

Minimum Complexity Pursuit

Shirin Jalali and Arian Maleki

Abstract—The fast growing field of compressed sensing is founded on the fact that if a signal is simple and has some ‘structure’, then it can be reconstructed accurately with far fewer samples than its ambient dimension. Many different plausible structures have been explored in this field, ranging from sparsity to low-rankness and to finite rate of innovation. However, there are important abstract questions that are yet to be answered. For instance, what are the general abstract meanings of structure and simplicity? Does there exist universal algorithms for recovering such simple structured objects from fewer samples than their ambient dimension? In this paper, we aim to address these two questions. Using algorithmic information theory tools such as Kolmogorov complexity, we provide a unified method of describing simplicity and structure. We then explore the performance of an algorithm motivated by Occams Razor (called MCP for minimum complexity pursuit) and show that it requires $O(k \log n)$ number of samples to recover a signal, where k and n represent its complexity and ambient dimension, respectively. Finally, we discuss more general classes of signals and provide guarantees on the performance of MCP.

I. INTRODUCTION

Compressed sensing (CS) refers to a body of techniques that undersample high-dimensional signals, and yet recover them accurately by exploiting their intrinsic ‘structure’ [1], [2]. This permits more efficient sensing systems that are proved to be valuable in many applications including magnetic resonance imaging (MRI) [3] and radar [4], to name a few. Some of the ‘structures’ that have been considered in the literature are as follows.

- i. Sparsity: A vector $x \in \mathbb{R}^n$ is called k -sparse if and only if $\|x\|_0 \triangleq \sum_{i=1}^n \mathbf{1}_{\{x_i \neq 0\}} \leq k$. Roughly speaking, according to compressed sensing a k -sparse signal x can be recovered from $d = O(k \log n)$ random linear measurements $y = Ax$.
- ii. Low rankness: If $X \in \mathbb{R}^{m \times n}$ is a low rank matrix with $\text{rank}(X) \leq k$, then $d = O(r(m+n) \log(mn))$ random linear measurements are sufficient for recovering X from its measurements accurately with high probability [5].
- iii. Model-based compressed sensing: [6] considers more structured signal models by assuming that from $\binom{n}{k}$ subspaces of k -sparse signals only m_k of them may occur. It is then proved that $O(\log(m_K))$ random linear measurements are sufficient for the accurate recovery of such signals. This class is a superset of some of the

other structures introduced in the literature such as the class of block-sparse signals [7]–[10].

- iv. Rate of innovation: [11] defines the rate of innovation of a signal as its “degrees of freedom”. Several important classes of functions such as the piecewise polynomial functions and sparse signals have clearly finite rate innovation. [11] suggests sampling schemes for several classes that recover the signal from $O(k)$ number of measurements, where k is the rate of innovation.

The above results seem to provide pieces of a bigger picture. Recently, [12] introduced the class of simple functions and atomic norm as a framework that unifies some of the above observations and extends them to some other signal classes. However, there is still an interesting conceptual question that needs to be addressed, i.e., what is the abstract meaning of ‘structure’ that allows fewer measurements than the ambient dimension of the signal? Given a simple signal, which scheme recovers the signal from an undersampled random linear set of measurements?

In the context of algorithmic information theory, Solomonoff [13] and Kolmogorov [14] suggested a universal notion of complexity for binary sequences, known as the Kolmogorov complexity. Given a binary sequence x , its Kolmogorov complexity $K(x)$ is defined as the length of the shortest computer program that prints x . In this paper, we extend the concept of Kolmogorov complexity to the real signals. Such extensions are straightforward and have been explored before [15]. Based on this notion of complexity, called Kolmogorov complexity of real signals, we show that Occams razor [16], i.e., finding the ‘simplest’ solution of the linear equations, correctly recovers the signal with much fewer measurements than the ambient dimension of the signal. Roughly speaking, we prove that the number of linear measurements required for recovering the correct solution is proportional to the complexity rather than the ambient dimension of the signal. We postpone the accurate exposition of our results to Section IV. We will further discuss the issue of model mismatch in the signal classes and will prove that the approach motivated by Occams razor is stable with respect to such non-idealities in the system.

Here is the organization of our paper. Section II defines the notation used throughout the paper. Section III defines Kolmogorov complexity of a real-valued signal. Section IV outlines our contribution. Section V calculates the Kolmogorov complexity of several classes that are popular in compressed sensing and clarifies the statements of our theorems on these classes. Section VI compares our work with other results in

S. Jalali is a postdoctoral scholar at the Center for Mathematics of Information, California Institute of Technology, Pasadena, CA, shirin@caltech.edu

A. Maleki is a postdoctoral scholar at Digital Signal Processing group, Rice University, Houston, TX, arian.maleki@rice.edu

the literature. Sections VII and VIII are devoted to the proofs of our main theorems.

II. DEFINITIONS

Calligraphic letters such as \mathcal{A} and \mathcal{B} denote sets. For a set \mathcal{A} , $|\mathcal{A}|$ and \mathcal{A}^c denote its size and its complement, respectively. For a sample space Ω and event set $\mathcal{A} \subseteq \Omega$, $\mathbf{1}_{\mathcal{A}}$ denotes the indicator function of the event \mathcal{A} .

