligned}$$

and the upper bound is achieved by a uniform scalar quantizer with

$$\lim_{R \rightarrow \infty} \lim_{(K, m, N) \rightarrow \infty} \sqrt{\frac{m}{2\pi e K}} 2^R \Delta_{u, SQ} = 1,$$

followed by Huffman encoding.

Proof: Given a quantization function, Huffman encoding gives an optimal prefix code that minimizes \bar{L} [15, Chapter 5]. Let $p_i = \Pr(Y : q(Y) = \omega_i)$ and let ℓ_i be the length of encoded codeword $f_{enc}(\omega_i)$. Let $H := \sum_{i=1}^M -p_i \log_2 p_i$. Then $H \leq \bar{L} = \sum_i p_i \ell_i \leq H + 1$. In addition, it is well known that the distortion of scalar quantization of a Gaussian source is lower bounded by $\frac{1}{12} 2^{2(h-H)} (1 + o_H(1))$, where h denotes the differential entropy of the source, and the lower bound is achieved by a uniform quantizer. Calculating h and interpreting H as a function of \bar{L} establish the claimed result. ■

As expected, for a given average description length, the average distortion of uniform scalar quantization and Huffman encoding is smaller than that of an optimal scalar quantizer with fixed length encoding.

C. Distortion of Vector Quantization

For the purpose of analyzing vector quantization schemes, we make the following assumptions.

Assumptions III:

- 1) Let $\Phi \in \mathbb{R}^{m \times N}$ be a matrix satisfying the RIP with parameter $\delta_K \in (0, 1)$.
- 2) Assume that $\mathbf{X} \in \mathbb{R}^n$ is exactly K -sparse, and that the nonzero entries of \mathbf{X} are i.i.d. standard Gaussian random variables.

Theorem 4: Suppose that Assumptions III hold. Then

$$(1 - \delta_K)(1 + o_K(1)) \leq \liminf_{R \rightarrow \infty} \frac{2^{2Rm/K}}{K} D_{VQ}^*(R) \quad (20)$$

$$\leq \limsup_{R \rightarrow \infty} \frac{2^{2R}}{m} D_{VQ}^*(R) \leq (1 + \delta_K)(1 + o_m(1)), \quad (21)$$

where $o_K(1) \xrightarrow{K \rightarrow \infty} 0$ and $o_m(1) \xrightarrow{m \rightarrow \infty} 0$. Another upper bound on $D_{VQ}^*(R)$ is given by

$$\limsup_{R \rightarrow \infty} \frac{2^{2R}}{K} D_{VQ}^*(R) \leq \frac{\pi\sqrt{3}}{2} \mu_2, \quad (22)$$

where μ_2 is as defined in (16).

Remark 2: The comparison of the two upper bounds in (21) and (22) depends on the ratio between m and K . Consider the case where $N = \beta K$, $m = \Theta(K \log(N/K)) = \alpha K$ for some $\alpha, \beta > 1$. The first upper bound becomes

$$\limsup_{R \rightarrow \infty} \frac{2^{2R}}{K} D_{VQ}^*(R) \leq \alpha(1 + \delta_K)(1 + o_m(1)).$$

It is smaller than the second upper bound if and only if

$$\delta_K < \frac{\pi\sqrt{3}}{2\alpha} \mu_2 - 1.$$

The upper bound (22) is obtained by using the Cartesian product of scalar quantizers and invoking the result in (17). The bounds (20) and (21) are proved in Appendix B. The basic ideas behind the proof are similar to those used for proving Theorem 2: the lower bound is obtained by averaging the distortions of optimal quantizers for every $T \in \binom{[N]}{K}$, while the upper bound is a uniform upper bound on the distortions of quantizers constructed for all $T \in \binom{[N]}{K}$.

Note that the lower bound in (20) is not achievable when $K < m$. The upper bounds (21) and (22) do not guarantee significant distortion reduction of vector quantization compared with scalar quantization. Due to their inherently high computational complexity, vector quantizers do not offer clear advantages that justify their use in practice.

D. CS Measurement Matrix Quantization Effects

In CS theory, the measurement matrix is generated either randomly or by some deterministic construction. Examples include Gaussian random matrices and the deterministic construction based on Vandermonde matrices [16], [17]. In both examples, the matrix entries typically have infinite precision, which is not the case in practice. It is therefore also plausible to study the effect of quantization of CS measurement matrix.

Consider Assumption I where the measurement matrix is randomly generated. Let us assume that every entry $\varphi_{i,j}$, $1 \leq i \leq m$ and $1 \leq j \leq N$, is quantized using a finite number of bits. Note that $\hat{\varphi}_{i,j} = \mathfrak{q}(\varphi_{i,j})$ is a bounded random variable and therefore Subgaussian distributed. The results in Theorem 1 are therefore automatically valid for quantized matrices as well.

Suppose that the measurement matrix is constructed deterministically and then quantized using a finite number of bits. The parameters μ_1 , μ_2 and δ_K of the quantized measurement matrix can be computed according to (15), (16) and (2), respectively. The results regarding scalar quantization and vector quantization described in Theorems 2 and 4 can be easily seen to hold in this case as well.

E. Reconstruction Distortion

Based on the results of the previous section, we are ready to quantify the reconstruction distortion of BP and SP methods introduced by quantization error.

It is well known from CS literature that the reconstruction distortion is dependent on the distortion in the measurements. Consider the quantized CS given by

$$\hat{\mathbf{Y}} = \mathbf{q}(\mathbf{Y}) = \Phi \mathbf{X} + \mathbf{E},$$

and where $\mathbf{E} \in \mathbb{R}^m$ denotes the quantization error. Let $\hat{\mathbf{X}}$ be the reconstructed signal based on the quantized measurements $\hat{\mathbf{Y}}$. Then the reconstruction distortion can be upper bounded by

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_2^2 \leq c^2 \|\mathbf{E}\|_2^2, \quad (23)$$

where the constant c differs for different reconstruction algorithms. The best bounding constant for the BP method was given in [7], and equals

$$c_{bp} = \frac{4}{\sqrt{3 - 3\delta_{4K} - \sqrt{1 + \delta_{4K}}}},$$

while for the SP algorithm, the constant was estimated in [5]

$$c_{sp} = \frac{1 + \delta_{3K} + \delta_{3K}^2}{\delta_{3K}(1 - \delta_{3K})}.$$

