

Ramsey theory reveals the conditions when sparse coding on subsampled data is unique

Christopher J. Hillar, Friedrich T. Sommer

Abstract—Sparse coding or dictionary learning has been widely used to reveal the sparse underlying structure of many kinds of sensory data. A related advance in signal processing is compressed sensing, a theory explaining how sparse data can be subsampled below the Nyquist-Shannon limit and then efficiently recovered from these subsamples. Here we study whether the conditions for recovery in compressed sensing are sufficient for dictionary learning to discover the original sparse causes of subsampled data. Using combinatorial Ramsey theory, we completely characterize when the learned dictionary matrix and sparse representations of subsampled data are unique (up to the natural equivalences of permutation and scaling). Surprisingly, uniqueness is shown to hold without any assumptions on the learned dictionaries or inferred sparse codes. Our result has implications for the learning of overcomplete dictionaries from subsampled data and has potential applications in data analysis and neuroscience. For instance, it identifies sparse coding as a possible learning mechanism for establishing lossless communication through severe bottlenecks, which might explain how different brain regions communicate through axonal fiber projections.

Index Terms—Dictionary learning, sparsity, compressed sensing, sparse coding, Ramsey theory

I. INTRODUCTION

INDPENDENT component analysis [1], [2] and dictionary learning with a sparse coding scheme [3], [4] have become important tools for revealing underlying structure in many different types of data [5], [6]. These algorithms share two main components: a *coding* step and a *reconstruction*. In the coding step of [3], for instance, a *sparse* code vector $\mathbf{b}(\mathbf{y})$ with a small number of nonzero coordinates is computed given a data point \mathbf{y} in a data set Y . The code vector $\mathbf{b}(\mathbf{y})$ is then used to linearly reconstruct \mathbf{y} as

$$\hat{\mathbf{y}} = B\mathbf{b}(\mathbf{y}), \quad (1)$$

using a matrix B (sometimes called a *dictionary* for Y). The code map $\mathbf{b}(\mathbf{y})$ and the matrix B are fit to training data using unsupervised learning.

If reconstruction succeeds and $\hat{\mathbf{y}} = \mathbf{y}$ for data in Y , then the representations $\mathbf{b}(\mathbf{y})$ and the dictionary B capture essential structure in the data. On the other hand, a nonzero difference $(\hat{\mathbf{y}} - \mathbf{y})$ gives an error signal to better tune the two steps of

C. Hillar is with the Mathematical Sciences Research Institute, Berkeley, CA, 94720 USA e-mail: chillar@msri.org and the Redwood Center for Theoretical Neuroscience, Berkeley, CA, 94720 USA. Hillar was partially supported by an NSF All-Institutes Postdoctoral Fellowship administered by the Mathematical Sciences Research Institute through its core grant DMS-0441170.

F. Sommer is with the Redwood Center for Theoretical Neuroscience, Berkeley, CA, 94720 USA e-mail: fsommer@berkeley.edu.

Adaptive Compressed Sampling

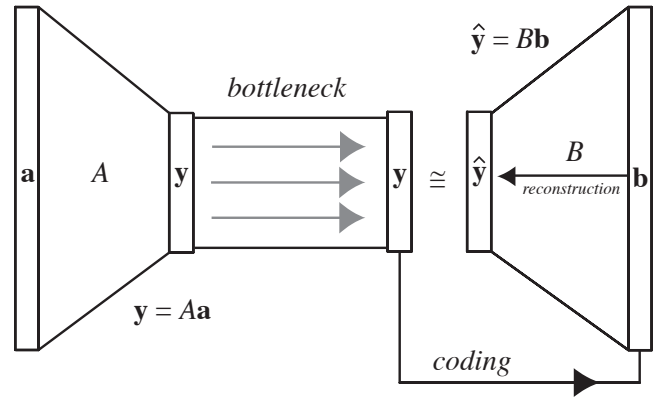


Fig. 1. In adaptive compressed sampling (ACS), signals with unknown sparse ambient structure \mathbf{a} are subsampled by an unknown compressed sensing (CS) matrix A to form $\mathbf{y} = A\mathbf{a}$. Unsupervised dictionary learning is then used to fit a sparse generative model $\hat{\mathbf{y}} = B\mathbf{b}(\mathbf{y})$ to the compressed data. When reconstruction succeeds and $\hat{\hat{\mathbf{y}}} = \mathbf{y}$, our main result (Theorem 1) implies that the resulting dictionary B and sparse code vectors \mathbf{b} are equal to the matrix A and sparse vectors \mathbf{a} up to a fixed permutation and scaling; see Eqn. (6). Such a scheme of sparse coding might enable brain regions to communicate sparse feature representations through axonal fiber projections with limited numbers of fibers [11], [12].

coding and reconstruction to the data. Such a control loop for optimizing a coding procedure is referred to as *predictive coding* or *self-supervised learning* in the literature [7]. The sparse vector $\mathbf{b} = \mathbf{b}(\mathbf{y})$ is also sometimes called an “efficient representation” of \mathbf{y} because it exposes the coefficients of the “independent components” or “causes” of data, thereby minimizing the redundancy of the description [8], [9].

Thus, the columns of a learned matrix B can be interpreted as structural primitives inherent in the data set Y , and the inferred vector $\mathbf{b}(\mathbf{y})$ as specifying a sparse weighted sum of these primitives which reconstruct \mathbf{y} . It has to be emphasized that sparse structure is a property empirically found in many types of sensory data, most notably natural images [3], as well as natural sounds and speech [4]. In the literature, predictive coding with a sparseness constraint is referred to as *sparse coding* or *sparse dictionary learning* (SDL). Importantly for applications, SDL can reveal *overcomplete* representations; that is, the dimension of \mathbf{b} can exceed the dimension of the data \mathbf{y} . Overcomplete representations can capture data mixtures that have been sparsely composed from a set of complete dictionaries [10].

A related advance in signal processing is the paradigm of compressed sensing (CS) [13], [14] (see also [15] for a recent

theoretical review). The theory of CS provides a collection of techniques to recover data vectors \mathbf{x} with sparse ambient structure after they have been linearly *subsampled* as $\mathbf{y} = \Phi\mathbf{x}$ by a (known) *compressive* matrix Φ (the number of rows n of Φ is significantly smaller than the number of its columns). Typically, the sparsity assumption enforced on \mathbf{x} is that it can be expressed $\mathbf{x} = \Psi\mathbf{a}$ for a fixed (known) dictionary matrix Ψ and an (unknown) m -dimensional vector \mathbf{a} with *at most* k nonzero components (i.e., entries). Such vectors \mathbf{a} are called *k-sparse*.

Surprisingly, under very mild CS conditions on the matrix $A = \Phi\Psi$, there are efficient and robust algorithms for recovering k -sparse m -dimensional \mathbf{a} (and thus \mathbf{x}) from the n -dimensional subsampled vector:

$$\mathbf{y} = A\mathbf{a}, \quad (2)$$

as long as the dimension of \mathbf{y} satisfies

$$n \geq Ck \log(m/k). \quad (3)$$

Here, C is a (small) universal constant independent of m , n , or k .¹ In other words, one can easily recover sparse high-dimensional vectors \mathbf{a} from “projections” (2) into spaces of dimension commensurate with number of active coordinates of those vectors.