Let $\{0, 1\}^*$ denote the set of all finite-length binary sequences, i.e., $\{0, 1\}^* \triangleq \cup_{n \geq 1} \{0, 1\}^n$. For a vector $x \in \mathbb{R}^n$, the ℓ_p norm of x is defined as $\|x\|_p \triangleq (\sum_{i=1}^n |x_i|^p)^{1/p}$. The ℓ_∞ norm of x is denoted by $\|x\|_\infty \triangleq \max_i |x_i|$.

For a real number $x \in [0, 1]$, let $[x]_m$ denote the m -bit approximation of x that results from taking the first m bits in the binary expansion of x . In other words, if $x = \sum_{i=1}^{\infty} 2^{-i}(x)_i$, where $(x)_i \in \{0, 1\}$ denotes the i^{th} bit in the binary expansion of x , then

$$[x]_m \triangleq \sum_{i=1}^m 2^{-i} x_i. \quad (1)$$

Similarly, for a vector $x^n \in [0, 1]^n$, define

$$[x^n]_m \triangleq ([x_1]_m, \dots, [x_n]_m). \quad (2)$$

For an integer $n \in \mathbb{N}$, let

$$\log^* n \triangleq \lceil \log_2 n \rceil + 2 \log_2 \max(\lceil \log_2 n \rceil, 1).$$

III. KOLMOGOROV COMPLEXITY

The Kolmogorov complexity of a finite-length sequence x with respect to a *universal computer* \mathcal{U} is defined as the minimum length over all programs that print x and halt.¹ For a universal computer \mathcal{U} and any computer \mathcal{A} , there exists a constant $c_{\mathcal{A}}$ such that $K_{\mathcal{U}}(x) \leq K_{\mathcal{A}}(x) + c_{\mathcal{A}}$, for all strings $x \in \{0, 1\}^*$ [17]. Hence, as suggested in [17], we drop the subscript \mathcal{U} , and let $K(x)$ denote the Kolmogorov complexity of the binary string x .

Similarly, the Kolmogorov complexity of an integer $n \in \mathbb{N}$, $K(n)$, is defined as the Kolmogorov complexity of its binary representation. It can be proved that

$$K(n) \leq \log^* n + c,$$

where c is a constant independent of n .

For $x = (x_1, x_2, \dots, x_n) \in [0, 1]^n$, define the Kolmogorov complexity of x at resolution m as

$$K^{[\cdot]m}(x) = K([x_1]_m, [x_2]_m, \dots, [x_n]_m). \quad (3)$$

Lemma 1: For $(x_1, x_2, \dots, x_n) \in [0, 1]^n$,

$$\limsup_{m \rightarrow \infty} \frac{K^{[\cdot]m}(x_1, x_2, \dots, x_n)}{m} \leq n.$$

The proof is very simple and is skipped.

Definition 1: The signal $x = (x_1, x_2, \dots, x_n)$ is called incompressible if and only if

$$\lim_{m \rightarrow \infty} \frac{K^{[\cdot]m}(x_1, x_2, \dots, x_n)}{m} = n.$$

¹Refer to Chapter 14 of [17] for the exact definition of a universal computer, and more details on the definition of the Kolmogorov complexity.

Proposition 1: Let $\{X_i\}_{i=1}^n \stackrel{iid}{\sim} U[0, 1]$. Then,

$$\frac{1}{m} K^{[\cdot]m}(X_1, X_2, \dots, X_n) \rightarrow n$$

in probability.

Proof: If $X_i = \sum_{j=1}^{\infty} (X_i)_j 2^{-j}$, where $(X_i)_j \in \{0, 1\}$, then $\{(X_i)_j\}_{j=1}^{\infty} \stackrel{iid}{\sim} \text{Bern}(1/2)$. Theorem 14.5.3 in [17] states that the normalized Kolmogorov's complexity of $([X_1]_m, \dots, [X_n]_m) = \{((X_i)_1, (X_i)_2, \dots, (X_i)_m)\}_{i=1}^n$, i.e.,

$$\frac{K(\{(X_i)_1, (X_i)_2, \dots, (X_i)_m\}_{i=1}^n | mn)}{mn} \rightarrow 1, \quad (4)$$

in probability. On the other hand,

$$\begin{aligned} & K(\{(X_i)_1, (X_i)_2, \dots, (X_i)_m\}_{i=1}^n | mn) \\ & \leq K(\{(X_i)_1, (X_i)_2, \dots, (X_i)_m\}_{i=1}^n) \\ & \leq K(\{(X_i)_1, (X_i)_2, \dots, (X_i)_m\}_{i=1}^n | mn) + \log^*(mn) + c, \end{aligned} \quad (5)$$

where c is a constant [17]. Hence, combining (4) and (5) proves the desired result. \blacksquare

IV. OUR CONTRIBUTION

Consider the problem of reconstructing a vector $x_o \in \mathbb{R}^n$ from d random linear measurements $y = Ax$ with $d < n$. We say a recovery algorithm is successful if as n grows the ℓ_2 -error between x_o and its reconstruction \hat{x}_o goes to zero, i.e., we want

$$\mathbb{P}(\|x_o^n - \hat{x}_o^n\|_2^2 > \epsilon) \rightarrow 0,$$

for any $\epsilon > 0$. Assuming that the signal is ‘structured’ in the sense that will be clarified later, we follow Ocam’s Razor and seek the simplest solution of $y = Ax$, i.e.,

$$\begin{aligned} & \arg \min \quad K^{[\cdot]m}(x_1, \dots, x_n) \\ & \text{s.t.} \quad Ax^n = y_o^n. \end{aligned} \quad (6)$$

We call this algorithm minimum complexity pursuit or MCP. The choice of m will be clarified later as well. Suppose that $A \in \mathbb{R}^{d \times n}$, where A_{ij} are iid $\mathcal{N}(0, 1/d)$, and assume that $y_o^n = Ax_o^n$. Let $\hat{x}_o^n = \hat{x}_o^n(y_o^n, A)$ denote the output of (6) to the inputs y_o^n and A .