A lower bound on the reconstruction distortion is given as follows. Suppose that the support set T of the sparse signal \mathbf{x} is perfectly reconstructed. The reconstructed signal $\hat{\mathbf{X}}$ is given by

$$\hat{\mathbf{X}} = (\Phi_T^* \Phi_T)^{-1} \Phi_T \hat{\mathbf{Y}},$$

and the reconstruction distortion is lower bounded by

$$\|\hat{\mathbf{X}} - \mathbf{X}\|_2^2 \geq \left(\frac{\sqrt{1 - \delta_K}}{1 + \delta_K} \right)^2 \|\hat{\mathbf{Y}} - \mathbf{Y}\|_2^2 = \frac{1 - \delta_K}{(1 + \delta_K)^2} \|\mathbf{E}\|_2^2. \quad (24)$$

For short, let

$$c_{lb} = \frac{\sqrt{1 - \delta_K}}{1 + \delta_K}.$$

Combining the bounds (23,24) and the results in Theorems 1-4, we summarize the asymptotic bounds on the reconstruction distortion as follows. Under Assumptions I, the reconstruction distortion of scalar quantization is bounded by

$$\begin{aligned} c_{lb}^2 \frac{\pi\sqrt{3}}{2} &\leq \lim_{R \rightarrow \infty} \lim_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E} \left[\|\hat{\mathbf{X}} - \mathbf{X}\|_2^2 \right] \\ &\leq \begin{cases} c_{sp}^2 \frac{\pi\sqrt{3}}{2} & \text{for subspace algorithm} \\ c_{bp}^2 \frac{\pi\sqrt{3}}{2} & \text{for basis pursuit algorithm} \end{cases}, \end{aligned}$$

and the reconstruction distortion of uniform scalar quantization is bounded by

$$\begin{aligned} c_{lb}^2 \frac{4 \log 2}{3} &\leq \lim_{R \rightarrow \infty} \lim_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{KR} \mathbb{E} \left[\|\hat{\mathbf{X}} - \mathbf{X}\|_2^2 \right] \\ &\leq \begin{cases} c_{sp}^2 \frac{4 \log 2}{3} & \text{for subspace algorithm} \\ c_{bp}^2 \frac{4 \log 2}{3} & \text{for basis pursuit algorithm} \end{cases}. \end{aligned}$$

Suppose that Assumption II holds. The reconstruction distortions for scalar quantization and uniform scalar quantization are

respectively bounded by

$$\begin{aligned} c_{lb}^2 \frac{\pi\sqrt{3}}{2} \mu_1 &\leq \liminf_{R \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right] \\ &\leq \limsup_{R \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right] \\ &\leq \begin{cases} c_{sp}^2 \frac{\pi\sqrt{3}}{2} \mu_2 & \text{for subspace algorithm} \\ c_{bp}^2 \frac{\pi\sqrt{3}}{2} \mu_2 & \text{for basis pursuit algorithm} \end{cases} \end{aligned}$$

and

$$c_{lb}^2 \frac{4 \log 2}{3} \mu_1 \leq \liminf_{R \rightarrow \infty} \frac{2^{2R}}{KR} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right].$$

Given the encoding rate R per measurement, the reconstruction distortion of the optimal scalar quantizer is bounded as

$$\begin{aligned} c_{lb}^2 \frac{\pi e}{6} &\leq \liminf_{R \rightarrow \infty} \liminf_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right] \\ &\leq \limsup_{R \rightarrow \infty} \limsup_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right] \\ &\leq \begin{cases} c_{sp}^2 \frac{\pi e}{3} & \text{for subspace algorithm} \\ c_{bp}^2 \frac{\pi e}{3} & \text{for basis pursuit algorithm} \end{cases}. \end{aligned}$$

The bounds for reconstruction distortion associated with vector quantization are given by

$$\begin{aligned} &c_{lb}^2 (1 - \delta_K) (1 + o_K(1)) \\ &\leq \liminf_{R \rightarrow \infty} \frac{2^{2Rm/K}}{K} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right] \\ &\leq \limsup_{R \rightarrow \infty} \frac{2^{2R}}{m} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right] \\ &\leq \begin{cases} c_{sp}^2 (1 + \delta_K) (1 + o_m(1)) & \text{for subspace algorithm} \\ c_{bp}^2 (1 + \delta_K) (1 + o_m(1)) & \text{for basis pursuit algorithm} \end{cases}, \end{aligned}$$

and

$$\begin{aligned} &\limsup_{R \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E} \left[\left\| \hat{\mathbf{X}} - \mathbf{X} \right\|_2^2 \right] \\ &\leq \begin{cases} c_{sp}^2 \frac{\pi\sqrt{3}}{2} \mu_2 & \text{for subspace algorithm} \\ c_{bp}^2 \frac{\pi\sqrt{3}}{2} \mu_2 & \text{for basis pursuit algorithm} \end{cases}. \end{aligned}$$

It is worth noting that the upper bound (23) on reconstruction distortion may not be tight. Empirical experiments show that this upper bound often significantly over-estimates the reconstruction distortion [5], [7].

IV. RECONSTRUCTION ALGORITHMS FOR QUANTIZED CS

We present next modifications of BP and SP algorithms that take into account quantization effects.

To describe these algorithms, we find the following notation useful. Let $\hat{\mathbf{Y}}$ be the quantized measurement vector. Given a vector $\hat{\mathbf{Y}}$, the corresponding quantization region can be easily identified: the quantization region of vector quantization $\mathcal{R}_{\hat{\mathbf{Y}}}$ is defined in (7); that of scalar quantization is given by the Cartesian product of the quantization regions for each coordinate, i.e., $\mathcal{R}_{\hat{\mathbf{Y}}} = \prod_{i=1}^m \mathcal{R}_{\hat{Y}_i}$ where $\mathcal{R}_{\hat{Y}_i}$ is the quantization region of \hat{Y}_i .

Similar to the standard BP method, the reconstruction problem can be now casted as

$$\min \|\mathbf{x}\|_1 \quad \text{subject to } \Phi \mathbf{x} \in \mathcal{R}_{\hat{\mathbf{Y}}}. \quad (25)$$

It can be verified that $\mathcal{R}_{\hat{\mathbf{Y}}}$ is a closed convex set and therefore (25) is a convex optimization problem and can be efficiently solved by linear programming techniques.