A common CS recovery condition on $A \in \mathbb{R}^{n \times m}$ is that it never maps two different sparse vectors to the same vector:²

$$A\mathbf{a}_1 = A\mathbf{a}_2 \text{ for } k\text{-sparse } \mathbf{a}_1, \mathbf{a}_2 \in \mathbb{R}^m \implies \mathbf{a}_1 = \mathbf{a}_2. \quad (4)$$

Note that a generic square matrix A is invertible; thus, (4) trivially holds for almost all matrices A whenever $n = m$. In the interesting regime of compressed sensing, however, the sample dimension n is significantly smaller than the original data dimension m . Thus, a condition such as (4) supplants “invertibility” of the matrix A with an “incoherence” among every $2k$ of its columns. Rather remarkably, even in the critical regimes close to equality of (3), condition (4) holds with very high probability for randomly generated $n \times m$ matrices A .

The theory of CS has implications for SDL. As the restriction (3) on the dimensions of A is necessary for successful coding in dictionary learning, it bounds the permitted degree of overcompleteness for SDL. Concretely, if (3) is violated so that overcompleteness exceeds $m = k \exp \frac{n}{Ck}$, the coding $\mathbf{b}(\mathbf{y})$ of \mathbf{y} in (1) cannot succeed and therefore dictionary learning is infeasible. Nonetheless, the theory does provide for an exponential degree of overcompleteness.

Here, we ask whether the necessary conditions of CS allow full recovery of original sparse representations using SDL.

Problem 1 (The ACS Problem): Let $Y = \{A\mathbf{a}_1, \dots, A\mathbf{a}_N\}$ be an n -dimensional data set generated linearly via (2) in which $A \in \mathbb{R}^{n \times m}$ satisfies compressed sensing conditions (3) and (4), but is unknown. Can sparse dictionary learning uncover the sparse causes $\mathbf{a}_1, \dots, \mathbf{a}_N$ and the matrix A ?

¹For a more detailed discussion of these facts (including proofs) and their relationship to approximation theory and concentration of measure phenomenon, we refer the reader to [16] and the references therein.

²Letting $c = \min\{p, 2k\}$, this condition is equivalent to the assertion that every c columns of A are linearly independent. This condition is sometimes expressed as $\sigma(A) > 2k$, where $\sigma(A)$ is the smallest number of linearly dependent columns of A (the *spark* of A).

By imposing CS conditions, Problem 1 combines dictionary learning with compressed sensing. We will refer to this combination as *adaptive compressed sampling* (ACS); see Fig. 1. Note that Problem 1 differs from *blind compressed sensing*, the problem of reconstructing the matrix Ψ from subsampled data, which is ill-posed without additional constraints [17].

Problem 1 is motivated by a basic question in theoretical neuroscience: How can synaptic learning enable communication between brain regions through axonal fiber tracts that form severe wiring bottlenecks? One theory of communication through these bottlenecks is sketched in [12]. There, data vectors \mathbf{a} correspond to firing patterns of m local neurons in a brain region and vectors \mathbf{y} generated through (2) represent the activity in another set of $n \ll m$ neurons in that brain region, all projecting their axonal fibers to one specific second region. Thus, in this model, brain regions communicate through bottlenecks by sending subsamples of their activity. For decoding of these signals, the theory of [12] posits that the second brain region recover the original activity using SDL. Thus, a positive answer to Problem 1 is crucial for this theory. See Section III below for more on this connection.

Since any permutation and (component-wise) scaling of a sparse vector is also a sparse vector, taken literally, Problem 1 is ill-posed [18]. More precisely, if P is a fixed permutation matrix³ and D is an invertible diagonal matrix, then

$$A\mathbf{a} = (AD^{-1}P^\top)(PDA)$$

for each sample $\mathbf{y} = A\mathbf{a}$. Thus, without access to A , one could not discriminate which of \mathbf{a} or PDA (resp. A or $AD^{-1}P^\top$) was the original sparse vector (resp. sampling matrix). Problem 1, therefore, should be interpreted as asking whether up to these natural transformations (permutation and scaling), recovery is possible with sparse dictionary learning.

In this note, we give a complete affirmative solution to Problem 1. Suppose equation (2) is used to generate a set of compressed data $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ from a (sufficiently large) set of N k -sparse vectors $\mathbf{a}_1, \dots, \mathbf{a}_N$. These sparse vectors and the dictionary A are unknown to the predictive coding unit, but A satisfies the conditions specified in Problem 1. If SDL is trained with subsamples, it tries to find a coding map $\mathbf{y} \mapsto \mathbf{b}(\mathbf{y})$ (with k -sparse output) and a dictionary B such that (1) reconstructs \mathbf{y} . We prove (in Theorem 1 below) that if the learning algorithm succeeds at predictive coding of the subsampled data:

$$\mathbf{y}_i = B\mathbf{b}_i, \quad i = 1, \dots, N; \quad (5)$$

then there is a fixed permutation matrix P and a fixed invertible diagonal matrix D such that

$$A = BPD \quad \text{and} \quad \mathbf{b}_i = PDA_i, \quad i = 1, \dots, N. \quad (6)$$

Thus, the matrix A and the uncompressed vectors are recovered (up to symmetry) whenever SDL of subsampled data succeeds.

³A *permutation matrix* P has binary entries and exactly one 1 in each column and exactly one 1 in each row (thus, $P\mathbf{v}$ for a column vector \mathbf{v} permutes its entries). Note that $PP^\top = P^\top P = I$, where I denotes the $p \times p$ identity matrix, and M^\top for a matrix M is its transpose. Thus, $P^{-1} = P^\top$.

To understand our result in the context of ACS, consider two receiver regions R_1 and R_2 , each having access to subsampled data Y . Suppose that both regions are able to reconstruct Y as in (5) using matrices B_1 and B_2 and code vectors $\mathbf{b}_i^{(1)}$ and $\mathbf{b}_i^{(2)}$, respectively. We remark that even if their coding paradigms are the same, regions R_1 and R_2 might have very different initial conditions for learning; thus, a priori, dictionaries and codes could be very different. Nonetheless, it follows from (6) that $A = B_1 P_1 D_1 = B_2 P_2 D_2$ for some permutation matrices P_1, P_2 and invertible diagonal matrices D_1, D_2 . Thus, the two dictionaries are related via

$$B_1 = B_2 P_2 D_2 D_1^{-1} P_1^\top = B_2 (P_2 P_1^\top) (P_1 D_2 D_1^{-1} P_1^\top).$$

The two matrices $P = P_2 P_1^\top$ and $D = P_1 D_2 D_1^{-1} P_1^\top$ are easily checked to be a permutation and diagonal, respectively. Therefore, B_1 and B_2 are a permutation scaling away from each other, and the same holds for $\mathbf{b}_i^{(1)}$ and $\mathbf{b}_i^{(2)}$.

The most general previous result on this problem is described in [18]. Under the additional assumption that B satisfies CS condition (4) and that certain *supports* (indices of nonzero components) of inferred vectors \mathbf{b}_i occur, the authors of [18] show that sufficiently many equations (5) imply (6).

