

On the Global-Local Dichotomy in Sparsity Modeling

Dmitry Batenkov, Yaniv Romano, and Michael Elad

Abstract The traditional sparse modeling approach, when applied to inverse problems with large data such as images, essentially assumes a sparse model for small overlapping data patches and processes these patches as if they were independent from each other. While producing state-of-the-art results, this methodology is suboptimal, as it does not attempt to model the entire global signal in any meaningful way—a nontrivial task by itself.

In this paper we propose a way to bridge this theoretical gap by constructing a global model from the bottom-up. Given local sparsity assumptions in a dictionary, we show that the global signal representation must satisfy a constrained underdetermined system of linear equations, which forces the patches to agree on the overlaps. Furthermore, we show that the corresponding global pursuit can be solved via local operations. We investigate conditions for unique and stable recovery and provide numerical evidence corroborating the theory.

Keywords Sparse representations · Inverse problems · Convolutional sparse coding

D. Batenkov (✉)

Department of Mathematics, Massachusetts Institute of Technology, Cambridge, MA 02139, USA
e-mail: batenkov@mit.edu

Y. Romano

Department of Electrical Engineering, Technion - Israel Institute of Technology, 32000 Haifa, Israel

e-mail: yromano@tx.technion.ac.il

M. Elad

Department of Computer Science, Technion - Israel Institute of Technology, 32000 Haifa, Israel
e-mail: elad@cs.technion.ac.il

© Springer International Publishing AG 2017

H. Boche et al. (eds.), *Compressed Sensing and its Applications*,

Applied and Numerical Harmonic Analysis,

https://doi.org/10.1007/978-3-319-69802-1_1

1 Introduction

1.1 *The Need for a New Local-Global Sparsity Theory*

The sparse representation model [17] provides a powerful approach to various inverse problems in image and signal processing such as denoising [18, 37], deblurring [14, 57], and super-resolution [47, 56], to name a few [38]. This model assumes that a signal can be represented as a sparse linear combination of a few columns (called atoms) taken from a matrix termed dictionary. Given a signal, the sparse recovery of its representation over a dictionary is called sparse coding or pursuit (such as the orthogonal matching pursuit, OMP, or basis pursuit, BP). Due to computational and theoretical aspects, when treating high-dimensional data, various existing sparsity-inspired methods utilize local patched-based representations rather than the global ones, i.e., they divide a signal into small overlapping blocks (patches), reconstruct these patches using standard sparse recovery techniques, and subsequently average the overlapping regions [11, 17]. While this approach leads to highly efficient algorithms producing state-of-the-art results, the global signal prior remains essentially unexploited, potentially resulting in suboptimal recovery.

As an attempt to tackle this flaw, methods based on the notion of *structured sparsity* [19, 29, 30, 32, 55] started to appear; for example, in [14, 37, 47] the observation that a patch may have similar neighbors in its surroundings (often termed the self-similarity property) is injected to the pursuit, leading to improved local estimations. Another possibility to consider the dependencies between patches is to exploit the multi-scale nature of the signals [36, 40, 53]. A different direction is suggested by the expected patch log likelihood (EPLL) method [40, 52, 60], which encourages the patches of the final estimate (i.e., after the application of the averaging step) to comply with the local prior. Also, a related work [45, 46] suggests promoting the local estimations to agree on their shared content (the overlap) as a way to achieve a coherent reconstruction of the signal.

Recently, an alternative to the traditional patch-based prior was suggested in the form of the convolutional, or shift-invariant, sparse coding (CSC) model [10, 25, 27, 28, 49, 54]. Rather than dividing the image into local patches and processing each of these independently, this approach imposes a specific structure on the global dictionary—a concatenation of banded circulant matrices—and applies a global pursuit. A thorough theoretical analysis of this model was proposed very recently in [41, 42], providing a clear understanding of its success.

The empirical success of the above algorithms indicates the great potential of reducing the inherent gap that exists between the independent local processing of patches and the global nature of the signal at hand. However, a key and highly desirable part is still missing—a theory which would suggest how to modify the basic sparse model to take into account the mutual dependencies between the patches, what approximation methods to use, and how to efficiently design and learn the corresponding structured dictionary.

1.2 *Content and Organization of the Paper*

In this paper we propose a systematic investigation of the signals which are implicitly defined by local sparsity assumptions. A major theme in what follows is that the presence of patch overlaps reduces the number of degrees of freedom, which, in turn, has theoretical and practical implications. In particular, this allows more accurate estimates for uniqueness and stability of local sparse representations, as well as better bounds on performance of existing sparse approximation algorithms. Moreover, the global point of view allows for development of new pursuit algorithms, which consist of local operation on one hand, while also taking into account the patch overlaps on the other hand. Some aspects of the offered theory are still incomplete, and several exciting research directions emerge as well.

The paper is organized as follows. In Section 2 we develop the basic framework for signals which are patch-sparse, building the global model from the “bottom-up,” and discuss some theoretical properties of the resulting model. In Section 3 we consider the questions of reconstructing the representation vector and of denoising a signal in this new framework. We describe “globalized” greedy pursuit algorithms [43] for these tasks, where the patch disagreements play a major role. We show that the frequently used local patch averaging (LPA) approach is in fact suboptimal in this case. In Section 4 and [Appendix E: Generative Models for Patch-Sparse Signals](#), we describe several instances/classes of the local-global model in some detail, exemplifying the preceding definitions and results. The examples include piecewise constant signals, signature-type (periodic) signals, and more general bottom-up models. In Section 5 we present results of some numerical experiments, where in particular we show that one of the new globalized pursuits, inspired by the ADMM algorithm [9, 23, 24, 33], turns out to have superior performance in all the cases considered. We conclude the paper in Section 6 by discussing possible research directions.

2 **Local-Global Sparsity**

We start with the local sparsity assumptions for every patch and subsequently provide two complimentary characterizations of the resulting global signal space. On one hand, we show that the signals of interest admit a global “sparse-like” representation with a dictionary of convolutional type and with additional linear constraints on the representation vector. On the other hand, the signal space is in fact a union of linear subspaces, where each subspace is a kernel of a certain linear map. To complement and connect these points of view, in [Appendix E: Generative Models for Patch-Sparse Signals](#), we show that the original local dictionary must carry a combinatorial structure, and based on this structure, we develop a generative model for patch-sparse signals. Concluding this section, we provide some theoretical analysis of the properties of the resulting model, in particular uniqueness and

stability of representation. For this task, we define certain measures of the dictionary, similar to the classical spark, coherence function, and the restricted isometry property, which take the additional dictionary structure into account. In general, this additional structure implies possibly better uniqueness as well as stability to perturbations; however, it is an open question to show they are provably better in certain cases.

2.1 Preliminaries

Let $[m]$ denote the set $\{1, 2, \dots, m\}$. If D is an $n \times m$ matrix and $S \subset [m]$ is an index set, then D_S denotes the submatrix of D consisting of the columns indexed by S .

Definition 1 (Spark of a Matrix). Given a dictionary $D \in \mathbb{R}^{n \times m}$, the *spark* of D is defined as the minimal number of columns which are linearly dependent:

$$\sigma(D) := \min \{j : \exists S \subset [m], |S| = j, \text{rank } D_S < j\}. \quad (1)$$

Clearly $\sigma(D) \leq n + 1$.

Definition 2. Given a vector $\alpha \in \mathbb{R}^m$, the ℓ_0 pseudo-norm is the number of nonzero elements in α :

$$\|\alpha\|_0 := \#\{j : \alpha_j \neq 0\}.$$

Definition 3. Let $D \in \mathbb{R}^{n \times m}$ be a dictionary with normalized atoms. The μ_1 coherence function (Tropp's Babel function) is defined as

$$\mu_1(s) := \max_{i \in [m]} \max_{S \subset [m] \setminus \{i\}, |S|=s} \sum_{j \in S} |\langle d_i, d_j \rangle|.$$

Definition 4. Given a dictionary D as above, the restricted isometry constant of order k is the smallest number δ_k such that

$$(1 - \delta_k) \|\alpha\|_2^2 \leq \|D\alpha\|_2^2 \leq (1 + \delta_k) \|\alpha\|_2^2$$

for every $\alpha \in \mathbb{R}^m$ with $\|\alpha\|_0 \leq k$.

For any matrix M , we denote by $\mathcal{R}(M)$ the column space (range) of M .

2.2 Globalized Local Model

In what follows we treat one-dimensional signals $x \in \mathbb{R}^N$ of length N , divided into $P = N$ overlapping patches of equal size n (so that the original signal is thought

to be periodically extended). The other natural choice is $P = N - n + 1$, but for simplicity of derivations, we consider only the periodic case.

Let $R_1 := [I_{n \times n} \ \mathbf{0} \ \mathbf{0} \ \dots \ \mathbf{0}] \in \mathbb{R}^{n \times N}$, and for each $i = 2, \dots, P$, we define $R_i \in \mathbb{R}^{n \times N}$ to be the circular column shift of R_1 by $n \cdot (i - 1)$ entries, i.e., this operator extracts the i -th patch from the signal in a circular fashion.

Definition 5. Given local dictionary $D \in \mathbb{R}^{n \times m}$, sparsity level $s < n$, signal length N , and the number of overlapping patches P , the *globalized local sparse* model is the set

$$\mathcal{M} = \mathcal{M}(D, s, P, N) := \{x \in \mathbb{R}^N, R_i x = D\alpha_i, \|\alpha_i\|_0 \leq s \ \forall i = 1, \dots, P\}. \quad (2)$$

This model suggests that each patch, $R_i x$ is assumed to have an s -sparse representation α_i , and this way we have characterized the global x by describing the local nature of its patches.

Next we derive a “global” characterization of \mathcal{M} . Starting with the equations

$$R_i x = D\alpha_i, \quad i = 1, \dots, P,$$

and using the equality $I_{N \times N} = \frac{1}{n} \sum_{i=1}^P R_i^T R_i$, we have a representation

$$x = \frac{1}{n} \sum_{i=1}^P R_i^T R_i x = \sum_{i=1}^P \left(\frac{1}{n} R_i^T D \right) \alpha_i.$$

Let the global “convolutional” dictionary D_G be defined as the horizontal concatenation of the (vertically) shifted versions of $\frac{1}{n}D$, i.e., (see Figure 1 on page 5)

$$D_G := \left[\left(\frac{1}{n} R_i^T D \right) \right]_{i=1 \dots P} \in \mathbb{R}^{N \times mP}. \quad (3)$$

Let $\Gamma \in \mathbb{R}^{mP}$ denote the concatenation of the local sparse codes, i.e.,

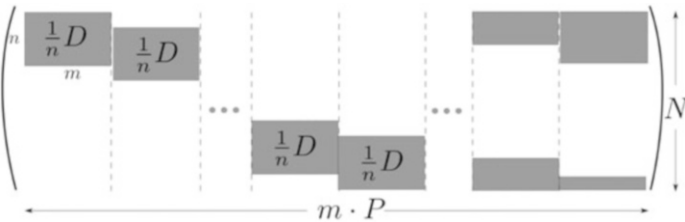


Fig. 1 The global dictionary D_G . After permuting the columns, the matrix becomes a union of circulant Toeplitz matrices, hence the term “convolutional”.

$$\Gamma := \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_P \end{bmatrix}.$$

Given a vector Γ as above, we will denote by \tilde{R}_i the operator of extracting its i -th portion,¹ i.e., $\tilde{R}_i \Gamma \equiv \alpha_i$.

Summarizing the above developments, we have the global convolutional representation for our signal as follows:

$$x = D_G \Gamma. \quad (4)$$

Next, applying R_i to both sides of (4) and using (2), we obtain

$$D\alpha_i = R_i x = R_i D_G \Gamma. \quad (5)$$

Let $\Omega_i := R_i D_G$ denote the i -th stripe from the global convolutional dictionary D_G . Thus (5) can be rewritten as

$$\underbrace{[\mathbf{0} \dots \mathbf{0} \ D \ \mathbf{0} \dots \mathbf{0}]}_{:=Q_i} \Gamma = \Omega_i \Gamma, \quad (6)$$

or $(Q_i - \Omega_i) \Gamma = 0$. Since this is true for all $i = 1, \dots, P$, we have shown that the vector Γ satisfies

$$\underbrace{\begin{bmatrix} Q_1 - \Omega_1 \\ \vdots \\ Q_P - \Omega_P \end{bmatrix}}_{:=M \in \mathbb{R}^{nP \times mP}} \Gamma = 0.$$

Thus, the condition that the patches $R_i x$ agree on the overlaps is equivalent to the global representation vector Γ residing in the null-space of the matrix M .

An easy computation provides the dimension of this null-space (see proof in [Appendix A: Proof of Lemma 1](#)), or in other words the overall number of degrees of freedom of admissible Γ .

¹Notice that while R_i extracts the i -th patch from the signal x , the operator \tilde{R}_i extracts the representation α_i of $R_i x$ from Γ .

Lemma 1. For any frame $D \in \mathbb{R}^{n \times m}$ (i.e., a full rank dictionary), we have

$$\dim \ker M = N(m - n + 1).$$

Note that in particular for $m = n$, we have $\dim \ker M = N$, and since in this case D is invertible, we have $R_i x = D \alpha_i$ where $\alpha_i = D^{-1} R_i x$, so that every signal admits a unique representation $x = D_G \Gamma$ with $\Gamma = (D^{-1} R_1 x, \dots, D^{-1} R_P x)^T$.

As we shall demonstrate now, the equation $M \Gamma = 0$ represents the requirement that the local sparse codes $\{\alpha_i\}$ are not independent but rather should be such that the corresponding patches $D \alpha_i$ agree on the overlaps.

Definition 6. Define the “extract from top/bottom” operators $S_T \in \mathbb{R}^{(n-1) \times n}$ and $S_B \in \mathbb{R}^{(n-1) \times n}$:

$$S_{T(op)} = [I_{n-1} \ \mathbf{0}], \quad S_{B(bottom)} = [\mathbf{0} \ I_{n-1}].$$

The following result is proved in [Appendix B: Proof of Lemma 2](#).