Theorem 1: Assume that $x_o = (x_{o,1}, x_{o,2}, \dots) \in [0, 1]^\infty$ is such that

$$\limsup_{n \rightarrow \infty} \frac{K^{[\cdot]m}(x_{o,1}, x_{o,2}, \dots, x_{o,n})}{m} \leq \kappa, \quad (7)$$

where $m = m_n = \lceil \log n \rceil$. Let $d = d_n = \lceil \kappa \log n \rceil$. Then, for any $\epsilon > 0$

$$\mathbb{P}(\|x_o^n - \hat{x}_o^n\|_2^2 > \epsilon) \rightarrow 0, \quad (8)$$

as n grows without bound.

This theorem indicates that when the Kolmogorov complexity of the signal is less than κ , then $O(\kappa \log n)$ linear measurements are sufficient for the successful recovery. Also, it provides an evidence for the success of Ocam’s Razor.

Although Theorem 1 is an asymptotic theorem, its proof provides information on the performance of MCP on finite length sequences as well.

Corollary 1: Assume that $x_o = (x_{o,1}, x_{o,2}, \dots, x_{o,n}) \in [0, 1]^n$ is such that

$$\frac{K^{[\cdot]m}(x_{o,1}, x_{o,2}, \dots, x_{o,n})}{m} \leq \kappa, \quad \forall m.$$

Let $m = m_n = \lceil \alpha \log n \rceil$ and $d = d_n = \lceil 2\alpha\kappa \log n \rceil$. Then, with probability $1 - n^{-\alpha\kappa}$

$$\|x_o^n - \hat{x}_o^n\|_2 \leq \frac{10n^{1/2-\alpha}}{\sqrt{\kappa} \log n}.$$

Now consider the following more general setting, where the original signal x_o^n to be recovered is not low-complexity, but is close to a low-complexity signal \tilde{x}^n , i.e., $\|x_o^n - \tilde{x}^n\|_2 \leq \epsilon_n$ with $\epsilon_n = o(1)$. Again, let $y_o^n = Ax_o^n$, and consider the following reconstruction algorithm for finding x_o^n from its linear measurements y_o^n :

$$\begin{aligned} \min \quad & K^{[\cdot]m}(x_1, \dots, x_n) \\ \text{s.t.} \quad & \|Ax^n - y_o^n\|_2 \leq \sigma_{\max}(A)\epsilon_n. \end{aligned}$$

Assume that $A \in \mathbb{R}^{d \times n}$ and A_{ij} are iid $\mathcal{N}(0, \frac{1}{d})$. Let $\hat{x}_o^n = \hat{x}_o^n(y_o^n, A)$.

Theorem 2: Assume that there exists \tilde{x}_o^n such that $\|x_o^n - \tilde{x}_o^n\|_2 \leq \epsilon_n$, and

$$\limsup_{m \rightarrow \infty} \frac{K^{[\cdot]m}(\tilde{x}_o^n)}{m} \leq \kappa_n. \quad (9)$$

Let $m = m_n = \lceil \log n \rceil$ and $d = d_n = \lceil \kappa_n \log n \rceil$. If $\epsilon_n = o(d_n/n)$, then for each $\epsilon > 0$,

$$\mathbb{P}(\|x_o^n - \hat{x}_o^n\|_2^2 > \epsilon) \rightarrow 0, \quad (10)$$

as n grows without bound.

In the next section we show that several popular classes of sequences studied in CS such as class of sparse signals and samples of piecewise smooth functions can be considered as special cases of the framework we introduced in this section and that Theorems 1 and 2 provide useful information about them.

V. APPLICATIONS

It is well-known that the Kolmogorov complexity is not computable. In fact, the only way to find the shortest program that generates a sequence is to run all the short programs and see if they generate the sequence or not. However, some short programs may not halt and there is no way to figure out if the program will halt or not. Hence, there is no effective way to calculate the Kolmogorov complexity. However, it is usually possible to find upper bounds for the Kolmogorov complexity. In this section, we consider several popular examples and provide upper bounds for their Kolmogorov complexity. Based on these upper bounds we use Theorems 1 and 2 to calculate the number of random linear measurements required by the MCP to recover these functions. This demonstrates the connection between the results of Section IV and the compressed sensing and finite

rate of innovation frameworks explained in Section I. It is straightforward to extend the results to the other classes we discussed in Section I.