In practical applications of SDL, however, one rarely has any guarantees on the learned dictionary matrix B or on the inferred sparse vectors $\mathbf{b}(\mathbf{y})$. It is also computationally intractable to verify a CS condition such as (4). Thus, the uniqueness guarantee of the sort (6) without assumptions on B or inferred vectors $\mathbf{b}(\mathbf{y})$ is ideal. As we shall see in the course of proving our main theorem, removing these assumptions is a technical challenge requiring methods from combinatorics, including the surprising use of a basic result in Ramsey theory [19]. This field of mathematics deals with proving statements of the following form: *In every group of 6 people, there are 3 people who all know each other or 3 people who all do not.*

The organization of this paper is as follows. Section II provides the precise mathematical formulation of the uniqueness result (6) alluded to above. Section III gives a short discussion of how our results fit into the general sparse dictionary learning literature and theoretical neuroscience. Finally, the proof of our main theorem appears in Section IV, where we first verify an instructive easier special case. (The Ramsey theory we need is proved in the Appendix.) To appeal to the widest audience possible, we have kept our mathematical arguments as self-contained and elementary as possible.

II. THE ACS RECONSTRUCTION THEOREM

Recall that a k -sparse vector \mathbf{a} is a column vector $\mathbf{a} \in \mathbb{R}^m$ with at most $k < m$ nonzero entries. Our main theorem is concerned with the recovery of k -sparse vectors \mathbf{a} after having received only subsampled versions $A\mathbf{a} \in \mathbb{R}^n$. Here, $A \in \mathbb{R}^{n \times m}$ is a matrix satisfying CS condition (4), and typically in applications $n \approx k \log(m/k)$ so that a significant dimension reduction takes place (although this is not used in our argument). We now state precisely our main theorem.

Theorem 1 (The ACS Theorem): Fix positive integers n and $k < m$. There are k -sparse $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ with the following property: if $A \in \mathbb{R}^{n \times m}$ satisfies (4) and $B \in \mathbb{R}^{n \times m}$

and k -sparse $\mathbf{b}_1, \dots, \mathbf{b}_N$ are such that (5) holds, then there exists an invertible diagonal matrix $D \in \mathbb{R}^{m \times m}$ and a permutation matrix $P \in \mathbb{R}^{m \times m}$ such that (6) is satisfied.

Remark 1: Equation $A = BPD$ and assumption (4) already imply the recovery result $\mathbf{b}_i = PD\mathbf{a}_i$. This follows since $A\mathbf{a}_i = AD^{-1}P^\top \mathbf{b}_i$ from (5) and thus $\mathbf{b}_i = PD\mathbf{a}_i$ from (4).

Remark 2: In applications, it is important to have bounds on the number N of equations of the form (5) guaranteeing uniqueness. We address this in Section III. We also note that our proof of Theorem 1 gives an algorithm to find P and D .

Remark 3: It is easy to see that the assumption (4) cannot be removed. For instance, trivially, there does not exist such a set of vectors $\mathbf{a}_1, \dots, \mathbf{a}_N$ when $A = 0$.

Theorem 1 is surprising and remarkably general: it says that any dictionary learning scheme producing sparse reconstructions in a subsampled space automatically gives faithful transmission of sparse signals (and dictionary) regardless of the CS sampling matrix A , the learned dictionary matrix B , or the inferred sparse vectors $\mathbf{b}(\mathbf{y})$.

To better understand some of its complexity, consider the statement of Theorem 1 when $k < n = m$, $B = I$ is the identity matrix, and $A \in \mathbb{R}^{n \times n}$ is invertible. In this case, the result implies that if $A\mathbf{a}$ is k -sparse for every k -sparse \mathbf{a} , it must be that $A = PD$ for some permutation matrix P and invertible diagonal matrix D . Since every such matrix $A = PD$ trivially satisfies this condition, Theorem 1 gives a complete characterization of all such matrices.

Corollary 1: Fix positive integers $k < n$. The set of invertible $n \times n$ matrices A having the property that $A\mathbf{a}$ is k -sparse for k -sparse \mathbf{a} is the set of matrices PD , where P and D run over permutation and invertible diagonal real matrices.

Remark 4: This corollary is also deducible from [18].

A surprising ingredient in our proof of Theorem 1 is a result from combinatorial Ramsey theory. We give here one limiting instance of the result we use (see Theorem 5 in the Appendix for the full statement and Figure 2 for an example with $s = 2$).

Theorem 2: Every coloring of the integer points \mathbb{Z}^2 in the plane with a finite set of colors has the following structural property. For each positive integer s , there are subsets $H_1, H_2 \subseteq \mathbb{Z}$ each containing s integers such that all the points in the grid $H_1 \times H_2$ possess the same color.

Experimental verification that a trained ACS unit robustly satisfies both implications (6) of Theorem 1 appears in [12]. These findings suggest that a noisy version of Theorem 1 or [18, Theorem 3] holds. This will be a focus of future work.

III. DISCUSSION

We have shown in this note that any learning scheme that converges a model of predictive coding (5) for a sufficient number of samples (2) of compressed data solves the ACS Problem uniquely as long as the conditions of compressed sensing are met (Theorem 1). For standard SDL, these conditions imply that the overcompleteness of a learned dictionary can reach but must not exceed $m \approx k \exp \frac{n}{Ck}$. Interestingly, our proof of Theorem 1 does not require any assumptions about the predictive model. In contrast, previous studies of uniqueness of dictionary learning were limited to complete

dictionaries, e.g. [20], or relied on additional assumptions for the predictive coding model. For example, [18] required B to fulfill CS condition (4) and put restrictions on the sparse codes \mathbf{b}_i that could occur. In general, it seems computationally challenging to enforce such requirements in a model of predictive coding that is dynamically evolving under the learning process.

The number of data samples N required in our proof of Theorem 1 has to be very large so that Ramsey theory can control the supports of inferred vectors. In the Appendix, we derive an upper bound given by inequality (23). In contrast, the uniqueness result in [18] requires far fewer samples: $N = (k+1)\binom{m}{k}$. It is an open problem how much smaller N can be made in Theorem 1 until additional assumptions on the predictive model become necessary. We suspect that a reduction to Ramsey theory might be necessary for proving Theorem 1; thus, without further assumptions, a large N might be unavoidable. In regimes of moderate overcompleteness or compression, theoretical [20], [21] and experimental [12] results suggest that the typical amount of data required for successful SDL is much smaller than the N used in our proof. One possible explanation for this discrepancy is that SDL coding maps $\mathbf{b}(\mathbf{y})$ are usually highly structured, whereas the code vectors \mathbf{b}_i in the hypothesis (6) of Theorem 1 are allowed to be arbitrary.

Another open problem left unaddressed in this work is to find conditions under which predictive coding (5) is guaranteed to converge. Although widely used in practice, all known (non-convex) sparse dictionary learning algorithms lack a mathematical proof of convergence, and finding such an argument is a major open problem in the field. The most significant progress on this problem appears in [22], [20], [21], where local convergence of SDL is established.⁴ What we have accomplished here with Theorem 1 is that whenever SDL converges, recovery of high-dimensional sparse codes is automatically guaranteed, independent of any assumptions on the learned dictionary or sparse code vectors.