Lemma 2. Let $\Gamma = [\alpha_1, \dots, \alpha_P]^T$. Under the above definitions, the following are equivalent:

1. $M \Gamma = 0$;
2. For each $i = 1, \dots, P$, we have $S_B D \alpha_i = S_T D \alpha_{i+1}$.

Definition 7. Given $\Gamma = [\alpha_1, \dots, \alpha_P]^T \in \mathbb{R}^{mP}$, the $\|\cdot\|_{0,\infty}$ pseudo-norm is defined by

$$\|\Gamma\|_{0,\infty} := \max_{i=1,\dots,P} \|\alpha_i\|_0.$$

Thus, every signal complying with the patch-sparse model, with sparsity s for each patch, admits the following representation.

Theorem 1. Given D, s, P , and N , the globalized local sparse model (2) is equivalent to

$$\begin{aligned} \mathcal{M} &= \{x \in \mathbb{R}^N : x = D_G \Gamma, M \Gamma = 0, \|\Gamma\|_{0,\infty} \leq s\} \\ &= \{x \in \mathbb{R}^N : x = D_G \Gamma, M_* \Gamma = 0, \|\Gamma\|_{0,\infty} \leq s\}, \end{aligned} \quad (7)$$

where the matrix $M_* \in \mathbb{R}^{(n-1)P \times mP}$ is defined as

$$M_* := \begin{bmatrix} S_B D & -S_T D & & & \\ & S_B D & -S_T D & & \\ & & & \ddots & \ddots \\ & & & & S_B D & -S_T D \end{bmatrix}.$$

Proof. If $x \in \mathcal{M}$ (according to (2)), then by the above construction x belongs to the set defined by the RHS of (7) (let's call it \mathcal{M}^* for the purposes of this proof only). In the other direction, assume that $x \in \mathcal{M}^*$. Now $R_i x = R_i D_G \Gamma = \Omega_i \Gamma$, and since $M\Gamma = 0$, we have $R_i x = Q_i \Gamma = D\tilde{R}_i \Gamma$. Denote $\alpha_i := \tilde{R}_i \Gamma$, and so we have that $R_i x = D\alpha_i$ with $\|\alpha_i\|_0 \leq s$, i.e., $x \in \mathcal{M}$ by definition. The second part follows from Lemma 2. \square

We say that α_i is a *minimal* representation of x_i if $x_i = D\alpha_i$ such that the matrix $D_{\text{supp } \alpha_i}$ has full rank—and therefore the atoms participating in the representation are linearly independent.²

Definition 8. Given a signal $x \in \mathcal{M}$, let us denote by $\rho(x)$ the set of all locally sparse and minimal representations of x :

$$\rho(x) := \left\{ \Gamma \in \mathbb{R}^{mP} : \|\Gamma\|_{0,\infty} \leq s, x = D_G \Gamma, M\Gamma = 0, D_{\text{supp } \tilde{R}_i \Gamma} \text{ is full rank} \right\}.$$

Let us now go back to the definition (2). Consider a signal $x \in \mathcal{M}$, and let $\Gamma \in \rho(x)$. Denote $S_i := \text{supp } \tilde{R}_i \Gamma$. Then we have $R_i x \in \mathcal{R}(D_{S_i})$, and therefore we can write $R_i x = P_{S_i} R_i x$, where P_{S_i} is the orthogonal projection operator onto $\mathcal{R}(D_{S_i})$. In fact, since D_{S_i} is full rank, we have $P_{S_i} = D_{S_i} D_{S_i}^\dagger$ where $D_{S_i}^\dagger = (D_{S_i}^T D_{S_i})^{-1} D_{S_i}^T$ is the Moore-Penrose pseudoinverse of D_{S_i} .

Definition 9. Given a support sequence $\mathcal{S} = (S_1, \dots, S_P)$, define the matrix $A_{\mathcal{S}}$ as follows:

$$A_{\mathcal{S}} := \begin{bmatrix} (I_n - P_{S_1}) R_1 \\ (I_n - P_{S_2}) R_2 \\ \vdots \\ (I_n - P_{S_P}) R_P \end{bmatrix} \in \mathbb{R}^{nP \times N}.$$

The map $A_{\mathcal{S}}$ measures the local patch discrepancies, i.e., how “far” is each local patch from the range of a particular subset of the columns of D .

Definition 10. Given a model \mathcal{M} , denote by $\Sigma_{\mathcal{M}}$ the set of all valid supports, i.e.,

$$\Sigma_{\mathcal{M}} := \{(S_1, \dots, S_P) : \exists x \in \mathcal{M}, \Gamma \in \rho(x) \text{ s.t. } \forall i = 1, \dots, P : S_i = \text{supp } \tilde{R}_i \Gamma\}.$$

With this notation in place, it is immediate to see that the global signal model is a union of subspaces.

Theorem 2. *The global model is equivalent to the union of subspaces*

$$\mathcal{M} = \bigcup_{\mathcal{S} \in \Sigma_{\mathcal{M}}} \ker A_{\mathcal{S}}.$$

²Notice that α_i might be a minimal representation but not a unique one with minimal sparsity. For discussion of uniqueness, see Subsection 2.3.

Remark 1. Contrary to the well-known union of subspaces model [7, 35], the subspaces $\{\ker A_{\mathcal{S}}\}$ do not have in general a sparse joint basis, and therefore our model is distinctly different from the well-known block-sparsity model [19, 20].

An important question of interest is to estimate $\dim \ker A_{\mathcal{S}}$ for a given $\mathcal{S} \in \Sigma_{\mathcal{M}}$. One possible solution is to investigate the “global” structure of the corresponding signals (as is done in Subsection 4.1 and Subsection 4.2), while another option is to utilize information about “local connections” (Appendix E: Generative Models for Patch-Sparse Signals).

2.3 Uniqueness and Stability

Given a signal $x \in \mathcal{M}$, it has a globalized representation $\Gamma \in \rho(x)$ according to Theorem 1. When is such a representation unique, and under what conditions can it be recovered when the signal is corrupted with noise?

In other words, we study the problem

$$\min \|\Gamma\|_{0,\infty} \quad \text{s.t. } D_G \Gamma = D_G \Gamma_0, M\Gamma = 0 \quad (P_{0,\infty})$$

and its noisy version

$$\min \|\Gamma\|_{0,\infty} \quad \text{s.t. } \|D_G \Gamma - D_G \Gamma_0\| \leq \varepsilon, M\Gamma = 0 \quad (P_{0,\infty}^\varepsilon).$$

For this task, we define certain measures of the dictionary, similar to the classical spark, coherence function, and the restricted isometry property, which take the additional dictionary structure into account. In general, the additional structure implies *possibly* better uniqueness as well as stability to perturbations; however, it is an open question to show they are *provably* better in certain cases.

The key observation is that the global model \mathcal{M} imposes a constraint on the allowed local supports.

Definition 11. Denote the set of allowed local supports by

$$\mathcal{T} := \{T : \exists (S_1, \dots, T, \dots, S_P) \in \Sigma_{\mathcal{M}}\}.$$

Recall the definition of the spark (1). Clearly $\sigma(D)$ can be equivalently rewritten as

$$\sigma(D) = \min \{j : \exists S_1, S_2 \subset [m], |S_1 \cup S_2| = j, \text{rank } D_{S_1 \cup S_2} < j\}. \quad (8)$$

Definition 12. The *globalized spark* $\sigma^*(D)$ is

$$\sigma^*(D) := \min \{j : \exists S_1, S_2 \in \mathcal{T}, |S_1 \cup S_2| = j, \text{rank } D_{S_1 \cup S_2} < j\}. \quad (9)$$

The following proposition is immediate by comparing (8) with (9).

Proposition 1. $\sigma^*(D) \geq \sigma(D)$.

The globalized spark provides a uniqueness result in the spirit of [15].

Theorem 3 (Uniqueness). *Let $x \in \mathcal{M}(D, s, N, P)$. If there exists $\Gamma \in \rho(x)$ for which $\|\Gamma\|_{0,\infty} < \frac{1}{2}\sigma^*(D)$ (i.e., it is a sufficiently sparse solution of $P_{0,\infty}$), then it is the unique solution (and so $\rho(x) = \{\Gamma\}$).*

Proof. Suppose that there exists $\Gamma_0 \in \rho(x)$ which is different from Γ . Put $\Gamma_1 := \Gamma - \Gamma_0$, then $\|\Gamma_1\|_{0,\infty} < \sigma^*(D)$, while $D_G\Gamma_1 = 0$ and $M\Gamma_1 = 0$. Denote $\beta_j := \tilde{R}_j\Gamma_1$. By assumption, there exists an index i for which $\beta_i \neq 0$, but we must have $D\beta_j = 0$ for every j , and therefore $D_{\text{supp}\beta_i}$ must be rank-deficient—contradicting the fact that $\|\beta_i\| < \sigma^*(D)$. \square

In classical sparsity, we have the bound

$$\sigma(D) \geq \min \{s : \mu_1(s-1) \geq 1\}, \quad (10)$$

where μ_1 is given by Definition 3. In a similar fashion, the globalized spark σ^* can be bounded by an appropriate analog of “coherence”—however, computing this new coherence appears to be in general intractable.

Definition 13. Given the model \mathcal{M} , we define the following globalized coherence function

$$\mu_1^*(s) := \max_{S \in \mathcal{T} \cup \mathcal{T}, |S|=s} \max_{j \in S} \sum_{k \in S \setminus \{j\}} |\langle d_j, d_k \rangle|,$$

where $\mathcal{T} \cup \mathcal{T} := \{S_1 \cup S_2 : S_1, S_2 \in \mathcal{T}\}$.

Theorem 4. *The globalized spark σ^* can be bounded by the globalized coherence as follows³:*

$$\sigma^*(D) \geq \min \{s : \mu_1^*(s) \geq 1\}.$$

Proof. Following closely the corresponding proof in [15], assume by contradiction that

$$\sigma^*(D) < \min \{s : \mu_1^*(s) \geq 1\}.$$

Let $S^* \in \mathcal{T} \cup \mathcal{T}$ with $|S^*| = \sigma^*(D)$ for which D_{S^*} is rank-deficient. Then the restricted Gram matrix $G := D_{S^*}^T D_{S^*}$ must be singular. On the other hand, $\mu_1^*(|S^*|) < 1$, and so in particular

$$\max_{j \in S^*} \sum_{k \in S^* \setminus \{j\}} |\langle d_j, d_k \rangle| < 1.$$

³In general $\min \{s : \mu_1^*(s-1) \geq 1\} \neq \max \{s : \mu_1^*(s) < 1\}$ because the function μ_1^* need not be monotonic.

But that means that G is diagonally dominant and therefore $\det G \neq 0$, a contradiction. \square

We see that $\mu_1^*(s+1) \leq \mu_1(s)$ since the outer maximization is done on a smaller set. Therefore, in general the bound of Theorem 4 appears to be sharper than (10).

A notion of globalized RIP can also be defined as follows.

Definition 14. The globalized RIP constant of order k associated to the model \mathcal{M} is the smallest number $\delta_{k,\mathcal{M}}$ such that

$$(1 - \delta_{k,\mathcal{M}}) \|\alpha\|_2^2 \leq \|D\alpha\|_2^2 \leq (1 + \delta_{k,\mathcal{M}}) \|\alpha\|_2^2$$

for every $\alpha \in \mathbb{R}^m$ with $\text{supp } \alpha \in \mathcal{T}$.

Immediately one can see the following (recall Definition 4).

Proposition 2. *The globalized RIP constant is upper bounded by the standard RIP constant:*

$$\delta_{k,\mathcal{M}} \leq \delta_k.$$

Definition 15. The generalized RIP constant of order k associated to signals of length N is the smallest number $\delta_k^{(N)}$ such that

$$(1 - \delta_k^{(N)}) \|\Gamma\|_2^2 \leq \|D_G \Gamma\|_2^2 \leq (1 + \delta_k^{(N)}) \|\Gamma\|_2^2$$

for every $\Gamma \in \mathbb{R}^{mN}$ satisfying $M\Gamma = 0$, $\|\Gamma\|_{0,\infty} \leq k$.