In order to adapt the SP algorithm to the quantization scenario at hand, we describe first a geometric interpretation of the projection operation in the SP algorithm. Given $\mathbf{y} \in \mathbb{R}^m$ and $\Phi_T \in \mathbb{R}^{m \times |T|}$, suppose that Φ_T has full column rank, in other words, suppose that the columns of Φ_T are linearly independent. The projection operation in (3) is equivalent to the optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^{|T|}} \|\mathbf{y} - \Phi_T \mathbf{x}\|_2^2. \quad (26)$$

Let \mathbf{x}^* be the solution of the quadratic optimization problem (26). Then functions (3-5) are equivalent to $\text{proj}(\mathbf{y}, \Phi_T) = \Phi_T \mathbf{x}^*$, $\text{resid}(\mathbf{y}, \Phi_T) = \mathbf{y} - \Phi_T \mathbf{x}^*$ and $\text{pcoeff}(\mathbf{y}, \Phi_T) = \mathbf{x}^*$.

The modified SP algorithm is based on the above geometric interpretation. More precisely, we use the following definition.

Definition 2: For given $\Phi_T \in \mathbb{R}^{m \times |T|}$, $\hat{\mathbf{Y}}$ and $\mathcal{R}_{\hat{\mathbf{Y}}}$, define

$$\mathcal{Q} := \left\{ (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{|T|} \times \mathcal{R}_{\hat{\mathbf{Y}}} : \|\mathbf{y} - \Phi_T \mathbf{x}\|_2 \leq \|\mathbf{y}' - \Phi_T \mathbf{x}'\|_2 \quad \forall (\mathbf{x}', \mathbf{y}') \in \mathbb{R}^{|T|} \times \mathcal{R}_{\hat{\mathbf{Y}}} \right\}, \quad (27)$$

and

$$(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \arg \min_{(\mathbf{x}, \mathbf{y}) \in \mathcal{Q}} \|\mathbf{y} - \hat{\mathbf{Y}}\|_2. \quad (28)$$

It can be verified that the pair $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ is well defined. See Appendix C for details.

This definition is introduced to identify the best approximation for $\hat{\mathbf{Y}}$ among multiple points in $\mathcal{R}_{\hat{\mathbf{Y}}}$ that minimize $\|\mathbf{y} - \Phi_T \mathbf{x}\|_2$. Based on this definition, we replace the resid and pcoeff functions in Algorithm 1 with new functions

$$\text{resid}^{(q)}(\hat{\mathbf{Y}}, \Phi_T) := \tilde{\mathbf{y}} - \Phi_T \tilde{\mathbf{x}}$$

and

$$\text{pcoeff}^{(q)}(\hat{\mathbf{Y}}, \Phi_T) := \tilde{\mathbf{x}},$$

where the superscript (q) emphasizes that these definitions are for the quantized case. This gives the modified SP algorithm.

The advantage of the modified algorithms are verified by the simulation results presented in the next section.

V. EMPIRICAL RESULTS

We performed extensive computer simulations in order to compare the performance of different quantizers and different reconstruction algorithms empirically. The parameters used in our simulations are $m = 128$, $N = 256$ and $K = 6$. Given these parameters, we generated realizations of $m \times N$ sampling matrices from the i.i.d. standard Gaussian ensemble and normalize the columns to have unit l_2 -norm. We also selected a support set T of size $|T| = K$ uniformly at random, generated the entries supported by T from the standard i.i.d. Gaussian distribution and set all other entries to zero. We let the quantization rates vary from two to six bits. For each quantization rate, we used Lloyd's algorithm (Section II-B) to obtain a nonuniform quantizer and employed brute-force search to find the optimal uniform quantizer. To test different quantizers and reconstruction algorithms, we randomly generated Φ and \mathbf{x} independently a thousand times. For each realization, we calculated the measurements \mathbf{Y} , the quantized measurements $\hat{\mathbf{Y}}$ and the reconstructed signal $\hat{\mathbf{X}}$.

Fig. 1 compares uniform and nonuniform quantizers with respect to measurement distortion. Though the quantization rates in our experiments are relatively small, the simulation results are consistent with the asymptotic results in Theorem 1: nonuniform quantization is better than uniform quantization and the gain increases with the quantization rate. Fig. 2a compares the reconstruction distortion of the standard BP and SP algorithms. The comparison of the modified algorithms is given in Fig. 2. The modified algorithms reduce the reconstruction distortion significantly. When the quantization rate is six bits, the reconstruction distortion of the modified algorithms is roughly one tenth of that of the standard algorithms. Furthermore, for both the standard and modified algorithms, the reconstruction distortion given by SP algorithms is much smaller than that of BP methods. Note that the computational complexity of the SP algorithms is also smaller than that of the BP methods, which shows

clear advantages for using SP algorithms in conjunction with quantized CS data. An interesting phenomenon occurs for the case of the modified BP method: although nonuniform quantization gives smaller measurement distortion, the corresponding reconstruction distortion is actually slightly larger than that of uniform quantization. We do not have solid analytical arguments to completely explain this somewhat counter-intuitive fact.

APPENDIX

A. Proof of Theorem 1

Let $T = \{1 \leq j \leq N : X_j \neq 0\}$ be the support set of \mathbf{x} , i.e., $x_i \neq 0$ for all $i \in T$ and $x_j = 0$ for all $j \notin T$. It is easy to show that for all $1 \leq i \leq m$ and $T \subset \{1, \dots, N\}$ such that $|T| = K$,

$$\mathbb{E} \left[\sum_{j \in T} A_{i,j} X_j \right] = 0$$

and

$$\mathbb{E} \left[\left(\sum_{j \in T} A_{i,j} X_j \right)^2 \right] = K.$$

According to the Central Limit Theorem, the distribution of $\frac{1}{\sqrt{K}} \sum_{j \in T} A_{i,j} X_j$ converges weakly to the standard Gaussian distribution as $K \rightarrow \infty$. This can be verified by the facts that $A_{i,j} X_j$ s are independent and identically distributed, and that the moment generating function of $A_{i,j} X_j$ is well defined. As a result, the distribution of $\sqrt{\frac{m}{K}} Y_i$ converges weakly to the standard Gaussian distribution as $K, m, N \rightarrow \infty$.