It has to be emphasized that ACS is compressed sensing with learning in the decoding stage. Specifically, the decoding stage has to infer or learn the product of the compressive matrix and the dictionary, both of which are made available to the decoder in standard compressed sensing. An earlier modification of compressed sensing proposed an altered encoding stage [23] (see also [24]). Rather than using a random sampling matrix, a learning algorithm called uncertain component analysis (UCA) was designed to optimize recovery in the decoder. In particular, for data that are not truly sparse (such as natural image patches), a learned sampling matrix can improve recovery quality considerably. It is an interesting question how ACS performs when the data are subsampled with a trained matrix rather than a random one.

Implications for theoretical neuroscience: An important consequence of our work is to provide the mathematical conditions under which a predictive coding model can be

⁴Roughly, local convergence means that given enough stochastically generated samples of the form (2), the learned matrix B and vectors \mathbf{b}_i in (5) used to reconstruct the data are local minima (with high probability) of a certain SDL objective function. See [21] for the most recent result of this kind.

trained in a compressed space (locally accessible by neurons in a receiver region) while accurately coding for data in some uncompressed space (neuronal firing in a sender region). Predictive coding has been widely proposed to describe representational learning in the brain, but it is often criticized as unrealistic to assume that the neural code is optimized to reproduce the full sensory signal. For instance, recovery of the dictionary Ψ for reconstructing the full data $\mathbf{x} = \Psi\mathbf{a}$ from subsamples (2) requires a factorization $A = \Phi\Psi$ which is ill-posed without additional constraints [17]. The theory of ACS suggests a scheme of representational learning that uses predictive coding without assuming that brain regions reconstruct full sensory input [12]. The theory predicts that learning recovers the sparse causes of the sensory signal, which might be useful for classification or other structural operations on input [5], [6]. Further, the learning is predicted to recruit only the number of neurons in the afferent region sufficient to represent the input (i.e., see Corollary 2 below).

ACS theory is consistent with old ideas of efficient coding in neuroscience [8], [9]. The difference is that the objective of efficient coding is to minimize the redundancy in the neural code whereas the objective of ACS is to optimize communication through a wiring bottleneck. However, as long as the conditions in Problem 1 are met, the appropriate learning algorithms and the resulting codes are indistinguishable. For example, when trained with randomly subsampled images, ACS produces V1-like receptive fields [12]. Interestingly, however, these receptive fields are not wired into the circuitry as in conventional efficient coding. The receptive fields are only reconstructable by an outside observer who has simultaneous access to the neural activity in the trained ACS circuit and the full images.

Implications for applications: A practical consequence of our main theorem is that it fully characterizes the feasible regime of overcompleteness in SDL. Another interesting result of this study is that obvious structure in the learned dictionary should not be the only criterion of success for SDL. For instance, a learned dictionary $B = \Phi\Psi D^{-1}P^T$ might not reveal clear structure even though learning has converged and the resulting sparse codes correctly represent the underlying sparse causes in the original data. It is an important question for the future how one can assess if SDL was successful when access to uncompressed data is impossible. For example, the performance of a classifier might be such a criterion [5], [6].

Finally, it would be interesting to explore whether compressing data with a random matrix before applying SDL is a form of regularization that reduces the number of free parameters in the model, thereby increasing the speed of learning.

IV. PROOF OF THEOREM 1

Our proof of Theorem 1 involves three main pieces: Theorem 5 from the Appendix, a combinatorial matrix theory result (Proposition 1), and a fact in basic linear algebra (Lemma 1). First, Ramsey theory is used to control the supports of \mathbf{b}_i relative to \mathbf{a}_i , and then Proposition 1 produces the permutation P and diagonal D (inductively) with the help of Lemma 1.

Before proving Theorem 1 in full generality, let us consider the simple case when $k = 1$. Set $\mathbf{a}_i = \mathbf{e}_i$ ($i = 1, \dots, m$) to

be the standard basis column vectors in \mathbb{R}^m .⁵ Assuming that (5) holds for some matrix B and 1-sparse \mathbf{b}_i , it follows that

$$A\mathbf{e}_i = B c_i \mathbf{e}_{\pi(i)}, \quad (7)$$

for some function $\pi : \{1, \dots, m\} \rightarrow \{1, \dots, m\}$ and $c_i \in \mathbb{R}$. Notice that if $c_i = 0$ for some i , then $A\mathbf{e}_i = 0$, contradicting assumption (4) for A . We show that π is necessarily injective (and thus is a permutation).⁶ Suppose $\pi(i) = \pi(j)$; then,

$$A\mathbf{e}_i = c_i B \mathbf{e}_{\pi(i)} = c_i B \mathbf{e}_{\pi(j)} = \frac{c_i}{c_j} B c_j \mathbf{e}_{\pi(j)} = \frac{c_i}{c_j} A \mathbf{e}_j.$$

Again by (4) this is only possible if $i = j$. Thus, π is injective.

Let P and D denote the permutation and diagonal matrices:

$$P = (\mathbf{e}_{\pi(1)} \cdots \mathbf{e}_{\pi(m)}), \quad D = \begin{pmatrix} c_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & c_m \end{pmatrix}. \quad (8)$$

The matrix formed by stacking left-to-right the column vectors on the right-hand side of (7) is easily seen to satisfy:

$$(B c_1 \mathbf{e}_{\pi(1)} \cdots B c_m \mathbf{e}_{\pi(m)}) = B P D.$$

On the other hand, the columns $A\mathbf{e}_i$ form the matrix A . Taken together, therefore, equations (7) imply that $A = B P D$. As discussed in Remark 1, the remainder of (6) follows directly from this matrix identity and (4).

This finishes the proof of Theorem 1 with sparsity $k = 1$. Unfortunately, the proof for a general k is considerably more difficult. The main trouble is that for $k > 1$, it is nontrivial to produce P and D as in (8) using only assumption (4) and (5); this is where methods from combinatorics are needed.

In what follows, we will use the notation $[m]$ for the set $\{1, \dots, m\}$, and $\binom{[m]}{k}$ for the set of k -element subsets of $[m]$. Also, recall that $\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_\ell\}$ for real vectors $\mathbf{v}_1, \dots, \mathbf{v}_\ell$ is the vector space consisting of their \mathbb{R} -linear span:

$$\text{Span}\{\mathbf{v}_1, \dots, \mathbf{v}_\ell\} = \left\{ \sum_{i=1}^{\ell} t_i \mathbf{v}_i : t_1, \dots, t_\ell \in \mathbb{R} \right\}.$$

Finally, for a set $S \subseteq [m]$ and a matrix A with columns $\{A_1, \dots, A_m\}$, we define

$$\text{Span}\{A_S\} = \text{Span}\{A_s : s \in S\}.$$

A main ingredient in the proof of Theorem 1 is the following fact in combinatorial matrix theory. For a concrete instance of this argument when $k = 2$, see Example 1 below.