Proposition 3. *We have*

$$\delta_k^{(N)} \leq \frac{\delta_{k,\mathcal{M}} + (n-1)}{n} \leq \frac{\delta_k + (n-1)}{n}.$$

Proof. Obviously it is enough to show only the leftmost inequality. If $\Gamma = (\alpha_i)_{i=1}^N$ and $\|\Gamma\|_{0,\infty} \leq k$, this gives $\|\alpha_i\|_0 \leq k$ for all $i = 1, \dots, N$. Further, setting $x := D_G \Gamma$ we clearly have $\Gamma \in \rho(x)$ and so $\text{supp } \Gamma \in \Sigma_{\mathcal{M}}$. Thus $\text{supp } \alpha_i \in \mathcal{T}$, and therefore

$$(1 - \delta_{k,\mathcal{M}}) \|\alpha_i\|_2^2 \leq \|D\alpha_i\|_2^2 \leq (1 + \delta_{k,\mathcal{M}}) \|\alpha_i\|_2^2.$$

By Corollary 3 we know that for every Γ satisfying $M\Gamma = 0$, we have

$$\|D_G \Gamma\|_2^2 = \frac{1}{n} \sum_{i=1}^N \|D\alpha_i\|_2^2.$$

Now for the lower bound,

$$\begin{aligned} \|D_G \Gamma\|_2^2 &\geq \frac{1 - \delta_{k, \mathcal{M}}}{n} \sum_{i=1}^N \|\alpha_i\|_2^2 = \left(1 - 1 + \frac{1 - \delta_{k, \mathcal{M}}}{n}\right) \|\Gamma\|_2^2 \\ &= \left(1 - \frac{\delta_{k, \mathcal{M}} + (n-1)}{n}\right) \|\Gamma\|_2^2. \end{aligned}$$

For the upper bound,

$$\begin{aligned} \|D_G \Gamma\|_2^2 &\leq \frac{1 + \delta_{k, \mathcal{M}}}{n} \sum_{i=1}^N \|\alpha_i\|_2^2 < \left(1 + \frac{\delta_{k, \mathcal{M}} + 1}{n}\right) \|\Gamma\|_2^2 \\ &\leq \left(1 + \frac{\delta_{k, \mathcal{M}} + (n-1)}{n}\right) \|\Gamma\|_2^2. \end{aligned}$$

□

Theorem 5 (Uniqueness and Stability of $P_{0, \infty}$ via RIP). *Suppose that $\delta_{2s}^{(N)} < 1$, and suppose further that $x = D_G \Gamma_0$ with $\|\Gamma_0\|_{0, \infty} = s$ and $\|D_G \Gamma_0 - x\|_2 \leq \varepsilon$. Then every solution $\hat{\Gamma}$ of the noise-constrained $P_{0, \infty}^\varepsilon$ problem*

$$\hat{\Gamma} \leftarrow \arg \min_{\Gamma} \|\Gamma\|_{0, \infty} \text{ s.t. } \|D_G \Gamma - x\| \leq \varepsilon, \quad M\Gamma = 0$$

satisfies

$$\|\hat{\Gamma} - \Gamma_0\|_2^2 \leq \frac{4\varepsilon^2}{1 - \delta_{2s}^{(N)}}.$$

In particular, Γ_0 is the unique solution of the noiseless $P_{0, \infty}$ problem.

Proof. Immediate using the definition of the globalized RIP:

$$\begin{aligned} \|\hat{\Gamma} - \Gamma_0\|_2^2 &< \frac{1}{1 - \delta_{2s}^{(N)}} \|D_G (\hat{\Gamma} - \Gamma_0)\|_2^2 \leq \frac{1}{1 - \delta_{2s}^{(N)}} \left(\|D_G \hat{\Gamma} - x\|_2 + \|D_G \Gamma_0 - x\|_2 \right)^2 \\ &\leq \frac{4\varepsilon^2}{1 - \delta_{2s}^{(N)}}. \end{aligned}$$

□

3 Pursuit Algorithms

In this section we consider the problem of efficient projection onto the model \mathcal{M} . First we treat the ‘‘oracle’’ setting, i.e., when the supports of the local patches (and therefore of the global vector Γ) are known. We show that the local patch averaging

(LPA) method is not a good projector; however, repeated application of it does achieve the desired result.

For the non-oracle setting, we consider “local” and “globalized” pursuits. The former type does not use any dependencies between the patches, and tries to reconstruct the supports α_i completely locally, using standard methods such as OMP—and as we demonstrate, it can be guaranteed to succeed in more cases than the standard analysis would imply. However a possibly better alternative exists, namely, a “globalized” approach with the patch disagreements as a major driving force.

3.1 *Global (Oracle) Projection, Local Patch Averaging (LPA) and the Local-Global Gap*

Here we briefly consider the question of efficient projection onto the subspace $\ker A_{\mathcal{S}}$, given \mathcal{S} .

As customary in the literature [12], the projector onto $\ker A_{\mathcal{S}}$ can be called *an oracle*. In effect, we would like to compute

$$x_G(y, \mathcal{S}) := \arg \min_x \|y - x\|_2^2 \quad \text{s.t. } A_{\mathcal{S}}x = 0, \quad (11)$$

given $y \in \mathbb{R}^N$.

To make things concrete, let us assume the standard Gaussian noise model:

$$y = x + \mathcal{N}(0, \sigma^2 I), \quad (12)$$

and let the mean squared error (MSE) of an estimator $f(y)$ of x be defined as usual, i.e., $MSE(f) := \mathbb{E}\|f(y) - x\|_2^2$. The following is well-known.

Proposition 4. *In the Gaussian noise model (12), the performance of the oracle estimator (11) is*

$$MSE(x_G) = (\dim \ker A_{\mathcal{S}}) \sigma^2.$$

Let us now turn to the local patch averaging (LPA) method. This approach suggests denoising an input signal by (i) breaking it into overlapping patches, (ii) denoising each patch independently, followed by (iii) averaging the local reconstructions to form the global signal estimate. The local denoising step is done by solving pursuit problems, estimating the local supports S_i , while the averaging step is the solution to the minimization problem:

$$\hat{x} = \arg \min_x \sum_{i=1}^P \|R_i x - P_{S_i} R_i y\|_2^2,$$

where y is the noisy signal. This has a closed-form solution:

$$\hat{x}_{LPA} = \left(\sum_i R_i^T R_i \right)^{-1} \left(\sum_i R_i^T P_{S_i} R_i \right) y = \underbrace{\left(\frac{1}{n} \sum_i R_i^T P_{S_i} R_i \right)}_{:=M_A} y. \quad (13)$$

Again, the following fact is well-established.

Proposition 5. *In the Gaussian noise model (12), the performance of the averaging estimator (13) is*

$$MSE(\hat{x}_{LPA}) = \sigma^2 \sum_{i=1}^N \lambda_i,$$

where $\{\lambda_1, \dots, \lambda_N\}$ are the eigenvalues of $M_A M_A^T$.

Thus, there exists a *local-global gap* in the oracle setting, illustrated in Figure 2 on page 14. In Subsection 4.1 we estimate this gap for a specific case of piecewise constant signals.

The following result is proved in [Appendix C: Proof of Theorem 6](#).

Theorem 6. *For any \mathcal{S} , we have*

$$\lim_{k \rightarrow \infty} M_A^k = P_{\ker A_{\mathcal{S}}},$$

where $P_{\ker A_{\mathcal{S}}}$ is the orthogonal projector onto $\ker A_{\mathcal{S}}$. Therefore for any y , iterations of (13) starting at y converge to $x_G(y)$ with a linear rate.

From the proof it is evident that the rate of convergence depends on the eigenvalues of M_A (which turn out to be related to the singular values of $A_{\mathcal{S}}$). Analyzing these

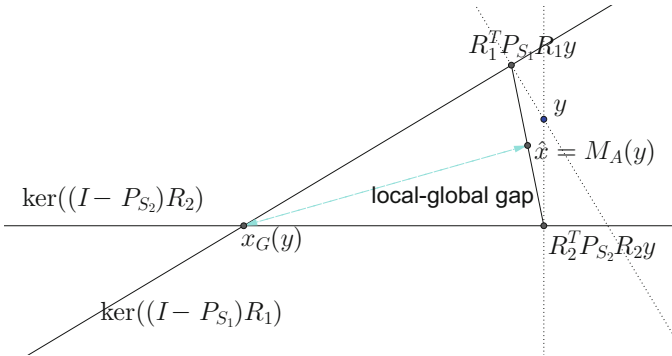


Fig. 2 The local-global gap, oracle setting. Illustration for the case $P = 2$. In details, the noisy signal y can be either projected onto $\ker A_{\mathcal{S}}$ (the point $x_G(y)$) or by applying the LPA (the point $\hat{x} = M_A(y)$). The difference between those two is the local-global gap, which can be significant.

eigenvalues (and therefore the convergence rate) appears to be a difficult problem for general \mathcal{M} and \mathcal{S} . In Theorem 8 we show one example where we consider the related problem of estimating the sum $\sum_{i=1}^N \lambda_i$ appearing in Proposition 5, in the case of the piecewise constant model (providing estimates for the local-global gap as well).

To conclude, we have shown that *the iterated LPA algorithm provides an efficient method for computing the global oracle projection x_G .*

3.2 Local Pursuit Guarantees

Now we turn to the question of projection onto the model \mathcal{M} when the support of Γ is not known.

Here we show that running OMP [13, 43] on each patch extracted from the signal in fact succeeds in more cases than can be predicted by the classical unconstrained sparse model for each patch. We use the modified coherence function (which is unfortunately intractable to compute):

$$\eta_1^*(s) := \max_{S \in \mathcal{T}} \left(\max_{j \in S} \sum_{k \in S \setminus \{j\}} |\langle d_k, d_j \rangle| + \max_{j \notin S} \sum_{k \in S} |\langle d_k, d_j \rangle| \right).$$

The proof of the following theorem is very similar to proving the guarantee for the standard OMP via the Babel function (Definition 3); see e.g., [22, Theorem 5.14]—and therefore we do not reproduce it here.

Theorem 7. *If $\eta_1^*(s) < 1$, then running OMP on each patch extracted from any $x \in \mathcal{M}$ will recover its true support.*

Since the modified coherence function takes the allowed local supports into consideration, one can readily conclude that

$$\eta_1^*(s) \leq \mu_1(s) + \mu_1(s-1),$$

and therefore Theorem 7 gives in general a possibly better guarantee than the one based on μ_1 .

3.3 Globalized Pursuits

We now turn to consider several pursuit algorithms, aiming at solving the $P_{0,\infty}/P_{0,\infty}^\epsilon$ problems, in the globalized model. The main question is how to project the patch supports onto the nonconvex set $\Sigma_{\mathcal{M}}$.

The core idea is to relax the constraint $M_*\Gamma = 0$, $\|\Gamma\|_{0,\infty} \leq s$ and allow for some patch disagreements, so that the term $\|M_*\Gamma_k\|$ is not exactly zero. Intuitive explanation is as follows: the disagreement term “drives” the pursuit, and the probability of success is higher because we only need to “jump-start” it with the first patch, and then by strengthening the weight of the penalty related to this constraint, the supports will “align” themselves correctly. Justifying this intuition, at least in some cases, is a future research goal.

3.3.1 Q-OMP

Given $\beta > 0$, we define

$$Q_\beta := \begin{bmatrix} D_G \\ \beta M_* \end{bmatrix}.$$

The main idea of the Q-OMP algorithm is to substitute the matrix Q_β as a proxy for the constraint $M_*\Gamma = 0$, by plugging it as a dictionary to the OMP algorithm. Then, given the obtained support \mathcal{S} , as a way to ensure that this constraint is met, one can construct the matrix $A_{\mathcal{S}}$ and project the signal onto the subspace $\ker A_{\mathcal{S}}$ (in Subsection 3.1 we show how such a projection can be done efficiently). The Q-OMP algorithm is detailed in Algorithm 1. Let us reemphasize the point that various values of β correspond to different weightings of the model constraint $M_*\Gamma = 0$ and this might possibly become useful when considering relaxed models (see Section 6).

Algorithm 1 The Q-OMP algorithm—a globalized pursuit

Given: noisy signal y , dictionary D , local sparsity s , parameter $\beta > 0$

1. Construct the matrix Q_β .
 2. Run the OMP algorithm on the vector $Y := \begin{bmatrix} y \\ \mathbf{0} \end{bmatrix}$, with the dictionary Q_β and sparsity sN . Obtain the global support vector $\hat{\Gamma}$ with $\text{supp } \hat{\Gamma} = \mathcal{S}$.
 3. Construct the matrix $A_{\mathcal{S}}$ and project y onto $\ker A_{\mathcal{S}}$.
-

3.3.2 ADMM-Inspired Approach

In what follows we extend the above idea and develop an ADMM-inspired pursuit [9, 23, 24, 33].

We start with the following global objective:

$$\hat{x} \leftarrow \arg \min_x \|y - x\|_2^2 \quad \text{s.t. } x = D_G\Gamma, M_*\Gamma = 0, \|\Gamma\|_{0,\infty} \leq K.$$

Clearly, it is equivalent to $\hat{x} = D_G \hat{\Gamma}$, where

$$\hat{\Gamma} \leftarrow \arg \min_{\Gamma} \|y - D_G \Gamma\|_2^2 \quad \text{s.t. } M_* \Gamma = 0, \|\Gamma\|_{0,\infty} \leq K. \quad (14)$$

Applying Corollary 3, we have the following result.

Proposition 6. *The following problem is equivalent to (14):*

$$\begin{aligned} \hat{\Gamma} \leftarrow \arg \min_{\{\alpha_i\}} \sum_{i=1}^P \|R_i y - D \alpha_i\|_2^2 \\ \text{s.t. } S_B D \alpha_i = S_T D \alpha_{i+1} \text{ and } \|\alpha_i\|_0 < K \text{ for } i = 1, \dots, P. \end{aligned} \quad (15)$$

We propose to approximate solution of the nonconvex problem (15) as follows. Define new variables z_i (which we would like to be equal to α_i eventually), and rewrite the problem by introducing the following variable splitting (here Z is the concatenation of all the z_i 's):

$$\{\hat{\Gamma}, \hat{Z}\} \leftarrow \arg \min_{\Gamma, Z} \sum_{i=1}^P \|R_i y - D \alpha_i\|_2^2 \quad \text{s.t. } S_B D \alpha_i = S_T D z_{i+1}, \alpha_i = z_i, \|\alpha_i\|_0 \leq K.$$

The constraints can be written in concise form

$$\underbrace{\begin{bmatrix} I \\ S_B D \end{bmatrix}}_{:=A} \alpha_i = \underbrace{\begin{bmatrix} I & 0 \\ 0 & S_T D \end{bmatrix}}_{:=B} \begin{pmatrix} z_i \\ z_{i+1} \end{pmatrix},$$

and so globally we would have the following structure (for $N = 3$)

$$\underbrace{\begin{bmatrix} A \\ A \\ A \end{bmatrix}}_{:=\tilde{A}} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{pmatrix} = \underbrace{\begin{bmatrix} I & & & \\ & S_T D & & \\ & & I & \\ & & & S_T D \\ S_T D & & & & I \end{bmatrix}}_{:=\tilde{B}} \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix}$$

Our ADMM-inspired method is defined in Algorithm 2.

Algorithm 2 The ADMM-inspired pursuit for $P_{0,\infty}^s$.