A. Sparsity

Let the signal $x_o = (x_{o,1}, x_{o,2}, \dots, x_{o,n})$ be k -sparse. Consider the following program for describing $[x_o^n]_m$. First, use a program of constant length to describe the structure of the signal as ‘sparse’ and the ordering of the rest of information. Then, spend $\log^* n + c$ bits to describe the length of the signal. Next, code the sparsity level k with $\log^* k$ bits, and spend $k(\log^* n + c)$ more bits to code the locations of the k non-zero elements. Finally, use km more bits to describe the quantized magnitudes of the non-zero coefficients. Therefore, we have

$$\begin{aligned} & \frac{K^{[\cdot]m}(x_{o,1}, x_{o,2}, \dots, x_{o,n})}{m} \\ & \leq k + \frac{(k+1)(\log^* n + c) + \log^* k + c}{m}. \end{aligned} \quad (11)$$

Plugging (11) into Theorem 1, we conclude that $\lceil (2k+1) \log n \rceil$ measurements are sufficient for the recovery of the k -sparse signals.

B. Piecewise polynomial

Let $(x_{o,1}, x_{o,2}, \dots, x_{o,n})$ be samples of a piecewise polynomial function $f(x)$ defined on $[0, 1]$ at locations $(0, 1/n, \dots, (n-1)/n)$. Further, assume that $0 \leq f(x) \leq 1$, for every x . Let Poly_N^Q represent the class of such functions which have at most Q singularities² and N is the maximum degree of each polynomial. Let $\{a_i^\ell\}_{i=0}^{N_\ell}$ denote the set of coefficients of the ℓ^{th} polynomial, where $N_\ell \leq N$ denotes its degree. For the notational simplicity, we assume that the coefficients of each polynomial belong to the $[0, 1]$ interval and that $\sum_{i=0}^{N_\ell} a_i^\ell < 1$ for every ℓ , where a_i^ℓ is the i^{th} coefficient of the ℓ^{th} polynomial. For a given length n , we derive an upper bound on the Kolmogorov complexity. Consider the following program for describing $[x_o^n]_m$. The code first specifies the model as ‘piecewise polynomial’ with parameters (n, Q, N) . This requires $\log^* n + \log^* N + \log^* k + c_1$ bits. Then, for each singularity point, the code first determines the largest sampling point i/n that is smaller than it. Since there are at most Q singularity points, describing this information requires at most $Q(\log^* n + c_2)$ bits. The next step is to describe the coefficients of each polynomial. Using an m' -bit quantizer for each coefficient, the induced error is bounded by

$$\begin{aligned} \left| \sum_{i=0}^{N_\ell} a_i^\ell t^n - \sum_{i=0}^{N_\ell} [a_i^\ell]_{m'} t^n \right| & \leq \sum_{i=0}^{N_\ell} |a_i^\ell - [a_i^\ell]_{m'}| \\ & \leq (N+1)2^{-m'}. \end{aligned} \quad (12)$$

To ensure that we are able to reconstruct the m -bit resolution of the samples from this description, $(N+1)2^{-m'} < 2^{-m}$. Therefore, describing the polynomials’ coefficients we need

²A singularity is a point at which the function is not infinitely differentiable.

$(Q+1)(N+1)(m + \lceil \log_2(N+1) \rceil)$ extra bits. Hence, overall, we conclude that

$$\begin{aligned} \frac{K^{[\cdot]m}(x_{o,1}, x_{o,2}, \dots, x_{o,n})}{m} &\leq (Q+1)(N+1) \\ &+ \frac{(Q+1)(N+1)\lceil \log_2(N+1) \rceil}{m} \\ &+ \frac{\log^* n + \log^* N + \log^* k + Q \log^* n + c_1 + c_2}{m}. \end{aligned} \quad (13)$$

It is straightforward to plug (13) into Theorem 2 and prove that, roughly speaking, for large values of n , $(QN + 2Q + 1) \log n$ measurements are sufficient for the successful recovery of the piecewise polynomial functions.

So far we have considered examples of low-complexity signals. However, in many applications the signals are not of low complexity but are rather close to low complexity signals. We present several examples here.

C. ℓ_p -constrained signals

While sparse signals have played an important role in the theory of compressed sensing, it is well-known that they do not occur in practice very often. More accurate models assume that either the magnitude of the signal follows a specific decay or the signal belongs to an ℓ_p ball with $p < 1$, i.e., $\|x_o\|_p \leq 1$ [1], [18]. For the signal $x_o \in \mathbb{R}^n$ with $\|x_o\|_p \leq 1$, let $(x_{o,(1)}, x_{o,(2)}, \dots, x_{o,(n)})$ denote the permuted version of x_o such that $x_{o,(1)} \geq x_{o,(2)} \geq \dots \geq x_{o,(n)}$. It is easy to show that $x_{o,(i)} \leq i^{-\frac{1}{p}}$. Therefore, if we just keep the k largest coefficients of this signal and set the rest to zero the resulting k -sparse vector \tilde{x}_o satisfies, $\|x_o - \tilde{x}_o\| \leq k^{-\frac{1}{p} + \frac{1}{2}}$. Setting the sparsity k to $n^{p/2}$, Theorem 2 proves that $d_n = n^{p/2} \log n$ samples are sufficient for asymptotically accurate recovery. It is interesting to note that as p decreases, the decay rate increases and the number of measurements required for the successful recovery decreases.