Proposition 1: Fix positive integers n and $k < m$. Let $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{n \times m}$ have columns $\{A_1, \dots, A_m\}$ and $\{B_1, \dots, B_m\}$. Suppose that A satisfies (4) and that

$$\alpha : \binom{[m]}{k} \rightarrow \binom{[m]}{k} \quad (9)$$

⁵The column vector \mathbf{e}_i has a 1 in its i th coordinate and zeroes elsewhere.

⁶*Injectivity* for a function $\pi : S \rightarrow T$ from a set S to a set T means that $\pi(s_1) = \pi(s_2)$ if and only if $s_1 = s_2$. The function π is *bijective* if in addition to being injective it is also *surjective*; that is, for each $t \in T$ there exists an $s \in S$ such that $\pi(s) = t$. A bijective function π has an inverse $\pi^{-1} : T \rightarrow S$ which satisfies $\pi(\pi^{-1}(t)) = t$ and $\pi^{-1}(\pi(s)) = s$ for all $s \in S$ and $t \in T$. When $S = T$ is a finite set, an injective function is also bijective; in this case, the function π is called a *permutation* of the set S .

is a map with the following property: for all $S \in \binom{[m]}{k}$,

$$\text{Span}\{A_S\} = \text{Span}\{B_{\alpha(S)}\}. \quad (10)$$

Then, there exists a permutation matrix $P \in \mathbb{R}^{m \times m}$ and an invertible diagonal matrix $D \in \mathbb{R}^{m \times m}$ such that $A = B P D$.

Proof: We shall induct on k , the base case $k = 1$ having already been worked out at the beginning of this section. We first prove that α is injective (and thus bijective). Suppose that $S_1, S_2 \in \binom{[m]}{k}$ are such that $\alpha(S_1) = \alpha(S_2)$; then by (10),

$$\text{Span}\{A_{S_1}\} = \text{Span}\{B_{\alpha(S_1)}\} = \text{Span}\{B_{\alpha(S_2)}\} = \text{Span}\{A_{S_2}\}.$$

In particular, using (13) from Lemma 1 below with $\ell = k$ and $M = A$, it follows that $S_1 = S_2$ and thus α is bijective. Moreover, from this bijectivity of α and the fact that every k columns of A are linearly independent, it follows that every k columns of B are also linearly independent.

We complete the proof, inductively, by producing a map:

$$\tau : \binom{[m]}{k-1} \rightarrow \binom{[m]}{k-1} \quad (11)$$

which satisfies

$$\text{Span}\{A_S\} = \text{Span}\{B_{\tau(S)}\} \text{ for } S \in \binom{[m]}{k-1}.$$

Let $\beta = \alpha^{-1}$ denote the inverse of α . Fix $S = \{i_1, \dots, i_{k-1}\} \in \binom{[m]}{k-1}$, and set $S_1 = S \cup \{r\}$ and $S_2 = S \cup \{s\}$ for some fixed $r, s \notin S$ with $r \neq s$ (so that $\beta(S_1) \neq \beta(S_2)$ by injectivity of β).⁷ Intersecting equations (10) with $S = \beta(S_1)$ and $S = \beta(S_2)$ and then applying identity (14) with $M = A$, it follows that

$$\text{Span}\{B_S, B_r\} \cap \text{Span}\{B_S, B_s\} = \text{Span}\{A_{\beta(S_1) \cap \beta(S_2)}\}. \quad (12)$$

Since the left-hand side of (12) is at least $k-1$ dimensional,⁸ the number of elements in the set $\beta(S_1) \cap \beta(S_2)$ is either $k-1$ or k . But $\beta(S_1) \neq \beta(S_2)$ so that $\beta(S_1) \cap \beta(S_2)$ consists of $k-1$ elements. Moreover,

$$\text{Span}\{B_S\} \subseteq \text{Span}\{A_{\beta(S_1) \cap \beta(S_2)}\}$$

implies that $\text{Span}\{B_S\} = \text{Span}\{A_{\beta(S_1) \cap \beta(S_2)}\}$.⁹

The association $S \mapsto \beta(S_1) \cap \beta(S_2)$ discussed above defines a function $\gamma : \binom{[m]}{k-1} \rightarrow \binom{[m]}{k-1}$ with the property that $\text{Span}\{B_S\} = \text{Span}\{A_{\gamma(S)}\}$. Finally, we show that γ is injective, which implies that $\tau = \gamma^{-1}$ is the map desired in (11) for the induction. If $\gamma(S) = \gamma(S')$, then $\text{Span}\{B_S\} = \text{Span}\{B_{S'}\}$. By (13) in Lemma 1 with $\ell = k-1$ and $M = B$, we have $S = S'$. Thus, γ is injective, finishing the proof. ■

The following elementary fact in linear algebra was used to prove Proposition 1 above.

Lemma 1: Let $M \in \mathbb{R}^{n \times m}$. If every set of $\ell+1$ columns of M are linearly independent, then for $S, S' \in \binom{[m]}{\ell}$,

$$\text{Span}\{M_S\} = \text{Span}\{M_{S'}\} \implies S = S'. \quad (13)$$

⁷Here we use the assumption that $k < m$ so that such a pair $r \neq s$ exists.

⁸Recall that we showed that every k columns of B is linearly independent.

⁹Use the following basic fact of linear algebra. If $U \subseteq V$ are two subspaces of a vector space W such that $\dim(U) = \dim(V)$ (i.e. they have the same vector-space dimension), then $U = V$.

If M satisfies condition (4) and $S_1, S_2 \in \binom{[m]}{k}$, then

$$\text{Span}\{M_{S_1 \cap S_2}\} = \text{Span}\{M_{S_1}\} \cap \text{Span}\{M_{S_2}\}. \quad (14)$$

Proof: To prove statement (13), suppose by way of contradiction that $S \neq S' \in \binom{[m]}{\ell}$ are such that $\text{Span}\{M_S\} = \text{Span}\{M_{S'}\}$. Then, without loss of generality, there is an $i \in S$ with $i \notin S'$. But $M_i \in \text{Span}\{M_{S'}\}$, which would imply that the ℓ columns of M determined by $S' \cup \{i\}$ are not linearly independent, a contradiction to the assumption on M .

We now prove (14). The inclusion \subseteq in (14) is trivial, so suppose $\mathbf{y} \in \text{Span}\{M_{S_1}\} \cap \text{Span}\{M_{S_2}\}$. Express \mathbf{y} as a linear combination of k columns of M indexed by S_1 and, separately, as a combination of k columns of M indexed by S_2 . By assumption (4), these linear combinations must be identical. In particular, \mathbf{y} was expressed as a linear combination of columns of M indexed by $S_1 \cap S_2$, and thus is in $\text{Span}\{M_{S_1 \cap S_2}\}$. ■

Example 1: We give a simple example of how the proof of Proposition 1 works in the case $n = m = 3$, $k = 2$. Suppose that $\alpha : \binom{[3]}{2} \rightarrow \binom{[3]}{2}$ is the (necessarily bijective) map:

$$\alpha(\{1, 2\}) = \{2, 3\}, \alpha(\{1, 3\}) = \{1, 2\}, \alpha(\{2, 3\}) = \{1, 3\}.$$

Following the proof of Proposition 1, one can check that

$$\gamma(\{1\}) = \{3\}, \gamma(\{2\}) = \{1\}, \gamma(\{3\}) = \{2\},$$

and thus we obtain the function $\tau = \gamma^{-1}$ as desired in (11). The resulting permutation P represents the cycle (123).