Given: noisy signal y , dictionary D , local sparsity s , parameter $\rho > 0$. The augmented Lagrangian is

$$L_\rho(\{\alpha_i\}, \{z_i\}, \{u_i\}) = \sum_{i=1}^P \|R_i y - D\alpha_i\|_2^2 + \frac{\rho}{2} \sum_{i=1}^P \|A\alpha_i - B \begin{pmatrix} z_i \\ z_{i+1} \end{pmatrix} + u_i\|_2^2.$$

1. Repeat until convergence:

a. Minimization wrt $\{\alpha_i\}$ is a batch-OMP:

$$\alpha_i^{k+1} \leftarrow \arg \min_{\alpha_i} \|R_i y - D\alpha_i\|_2^2 + \frac{\rho}{2} \|A\alpha_i - B \begin{pmatrix} z_i^k \\ z_{i+1}^k \end{pmatrix} + u_i^k\|_2^2, \quad s.t. \|\alpha_i\|_0 \leq K$$

$$\alpha_i^{k+1} \leftarrow OMP \left(\tilde{D} = \begin{bmatrix} D \\ \sqrt{\frac{\rho}{2}} A \end{bmatrix}, \tilde{y}_i^k = \begin{pmatrix} R_i y \\ \sqrt{\frac{\rho}{2}} \left(B \begin{pmatrix} z_i^k \\ z_{i+1}^k \end{pmatrix} - u_i^k \right) \end{pmatrix}, K \right).$$

b. Minimization wrt z is a least squares problem with a sparse matrix, which can be implemented efficiently:

$$z^{k+1} \leftarrow \arg \min_z \|\tilde{A}\Gamma^{k+1} + U^k - \tilde{B}Z\|_2^2$$

c. Dual update:

$$U^{k+1} \leftarrow \tilde{A}\Gamma^{k+1} - \tilde{B}Z + U^k.$$

2. Compute $\hat{y} := D_G \hat{\Gamma}$.

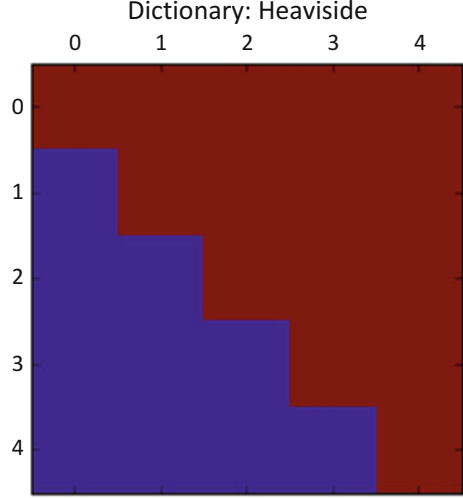
4 Examples

We now turn to present several classes of signals that belong to the proposed globalized model, where each of these is obtained by imposing a special structure on the local dictionary. Then, we demonstrate how one can sample from \mathcal{M} and generate such signals. Additional examples are given in [Appendix E: Generative Models for Patch-Sparse Signals](#).

4.1 Piecewise Constant (PWC) Signals

The (unnormalized) Heaviside $n \times n$ dictionary H_n is the upper triangular matrix with 1's in the upper part (see Figure 3 on page 19). Formally, each local atom d_i of length n is expressed as a step function, given by $d_i^T = [\mathbf{1}_i, \mathbf{0}_{n-i}]^T$, $1 \leq i \leq n$, where $\mathbf{1}_i$ is a vector of ones of length i . Similarly, $\mathbf{0}_{n-i}$ is a zero vector of length $n-i$. The following property is verified by noticing that H_n^{-1} is the discrete difference operator.

Fig. 3 Heaviside dictionary
 H_4 . Red is 1, blue is 0.



Proposition 7. *If a patch $x_i \in \mathbb{R}^n$ has $L - 1$ steps, then its (unique) representation in the Heaviside dictionary H_n has at most L nonzeros.*

Corollary 1. *Let $x \in \mathbb{R}^N$ be a piecewise constant signal with at most $L - 1$ steps per each segment of length n (in the periodic sense). Then*

$$x \in \mathcal{M}(H_n, L, N, P = N).$$

Remark 2. The model $\mathcal{M}(H_n, L, N, P = N)$ contains also some signals having exactly L steps in a particular patch, but those patches must have their last segment with zero height.

As an example, one might synthesize signals with sparsity $\|\Gamma\|_{0,\infty} \leq 2$ according to the following scheme:

1. Draw at random the support of Γ with the requirement that the distance between the jumps within the signal will be at least the length of a patch (this allows at most two nonzeros per patch, one for the step and the second for the bias/DC).
2. Multiply each step by a random number.

The global subspace $A_{\mathcal{S}}$ and the corresponding global oracle denoiser x_G (11) in the PWC model can be explicitly described.

Proposition 8. *Let $x \in \mathbb{R}^N$ consist of s constant segments with lengths ℓ_r , $r = 1, \dots, s$, and let Γ be the (unique) global representation of x in \mathcal{M} (i.e., $\rho(x) = \{\Gamma\}$). Denote $B := \text{diag}(B_r)_{r=1}^s$, where $B_r = \frac{1}{\ell_r} \mathbf{1}_{\ell_r \times \ell_r}$. Then*

1. We have

$$\ker A_{\text{supp } \Gamma} = \ker (I_N - B), \quad (16)$$

and therefore $\dim \ker A_{\text{supp } \Gamma} = s$ and $MSE(\hat{x}_G) = s\sigma^2$ under the Gaussian noise model (12).

2. Furthermore, the global oracle estimator x_G is given by

$$x_G(y, \text{supp } \Gamma) = By, \quad (17)$$

i.e., the global oracle is the averaging operator within the constant segments of the signal.

Proof. Every signal $y \in \ker A_{\text{supp } \Gamma}$ has the same “local jump pattern” as x , and therefore it also has the same *global* jump pattern. That is, every such y consists of s constant segments with lengths ℓ_r . It is an easy observation that such signals satisfy $y = By$, which proves (16). It is easy to see that $\dim \ker (I_{\ell_r} - B_r) = 1$, and therefore

$$\dim \ker (I_N - \text{diag}(B_r)_{r=1}^s) = s.$$

The proof of 1) is finished by invoking Proposition 4.

To prove (17), notice that by the previous discussion the null-space of $A_{\text{supp } \Gamma}$ is spanned by the orthogonal set $e_r = \frac{1}{\sqrt{\ell_r}} \begin{bmatrix} 0, \dots, 0, \underbrace{1, 1, \dots, 1}_{\ell_r}, 0, \dots, 0 \end{bmatrix}^T$, $r = 1, \dots, s$. Let $K = [e_1, \dots, e_s]$, then $x_G = KK^\dagger = KK^T$. It can be easily verified by direct computation that $KK^T = B$. \square

It turns out that the LPA performance (and the local-global gap) can be accurately described by the following result. We provide an outline of proof in [Appendix D: Proof of Theorem 8](#).

Theorem 8. Let $x \in \mathbb{R}^N$ consist of s constant segments with lengths ℓ_r , $r = 1, \dots, s$, and assume the Gaussian noise model (12). Then

1. There exists a function $R(n, \alpha) : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{R}^+$, with $R(n, \alpha) > 1$, such that

$$MSE(\hat{x}_{LPA}) = \sigma^2 \sum_{r=1}^s R(n, \ell_r).$$

2. The function $R(n, \alpha)$ satisfies:

- a. $R(n, \alpha) = 1 + \frac{\alpha(2\alpha H_\alpha^{(2)} - 3\alpha + 2) - 1}{n^2}$ if $n \geq \alpha$, where $H_\alpha^{(2)} = \sum_{k=1}^{\alpha} \frac{1}{k^2}$;
- b. $R(n, \alpha) = \frac{11}{18} + \frac{2\alpha}{3n} + \frac{6\alpha - 11}{18n^2}$ if $n \leq \frac{\alpha}{2}$.

Corollary 2. *The function $R(n, \alpha)$ is monotonically increasing in α (with n fixed) and monotonically decreasing in n (with α fixed). Furthermore,*

1. $\lim_{n \rightarrow \infty} R(n, n) = \frac{\pi^2}{3} - 2 \approx 1.29$;
2. $\lim_{n \rightarrow \infty} R(n, 2n) = \frac{35}{18} \approx 1.94$.

Thus, for reasonable choices of the patch size, the local-global gap is roughly a constant multiple of the number of segments, reflecting the global complexity of the signal.

For numerical examples of reconstructing the PWC signals using our local-global framework, see Subsection 5.2.

4.2 Signature-Type Dictionaries

Another type of signals that comply with our model are those represented via a signature dictionary, which has been shown to be effective for image restoration [3]. This dictionary is constructed from a small signal, $x \in \mathbb{R}^m$, such that its every patch (in varying location, extracted in a cyclic fashion), $R_i x \in \mathbb{R}^n$, is a possible atom in the representation, namely, $d_i = R_i x$. As such, every consecutive pair of atoms $(i, i + 1)$ is essentially a pair of overlapping patches that satisfy $S_B d_i = S_T d_{i+1}$ (before normalization). The complete algorithm is presented for convenience in Algorithm 3.

Algorithm 3 Constructing the signature dictionary

1. Choose the base signal $x \in \mathbb{R}^m$.
 2. Compute $D(x) = [R_1 x, R_2 x, \dots, R_m x]$, where R_i extracts the i -th patch of size n in a cyclic fashion.
 3. Normalization: $\tilde{D}(x) = [d_1, \dots, d_m]$, where $d_i = \frac{R_i x}{\|R_i x\|_2}$.
-

Given D as above, one can generate signals $y \in \mathbb{R}^N$, where N is an integer multiple of m , with s nonzeros per patch, by the easy procedure outlined below.

1. Init: Construct a base signal $b \in \mathbb{R}^N$ by replicating $x \in \mathbb{R}^m$ N/m times (note that b is therefore periodic). Set $y = 0$.
2. Repeat for $j = 1, \dots, s$:
 - a. Shift: Circularly shift the base signal by t_j positions, denoted by $\text{shift}(b, t_j)$, for some $t_j = 0, 1, \dots, m - 1$ (drawn at random).
 - b. Aggregate: $y = y + \omega_j \cdot \text{shift}(b, t_j)$, where ω is an arbitrary random scalar.

Notice that a signal constructed in this way must be periodic, as it is easily seen that

$$\ker A_{\mathcal{S}} = \text{span} \{ \text{shift}(b, t_i) \}_{i=1}^s,$$

while the support sequence \mathcal{S} is

$$\mathcal{S} = ([t_1, t_2, \dots, t_s], [t_1, t_2, \dots, t_s] + 1, \dots, [t_1, t_2, \dots, t_s] + N) \pmod{m}.$$

Assuming that there are no additional relations between the single atoms of D except those from the above construction, all $\mathcal{S} \in \Sigma_{\mathcal{M}}$ are easily seen to be of the above form.

In Figure 4 on page 23, we give an example of a signature-type dictionary D for $(n, m) = (6, 10)$ and a signal x with $N = P = 30$ together with its corresponding sparse representation Γ .

Remark 3. It might seem that every $n \times m$ Hankel matrix such as the one shown in Figure 4 on page 23 produces a signature-type dictionary with a nonempty signal space \mathcal{M} . However this is not the case, because such a dictionary will usually fail to generate signals of length larger than $n + m - 1$.

4.2.1 Multi-Signature Dictionaries

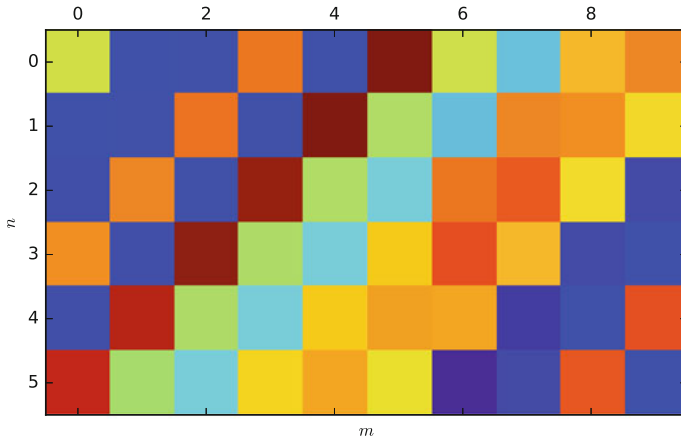
One can generalize the construction of Subsection 4.2 and consider k -tuples of initial base signals x_1, \dots, x_k , instead of a single x . The desired dictionary D will consist of corresponding k -tuples of atoms, which are constructed from those base signals. In order to avoid ending up with the same structure as the case $k = 1$, we also require a “mixing” of the atoms. The complete procedure is outlined in Algorithm 4.