D. Smooth functions

Suppose that x_1, x_2, \dots, x_n are equispaced samples of a smooth function $f : [0, 1] \rightarrow \mathbb{R}$ with $0 \leq f(x) \leq 1$. Let the function be $\beta + 1$ times differentiable and $\|f^{(\beta+1)}\|_\infty \leq \gamma$. For the notational simplicity we assume that $|f^{(m)}(x)| \leq 1$ for every $m \leq \beta + 1$. This function is not necessarily a low-complexity signal, but it can be well approximated with a piecewise polynomial function. To show this, consider partitioning the $[0, 1]$ interval into subintervals of size r_n , and approximating the function f with a polynomial of degree β in each subinterval. Let $f_\beta(x)$ denote the resulting piecewise polynomial function. It is easy to prove that $\|f - f_\beta\|_\infty \leq \gamma r_n^{\beta+1}$. Hence, if x and x_o denote vectors consisting of the equispaced samples of the original signal and its piecewise polynomial approximation, respectively, it follows that $\|x - x_o\|_2 \leq \gamma \sqrt{n} r_n^{\beta+1}$.

On the other hand the complexity of the piecewise polynomial signal is essentially proportional to β/r_n . Setting $r_n = n^{\frac{2}{2\beta}}$, Theorem 2 proves that $d_n = O(n^{1/\beta} \log n)$ is enough for the accurate recovery of the samples of such signals. Clearly, for $\beta < 1$, this bound indicates that the

number of samples we need is at the same order as the ambient dimension. However, as β increases fewer number of samples are required.

Similar results hold for the piecewise smooth functions, which are very popular in image and signal processing.

VI. RELATED WORK

Our work is inspired by [19] and [20]. [19] considers the well studied problem of estimation, where the goal is to recover a vector θ from its noisy observations $s = \theta + z$, where z represents the noise in the system. It then suggests using the *minimum Kolmogorov complexity estimation* (MKCE) approach and proves that if $\theta_i \stackrel{iid}{\sim} \pi$, under several scenarios for the signal and noise, the average marginal distribution of the estimate of MKCE tends to the actual posterior distribution. On the other hand, [20] considers the problem of compressed sensing over binary sequences. Consider the set of all the binary sequences with Kolmogorov complexity less than or equal to k_0 , i.e.,

$$\mathcal{S}(k_0) \triangleq \{\mathbf{x} : K(\mathbf{x}) \leq k_0\}.$$

Let A denote a $d \times n$ binary matrix, $\mathbf{x}_o = (x_1, x_2, \dots, x_n)^T$, $\mathbf{y}_o = A\mathbf{x}_o$. Consider the following algorithm for reconstructing signal \mathbf{x}_o from its linear measurements \mathbf{y}_o :

$$\hat{\mathbf{x}}(\mathbf{y}_o, A) \triangleq \arg \min_{\mathbf{y}_o = A\mathbf{x}} K(\mathbf{x}). \quad (14)$$

[20] considers this scheme and proves that $2k$ random linear binary measurements are sufficient for recovering the binary sequences in $\mathcal{S}(k_0)$ with, high probability. This result does not provide any information on the successful recovery of real signals and it does not consider the non-idealities in the signals either. Our paper settles both questions.

As mentioned in Section I the problem we discuss in this paper is a central problem in the field of compressed sensing [1], [2]. Several papers have considered different generalization of sparsity [5], [6], [11], [12]. As mentioned before, all these models can be considered as subclasses of the general model we consider here. However, it is worth noting that even though the recovery approach proposed in our paper is universal, since Kolmogorov complexity is not computable, it is not useful for practical purposes.

In this paper, we considered deterministic models for the signals. Similar extensions have been considered in the random settings as well. For instance, [21] considers the problem of recovering a memoryless process from a linear set of measurements and proves the connection between the number of measurements required and the Renyi entropy. Also, our work is in the same spirit with the minimum entropy decoder proposed by Csiszar in [22]. He suggests a universal minimum entropy decoder, for reconstructing an iid signal from its linear measurements at a rate determined by the entropy of the source.

VII. PROOF OF THEOREM 1

The following Lemma will be used in the proof of the main theorem.

Lemma 2 (Chi-square concentration): Fix $\tau > 0$ and $x \in \mathbb{R}^n$. Assume that $\|x\|_2^2 = 1$. Let $Z_i \triangleq \sum_{j=1}^n A_{ij}x_j$, $i = 1, 2, \dots, d$. We then have,

$$\mathbb{P} \left(\sum_{i=1}^d Z_i^2 - 1 < -\tau \right) \leq e^{\frac{d}{2}(\tau + \log(1-\tau))}. \quad (15)$$

Proof: Note that $\{Z_i\}_{i=1}^d$ are iid $\mathcal{N}(0, 1/d)$. By Markov inequality, for any $\lambda > 0$, we have

$$\begin{aligned} \mathbb{P} \left(\sum_{i=1}^d Z_i^2 - 1 < -\tau \right) &= \mathbb{P} \left(-\sum_i Z_i^2 + 1 > \tau \right) \\ &\leq e^{-\lambda\tau} \mathbb{E} \left[e^{\lambda(1 - \sum Z_i^2)} \right] \\ &= e^{-\lambda\tau + \lambda} \left(\mathbb{E} [e^{-\lambda Z_1^2}] \right)^d \\ &= e^{-\lambda\tau + \lambda} \left(1 + \frac{2\lambda}{d} \right)^{-d/2}. \end{aligned} \quad (16)$$