We now prove Theorem 1 by combining Theorem 5 from the Appendix with Proposition 1.

Proof of Theorem 1: It is sufficient to construct a map

$$\alpha : \binom{[m]}{k} \rightarrow \binom{[m]}{k}, \quad \{i_1, \dots, i_k\} \mapsto \{r_1, \dots, r_k\} \quad (15)$$

satisfying hypothesis (10) in Proposition 1.

Fix a finite subset $T \subset \mathbb{R}$. Also, fix $\{i_1, \dots, i_k\} \subseteq [m]$, and let $\mathbf{a} = t_1 \mathbf{e}_{i_1} + \dots + t_k \mathbf{e}_{i_k}$ for $t_i \in T$. Suppose that \mathbf{b} is k -sparse with $\mathbf{Aa} = \mathbf{Bb}$; then, $\mathbf{b} \in \text{Span}\{\mathbf{e}_{j_1}, \dots, \mathbf{e}_{j_k}\}$ for some $\{j_1, \dots, j_k\} \in \binom{[m]}{k}$. Viewing each k -element subset of $[m]$ as a *color*, this map

$$f : T^k \rightarrow \binom{[m]}{k}, \quad (t_1, \dots, t_k) \mapsto \{j_1, \dots, j_k\}$$

is a coloring of the finite set T^k with colors in $C = \binom{[m]}{k}$.

We now apply Ramsey theory. For sufficiently large sets T , Theorem 5 below guarantees that *regardless of the map* f there are subsets $H_1, \dots, H_k \subseteq T$ each having $s = 2$ elements and $\{r_1, \dots, r_k\} \in \binom{[m]}{k}$ such that $f(t_1, \dots, t_k) = \{r_1, \dots, r_k\}$ is constant for all $(t_1, \dots, t_k) \in H_1 \times \dots \times H_k$.

We claim that defining $\alpha(\{i_1, \dots, i_k\}) = \{r_1, \dots, r_k\}$ according to the above recipe gives a map satisfying (10). To verify the claim, we shall prove that

$$\text{Span}\{\mathbf{Ae}_{i_1}, \dots, \mathbf{Ae}_{i_k}\} \subseteq \text{Span}\{\mathbf{Be}_{r_1}, \dots, \mathbf{Be}_{r_k}\}, \quad (16)$$

which is inclusion \subseteq of (10). To see how the other inclusion follows from this, observe that the left-hand vector space in (16) is k -dimensional, while the right-hand space is at most k -dimensional; thus, equality holds.

To verify (16), we need to show that each \mathbf{Ae}_{i_ℓ} is in the right-hand span. To see this, consider a pair of elements in

$H_1 \times \dots \times H_k \subseteq \mathbb{R}^k$ which differ only in the ℓ th coordinate. By construction, the vectors $\mathbf{Aa}_1, \mathbf{Aa}_2$ corresponding to these two points have difference $\mathbf{A}(\mathbf{a}_1 - \mathbf{a}_2)$ that is a nonzero scalar multiple of \mathbf{Ae}_{i_ℓ} and is also in the right-hand span of (16). ■

We close this section by stating a generalization of Theorem 1 that allows for arbitrary receiver dimensions $p > m$. We omit the similar proof. For this generalization, we define a *column permutation matrix* to be a binary matrix having exactly one 1 in each column and at most one 1 in each row.

Theorem 3 (Overcomplete ACS Theorem): Fix positive integers n and $k < m < p$. There are k -sparse $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ with the following property: if $A \in \mathbb{R}^{n \times m}$ satisfies (4) and $B \in \mathbb{R}^{n \times p}$ and k -sparse $\mathbf{b}_1, \dots, \mathbf{b}_N \in \mathbb{R}^p$ are such that (5) hold, then there is an invertible diagonal matrix $D \in \mathbb{R}^{m \times m}$ and a column permutation $P \in \mathbb{R}^{p \times m}$ such that $A = BPD$.

Moreover, if the coding map $\mathbf{b}(\mathbf{y})$ satisfies $\mathbf{b}(t\mathbf{a}) \in \text{Span}\{\mathbf{b}(\mathbf{a})\}$ for $t \in \mathbb{R}$ and 1-sparse \mathbf{a} , then $\mathbf{b}_i = P\mathbf{D}\mathbf{a}_i$ and $PP^\top \mathbf{b}_i = \mathbf{b}_i$ for all $i = 1, \dots, N$.

Corollary 2 (ACS Efficiency): Under the assumptions of Theorem 3, there is a fixed set of $p - m$ coordinates of inferred sparse vectors \mathbf{b}_i which are always zero.

Proof: Let P be the column permutation matrix from the conclusion of Theorem 3. The number of rows of P that consist of all zeroes is $p - m$ since there are only m entries of P equal to 1. But then left multiplication by the matrix P on vectors $P^\top \mathbf{b}_i$ produces vectors with a fixed set of at least $p - m$ coordinates that are zero. ■

As an application of Theorem 3, consider its interpretation in the context of sender and receiver populations of neurons. Suppose that sparse dictionary learning has been successfully applied in a receiver region to decode a sender's signals. Then, Theorem 3 says that the receiver will recover the original codes up to natural equivalences, while Corollary 2 predicts that any extra neurons will become decoupled from the input stream and thus be free to be utilized for other purposes.

APPENDIX RAMSEY THEORY

We now explain how to prove Theorem 5 below; it was a crucial ingredient in the proof of Theorem 1. Its statement is very similar to a basic result in Ramsey theory [19, Theorem A]. For a recent survey of the field of Ramsey theory, see the article [25], and for a compilation of several applications to computer science and mathematics, see the paper [26].

Given positive integers c and s_1, \dots, s_c , the *Ramsey number* $R(s_1, \dots, s_c)$ is defined as the least integer R (if it exists) such that if the edges of the complete graph K_R on R vertices are colored with c colors, there is an $i \in [c]$ and a subset of s_i vertices of K_R all of whose pair-wise edges are the same color i . Ramsey's Theorem is then the statement that a finite $R(s_1, \dots, s_c)$ always exists. For instance, as pointed out near the end of the introduction, we have $R(3, 3) = 6$. To prove Theorem 1, however, we need the following modification of Ramsey's result.