Algorithm 4 Constructing the multi-signature dictionary

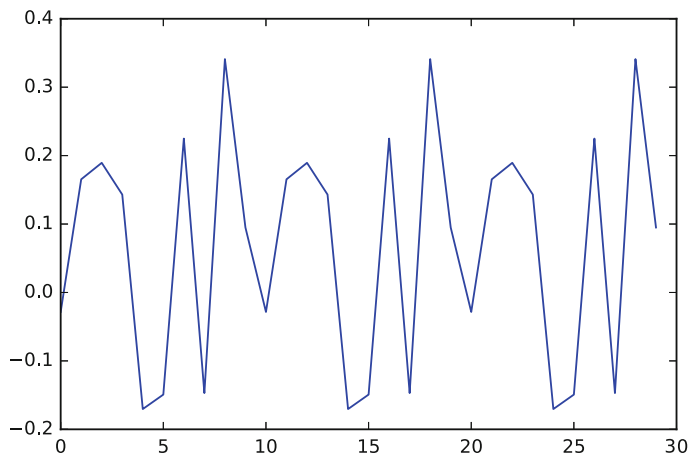
1. Input: n, m, k such that k divides m . Put $r := \frac{m}{k}$.
 2. Select a signal basis matrix $X \in \mathbb{R}^{r \times k}$ and r nonsingular transfer matrices $M_i \in \mathbb{R}^{k \times k}$, $i = 1, \dots, r$.
 3. Repeat for $i = 1, \dots, r$:
 - a. Let $Y_i = [y_{i,1}, \dots, y_{i,k}] \in \mathbb{R}^{n \times k}$, where each $y_{i,j}$ is the i -th patch (of length n) of the signal x_j .
 - b. Put the k -tuple $[d_{i,1}, \dots, d_{i,k}] = Y_i \times M_i$ as the next k atoms in D .
-

In order to generate a signal of length N from \mathcal{M} , one can follow these steps (again we assume that m divides N):

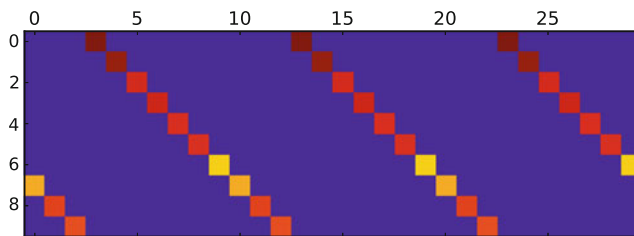
1. Create a base signal matrix $X^G \in \mathbb{R}^{N \times k}$ by stacking $k \frac{N}{m}$ copies of the original basis matrix X . Set $y = 0$.
2. Repeat for $j = 1, \dots, s$:
 - a. Select a base signal $b_j \in \mathcal{R}(X^G)$ and shift it (in a circular fashion) by some $t_j = 0, 1, \dots, R - 1$.
 - b. Aggregate: $y = y + \text{shift}(b_j, t_j)$ (note that here we do not need to multiply by a random scalar).



(a) The dictionary matrix D



(b) The signal $x \in \ker A_{\mathcal{S}}$ for \mathcal{S} generated by $t_1 = 6$ and $s = 1$, with $P = N = 30$.



(c) The coefficient matrix Γ corresponding to the signal x in (c)

Fig. 4 An example of the signature dictionary with $n = 6$, $m = 10$. See Remark 3.

This procedure will produce a signal y of local sparsity $k \cdot s$. The corresponding support sequence can be written as

$$\mathcal{S} = (s_1, s_2, \dots, s_N),$$

where $s_i = s_1 + i \pmod{m}$ and

$$s_1 = [(t_1, 1), (t_1, 2), \dots, (t_1, k), \dots, (t_s, 1), (t_s, 2), \dots, (t_s, k)].$$

Here (t_j, i) denotes the atom $d_{t_j, i}$ in the notation of Algorithm 4. The corresponding signal space is

$$\ker A_{\mathcal{S}} = \text{span} \left\{ \text{shift}(X^G, t_j) \right\}_{j=1}^s,$$

and it is of dimension $k \cdot s$.

An example of a multi-signature dictionary and corresponding signals may be seen in Figure 5 on page 25.

4.3 Convolutional Dictionaries

An important class of signals is the *sparse convolution model*, where each signal $x \in \mathbb{R}^N$ can be written as a linear combination of shifted “waveforms” $\mathbf{d}_i \in \mathbb{R}^n$, each \mathbf{d}_i being a column in the local dictionary $D' \in \mathbb{R}^{n \times m}$. More conveniently, any such x can be represented as a circular convolution of \mathbf{d}_i with a (sparse) “feature map” $\psi_i \in \mathbb{R}^N$:

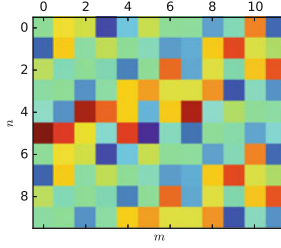
$$x = \sum_{i=1}^m \mathbf{d}_i *_{\mathbb{N}} \psi_i. \quad (18)$$

Such signals arise in various applications, such as audio classification [6, 26, 50], neural coding [16, 44], and mid-level image representation and denoising [31, 58, 59].

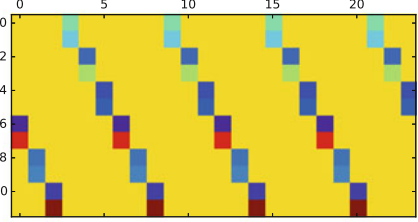
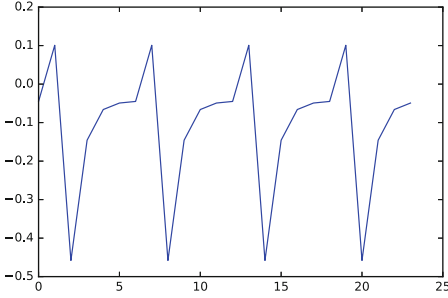
Formally, the convolutional class can be recast into the patch-sparse model of this paper as follows. First, we can rewrite (18) as

$$x = \underbrace{[\mathbf{C}_1 \ \mathbf{C}_2 \ \dots \ \mathbf{C}_m]}_{:=\mathbf{E}} \Psi,$$

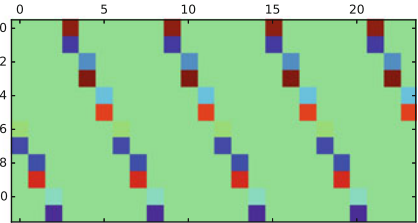
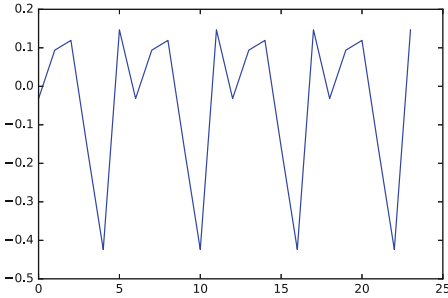
where each $\mathbf{C}_i \in \mathbb{R}^{N \times N}$ is a banded circulant matrix with its first column being equal to \mathbf{d}_i and $\Psi \in \mathbb{R}^{Nm}$ is the concatenation of the ψ_i 's. It is easy to see that by permuting the columns of \mathbf{E} , one obtains precisely the global convolutional dictionary nD_G based on the local dictionary D' (recall (3)). Therefore we obtain



(a) The dictionary D



(b) The first signal and its sparse representation in $\ker A_{\mathcal{S}}$ with $N = 24$, $s = 1$ and $t_1 = 5$.



(c) The second signal and its sparse representation in $\ker A_{\mathcal{S}}$.

Fig. 5 Example of multi-signature dictionary with $n = 10$, $m = 12$, and $k = 2$.

$$x = \underbrace{D_G(D')}_{:=D'_G} \Gamma'. \tag{19}$$

While it is tempting to conclude from comparing (19) and (4) that the convolutional model is equivalent to the patch-sparse model, an essential ingredient is missing, namely, the requirement of equality on overlaps, $M\Gamma' = 0$. Indeed, nothing in the definition of the convolutional model restricts the representation Ψ (and therefore Γ'); therefore, in principle the number of degrees of freedom remains Nm , as compared to $N(m - n + 1)$ from Proposition 15.

To fix this, following [42], we apply R_i to (19) and obtain $R_i x = R_i D'_G \mathbf{F}'$. The “stripe” $\Omega'_i = R_i D'_G$ has only $(2n - 1)m$ nonzero consecutive columns, and in fact the nonzero portion of Ω'_i is equal for all i . This implies that every x_i has a representation $x_i = \Theta \mathbf{y}_i$ in the “pseudo-local” dictionary

$$\Theta(D') := \left[Z_B^{(n-1)} D' \dots D' \dots Z_T^{(n-1)} D' \right] \in \mathbb{R}^{n \times (2n-1)m},$$

where the operators $Z_B^{(k)}$ and $Z_T^{(k)}$ are given by Definition 6 in Appendix B: Proof of Lemma 2. If we now assume that our convolutional signals satisfy

$$\|\mathbf{y}_i\|_0 \leq s \quad \forall i,$$

then we have shown that they belong to $\mathcal{M}(\Theta(D'), s, P, N)$ and thus can be formally treated by the framework we have developed.

It turns out that this direct approach is quite naive, as the dictionary $\Theta(D')$ is extremely ill-equipped for sparse reconstruction (e.g., it has repeated atoms, and therefore $\mu(\Theta(D')) = 1$). To tackle this problem, a convolutional sparse coding framework was recently developed in [42], where the explicit dependencies between the sparse representation vectors \mathbf{y}_i (and therefore the special structure of the corresponding constraint $M(D') \mathbf{F}' = 0$) were exploited quite extensively, resulting in efficient recovery algorithms and nontrivial theoretical guarantees. We refer the reader to [42] for further details and examples.

5 Numerical Experiments

In this section, we test the effectiveness of the globalized model for recovering the signals from Section 4, both in the noiseless and noisy cases. For the PWC, we show a real-world example. These results are also compared to several other approaches such as the LPA, total variation denoising (for the PWC), and a global pursuit based on OMP.

5.1 Signature-Type Signals

In this section we investigate the performance of the pursuit algorithms on signals complying with the signature dictionary model elaborated in Subsection 4.2, constructed from one or two base signals ($k = 1, 2$), and allowing for varying values of s . We compare the results to both LPA and a global pursuit, which uses the dictionary explicitly constructed from the signature model. In detail, the global dictionary D^* is an $N \times (km)$ matrix consisting of the base signal matrix X^G and all its shifts, i.e. (recall the definitions in Subsection 4.2.1)

$$D^* = \left[\text{shift}(X^G, i)_{i=0}^{m-1} \right].$$

Given that, the global OMP algorithm is defined to run for $k \cdot s$ steps on D^* .

5.1.1 Constructing the Dictionary

In the context of the LPA algorithm, the condition for its success in recovering the representation is a function of the mutual coherence of the local dictionary—the smaller this measure, the larger the number of nonzeros that are guaranteed to be recovered. Leveraging this, we aim at constructing $D \in \mathbb{R}^{n \times m}$ of a signature type that has a small coherence. This can be cast as an optimization problem

$$D = \tilde{D}(x_0), \quad x_0 = \arg \min_{x \in \mathbb{R}^m} \mu(\tilde{D}(x)),$$

where $\tilde{D}(x)$ is computed by Algorithm 3 (or Algorithm 4) and μ is the (normalized) coherence function.

In our experiments, we choose $(n, m) = (15, 20)$ for $k = 1$ and $(n, m) = (10, 20)$ for $k = 2$. We minimize the above loss function via gradient descent, resulting in $\mu(\tilde{D}(x)) = 0.20$ for $k = 1$ and $\mu = 0.26$ for $k = 2$. We used the TensorFlow open source package [1]. As a comparison, the coherence of a random signature dictionary is about 0.5.

5.1.2 Noiseless Case

In this setting, we test the ability of the globalized OMP (Subsection 3.3.1) to perfectly recover the sparse representation of clean signature-type signals. Figure 6 compares the proposed algorithm (for different choices of $\beta \in \{0.25, 0.5, 1, 2, 5\}$) with the LPA by providing their probability of success in recovering the true sparse vectors, averaged over 10^3 randomly generated signals of length $N = 100$. For brevity we show only the results for $k = 1$ here.

From a theoretical perspective, since $\mu(D) = 0.20$, the LPA algorithm is guaranteed to recover the representation when $\|\Gamma\|_{0,\infty} \leq 3$, as indeed it does. Comparing the LPA approach to the globalized OMP, one can observe that for $\beta \geq 1$ the latter consistently outperforms the former, having a perfect recovery for $\|\Gamma\|_{0,\infty} \leq 4$. Another interesting insight of this experiment is the effect of β on the performance; roughly speaking, a relatively large value of this parameter results in a better success rate than the very small ones, thereby emphasizing importance of the constraint $M_* \Gamma = 0$. On the other hand, β should not be too large since the importance of the signal is reduced compared to the constraint, which might lead to deterioration in the success rate (see the curve that corresponds to $\beta = 5$ in Figure 6).

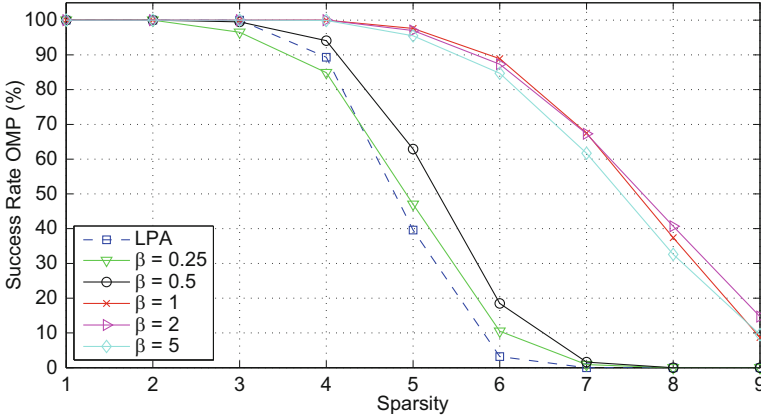
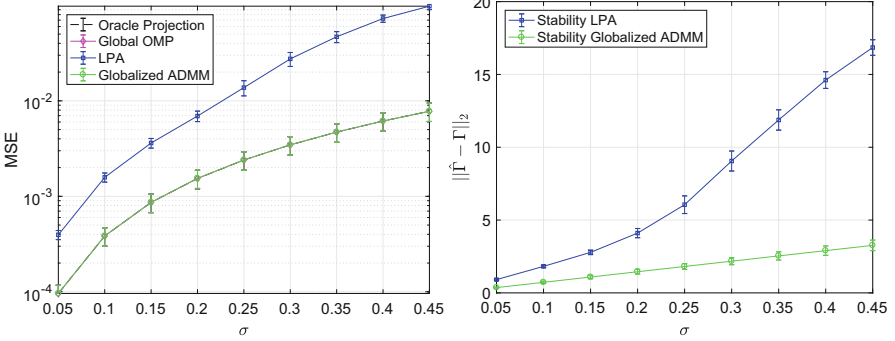


Fig. 6 Probability of the success (%) of the globalized OMP (for various values of β) and the LPA algorithms to perfectly recover the sparse representations of test signals from the signature dictionary model, averaged over 10^3 realizations, as a function of sparsity per patch.