We optimize over λ to obtain

$$\lambda^* = \frac{d\tau}{2(1-\tau)}. \quad (17)$$

If we plug (17) into (16) we obtain (15). \blacksquare

Proof: [Proof of Theorem 1] Let $e_m^n = x_o^n - [x_o^n]_m$ and $\hat{e}_m^n = \hat{x}_o^n - [\hat{x}_o^n]_m$ denote the quantization errors of the original and the reconstructed signals, respectively. Since both $Ax_o^n = y_o$ and $A\hat{x}_o^n = y_o$, it follows that

$$A([x_o^n]_m + e_m^n) = A([\hat{x}_o^n]_m + \hat{e}_m^n)$$

and

$$A([x_o^n]_m - [\hat{x}_o^n]_m) = A(\hat{e}_m^n - e_m^n). \quad (18)$$

On the other hand, since $|y - [y]_m| \leq 2^{-m}$, for each $y \in [0, 1]$, we have

$$\|\hat{e}_m^n - e_m^n\|_2^2 \leq n2^{-2m+1}.$$

Hence,

$$\begin{aligned} \|A([x_o^n]_m - [\hat{x}_o^n]_m)\|_2 &= \|A(\hat{e}_m^n - e_m^n)\|_2 \\ &\leq \sigma_{\max}(A) \sqrt{n2^{-2m+1}}. \end{aligned} \quad (19)$$

Since, by assumption, (7) holds for x_o , for each $\delta > 0$, there exists N_δ , such that for any $n > N_\delta$,

$$\frac{K^{[\cdot]m}(x_o^n)}{m} \leq \kappa + \delta \quad (20)$$

Since \hat{x}_o^n is the solution of (6),

$$K^{[\cdot]m}(\hat{x}_o^n) \leq K^{[\cdot]m}(x_o^n). \quad (21)$$

Moreover,

$$K([x_o^n]_m - [\hat{x}_o^n]_m) \leq K^{[\cdot]m}(x_o^n) + K^{[\cdot]m}(\hat{x}_o^n) + C, \quad (22)$$

where C is a constant independent of all the other variables in the problem [17]. Combining (20), (21) and (22) yields

$$K([x_o^n]_m - [\hat{x}_o^n]_m) \leq 2(\kappa + \delta)m + C. \quad (23)$$

If for each sequence y^n with $K^{[\cdot]m}(y^n) \leq 2(\kappa + \delta)m + C$, $\|A[y^n]_m\|_2 \geq \tau\|y^n\|_2$, for some fixed $\tau > 0$, then from (19)

$$\begin{aligned} \|x_o^n - \hat{x}_o^n\|_2 &= \|[x_o^n]_m + e_m^n - [\hat{x}_o^n]_m - \hat{e}_m^n\|_2 \\ &\leq \|[x_o^n]_m - [\hat{x}_o^n]_m\|_2 + \|e_m^n - \hat{e}_m^n\|_2 \\ &\leq \tau^{-1} \sigma_{\max}(A) \sqrt{n2^{-2m+1}} + \sqrt{n2^{-2m+1}} \\ &\leq (\tau^{-1} \sigma_{\max}(A) + 1) \sqrt{n2^{-2m+1}}. \end{aligned} \quad (24)$$

Define the events $\mathcal{E}_1^{(n)}$ and $\mathcal{E}_2^{(n)}$ as

$$\begin{aligned} \mathcal{E}_1^{(n)} &\triangleq \{A_{d \times n} : \\ &\nexists y^n; K^{[\cdot]m}(y^n) \leq 2(\kappa + \delta)m + C, \|Ay^n\|_2 < \tau\|y^n\|_2\}, \end{aligned} \quad (25)$$

and

$$\mathcal{E}_2^{(n)} \triangleq \left\{ A_{d \times n} : \sigma_{\max}(A) - 1 - \sqrt{\frac{n}{d}} < t \right\}, \quad (26)$$

for some $t > 0$.

Using these definitions plus the union bound, it follows that

$$\begin{aligned} \mathbb{P}(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon) &= \mathbb{P} \left(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon, \mathcal{E}_1^{(n)} \cap \mathcal{E}_2^{(n)} \right) \\ &\quad + \mathbb{P} \left(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon, (\mathcal{E}_1^{(n)})^c \cap \mathcal{E}_2^{(n)} \right) \\ &\leq \mathbb{P} \left(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon, \mathcal{E}_1^{(n)} \cap \mathcal{E}_2^{(n)} \right) \\ &\quad + \mathbb{P} \left((\mathcal{E}_1^{(n)})^c \cap \mathcal{E}_2^{(n)} \right) \\ &\leq \mathbb{P} \left(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon, \mathcal{E}_1^{(n)} \cap \mathcal{E}_2^{(n)} \right) \\ &\quad + \mathbb{P} \left(\mathcal{E}_1^{(n),c} \right) + \mathbb{P} \left(\mathcal{E}_2^{(n),c} \right). \end{aligned} \quad (27)$$

If $A \in \mathcal{E}_1^{(n)} \cap \mathcal{E}_2^{(n)}$, then from (19)

$$\|x_o^n - \hat{x}_o^n\|_2 \leq \left(\tau^{-1} \left(\sqrt{\frac{n}{d}} + 1 + t \right) + 1 \right) \sqrt{n2^{-2m+1}}. \quad (28)$$