Theorem 4 (Infinite version): Fix a finite set C of c colors and positive integers k, s . For all sufficiently large sets T_1, \dots, T_k , every coloring

$$f : T_1 \times \dots \times T_k \rightarrow C$$

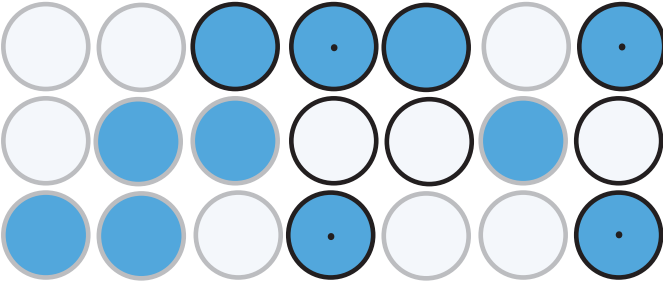


Fig. 2. **Iterated pigeon-hole principle.** How to find a 2×2 submatrix of the same color in a 2-coloring of a 3×7 grid (this is the $s = 2$, $c = 2$, $k = 2$ case of Theorem 5). In the top row of the figure above, there are 4 entries with the same color blue (they have black borders). Of the 4 entries directly below the blue from the top row, there are 3 of them that are colored white. Finally, below these lie 2 that are blue. The entries comprising the resulting 2×2 submatrix with the same color are shown with a black dot in their centers. Our example also shows that such a 2×2 submatrix need not exist in a 3×6 grid. Note that a direct application of Ramsey's theorem, as suggested by Remark 5, requires a grid of size 18×18 since $R(4, 4) = 18$. To give the reader some sense of the complexity of these questions, we remark that the number $R(5, 5)$ is not known although it is between 43 and 49 [28].

of the points in $T_1 \times \dots \times T_k$ obeys the following structural property: there are subsets $H_i \subseteq T_i$ of size s such that all points in $H_1 \times \dots \times H_k$ have the same color.

Remark 5: To see how Theorem 4 is related to Ramsey's theorem, consider when $k = 2$ and $c = 2$ (see Figure 2). In this case, one easily verifies that if $|T_1|, |T_2| \geq R(2s, 2s)$, there are subsets H_1, H_2 of size s as in the theorem statement.

Figure 2 shows how "sufficiently large" in Theorem 4 may be taken to be $|T_1| \geq 3$ and $|T_2| \geq 7$ in the case $c = k = s = 2$. Thus, defining $T(c, k, s)$ to be the smallest such product $|T_1| \dots |T_k|$ among sizes of sets T_i guaranteeing the conclusion of Theorem 4, we have $T(2, 2, 2) \leq 21$. More generally, the number $T(c, k, 2)$ will be important when we address below the question posed in Remark 2.

As Melody Chan points out [27], Theorem 4 can be deduced in the same manner as the standard inductive proof of Ramsey's Theorem. However, to best answer Remark 2, we desire more effective bounds than those offered by this argument.

Given positive integers k, s , and c , define numbers d_0, d_1, \dots, d_k recursively as follows:

$$d_i = s \cdot c^{2d_{i-1}}, \quad i = 1, \dots, k; \quad d_0 = \frac{1}{2}. \quad (17)$$

Notice that these (towers of) numbers grow very rapidly:

$$d_1 = sc, \quad d_2 = s \cdot c^{2sc}, \quad \text{and} \quad d_3 = s \cdot c^{2s \cdot c^{2sc}}.$$

A version of the following fact is likely known, but we include a proof for completeness.

Theorem 5 (Effective version): Fix a finite set C of c colors and positive integers k, s . If T_1, \dots, T_k are finite sets with sizes $|T_i| \geq d_i$, then for every coloring

$$f : T_1 \times \dots \times T_k \rightarrow C$$

of the points in $T_1 \times \dots \times T_k$, there are $H_i \subseteq T_i$ of size s such that all points in $H_1 \times \dots \times H_k$ have the same color.

Proof: First note that it is enough to prove the theorem for set sizes $|T_i| = d_i$ (by removing points as necessary). When $k = 1$, the proof boils down to the pigeon-hole principle: since

the range of f is finite of size c and since $|T_1| = d_1 = sc$, there must be at least $s = \frac{d_1}{c}$ points in $|T_1|$ which map to the same element of C . For expositional clarity, we only sketch the details of the proof when $k = 3$, as the ideas are the same in general. Since $c = 1$ is trivial, we shall assume $c > 1$.

Enumerate the elements of the three given sets as:

$$\begin{aligned} T_1 &= \{t_{11}, \dots, t_{1d_1}\}, \quad T_2 = \{t_{21}, \dots, t_{2d_2}\}, \\ T_3 &= \{t_{31}, \dots, t_{3d_3}\}. \end{aligned} \quad (18)$$

Consider the tree of height 2 with root t_{11} and children t_{21}, \dots, t_{2d_2} , each of whom have children t_{31}, \dots, t_{3d_3} . Since the number of children of t_{21} is d_3 , there is a subset $T'_3 \subseteq T_3$ with $\frac{d_3}{c}$ elements such that

$$f(t_{11}, t_{21}, T'_3) := \{f(t_{11}, t_{21}, t') : t' \in T'_3\} \subseteq C$$

is a single color in C . Consider now those children of t_{22} which are also members of T'_3 . As before, there is a subset $T''_3 \subseteq T'_3$ of size $\frac{d_3}{c^2}$ satisfying $|f(t_{11}, t_{22}, T''_3)| = 1$. Continuing in this manner a total of d_2 times, one produces a subset $T_3^{(1)} \subseteq T_3$ of size $\frac{d_3}{c^{d_2}}$ with the property that $|f(t_{11}, t_{2i}, T_3^{(1)})| = 1$ for all $i = 1, \dots, d_2$. Examine now the images $f(t_{11}, t_{2i}, T_3^{(1)})$ as i varies. As is easily seen, there is a subset $T_2^{(1)} \subseteq T_2$ of size $\frac{d_2}{c}$ such that $f(t_{11}, T_2^{(1)}, T_3^{(1)})$ consists of only 1 element.

The procedure in the previous paragraph constitutes the $j = 1^{\text{st}}$ round in a process of d_1 rounds that will produce the desired subsets H_1, H_2, H_3 . Set $T_i^{(0)} = T_i$ for each i . To move generally from round j to round $j + 1$, the idea is to consider the same tree as before but with root t_{1j} having children $T_2^{(j)}$ (each of which have children $T_3^{(j)}$), and to produce new subsets $T_3^{(j+1)} \subseteq T_3^{(j)}$ and $T_2^{(j+1)} \subseteq T_2^{(j)}$ with the same procedure as above. At the end of d_1 such rounds, we will have produced nested subsets

$$T_3^{(d_1)} \subseteq \dots \subseteq T_3^{(1)} \subseteq T_3^{(0)} \quad \text{and} \quad T_2^{(d_1)} \subseteq \dots \subseteq T_2^{(1)} \subseteq T_2^{(0)} \quad (19)$$

such that

$$|f(t_{1j}, T_2^{(j)}, T_3^{(j)})| = 1, \quad \text{for each } j = 1, \dots, d_1. \quad (20)$$

It is easy to check that after j such rounds of culling,

$$|T_2^{(j)}| = \frac{|T_2^{(j-1)}|}{c} \quad \text{and} \quad |T_3^{(j)}| = \frac{|T_3^{(j-1)}|}{c^{|T_2^{(j-1)}|}}. \quad (21)$$

A straightforward induction shows that eqs. (21) have solution:

$$|T_2^{(j)}| = \frac{d_2}{c^j} \quad \text{and} \quad |T_3^{(j)}| = \frac{d_3}{c^{d_2 \sum_{\ell=0}^{j-1} \frac{1}{c^\ell}}}. \quad (22)$$

Set $H_2 = T_2^{(d_1)}$ and $H_3 = T_3^{(d_1)}$. Since $d_1 = sc$, we know from (19) and (20) that there is a subset $H_1 \subseteq T_1$ of s elements such that $f(H_1, H_2, H_3)$ consists of only one color from C .