5.1.3 Noisy Case

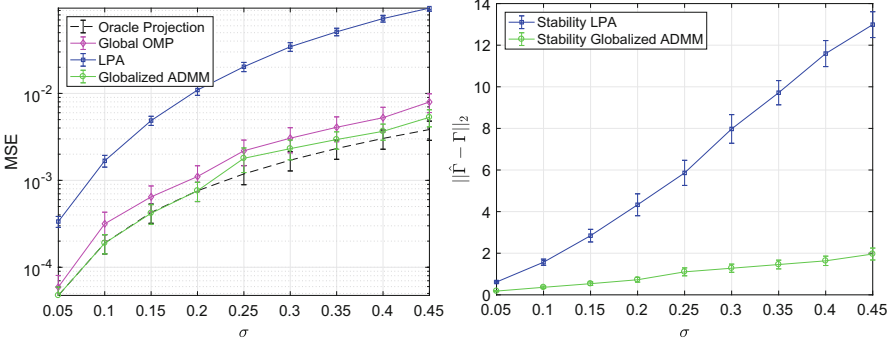
In what follows, the stability of the proposed globalized ADMM-inspired pursuit is tested and compared to the traditional LPA algorithm, as well as to the global OMP. In addition to the above, we provide the restoration performance of the oracle estimator, serving as an indication for the best possible denoising that can be achieved. In this case, the oracle projection matrix A_S is constructed according to the ground-truth support S .

We generate ten random signature-type signals, where each of these is corrupted by white additive Gaussian noise with standard deviation σ , ranging from 0.05 up to 0.5. The global number of nonzeros is injected to the global OMP, and the information regarding the local sparsity is utilized both by the LPA algorithm as well as by our ADMM-inspired pursuit (which is based on local sparse recovery operations). Following Figure 7 parts (a, c), which plot the mean squared error (MSE) of the estimation as a function of the noise level, the ADMM-inspired pursuit achieves the best denoising performance, having similar results to the oracle estimator for all noise levels and sparsity factors. The source of superiority of the ADMM-inspired pursuit might be its inherent ability to obtain an estimation that perfectly fits to the globalized model. The second best algorithm is the global OMP; using complete global information about the signal space, this fact is to be expected. The LPA algorithm is the least accurate; it shows that for our signals the assumption of patch independence severely degrades performance. This sheds light on the difficulty of finding the true supports, the nontrivial solution of this problem, and the great advantage of the proposed globalized model.



(a) $k = 1, s = 3$, denoising. The globalized ADMM-inspired (green curve) and oracle projection (black curve) coincide with Global OMP (magenta curve).

(b) $k = 1, s = 3$, stability



(c) $k = 2, s = 1$, denoising

(d) $k = 2, s = 1$, stability

Fig. 7 (a, c) Denoising performance of the global OMP, ADMM-inspired pursuit, and LPA algorithm for signals from the signature model with (a) $k = 1, s = 5$ and (c) $k = 2, s = 1$. The performance of the oracle estimator is provided as well, demonstrating the best possible restoration that can be achieved. (b, d) Stability of the ADMM-inspired pursuit and LPA algorithm for (b) $k = 1, s = 5$ and (d) $k = 2, s = 1$. For (a, b) the signal size was $N = 100$, while for (c, d) it was $N = 80$.

Similar conclusion holds for the stable recovery of the sparse representations. Per each pursuit algorithm, Figure 7 parts (b, d) illustrate the ℓ_2 distance between the original sparse vector Γ and its estimation $\hat{\Gamma}$, averaged over the different noise realizations. As can be seen, the ADMM-inspired pursuit achieves the most stable recovery, outperforming the LPA algorithm especially in the challenging cases of high noise levels and/or large sparsity factors.

5.2 Denoising PWC Signals

5.2.1 Synthetic Data

In this scenario, we test the ability of the globalized ADMM-inspired pursuit to restore corrupted PWC signals and compare these to the outcome of the LPA algorithm.

In addition, we run the total variation (TV) denoising [48] on the signals, which is known to perform well on PWC. We chose the regularization parameter in the TV by running an exhaustive search over a wide range of values per input signal and picked the one that minimizes the MSE between the estimated and the true signal. Notice that this results in the best possible denoising performance that can be obtained by the TV.

The projected versions of both ADMM-inspired pursuit and LPA are provided along with the one of the oracle estimator. Following the description in Section 4.1, we generate a signal of length $N = 200$, composed of patches of size $n = m = 20$ with a local sparsity of at most 2 nonzeros in the $\ell_{0,\infty}$ sense. These signals are then contaminated by a white additive Gaussian noise with σ in the range of 0.1 to 0.9.

The restoration performance (in terms of MSE) of the abovementioned algorithms and their stability are illustrated in Figure 8, where the results are averaged over 10 noise realizations. As can be seen, the globalized approach significantly outperforms the LPA algorithm for all noise levels. Furthermore, when $\sigma \leq 0.5$, the ADMM-inspired pursuit performs similarly to the oracle estimator. One can also notice that the ADMM-inspired pursuit and its projected version result in the very same estimation, i.e., this algorithm forces the signal to conform with the patch-sparse model globally. On the other hand, following the visual illustration given in Figure 9, the projected version of the LPA algorithm has only two nonzero segments,

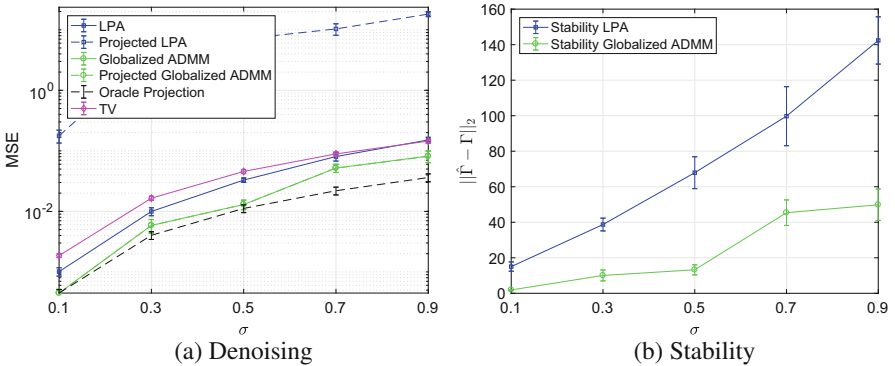
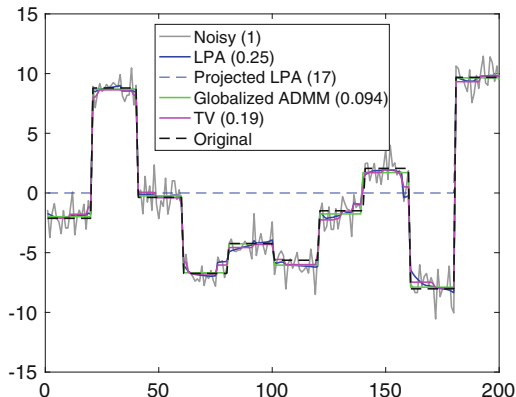


Fig. 8 (a) Denoising performance and (b) stability for various noise levels, tested for signals from the piecewise constant model with $\|\Gamma\|_{0,\infty} \leq 2$.

Fig. 9 Denoising of a PWC signal contaminated with additive Gaussian noise ($\sigma = 1.1$) via several pursuit algorithms: input noisy signal (MSE = 1.0), LPA algorithm (MSE = 0.25), projected LPA (MSE = 17), ADMM-inspired pursuit (MSE = 0.094), and TV (MSE=0.19). Projected ADMM is identical to the ADMM-inspired pursuit.



which are the consequence of almost complete disagreement in the support (local inconsistency). This is also reflected in Figure 8a, illustrating that even for a very small noise level ($\sigma = 0.1$), the projected version of the LPA algorithm has a very large estimation error (MSE ≈ 0.18) compared to the one of the ADMM-inspired pursuit (MSE ≈ 0.0004), indicating that the former fails in obtaining a consistent representation of the signal. The TV method is unable to take into account the local information, resulting in reconstruction of lesser quality than both the ADMM-inspired and the LPA.

5.2.2 Real-World Data

Here we apply the globalized ADMM for the PWC model on a real-world DNA copy number data from [51]. The data (see also [34]) come from a single experiment on 15 fibroblast cell lines with each array containing over 2000 (mapped) BACs (bacterial artificial chromosomes) spotted in triplicate. The results (see Figure 10) appear to be reasonably significant.

6 Discussion

In this work we have presented an extension of the classical theory of sparse representations to signals which are locally sparse, together with novel pursuit algorithms. We envision several promising research directions which might emerge from this work.

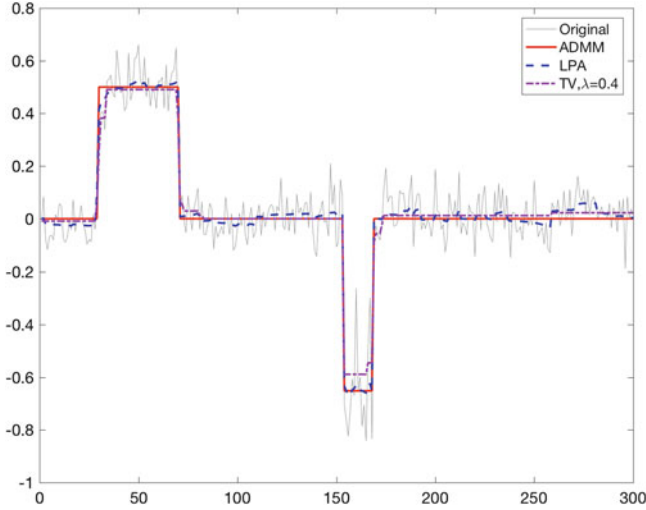


Fig. 10 Applying the PWC reconstruction to a single fibroblast cell line, as described in [51]. The value of λ in TV was chosen empirically based on visual quality. For the ADMM, we chose $n = 40$ and $k = 2$. The ordinate is the normalized average of the log base 2 test over reference ratio of the cell line.

6.1 Relation to Other Models

Viewed globally, the resulting signal model can be considered a sort of “structured sparse” model; however, in contrast to other such constructions ([29, 30, 32, 55] and others), our model incorporates both structure in the representation coefficients and a structured dictionary.

The recently developed framework of convolutional sparse coding (CSC) [41, 42] bears some similarities to our work, in that it, too, has a convolutional representation of the signal via a dictionary identical in structure to D_G . However, the underlying local sparsity assumptions are drastically different in the two models, resulting in very different guarantees and algorithms. That said, we believe that it would be important to provide precise connections between the results, possibly leading to their deeper understanding. First steps in this direction are outlined in Subsection 4.3.

6.2 Further Extensions

The decomposition of the global signal $x \in \mathbb{R}^N$ into its patches,

$$x \mapsto (R_i x)_{i=1}^P, \quad (20)$$

is a special case of a more general decomposition, namely,

$$x \mapsto (w_i \mathcal{P}_i x)_{i=1}^P, \quad (21)$$

where \mathcal{P}_i is the (orthogonal) projection onto a subspace W_i of \mathbb{R}^N and w_i are some weights. This observation naturally places our theory, at least partially, into the framework of *fusion frames*, a topic which is generating much interest recently in the applied harmonic analysis community [21, Chapter 13]. In fusion frame theory, which is motivated by applications such as distributed sensor networks, the starting point is precisely the decomposition (21). Instead of the reconstruction formula $x = \sum_i \frac{1}{n} R_i^T R_i x$, in fusion frame theory we have

$$x = \sum_i w_i^2 S_{\mathcal{W}}^{-1} (\mathcal{P}_i x),$$

where $S_{\mathcal{W}}$ is the associated *fusion frame operator*. The natural extension of our work to this setting would seek to enforce some sparsity of the projections $\mathcal{P}_i x$. Perhaps the most immediate variant of (20) in this respect would be to drop the periodicity requirement, resulting in a slightly modified R_i operators near the endpoints of the signal. We would like to mention some recent works which investigate different notions of fusion frame sparsity [2, 4, 8].

Another intriguing possible extension of our work is to relax the complete overlap requirement between patches and consider an ‘‘approximate patch sparsity’’ model, where the patch disagreement vector $M\Gamma$ is not zero but ‘‘small.’’ In some sense, one can imagine a full ‘‘spectrum’’ of such models, ranging from a complete agreement (this work) to an arbitrary disagreement (such as in the CSC framework mentioned above).

6.3 Learning Models from Data

The last point above brings us to the question of how to obtain ‘‘good’’ models, reflecting the structure of the signals at hand (such as speech/images, etc.). We hope that one might use the ideas presented here in order to create novel learning algorithms. In this regard, the main difficulty is how to parametrize the space of allowed models in an efficient way. While we presented some initial ideas in [Appendix E: Generative Models for Patch-Sparse Signals](#), in the most general case (incorporating the approximate sparsity direction above), the problem remains widely open.

Acknowledgments The research leading to these results has received funding from the European Research Council under European Union’s Seventh Framework Programme, ERC Grant agreement no. 320649. The authors would also like to thank Jeremias Sulam, Vardan Papayan, Raja Giryes, and Gitta Kutinyok for inspiring discussions.

Appendix A: Proof of Lemma 1

Proof. Denote $Z := \ker M$ and consider the linear map $A : Z \rightarrow \mathbb{R}^N$ given by the restriction of the “averaging map” $D_G : \mathbb{R}^{mP} \rightarrow \mathbb{R}^N$ to Z .