Since, by assumption, $m = m_n = \lceil \log n \rceil$ and $d = d_n = \lceil \kappa \log n \rceil$, if n large enough,

$$\left(\tau^{-1} \left(\sqrt{\frac{n}{d}} + 1 + t \right) + 1 \right) \sqrt{n2^{-2m+1}} < \epsilon. \quad (29)$$

Hence, for n large enough

$$\mathbb{P} \left(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon, \mathcal{E}_1^{(n)} \cap \mathcal{E}_2^{(n)} \right) = 0. \quad (30)$$

On the other hand, by Lemma 2, for each sequence $x^n \in \mathbb{R}^n$,

$$\begin{aligned} \mathbb{P} \{ \|Ax^n\|_2^2 \leq \tau \|x^n\|_2^2 \} &= \mathbb{P} \left\{ \left\| A \frac{x^n}{\|x^n\|_2} \right\|_2^2 \leq \tau^2 \right\} \\ &\leq e^{\frac{d}{2}(1-\tau^2+2\log\tau)}. \end{aligned} \quad (31)$$

Therefore,

$$\begin{aligned} \mathbb{P} \left(\mathcal{E}_1^{(n),c} \right) &= \\ \mathbb{P} \left\{ \exists y^n : K^{[\cdot]m}(y^n) \leq 2(\kappa + \delta)m + C, \|Ay^n\|_2^2 < \tau \|y^n\|_2^2 \right\} \\ &\leq 2^{2(\kappa + \delta)m + C} e^{-\frac{d}{2}(1-\tau^2+2\log\tau)}. \end{aligned} \quad (32)$$

If we set $\tau = 0.04$ and $d = \lceil \kappa \log n \rceil$ it is simple to see that this probability goes to zero. Finally, we can use the concentration of Lipschitz function of a Gaussian random vector to prove [23]

$$\begin{aligned} \mathbb{P}\left(\mathcal{E}_2^{(n),c}\right) &= \mathbb{P}\left(\sigma_{\max}(A) - 1 - \sqrt{\frac{n}{d}} > t\right) \\ &\leq e^{-dt^2/2}. \end{aligned} \quad (33)$$

Setting t to a constant and $d = \lceil \kappa \log n \rceil$ proves that this probability also goes to zero. ■

VIII. PROOF OF THEOREM 2

Let $x_o^n = [x_o^n]_m + e_m^n$, $\tilde{x}_o^n = [\tilde{x}_o^n]_m + \tilde{e}_m^n$, and $\hat{x}_o^n = [\hat{x}_o^n]_m + \hat{e}_m^n$.

Note that since $\|A\tilde{x}_o^n - y_o^n\|_2 = \|A(\tilde{x}_o^n - x_o^n)\|_2 \leq \sigma_{\max}(A)\epsilon_n$, \tilde{x}_o^n is also a feasible solution. Therefore, since \tilde{x}_o^n and \hat{x}_o^n are both feasible, by triangle inequality,

$$\begin{aligned} \|A\tilde{x}_o^n - A\hat{x}_o^n\|_2 &= \|A\tilde{x}_o^n - y_o^n - (A\hat{x}_o^n - y_o^n)\|_2 \\ &\leq 2\sigma_{\max}(A)\epsilon_n. \end{aligned} \quad (34)$$

Again, by triangle inequality,

$$\begin{aligned} \|A\tilde{x}_o^n - A\hat{x}_o^n\|_2 &= \|A([\tilde{x}_o^n]_m + \tilde{e}_m^n) - A([\hat{x}_o^n]_m + \hat{e}_m^n)\|_2 \\ &\geq \|A([\tilde{x}_o^n]_m - [\hat{x}_o^n]_m)\|_2 - \|A([\tilde{e}_m^n]_m - [\hat{e}_m^n]_m)\|_2 \\ &\geq \|A([\tilde{x}_o^n]_m - [\hat{x}_o^n]_m)\|_2 - \sigma_{\max}(A)\|[\tilde{e}_m^n]_m - [\hat{e}_m^n]_m\|_2 \\ &\geq \|A([\tilde{x}_o^n]_m - [\hat{x}_o^n]_m)\|_2 - \sigma_{\max}(A)\sqrt{n}2^{-2m+1}. \end{aligned} \quad (35)$$

Combining (34) and (35), it follows

$$\|A([\tilde{x}_o^n]_m - [\hat{x}_o^n]_m)\|_2 \leq \sigma_{\max}(A)\sqrt{n}2^{-2m+1} + 2\sigma_{\max}(A)\epsilon_n. \quad (36)$$

Since both \tilde{x}_o^n and \hat{x}_o^n are feasible, and \hat{x}_o^n is the optimizer of (9), we have

$$K^{[\cdot]m}(\hat{x}_o^n) \leq K^{[\cdot]m}(\tilde{x}_o^n) \leq m(\kappa_n + \delta), \quad (37)$$

and therefore

$$K^{[\cdot]m}(\hat{x}_o^n - \tilde{x}_o^n) \leq m2(\kappa_n + \delta) + C, \quad (38)$$

where C is a constant independent of m and n .