We are done, therefore, as long as H_2, H_3 have size at least s . For H_2 this is clear, while for H_3 we compute using (22):

$$|H_3| = |T_3^{(d_1)}| \geq \frac{d_3}{c^{d_2 \sum_{\ell=0}^{\infty} \frac{1}{c^\ell}}} = \frac{d_3}{c^{d_2 \frac{c}{c-1}}} \geq \frac{d_3}{c^{2d_2}} = s.$$

■

Remark 6: Clearly, the tower of numbers d_i defined by (17) can be decreased in Theorem 5 (see Figure 2), but we do not know by how much.

Corollary 3: We have the following bound:

$$T(c, k, s) \leq \prod_{i=1}^k d_i.$$

In the application of Theorem 5 for Theorem 1, we have $c = \binom{m}{k}$ and $s = 2$. Thus, the following is an upper bound on the the number N of \mathbf{a}_i used in our proof of Theorem 1:

$$N \leq \binom{m}{k} \cdot T\left(\binom{m}{k}, k, 2\right). \quad (23)$$

As a final remark, we note that the computational decision problem associated with finding $H_1 \times \cdots \times H_k$ all of the same color as in Theorem 5 is NP-complete when $c > 1$ (reduce to BICLIQUE and then apply the results of [29]).

ACKNOWLEDGMENT

The authors would like to thank the following people for helpful discussions: Charles Cadieu, Jack Culpepper, Mike DeWeese, Guy Isely, Amir Khosrowshahi, and Chris Rozell. We also thank Melody Chan for discussions on Ramsey theory and Matthias Mnich for explaining the NP-completeness lurking inside of Theorem 5.

REFERENCES

- [1] P. Comon, "Independent component analysis, a new concept?" *Signal processing*, vol. 36, no. 3, pp. 287–314, 1994.
- [2] A. Bell and T. Sejnowski, "Learning the higher-order structure of a natural sound," *Network: Computation in Neural Systems*, vol. 7, no. 2, pp. 261–266, 1996.
- [3] B. Olshausen and D. Field, "Emergence of simple-cell receptive field properties by learning a sparse code for natural images," *Nature*, vol. 381, no. 6583, pp. 607–609, 1996.
- [4] E. Smith and M. Lewicki, "Efficient auditory coding," *Nature*, vol. 439, no. 7079, pp. 978–982, 2006.
- [5] J. Yang, K. Yu, Y. Gong, and T. Huang, "Linear spatial pyramid matching using sparse coding for image classification," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009.
- [6] Y. Boureau, F. Bach, Y. LeCun, and J. Ponce, "Learning mid-level features for recognition," *Proc. International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [7] G. Hinton, "Connectionist learning procedures," *Artificial intelligence*, vol. 40, no. 1-3, pp. 185–234, 1989.
- [8] F. Attneave, "Informational aspects of visual perception," *Psychol. Rev.*, vol. 61, pp. 183–93, 1954.
- [9] H. B. Barlow, "Possible principles underlying the transformation of sensory messages," 1961.
- [10] S. Chen, D. Donoho, and M. Saunders, "Atomic decomposition by basis pursuit," *SIAM review*, vol. 43, p. 129, 2001.
- [11] W. Coulter, C. Hillar, G. Isley, and F. Sommer, "Adaptive compressed sensing: A new class of self-organizing coding models for neuroscience," in *Acoustics Speech and Signal Processing (ICASSP), 2010 IEEE International Conference on*. IEEE, pp. 5494–5497.
- [12] G. Isely, C. Hillar, and F. Sommer, "Deciphering subsampled data: adaptive compressive sampling as a principle of brain communication," *Advances in neural information processing systems*, 2010.
- [13] D. Donoho, "Compressed sensing," *Information Theory, IEEE Transactions on*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [14] E. Candes and T. Tao, "Decoding by linear programming," *Information Theory, IEEE Transactions on*, vol. 51, no. 12, pp. 4203–4215, 2005.
- [15] J. D. Blanchard, C. Cartis, and J. Tanner, "Compressed sensing: How sharp is the restricted isometry property?" *SIAM Review*, vol. 53, no. 1, pp. 105–125, 2011. [Online]. Available: <http://link.aip.org/link/?SIR/53/105/1>
- [16] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin, "A simple proof of the restricted isometry property for random matrices," *Constructive Approximation*, vol. 28, no. 3, pp. 253–263, 2008.
- [17] S. Gleichman and Y. Eldar, "Blind compressed sensing," *Information Theory, IEEE Transactions on*, vol. 57, no. 10, pp. 6958–6975, 2011.
- [18] M. Aharon, M. Elad, and A. Bruckstein, "On the uniqueness of overcomplete dictionaries, and a practical way to retrieve them," *Linear algebra and its applications*, vol. 416, no. 1, pp. 48–67, 2006.
- [19] F. P. Ramsey, "On a problem of formal logic," *Proceedings of the London Mathematical Society*, vol. s2-30, no. 1, pp. 264–286, 1930. [Online]. Available: <http://plms.oxfordjournals.org/content/s2-30/1/264.short>
- [20] R. Gribonval and K. Schnass, "Dictionary Identification–Sparse Matrix-Factorization via ℓ_1 -Minimization," *Information Theory, IEEE Transactions on*, vol. 56, no. 7, pp. 3523–3539, 2010.
- [21] Q. Geng, H. Wang, and J. Wright, "On the local correctness of l_1 -minimization for dictionary learning," *CoRR*, vol. abs/1101.5672, 2011.
- [22] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online learning for matrix factorization and sparse coding," *The Journal of Machine Learning Research*, vol. 11, pp. 19–60, 2010.
- [23] Y. Weiss, H. Chang, and W. Freeman, "Learning compressed sensing," in *Snowbird Learning Workshop, Allerton, CA*. Citeseer, 2007.
- [24] V. Abolghasemi, D. Jarchi, and S. Sanei, "A robust approach for optimization of the measurement matrix in compressed sensing," in *Cognitive Information Processing (CIP), 2010 2nd International Workshop on*. IEEE, pp. 388–392.
- [25] R. Graham, "Old and new problems and results in ramsey theory," *Horizons of Combinatorics*, pp. 105–118, 2008.
- [26] V. Rosta, "Ramsey theory applications," *the electronic journal of combinatorics*, pp. 1–43.
- [27] M. Chan, *Private communication*, 2011.
- [28] S. Radziszowski *et al.*, "Small ramsey numbers," *Electronic Journal of Combinatorics*, vol. 1, p. 28, 1994.
- [29] M. Dawande, P. Keskinocak, J. Swaminathan, and S. Tayur, "On bipartite and multipartite clique problems," *Journal of Algorithms*, vol. 41, no. 2, pp. 388–403, 2001.