1. Let us see first that $\text{im}(A) = \mathbb{R}^N$. Indeed, for every $x \in \mathbb{R}^N$, consider its patches $x_i = R_i x$. Since D is full rank, there exist $\{\alpha_i\}$ for which $D\alpha_i = x_i$. Then setting $\Gamma := (\alpha_1, \dots, \alpha_p)$, we have both $D_G \Gamma = x$ and $M\Gamma = 0$ (by construction, see Section 2), i.e., $\Gamma \in Z$ and the claim follows.
2. Define

$$J := \ker D \times \ker D \times \dots \times \ker D \subset \mathbb{R}^{mP}.$$

We claim that $J = \ker A$.

- a. In one direction, let $\Gamma = (\alpha_1, \dots, \alpha_p) \in \ker A$, i.e., $M\Gamma = 0$ and $D_G \Gamma = 0$. Immediately we see that $\frac{1}{n} D\alpha_i = 0$ for all i , and therefore $\alpha_i \in \ker D$ for all i , thus $\Gamma \in J$.
 - b. In the other direction, let $\Gamma = (\alpha_1, \dots, \alpha_p) \in J$, i.e., $D\alpha_i = 0$. Then the local representations agree, i.e., $M\Gamma = 0$, thus $\Gamma \in Z$. Furthermore, $D_G \Gamma = 0$ and therefore $\Gamma \in \ker A$.
3. By the fundamental theorem of linear algebra, we conclude

$$\begin{aligned} \dim Z &= \dim \text{im}(A) + \dim \ker A = N + \dim J \\ &= N + (m - n)N = N(m - n + 1). \end{aligned}$$

□

Appendix B: Proof of Lemma 2

We start with an easy observation.

Proposition 9. *For any vector $\rho \in \mathbb{R}^N$, we have*

$$\|\rho\|_2^2 = \frac{1}{n} \sum_{j=1}^N \|R_j \rho\|_2^2.$$

Proof. Since

$$\|\rho\|_2^2 = \sum_{j=1}^N \rho_j^2 = \frac{1}{n} \sum_{j=1}^N n \rho_j^2 = \frac{1}{n} \sum_{j=1}^N \sum_{k=1}^n \rho_j^2,$$

we can rearrange the sum and get

$$\begin{aligned}\|\rho\|_2^2 &= \frac{1}{n} \sum_{k=1}^n \sum_{j=1}^N \rho_j^2 = \frac{1}{n} \sum_{k=1}^n \sum_{j=1}^N \rho_{(j+k) \bmod N}^2 = \frac{1}{n} \sum_{j=1}^N \sum_{k=1}^n \rho_{(j+k) \bmod N}^2 \\ &= \frac{1}{n} \sum_{j=1}^N \|\mathcal{R}_j \rho\|_2^2.\end{aligned}$$

□

Corollary 3. *Given $M\Gamma = 0$, we have*

$$\|y - D_G \Gamma\|_2^2 = \frac{1}{n} \sum_{j=1}^N \|\mathcal{R}_j y - D \alpha_j\|_2^2.$$

Proof. Using Proposition 9, we get

$$\|y - D_G \Gamma\|_2^2 = \frac{1}{n} \sum_{j=1}^N \|\mathcal{R}_j y - \mathcal{R}_j D_G \Gamma\|_2^2 = \frac{1}{n} \sum_{j=1}^N \|\mathcal{R}_j y - \Omega_j \Gamma\|_2^2.$$

Now since $M\Gamma = 0$, then by definition of M , we have $\Omega_j \Gamma = D \alpha_j$ (see (6)), and this completes the proof. □

Recall Definition 6. Multiplying the corresponding matrices gives

Proposition 10. *We have the following equality for all $i = 1, \dots, P$:*

$$S_B \mathcal{R}_i = S_T \mathcal{R}_{i+1}. \quad (22)$$

To facilitate the proof, we introduce extension of Definition 6 to multiple shifts as follows.

Definition 16. Let n be fixed. For $k = 0, \dots, n-1$ let

1. $S_T^{(k)} := [I_{n-k} \ \mathbf{0}]$ and $S_B^{(k)} := [\mathbf{0} \ I_{n-k}]$ denote the operators extracting the top (resp. bottom) $n-k$ entries from a vector of length n ; the matrices have dimension $(n-k) \times n$.
2. $Z_B^{(k)} := \begin{bmatrix} S_B^{(k)} \\ \mathbf{0}_{k \times n} \end{bmatrix}$ and $Z_T^{(k)} := \begin{bmatrix} \mathbf{0}_{k \times n} \\ S_T^{(k)} \end{bmatrix}$.
3. $W_B^{(k)} := \begin{bmatrix} \mathbf{0}_{k \times n} \\ S_B^{(k)} \end{bmatrix}$ and $W_T^{(k)} := \begin{bmatrix} S_T^{(k)} \\ \mathbf{0}_{k \times n} \end{bmatrix}$.

Note that $S_B = S_B^{(1)}$ and $S_T = S_T^{(1)}$. We have several useful consequences of the above definitions. The proofs are carried out via elementary matrix identities and are left to the reader.

Proposition 11. *For any $n \in \mathbb{N}$, the following hold:*

1. $Z_T^{(k)} = \left(Z_T^{(1)}\right)^k$ and $Z_B^{(k)} = \left(Z_B^{(1)}\right)^k$ for $k = 0, \dots, n-1$;
2. $W_T^{(k)} W_T^{(k)} = W_T^{(k)}$ and $W_B^{(k)} W_B^{(k)} = W_B^{(k)}$ for $k = 0, \dots, n-1$;
3. $W_T^{(k)} W_B^{(j)} = W_B^{(j)} W_T^{(k)}$ for $j, k = 0, \dots, n-1$;
4. $Z_B^{(k)} = Z_B^{(k)} W_B^{(k)}$ and $Z_T^{(k)} = Z_T^{(k)} W_T^{(k)}$ for $k = 0, \dots, n-1$;
5. $W_B^{(k)} = Z_T^{(1)} W_B^{(k-1)} Z_B^{(1)}$ and $W_T^{(k)} = Z_B^{(1)} W_T^{(k-1)} Z_T^{(1)}$ for $k = 1, \dots, n-1$;
6. $Z_B^{(k)} Z_T^{(k)} = W_T^{(k)}$ and $Z_T^{(k)} Z_B^{(k)} = W_B^{(k)}$ for $k = 0, \dots, n-1$;
7. $(n-1) I_{n \times n} = \sum_{k=1}^{n-1} \left(W_B^{(k)} + W_T^{(k)}\right)$.

Proposition 12. *If the vectors $u_1, \dots, u_N \in \mathbb{R}^n$ satisfy pairwise*

$$S_B u_i = S_T u_{i+1},$$

then they also satisfy for each $k = 0, \dots, n-1$ the following:

$$W_B^{(k)} u_i = Z_T^{(k)} u_{i+k}, \quad (23)$$

$$Z_B^{(k)} u_i = W_T^{(k)} u_{i+k}. \quad (24)$$

Proof. It is easy to see that the condition $S_B u_i = S_T u_{i+1}$ directly implies

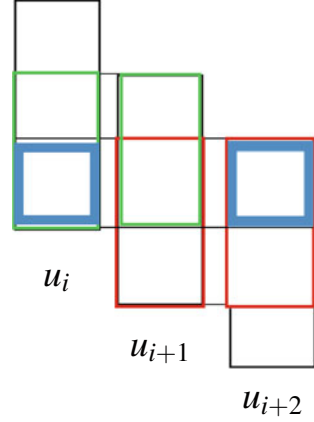
$$Z_B^{(1)} u_i = W_T^{(1)} u_{i+1}, \quad W_B^{(1)} u_i = Z_T^{(1)} u_{i+1} \quad \forall i. \quad (25)$$

Let us first prove (23) by induction on k . The base case $k = 1$ is precisely (25). Assuming validity for $k-1$ and $\forall i$, we have

$$\begin{aligned} W_B^{(k)} u_i &= Z_T^{(1)} W_B^{(k-1)} Z_B^{(1)} u_i && \text{(by Proposition 11, item 5)} \\ &= Z_T^{(1)} W_B^{(k-1)} W_T^{(1)} u_{i+1} && \text{(by (25))} \\ &= Z_T^{(1)} W_T^{(1)} W_B^{(k-1)} u_{i+1} && \text{(by Proposition 11, item 3)} \\ &= Z_T^{(1)} W_T^{(1)} Z_T^{(k-1)} u_{i+k} && \text{(by the induction hypothesis)} \\ &= Z_T^{(1)} Z_T^{(k-1)} u_{i+k} && \text{(by Proposition 11, item 4)} \\ &= Z_T^{(k)} u_{i+k}. && \text{(by Proposition 11, item 1)} \end{aligned}$$

To prove (24) we proceed as follows:

Fig. 11 Illustration to the proof of Proposition 12. The green pair is equal, as well as the red pair. It follows that the blue elements are equal as well.



$$\begin{aligned}
 Z_B^{(k)} u_i &= Z_B^{(k)} W_B^{(k)} u_i && \text{(by Proposition 11, item 4)} \\
 &= Z_B^{(k)} Z_T^{(k)} u_{i+k} && \text{(by (23) which is already proved)} \\
 &= W_T^{(k)} u_{i+k}. && \text{(by Proposition 11, item 6)}
 \end{aligned}$$

This finishes the proof of Proposition 12. □

Example 1. To help the reader understand the claim of Proposition 12, consider the case $k = 2$, and take some three vectors u_i, u_{i+1}, u_{i+2} . We have $S_B u_i = S_T u_{i+1}$ and also $S_B u_{i+1} = S_T u_{i+2}$. Then clearly $S_B^{(2)} u_i = S_T^{(2)} u_{i+2}$ (see Figure 11 on page 37) and therefore $W_B^{(2)} u_i = Z_T^{(2)} u_{i+2}$.

Let us now present the proof of Lemma 2.

Proof. We show equivalence in two directions.

- (1) \implies (2): Let $M\Gamma = 0$. Define $x := D_G\Gamma$, and then further denote $x_i := R_i x$. Then on the one hand:

$$\begin{aligned}
 x_i &= R_i D_G \Gamma \\
 &= \Omega_i \Gamma && \text{(definition of } \Omega_i) \\
 &= D\alpha_i. && (M\Gamma = 0)
 \end{aligned}$$

On the other hand, because of (22) we have $S_B R_i x = S_T R_{i+1} x$, and by combining the two, we conclude that $S_B D\alpha_i = S_T D\alpha_{i+1}$.

- (2) \implies (1): In the other direction, suppose that $S_B D\alpha_i = S_T D\alpha_{i+1}$. Denote $u_i := D\alpha_i$. Now consider the product $\Omega_i \Gamma$ where $\Omega_i = R_i D_G$. One can easily be convinced that in fact

$$\Omega_i \Gamma = \frac{1}{n} \left(\sum_{k=1}^{n-1} \left(Z_B^{(k)} u_{i-k} + Z_T^{(k)} u_{i+k} \right) + u_i \right).$$

Therefore

$$\begin{aligned} (\Omega_i - Q_i) \Gamma &= \frac{1}{n} \left(u_i + \sum_{k=1}^{n-1} \left(Z_B^{(k)} u_{i-k} + Z_T^{(k)} u_{i+k} \right) \right) - u_i \\ &= \frac{1}{n} \left(\sum_{k=1}^{n-1} \left(W_T^{(k)} u_i + W_B^k u_i \right) - (n-1) u_i \right) \quad (\text{by Proposition 12}) \\ &= 0. \quad (\text{by Proposition 11, item 7}) \end{aligned}$$

Since this holds for all i , we have shown that $M\Gamma = 0$.

□

Appendix C: Proof of Theorem 6

Recall that $M_A = \frac{1}{n} \sum_i R_i^T P_{s_i} R_i$. We first show that M_A is a contraction.

Proposition 13. $\|M_A\|_2 \leq 1$.

Proof. Closely following a similar proof in [45], divide the index set $\{1, \dots, N\}$ into n groups representing *non-overlapping* patches: for $i = 1, \dots, n$ let

$$K(i) := \left\{ i, i+n, \dots, i + \left(\left\lfloor \frac{N}{n} \right\rfloor - 1 \right) n \right\} \pmod{N}.$$

Now

$$\begin{aligned} \|M_A x\|_2 &= \frac{1}{n} \left\| \sum_{i=1}^N R_i^T P_{s_i} R_i x \right\|_2 \\ &= \frac{1}{n} \left\| \sum_{i=1}^n \sum_{j \in K(i)} R_j^T P_{s_j} R_j x \right\|_2 \\ &\leq \frac{1}{n} \sum_{i=1}^n \left\| \sum_{j \in K(i)} R_j^T P_j R_j x \right\|_2. \end{aligned}$$

By construction, $R_j R_k^T = \mathbf{0}_{n \times n}$ for $j, k \in K(i)$ and $j \neq k$. Therefore for all $i = 1, \dots, n$ we have

$$\begin{aligned} \left\| \sum_{j \in K(i)} R_j^T P_{s_j} R_j x \right\|_2^2 &= \sum_{j \in K(i)} \|R_j^T P_{s_j} R_j x\|_2^2 \\ &\leq \sum_{j \in K(i)} \|R_j x\|_2^2 \leq \|x\|_2^2. \end{aligned}$$

Substituting in back into the preceding inequality finally gives

$$\|M_A x\|_2 \leq \frac{1}{n} \sum_{i=1}^n \|x\|_2 = \|x\|_2.$$

□

Now let us move on to prove Theorem 6.