We apply a scalar quantizer with 2^R levels to the random variable $\sqrt{\frac{m}{K}} Y_i$. In this case, one has

$$\begin{aligned} & \frac{1}{K} \mathbb{E} \left[\left\| \hat{\mathbf{Y}} - \mathbf{Y} \right\|_2^2 \right] \\ &= \frac{1}{m} \frac{m}{K} \mathbb{E} \left[\sum_{i=1}^m (\hat{Y}_i - Y_i)^2 \right] \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left[\left(\sqrt{\frac{m}{K}} \hat{Y}_i - \sqrt{\frac{m}{K}} Y_i \right)^2 \right] \\ &= \mathbb{E} \left[\left(\sqrt{\frac{m}{K}} \hat{Y}_i - \sqrt{\frac{m}{K}} Y_i \right)^2 \right], \end{aligned} \quad (29)$$

where the last line represents the distortion of quantizing $\sqrt{\frac{m}{K}} Y_i$. Note that the distortion-rate function for scalar quantization of a Gaussian random variable is given by

$$\lim_{R \rightarrow \infty} 2^{2R} D_g^*(R) = \frac{\pi\sqrt{3}}{2} \sigma^2, \quad (30)$$

where σ^2 is the variance of the underlying Gaussian source (see [18] for a detailed proof of this result). We then have

$$\lim_{R \rightarrow \infty} \lim_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{K} D^*(R) = \lim_{R \rightarrow \infty} 2^{2R} D_g^*(R) = \frac{\pi\sqrt{3}}{2},$$

which completes the proof of (13).

Consider a uniform quantizer with codebook \mathcal{C}_u , such that $|\mathcal{C}_u| = 2^R$, and apply the corresponding uniform quantizer to the random variable $\sqrt{\frac{m}{K}} Y_i$. It was shown in [19] that the distortion-rate function of uniform scalar quantizers of a Gaussian random variable equals

$$\lim_{R \rightarrow \infty} \frac{2^{2R}}{R} D_{u,g}^*(R) = \frac{4}{3} \sigma^2 \log 2. \quad (31)$$

It is clear that

$$\lim_{R \rightarrow \infty} \lim_{(K, m, N) \rightarrow \infty} \frac{2^{2R}}{KR} D_u^*(R) = \lim_{R \rightarrow \infty} \frac{2^{2R}}{R} D_{u,g}^*(R) = \frac{4}{3} \log 2,$$

This proves Theorem 1.

B. Proof of Theorems 2 and 4

For completeness, let us first briefly review the key results used for deriving the asymptotic distortion-rate function for CS vector quantization. Suppose the source $\mathbf{Y} \in \mathbb{R}^k$ has probability density function $f(\mathbf{y})$. Let $\mathcal{R} \subset \mathbb{R}^k$ be a quantization region and $\boldsymbol{\omega} \in \mathcal{C}$ be the corresponding quantization level. The corresponding normalized moment of inertia (NMI) is defined as

$$m(\mathcal{R}) = \frac{\frac{1}{k} \int_{\mathcal{R}} \|\mathbf{y} - \boldsymbol{\omega}\|_2^2 f(\mathbf{y}) d\mathbf{y}}{\left(\int_{\mathcal{R}} d\mathbf{y}\right)^{1+2/k}}.$$

The optimal NMI equals

$$m_k^* = \inf_{\mathcal{R} \subset \mathbb{R}^k} m(\mathcal{R}),$$

only depends on the number of dimensions: $m_k^* = c_k$ with $c_k = \frac{1}{12}$ when $k = 1$ and $c_k \rightarrow \frac{1}{2\pi e}$ when $k \rightarrow \infty$. Thus the distortion rate function satisfies

$$\lim_{R \rightarrow \infty} \frac{2^R}{k} D(R) = \int \frac{f(\mathbf{y})}{\lambda_k^{2/k}(\mathbf{y})} m_k^* d\mathbf{y}, \quad (32)$$

where R is the quantization rate per dimension, and $\lambda_k(\mathbf{y})$ denotes the point density function. In this case, the integral

$$\int_{\mathcal{M}} \lambda_k(\mathbf{y}) d\mathbf{y}$$

gives the fraction of quantization levels belonging to \mathcal{M} for all measurable sets $\mathcal{M} \subset \mathbb{R}^k$. For simplicity, we have assumed that $\lambda_k(\mathbf{y})$ is continuous on \mathbb{R}^k . For fixed m_k^* , the problem of designing an asymptotically optimal quantizer can be reduced to the problem of finding the point density function $\lambda_k^*(\mathbf{y})$ that minimizes (32). By Hölder's inequality, the optimal point density function is given by

$$\lambda_k^*(\mathbf{y}) = \frac{f^{k/(k+2)}(\mathbf{y})}{\int f^{k/(k+2)}(\mathbf{y}) \cdot d\mathbf{y}},$$

and the asymptotic distortion rate function is therefore

$$\lim_{R \rightarrow \infty} \frac{2^R}{k} D^*(R) = c_k \left(\int f^{k/(k+2)}(\mathbf{y}) \cdot d\mathbf{y} \right)^{\frac{k+2}{k}}. \quad (33)$$

If the source \mathbf{Y} is Gaussian distributed with covariance matrix $\boldsymbol{\Sigma} > 0$, then the asymptotic distortion rate function (33) can be explicitly evaluated as

$$\begin{aligned} \lim_{R \rightarrow \infty} \frac{2^R}{k} D^*(R) &= c_k |2\pi\boldsymbol{\Sigma}|^{\frac{1}{k}} \left(\frac{k+2}{k} \right)^{\frac{k+2}{2}} \\ &= |\boldsymbol{\Sigma}|^{\frac{1}{k}} (1 + o_K(1)), \end{aligned} \quad (34)$$

where $o_K(1) \rightarrow 0$ as $K \rightarrow \infty$, and the last equality follows from the fact that $c_k \rightarrow \frac{1}{2\pi e}$ and $\left(\frac{k+2}{2}\right)^{\frac{k+2}{2}} \rightarrow e$ as $k \rightarrow \infty$.

We present next the key results used for proving the upper bounds in (17) and (21).