Consider defining the events \mathcal{E}_1 and \mathcal{E}_2 as done in (25) and (26), in the proof of Theorem 1. Then, using the same argument used in that proof,

$$\begin{aligned} \mathbb{P}(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon) &\leq \mathbb{P}\left(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon, \mathcal{E}_1^{(n)} \cap \mathcal{E}_2^{(n)}\right) \\ &\quad + \mathbb{P}\left(\mathcal{E}_1^{(n),c}\right) + \mathbb{P}\left(\mathcal{E}_2^{(n),c}\right). \end{aligned} \quad (39)$$

However, our choice of parameters guarantees that for large enough n , $\mathbb{P}(\|x_o^n - \hat{x}_o^n\|_2 > \epsilon, \mathcal{E}_1^{(n)} \cap \mathcal{E}_2^{(n)}) = 0$, and moreover, $\mathbb{P}(\mathcal{E}_1^{(n),c})$ and $\mathbb{P}(\mathcal{E}_2^{(n),c})$ both go to 0 as n grows to infinity.

IX. CONCLUSION

In this paper, we consider the problem of recovering structured signals from their linear measurements. We use the Komogorov complexity of the quantized signal as a universal measure of complexity that covers many different examples explored in compressed sensing literature and related areas. We then show that, if we consider low-complexity signals, the minimum complexity pursuit scheme inspired by the Occam's razor recovers the simplest solution of a set of random linear measurements. In fact, we prove that the number of measurements required is proportional to the complexity and logarithmically to the ambient dimension of the signal. We also consider more practical scenarios where the signal is not 'simple' but is 'close' to a low complexity signal. We show that even in such cases following minimum complexity pursuit algorithm provides a good estimate of the signal from much fewer samples than the ambient dimension of the signal.

As mentioned in the paper, Kolmogorov complexity of a sequence is not computable. However, currently we are working on deriving implementable schemes by replacing Kolmogorov complexity by computable measures such as minimum description length [24].

REFERENCES

- [1] D. L. Donoho. Compressed sensing. *IEEE Transactions on Information Theory*, 52(4):489–509, April 2006.
- [2] E. Candès, J. Romberg, and T. Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *IEEE Transactions on Information Theory*, 52(2):489–509, February 2006.
- [3] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly. Compressed sensing MRI. *IEEE Signal Processing Magazine*, 25(2):72–82, March 2008.
- [4] M. A. Herman and T. Strohmer. High-resolution radar via compressed sensing. *IEEE Transactions on Signal Processing*, 57(6):2275–2284, June 2009.
- [5] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum rank solutions to linear matrix equations via nuclear norm minimization. *Siam Review*, 52(3):471–501, April 2010.
- [6] R.G. Baraniuk, V. Cevher, M.F. Duarte, and C. Hegde. Model-based compressive sensing. *IEEE Transactions on Information Theory*, 56(4):1982–2001, April 2010.
- [7] Y.C. Eldar, P. Kuppinger, and H. Bolcskei. Block-sparse signals: Uncertainty relations and efficient recovery. *IEEE Transactions on Signal Processing*, 58(6):3042–3054, 2010.
- [8] M. Stojnic, F. Parvaresh, and B. Hassibi. On the reconstruction of block-sparse signals with an optimal number of measurements. *IEEE Transactions on Signal Processing*, 57(8):3075–3085, 2009.
- [9] M. Stojnic. Block-length dependent thresholds in block-sparse compressed sensing. *Arxiv preprint arXiv:0907.3679*, 2009.
- [10] D. Malioutov, M. Cetin, and A.S. Willsky. A sparse signal reconstruction perspective for source localization with sensor arrays. *IEEE Transactions on Signal Processing*, 53(8):3010 – 3022, August 2005.
- [11] M. Vetterli, P. Marziliano, and T. Blu. Sampling signals with finite rate of innovation. *IEEE Transactions on Signal Processing*, 50(6):1417–1428, June 2002.
- [12] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. Willsky. The convex geometry of linear inverse problems. *Preprint*, 2010.
- [13] R. J. Solomonoff. A formal theory of inductive inference. *Inform. Contr.*, 7:224–254, 1964.
- [14] A. N. Kolmogorov. Logical basis for information theory and probability theory. *IEEE Transactions on Information Theory*, 14:662–664, 1968.
- [15] L. Staiger. The kolmogorov complexity of real numbers. *Theoretical Computer Science*, 284(2):455 – 466, 2002.

- [16] S. C. Tormay. *Ockham: studies and selections*. Open court publishers, La Salle, IL, 1938.
- [17] T. Cover and J. Thomas. *Elements of Information Theory*. Wiley, New York, 2nd edition, 2006.
- [18] A. Maleki. Approximate message passing algorithm for compressed sensing. *Stanford University PhD Thesis*, 2010.
- [19] D. L. Donoho. Kolmogorov sampler. *Preprint*, 2002.
- [20] D. L. Donoho, H. Kakavand, and J. Mammen. The simplest solution to an underdetermined system of linear equations. In *2006 IEEE International Symposium on Information Theory*, pages 1924 –1928, July 2006.
- [21] Y. Wu and S. Verdú. Renyi information dimension: Fundamental limits of almost lossless analog compression. *Information Theory, IEEE Transactions on*, 56(8):3721 –3748, August 2010.
- [22] I. Csiszar. Linear codes for sources and source networks: Error exponents, universal coding. *IEEE Transaction on Information Theory*, 28:585 – 592, 1982.
- [23] E. Candès, J. Romberg, , and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203 – 4215, Dec. 2005.
- [24] J. Rissanen. Stochastic complexity and modeling. *Annals of Statistics*, 14:1080–1100, 1986.