Proposition 1: Let $\mathbf{Y}_0 \in \mathbb{R}^k$ be a Gaussian random vector with zero mean and covariance matrix $\boldsymbol{\Sigma}_0$. Let $\{\mathbf{q}_R(\cdot)\}$, where the subscript R denotes the quantization rate, be a sequence of quantizers designed to achieve the asymptotic distortion rate function for Gaussian source $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_1)$ with $\mathbf{0} < \boldsymbol{\Sigma}_1 \in \mathbb{R}^{k \times k}$. Apply $\mathbf{q}_R(\cdot)$ to \mathbf{Y}_0 . If $\boldsymbol{\Sigma}_0 < \boldsymbol{\Sigma}_1$, then

$$\begin{aligned} \lim_{R \rightarrow \infty} \frac{2^{2R}}{k} \mathbb{E}_{Y_0} \left[\|\mathbf{Y}_0 - \mathbf{q}_R(\mathbf{Y}_0)\|_2^2 \right] \\ \leq c_k (2\pi\boldsymbol{\Sigma}_1)^{\frac{1}{k}} \left(\frac{k+2}{k} \right)^{\frac{k+2}{2}}. \end{aligned} \quad (35)$$

Proof: First assume that $\mathbf{0} < \boldsymbol{\Sigma}_0$. Let $f_0(\mathbf{y})$ and $f_1(\mathbf{y})$ be the probability density functions for \mathbf{Y}_0 and \mathbf{Y}_1 , respectively.

Denote $E_{Y_0} \left[\|\mathbf{Y}_0 - \mathbf{q}_R(\mathbf{Y}_0)\|_2^2 \right]$ by $D(R)$. It is clear that

$$\begin{aligned}
& \lim_{R \rightarrow \infty} \frac{2^R}{k} D(R) \\
&= c_k \int \frac{f_0(\mathbf{y})}{\left(\lambda_{k,1}^*(\mathbf{y})\right)^{2/k}} d\mathbf{y} \\
&= c_k \int \frac{f_0(\mathbf{y})}{f_1^{2/(k+2)}(\mathbf{y})} d\mathbf{y} \cdot \left(\int f_1^{k/(k+2)}(\mathbf{y}) d\mathbf{y} \right)^{\frac{2}{k}}.
\end{aligned} \tag{36}$$

We upper bound the first integral as follows

$$\begin{aligned}
& \int \frac{f_0(\mathbf{y})}{f_1^{2/(k+2)}(\mathbf{y})} d\mathbf{y} \\
&= \frac{|2\pi \boldsymbol{\Sigma}_1|^{\frac{1}{k+2}}}{|2\pi \boldsymbol{\Sigma}_0|^{1/2}} \int \exp \left\{ -\frac{1}{2} \mathbf{y}^* \left(\boldsymbol{\Sigma}_0^{-1} - \frac{2}{k+2} \boldsymbol{\Sigma}_1^{-1} \right) \mathbf{y} \right\} d\mathbf{y} \\
&\stackrel{(a)}{=} \frac{|2\pi \boldsymbol{\Sigma}_1|^{\frac{1}{k+2}}}{|2\pi \boldsymbol{\Sigma}_0|^{1/2}} \frac{|2\pi \boldsymbol{\Sigma}_0|^{1/2}}{\left| \mathbf{I}_k - \frac{2}{k+2} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_1^{-1} \right|^{1/2}} \\
&\stackrel{(b)}{\leq} |2\pi \boldsymbol{\Sigma}_1|^{\frac{1}{k+2}} \left(\frac{k+2}{k} \right)^{\frac{k}{2}} \\
&= \int f_1^{\frac{k}{k+2}}(\mathbf{x}) d\mathbf{x},
\end{aligned} \tag{37}$$

where (a) holds because

$$\begin{aligned}
& \boldsymbol{\Sigma}_0^{-1} - \frac{2}{k+2} \boldsymbol{\Sigma}_1^{-1} \\
&= \boldsymbol{\Sigma}_0^{-1} \left(\mathbf{I}_k - \frac{2}{k+2} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_1^{-1} \right) \\
&= \left[\left(\mathbf{I}_k - \frac{2}{k+2} \boldsymbol{\Sigma}_0 \boldsymbol{\Sigma}_1^{-1} \right)^{-1} \boldsymbol{\Sigma}_0 \right]^{-1},
\end{aligned}$$

and (b) follows from the assumption $\boldsymbol{\Sigma}_0 < \boldsymbol{\Sigma}_1$. Substituting (37) into (36), one obtains

$$\begin{aligned}
& \lim_{R \rightarrow \infty} \frac{2^R}{k} D(R) \\
&\leq c_k \left(\int f_1^{k/(k+2)}(\mathbf{y}) d\mathbf{y} \right)^{\frac{k+2}{k}} \\
&= c_k |2\pi \boldsymbol{\Sigma}_1|^{\frac{1}{k}} \left(\frac{k+2}{k} \right)^{\frac{k+2}{2}},
\end{aligned}$$

which will be used to prove the upper bounds in (17) and (21).

Suppose that $|\boldsymbol{\Sigma}_0| = 0$ (some of the eigenvalues of $\boldsymbol{\Sigma}_0$ are zero). Since $\boldsymbol{\Sigma}_0 < \boldsymbol{\Sigma}_1$, when $\epsilon > 0$ is sufficiently small, we have $\mathbf{0} < \boldsymbol{\Sigma}_\epsilon := \boldsymbol{\Sigma}_0 + \epsilon \mathbf{I} < \boldsymbol{\Sigma}_1$. Let $f_\epsilon(\mathbf{y})$ be the probability density function of Gaussian vector with zero mean and variance $\boldsymbol{\Sigma}_\epsilon$.

Then,

$$\begin{aligned}
& \lim_{R \rightarrow \infty} \frac{2^R}{k} D(R) \\
&= c_k \int \frac{f_0(\mathbf{y})}{\left(\lambda_{k,1}^*(\mathbf{y})\right)^{2/k}} d\mathbf{y} \\
&= c_k \int \frac{\lim_{\epsilon \rightarrow 0} f_\epsilon(\mathbf{y})}{\left(\lambda_{k,1}^*(\mathbf{y})\right)^{2/k}} d\mathbf{y} \\
&\stackrel{(c)}{\leq} c_k \liminf_{\epsilon \rightarrow 0} \int \frac{f_\epsilon(\mathbf{y})}{\left(\lambda_{k,1}^*(\mathbf{y})\right)^{2/k}} d\mathbf{y} \\
&\stackrel{(d)}{\leq} c_k |2\pi \Sigma_1|^{\frac{1}{k}} \left(\frac{k+2}{k}\right)^{\frac{k+2}{2}},
\end{aligned}$$

where (c) follows from Fatou's lemma [20], and (d) follows from the first part of this proof. This proves the proposition. ■

1) *Lower Bounds for Scalar Quantization:*

We prove the lower bound in (17). Given Assumptions II, each Y_i , $1 \leq i \leq m$, is a linear combination of Gaussian random variables, and therefore each Y_i is a Gaussian random variable itself. For a given i and a given T , the mean and the variance of Y_i are $E[Y_i] = 0$ and $\sigma_{i,T}^2 = E[Y_i^2] = \sum_{j \in T} \varphi_{i,j}^2$, respectively. The variance depends on the row index i and the support set T . We calculate the average variance across all rows and all support sets as

$$\begin{aligned}
\bar{\sigma}^2 &= \frac{1}{m} \sum_{i=1}^m \left(\frac{1}{\binom{N}{K}} \sum_T \sum_{j \in T} \varphi_{i,j}^2 \right) \\
&= \frac{1}{m} \frac{1}{\binom{N}{K}} \sum_T \sum_{j \in T} \left(\sum_{i=1}^m \varphi_{i,j}^2 \right) \\
&\stackrel{(a)}{=} \frac{1}{m} \frac{1}{\binom{N}{K}} \sum_{j=1}^N \left(\sum_{T: j \in T} \|\varphi_j\|_2^2 \right) \\
&\stackrel{(b)}{=} \frac{1}{m} \frac{1}{\binom{N}{K}} \sum_{j=1}^N \binom{N-1}{K-1} \|\varphi_j\|_2^2 \\
&\stackrel{(c)}{=} \frac{K}{m} \frac{1}{N} \sum_{j=1}^N \|\varphi_j\|_2^2 \\
&\stackrel{(d)}{=} \frac{K}{m} \mu_1,
\end{aligned} \tag{38}$$

where

- (a) is obtained by exchanging the sums over T and j ,
- (b) holds because for any given $1 \leq j \leq N$, there are $\binom{N-1}{K-1}$ many subsets T containing the index j ,
- (c) is due to the fact that $\binom{N-1}{K-1} / \binom{N}{K} = K/N$,
- (d) follows from the definition (15).

Suppose that one deals with the ideal case: the support set T is known before taking the measurements; and for different values of i and T , we are allowed to use different quantizers. Given i and T , we apply the optimal quantizer for the Gaussian random variable $\sqrt{\frac{m}{K}} Y_i$, so that the quantization distortion of Y_i satisfies

$$\lim_{R \rightarrow \infty} 2^{2R} D_{i,T}^*(R) = \frac{\pi \left(\frac{m}{K} \sigma_{i,T}^2\right)}{2} \sqrt{3},$$

which is a direct application of (33) with $k = 1$. Taking the average over all i and all T gives

$$\begin{aligned}
& \lim_{R \rightarrow \infty} \frac{1}{m} \sum_{i=1}^m \mathbb{E}_T [2^{2R} D_{i,T}^*(R)] \\
&= \frac{1}{m} \sum_{i=1}^m \frac{1}{\binom{T}{K}} \sum_T \left(\lim_{R \rightarrow \infty} 2^{2R} D_{i,T}^*(R) \right) \\
&= \frac{1}{m} \frac{1}{\binom{T}{K}} \sum_{i=1}^m \sum_T \left(\frac{\pi \left(\frac{m}{K} \sigma_{i,T}^2 \right)}{2} \sqrt{3} \right) \\
&= \frac{\pi \mu_1}{2} \sqrt{3},
\end{aligned}$$

where the last equality follows from (38).

However, the support set T is unknown before taking the measurements. Furthermore, the same quantizer has to be employed for different choices of i and T . Thus, for every R , i and T , $\mathbb{E}_{Y_i} \left[\frac{m}{K} |Y_i - \hat{Y}_i|^2 \right] \geq D_{i,T}^*(R)$. As a result,

$$\begin{aligned}
& \liminf_{R \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E}_T \left[\mathbb{E}_Y \left[\left\| \hat{\mathbf{Y}} - \mathbf{Y} \right\|_2^2 \right] \right] \\
&= \liminf_{R \rightarrow \infty} \frac{2^{2R}}{\binom{N}{T}} \sum_T \frac{1}{m} \sum_{i=1}^m \frac{m}{K} \mathbb{E}_Y \left[(\hat{y}_i - y_i)^2 \right] \\
&\geq \liminf_{R \rightarrow \infty} \frac{2^{2R}}{\binom{N}{T}} \sum_T \frac{1}{m} \sum_{i=1}^m D_{i,T}^*(R) \\
&= \frac{\pi \mu_1}{2} \sqrt{3}.
\end{aligned}$$

Since the above derivation is valid for all K , m and N , the claim in (17) holds.

The result in (18) for uniform quantizers can be proved using similar arguments. For the ideal case, given i and T , apply the optimal *uniform* quantizer for the standard Gaussian random variable to $\sqrt{\frac{m}{K}} y_i$. The corresponding distortion rate function for this case was characterized in [19] and s given by

$$\lim_{R \rightarrow \infty} 2^{2R} D_{u,i,T}^*(R) = \frac{4}{3} \sigma_{i,T}^2 \ln 2.$$

Therefore,

$$\begin{aligned}
& \liminf_{R \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E}_T \left[\mathbb{E}_Y \left[\left\| \hat{\mathbf{Y}} - \mathbf{Y} \right\|_2^2 \right] \right] \\
&\geq \frac{4}{3} \mu_1 \ln 2,
\end{aligned}$$

which completes the proof of (18).

2) The Upper Bound for Scalar Quantization:

By the definition of μ_2 in (16), the variance of the Gaussian random variable $\sqrt{\frac{m}{K}} Y_i$ is upper bounded by μ_2 uniformly for all i and all T . For each quantization rate R , we design the optimal quantizer for a Gaussian source with variance μ_2 and apply this quantizer to quantize all components of \mathbf{Y} . Using (35), one can show that the quantization distortion for all i and T satisfies

$$\begin{aligned}
& \limsup_{R \rightarrow \infty} \frac{2^{2R}}{K} \mathbb{E}_T \mathbb{E}_Y \left[\left\| \hat{\mathbf{Y}} - \mathbf{Y} \right\|_2^2 \right] \\
&\leq \frac{\pi}{2} \mu_2 \sqrt{3},
\end{aligned}$$

which proves the upper bound in (17).

3) The Lower Bound for Vector Quantization:

The basic idea for proving the lower bound in (20) is similar to that behind (17). For each T , a lower bound on the minimum achievable distortion is derived. The average distortion taken over all the sets T serves as a lower bound of the

overall distortion-rate function.

Suppose the ideal case where we have prior knowledge of $T \in \binom{[N]}{K}$. We study the distortion rate function for every given T . The measurement vector \mathbf{Y} is Gaussian distributed with zero mean and covariance matrix $\Phi_T \Phi_T^*$, where Φ_T consists of the columns of Φ indexed by T . The singular value decomposition of $\Phi_T \Phi_T^*$ gives $\mathbf{U}_T \Lambda_T \mathbf{U}_T^*$, where $\mathbf{U}_T \in \mathbb{R}^{m \times m}$ has orthonormal columns and $\Lambda_T = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$ is the diagonal matrix formed by the singular values $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$. Note that $\lambda_i(\Phi_T^* \Phi_T) = \lambda_i(\Phi_T \Phi_T^*)$ for $1 \leq i \leq K$. According to Assumption III.1, the measurement matrix Φ satisfies the RIP with constant parameter δ_K , which implies that $1 - \delta_K \leq \lambda_i(\Phi_T^* \Phi_T) \leq 1 + \delta_K$ for all $1 \leq i \leq K$. It can be concluded that $1 - \delta_K \leq \lambda_i \leq 1 + \delta_K$ for $1 \leq i \leq K$ and $\lambda_i = 0$ for $K + 1 \leq i \leq m$. As a result, $\Phi_T \Phi_T^* = \mathbf{U}_{T,K} \Lambda_{T,K} \mathbf{U}_{T,K}^*$ where $\mathbf{U}_{T,K} \in \mathbb{R}^{m \times K}$ contains the first K columns of \mathbf{U}_T and $\Lambda_{T,K} \in \mathbb{R}^{K \times K}$ is the diagonal matrix formed by the K largest singular values. Denote the matrix formed by the last $m - K$ columns of \mathbf{U} by $\mathbf{U}_{T,K}^\perp$: clearly, $\mathbf{U}_T = [\mathbf{U}_{T,K} | \mathbf{U}_{T,K}^\perp]$.

The best quantization strategy is to quantize $\tilde{\mathbf{Y}} = \mathbf{U}_{T,K}^* \mathbf{Y}$ so that no quantization bit is used for the ‘‘trivial signal’’ $(\mathbf{U}_{T,K}^\perp)^* \mathbf{Y}$. It is clear that $\tilde{\mathbf{Y}} \sim \mathcal{N}(\mathbf{0}, \Lambda_{T,K})$ and $\mathbf{0} < \Lambda_{T,K}$. The corresponding asymptotic distortion rate function is therefore

$$\begin{aligned} \lim_{R \rightarrow \infty} \frac{2^{2mR/K}}{K} D_T^*(R) &\stackrel{(34)}{=} c_K (2\pi \Lambda_{T,K})^{\frac{1}{K}} \left(\frac{K+2}{K} \right)^{\frac{K+2}{2}} \\ &\geq (1 - \delta_K) (1 + o_K(1)), \end{aligned}$$

where the $2^{2mR/K}$ term comes from the fact that the total quantization rate mR is used to quantize a K -dimensional signal. Since this lower bound is valid for all $T \in \binom{[N]}{K}$, we have proved the lower bound in (20).

4) The Upper Bound for Vector Quantization:

Let $\epsilon > 0$ be a small constant. Let $\{q_R(\cdot)\}$ be a sequence of quantizers that approaches the asymptotic distortion rate function for quantizing $\tilde{\mathbf{Y}} \sim \mathcal{N}(\mathbf{0}, (1 + \delta_K + \epsilon) \mathbf{I}_m)$. To prove the upper bound in (21), apply the quantizer sequence $\{q_R(\cdot)\}$ to \mathbf{Y} . For every $T \in \binom{[N]}{K}$, $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \Phi_T \Phi_T^*)$. According to the Assumption III.1, $\Phi_T \Phi_T^* < (1 + \delta_K + \epsilon) \mathbf{I}_m$. Applying Proposition 1, we have

$$\begin{aligned} \lim_{R \rightarrow \infty} \frac{2^{2R}}{m} \mathbb{E}_{\mathbf{Y}} \left[\|\mathbf{Y} - q_R(\mathbf{Y})\|_2^2 \right] \\ \leq (1 + \delta_K + \epsilon) (1 + o_M(1)). \end{aligned}$$

The upper bound in (21) is proved by taking the limit $\epsilon \downarrow 0$.

C. The Existence and Uniqueness of $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ in Equation (28)

Consider the optimization problem

$$\min_{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{|T|} \times \mathcal{R}_{\tilde{\mathbf{y}}}} \|\mathbf{y} - \Phi_T \mathbf{x}\|_2, \quad (39)$$

which is equivalent to

$$\min_{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{|T|} \times \mathcal{R}_{\tilde{\mathbf{y}}}} \left\| \begin{bmatrix} -\Phi_T & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \right\|_2^2. \quad (40)$$

Note that the objective function is convex and the constraint set is convex and closed. The optimization problem (40) has at least one solution. Note that the matrix $[-\Phi_T \ \mathbf{I}]$ does not have full row-rank. Hence, the solution may not be unique: the set \mathcal{Q} defined in (27) gives all the possible solutions, and is convex and closed.

Let \mathfrak{P} be the projection function from $\mathbb{R}^{|T|} \times \mathbb{R}^m$ to \mathbb{R}^m , i.e., $\mathfrak{P}((\mathbf{x}, \mathbf{y})) = \mathbf{y}$. Since the set \mathcal{Q} is convex, the set $\mathfrak{P}(\mathcal{Q})$ is also convex. The quadratic optimization problem

$$\min_{\mathbf{y} \in \mathfrak{P}(\mathcal{Q})} \|\hat{\mathbf{Y}} - \mathbf{y}\|_2$$

has a unique solution. Denote this unique solution by $\tilde{\mathbf{y}}$. Furthermore, recall our assumption that Φ_T has full column rank.

For any given $\mathbf{y} \in \mathbb{R}^m$, the solution of

$$\min_{\mathbf{x} \in \mathbb{R}^{|T|}} \|\mathbf{y} - \Phi_T \mathbf{x}\|_2$$

is therefore unique. As a result, there exists a unique $\tilde{\mathbf{x}} \in \mathbb{R}^{|T|}$ such that $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \in \mathcal{Q}$. This establishes the existence and uniqueness of the point $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$.

REFERENCES

- [1] D. Donoho, "Compressed sensing," *IEEE Trans. Inform. Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [2] E. Candès and T. Tao, "Decoding by linear programming," *Information Theory, IEEE Transactions on*, vol. 51, no. 12, pp. 4203–4215, 2005.
- [3] E. Candès, M. Rudelson, T. Tao, and R. Vershynin, "Error correction via linear programming," in *IEEE Symposium on Foundations of Computer Science (FOCS)*, pp. 295 – 308, 2005.
- [4] E. J. Candès and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?," *IEEE Trans. Inform. Theory*, vol. 52, no. 12, pp. 5406–5425, 2006.
- [5] W. Dai and O. Milenkovic, "Subspace pursuit for compressive sensing signal reconstruction," *IEEE Trans. Inform. Theory*, accepted, 2008.
- [6] D. Needell and J. A. Tropp, "CoSaMP: Iterative signal recovery from incomplete and inaccurate samples," *Appl. Comp. Harmonic Anal.*, accepted, 2008.
- [7] E. J. Candès, J. K. Romberg, and T. Tao, "Stable signal recovery from incomplete and inaccurate measurements," *Comm. Pure Appl. Math.*, vol. 59, no. 8, pp. 1207–1223, 2006.
- [8] E. Candès and J. Romberg, "Encoding the ℓ_p ball from limited measurements," *Data Compression Conference*, pp. 33–42, March 2006.
- [9] P. Boufounos and R. Baraniuk, "Quantization of sparse representations," *Preprint*, 2008.
- [10] P. Boufounos and R. G. Baraniuk, "1-bit compressive sensing," in *Conf. on Info. Sciences and Systems (CISS)*, (Princeton, NJ), pp. 16–21, March 2008.
- [11] V. Goyal, A. Fletcher, and S. Rangan, "Compressive sampling and lossy compression," *IEEE Signal Processing Magazine*, vol. 25, pp. 48–56, March 2008.
- [12] I. E. Nesterov, A. Nemirovskii, and Y. Nesterov, *Interior-Point Polynomial Algorithms in Convex Programming*. SIAM, 1994.
- [13] S. Lloyd, "Least squares quantization in pcm," *Information Theory, IEEE Transactions on*, vol. 28, pp. 129–137, Mar 1982.
- [14] K. Sayood, *Introduction to Data Compression*. Morgan Kaufmann, 3rd edition ed., 2005.
- [15] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. John Wiley & Sons, 1st edition ed., 1991.
- [16] M. Akcakaya and V. Tarokh, "On sparsity, redundancy and quality of frame representations," pp. 951–955, June 2007.
- [17] E. Ardestanizadeh, M. Cheraghchi, and A. Shokrollahi, "Bit precision analysis for compressed sensing," *Preprint*, 2009.
- [18] P. Zador, *Development and evaluation of procedures for quantizing multivariate distributions*. PhD thesis, Stanford University, Stanford, CA, 1964.
- [19] D. Hui and D. Neuhoff, "Asymptotic analysis of optimal fixed-rate uniform scalar quantization," *IEEE Trans. Inform. Theory*, vol. 47, pp. 957–977, Mar 2001.
- [20] H. Royden, *Real Analysis*. Prentice Hall, 3 edition ed., 1988.

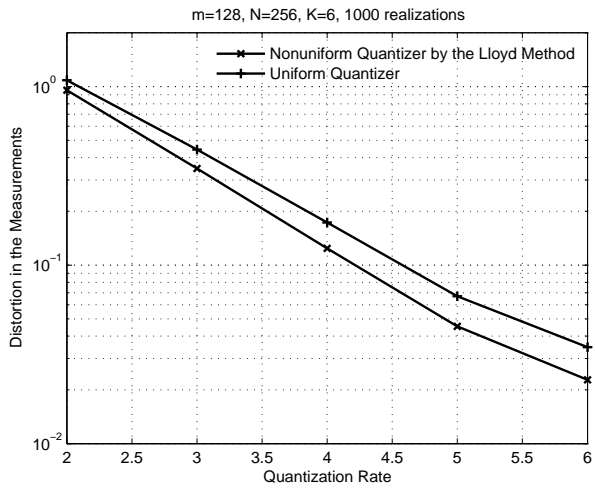
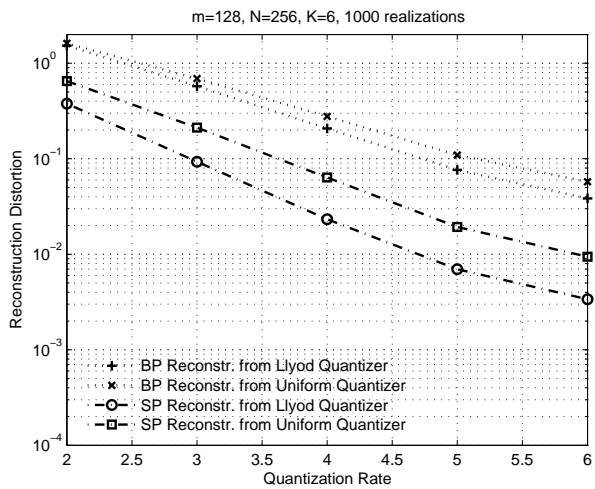
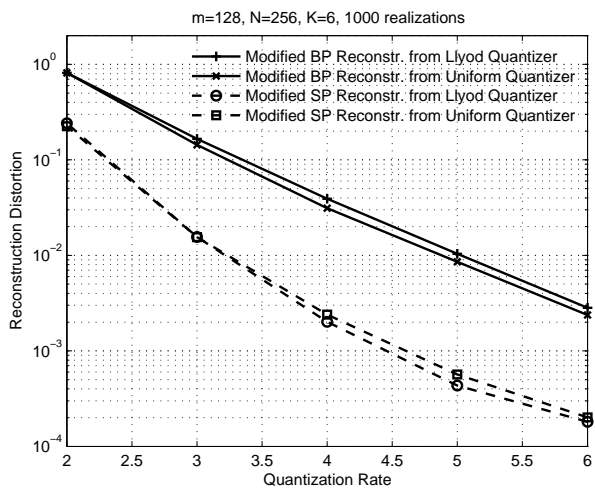


Figure 1: Distortion in the measurements.



(a) By standard reconstruction algorithms



(b) By modified reconstruction algorithms

Figure 2: Distortion in the reconstruction signals.