Proof. Define

$$\hat{P}_i := (I - P_{s_i}) R_i.$$

It is easy to see that

$$\sum_i \hat{P}_i^T \hat{P}_i = A_{\mathcal{J}}^T A_{\mathcal{J}}.$$

Let the SVD of $A_{\mathcal{J}}$ be

$$A_{\mathcal{J}} = U \Sigma V^T.$$

Now

$$\begin{aligned} V \Sigma^2 V^T &= A_{\mathcal{J}}^T A_{\mathcal{J}} = \sum_i \hat{P}_i^T \hat{P}_i = \sum_i R_i^T R_i - \underbrace{\sum_i R_i^T P_{s_i} R_i}_{:=T} \\ &= nI - T. \end{aligned}$$

Therefore $T = nI - V \Sigma^2 V^T$, and

$$M_A = \frac{1}{n} T = I - \frac{1}{n} V \Sigma^2 V^T = V \left(I - \frac{\Sigma^2}{n} \right) V^T.$$

This shows that the eigenvalues of M_A are $\tau_i = 1 - \frac{\sigma_i^2}{n}$ where $\{\sigma_i\}$ are the singular values of $A_{\mathcal{J}}$. Thus we obtain

$$M_A^k = V \text{diag} \{ \tau_i^k \} V^T.$$

If $\sigma_i = 0$ then $\tau_i = 1$, and in any case, by Proposition 13, we have $|\tau_i| \leq 1$. Let the columns of the matrix W consist of the singular vectors of $A_{\mathcal{S}}$ corresponding to $\sigma_i = 0$ (and so $\text{span } W = \mathcal{N}(A_{\mathcal{S}})$), then

$$\lim_{k \rightarrow \infty} M_A^k = WW^T.$$

Thus, as $k \rightarrow \infty$, M_A^k tends to the orthogonal projector onto $\mathcal{N}(A_{\mathcal{S}})$. The convergence is evidently linear, the constant being dependent upon $\{\tau_i\}$. \square

Appendix D: Proof of Theorem 8

Recall that the signal consists of s constant segments of corresponding lengths ℓ_1, \dots, ℓ_s . We would like to compute the MSE for every pixel within every such segment of length $\alpha := \ell_r$. For each patch, the oracle provides the locations of the jump points within the patch.

Let us calculate the MSE for pixel with index 0 inside a constant (**nonzero**) segment $[-k, \alpha - k - 1]$ with value v (Figure 12 on page 41 might be useful). The oracle estimator has the explicit formula

$$\hat{x}_A^{r,k} = \frac{1}{n} \sum_{j=1}^n \frac{1}{b_j - a_j + 1} \sum_{i=a_j}^{b_j} (v + z_i), \quad (26)$$

where $j = 1, \dots, n$ corresponds to the index of the overlapping patch containing the pixel, intersecting the constant segment on $[a_j, b_j]$, so that

$$\begin{aligned} a_j &= -\min(k, n - j), \\ b_j &= \min(\alpha - k - 1, j - 1). \end{aligned}$$

Now, the oracle error for the pixel is

$$\begin{aligned} \hat{x}_A^{r,k} - v &= \frac{1}{n} \sum_{j=1}^n \frac{1}{b_j - a_j + 1} \sum_{i=a_j}^{b_j} z_i \\ &= \sum_{i=-k}^{\alpha-k-1} c_{i,\alpha,n,k} z_i, \end{aligned}$$

where the coefficients $c_{i,\alpha,n,k}$ are some *positive* rational numbers depending only on i, α, n and k . It is easy to check by rearranging the above expression that

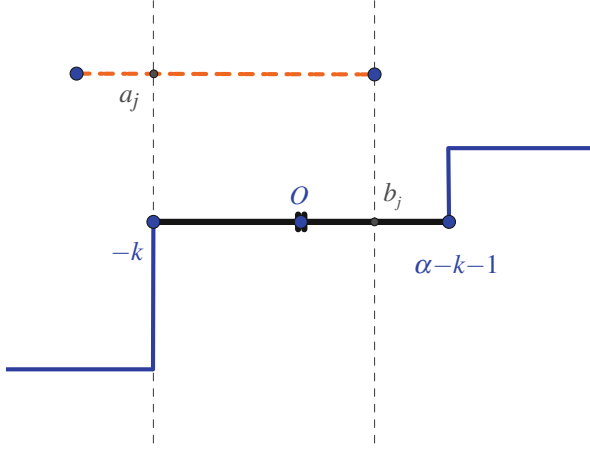


Fig. 12 The oracle estimator for the pixel O in the segment (black). The orange line is patch number $j = 1, \dots, n$, and the relevant pixels are between a_j and b_j . The signal itself is shown to extend beyond the segment (blue line).

$$\sum_{i=-k}^{\alpha-k-1} c_{i,\alpha,n,k} = 1, \tag{27}$$

and furthermore, denoting $d_i := c_{i,\alpha,n,k}$ for fixed α, n, k , we also have that

$$d_{-k} < d_{-k+1} < \dots < d_0 > d_1 > \dots > d_{\alpha-k-1}. \tag{28}$$

Example 2. $n = 4, \alpha = 3$

- For $k = 1$:

$$\begin{aligned} \hat{x}_A^{r,k} - v &= \frac{1}{4} \left(\frac{1}{2} + \frac{1}{3} + \frac{1}{3} + \frac{1}{2} \right) z_0 + \frac{1}{4} \left(\frac{1}{2} + \frac{1}{3} + \frac{1}{3} \right) z_{-1} + \frac{1}{4} \left(\frac{1}{3} + \frac{1}{3} + \frac{1}{2} \right) z_1 \\ &= \underbrace{\frac{7}{24}}_{d_{-1}} z_{-1} + \underbrace{\frac{5}{12}}_{d_0} z_0 + \underbrace{\frac{7}{24}}_{d_1} z_1 \end{aligned}$$

- For $k = 2$:

$$\begin{aligned} \hat{x}_A^{r,k} - v &= \frac{1}{4} \left(\frac{1}{3} + \frac{1}{3} + \frac{1}{2} + 1 \right) z_0 + \frac{1}{4} \left(\frac{1}{3} + \frac{1}{3} + \frac{1}{2} \right) z_{-1} + \frac{1}{4} \left(\frac{1}{3} + \frac{1}{3} \right) z_{-2} \\ &= \frac{13}{24} z_0 + \frac{7}{24} z_{-1} + \frac{1}{6} z_{-2} \end{aligned}$$

Now consider the optimization problem

$$\min_{c \in \mathbb{R}^\alpha} c^T c \quad \text{s.t. } \mathbf{1}^T c = 1.$$

It can be easily verified that it has the optimal value $\frac{1}{\alpha}$, attained at $c^* = \alpha \mathbf{1}$. From this, (27) and (28), it follows that

$$\sum_{i=-k}^{\alpha-k-1} c_{i,\alpha,n,k}^2 > \frac{1}{\alpha}.$$

Since the z_i are i.i.d., we have

$$\mathbb{E} \left(\hat{x}_A^{r,k} - v \right)^2 = \sigma^2 \sum_{i=-k}^{\alpha-k-1} c_{i,\alpha,n,k}^2,$$

while for the entire nonzero segment of length $\alpha = \ell_r$

$$E_r := \mathbb{E} \left(\sum_{k=0}^{\alpha-1} \left(\hat{x}_A^{r,k} - v \right)^2 \right) = \sum_{k=0}^{\alpha-1} \mathbb{E} \left(\hat{x}_A^{r,k} - v \right)^2 = \sigma^2 \sum_{k=0}^{\alpha-1} \sum_{i=-k}^{\alpha-k-1} c_{i,\alpha,n,k}^2.$$

Defining

$$R(n, \alpha) := \sum_{k=0}^{\alpha-1} \sum_{i=-k}^{\alpha-k-1} c_{i,\alpha,n,k}^2,$$

we obtain that $R(n, \alpha) > 1$ and furthermore

$$\mathbb{E} \|\hat{x}_A - x\|^2 = \sum_{r=1}^s E_r = \sigma^2 \sum_{r=1}^s R(n, \ell_r) > s\sigma^2.$$

This proves item (1) of Theorem 8. For showing the explicit formulas for $R(n, \alpha)$ in item (2), we have used automatic symbolic simplification software MAPLE [39].

By construction (26), it is not difficult to see that if $n \geq \alpha$ then

$$\begin{aligned} R(n, \alpha) &= \frac{1}{n^2} \sum_{k=0}^{\alpha-1} \left(\sum_{j=0}^k (2H_{\alpha-1} - H_k + \frac{n-\alpha+1}{\alpha} - H_{\alpha-1-j})^2 \right. \\ &\quad \left. + \sum_{j=k+1}^{\alpha-1} (2H_{\alpha-1} - H_{\alpha-k-1} + \frac{n-\alpha+1}{\alpha} - H_j)^2 \right), \end{aligned}$$

where $H_k := \sum_{i=1}^k \frac{1}{i}$ is the k -th harmonic number. This simplifies to

$$R(n, \alpha) = 1 + \frac{\alpha(2\alpha H_\alpha^{(2)} + 2 - 3\alpha) - 1}{n^2},$$

where $H_k^{(2)} = \sum_{i=1}^k \frac{1}{i^2}$ is the k -th harmonic number of the second kind.

On the other hand, for $n \leq \frac{\alpha}{2}$ we have

$$R(n, \alpha) = \sum_{k=0}^{n-2} c_{n,k}^{(1)} + \sum_{k=n-1}^{\alpha-n} c_{n,k}^{(2)} + \sum_{k=\alpha-n+1}^{\alpha-1} c_{n,\alpha-1-k}^{(1)},$$

where

$$c_{n,k}^{(1)} = \frac{1}{n^2} \left(\sum_{j=k}^{n-1} \left(H_{n-1} - H_j + \frac{k+1}{n} \right)^2 + \sum_{i=n-k}^{n-1} \left(\frac{n-i}{n} \right)^2 + \sum_{i=0}^{k-1} \left(H_{n-1} - H_k + \frac{k-i}{n} \right)^2 \right)$$

and

$$c_{n,k}^{(2)} = \frac{1}{n^2} \left(\sum_{j=k-n+1}^k \left(\frac{j-k+n}{n} \right)^2 + \sum_{j=k+1}^{k+n-1} \left(\frac{k+n-j}{n} \right)^2 \right).$$

Automatic symbolic simplification of the above gives

$$R(n, \alpha) = \frac{11}{18} + \frac{2\alpha}{3n} - \frac{5}{18n^2} + \frac{\alpha-1}{3n^3}.$$

Appendix E: Generative Models for Patch-Sparse Signals

In this section we propose a general framework aimed at generating signals from the patch-sparse model. Our approach is to construct a graph-based model for the dictionary and subsequently use this model to generate dictionaries and signals which turn out to be much richer than those considered in Section 4.

Local Support Dependencies

We start by highlighting the importance of the local connections (recall Lemma 2) between the neighboring patches of the signal and therefore between the corresponding subspaces containing those patches. This in turn allows to characterize $\Sigma_{\mathcal{M}}$ as the set of all “realizable” paths in a certain dependency graph derived from the dictionary D . This point of view allows to describe the model \mathcal{M} using only the intrinsic properties of the dictionary, in contrast to Theorem 2.

Proposition 14. *Let $0 \neq x \in \mathcal{M}$ and Γ a gamma $\in \rho(x)$ with $\text{supp } \Gamma = (S_1, \dots, S_P)$. Then for $i = 1, \dots, P$*

$$\text{rank} [S_B D_{S_i} - S_T D_{S_{i+1}}] < |S_i| + |S_{i+1}| \leq 2s, \quad (29)$$

where by convention $\text{rank } \emptyset = -\infty$.

Proof. $x \in \mathcal{M}$ implies by Lemma 2 that for every $i = 1, \dots, P$

$$[S_B D \quad - S_T D] \begin{bmatrix} \alpha_i \\ \alpha_{i+1} \end{bmatrix} = 0.$$

But

$$[S_B D \quad - S_T D] \begin{bmatrix} \alpha_i \\ \alpha_{i+1} \end{bmatrix} = [S_B D_{S_i} \quad - S_T D_{S_{i+1}}] \begin{bmatrix} \alpha_i |S_i| \\ \alpha_{i+1} |S_{i+1}| \end{bmatrix} = 0,$$

and therefore the matrix $[S_B D_{S_i} - S_T D_{S_{i+1}}]$ must be rank-deficient. Note in particular that the conclusion still holds if one (or both) of the $\{s_i, s_{i+1}\}$ is empty. \square

The preceding result suggests a way to describe all the supports in $\Sigma_{\mathcal{M}}$.

Definition 17. Given a dictionary D , we define an abstract directed graph $\mathcal{G}_{D,s} = (V, E)$, with the vertex set

$$V = \{(i_1, \dots, i_k) \subset \{1, \dots, m\} : \text{rank } D_{i_1, \dots, i_k} = k < n\},$$

and the edge set

$$E = \left\{ (S_1, S_2) \in V \times V : \text{rank} [S_B D_{S_1} \quad - S_T D_{S_2}] < \min \{n - 1, |S_1| + |S_2|\} \right\}.$$

In particular, $\emptyset \in V$ and $(\emptyset, \emptyset) \in E$ with $\text{rank} [\emptyset] := -\infty$.

Remark 4. It might be impossible to compute $\mathcal{G}_{D,s}$ in practice. However we set this issue aside for now and only explore the theoretical ramifications of its properties.

Definition 18. The set of all directed paths of length P in $\mathcal{G}_{D,s}$, not including the self-loop $\underbrace{(\emptyset, \emptyset, \dots, \emptyset)}_{\times P}$, is denoted by $\mathcal{C}_{\mathcal{G}}(P)$.

Definition 19. A path $\mathcal{S} \in \mathcal{C}_{\mathcal{G}}(P)$ is called *realizable* if $\dim \ker A_{\mathcal{S}} > 0$. The set of all realizable paths in $\mathcal{C}_{\mathcal{G}}(P)$ is denoted by $\mathcal{R}_{\mathcal{G}}(P)$.

Thus we have the following result.

Theorem 9. Suppose $0 \neq x \in \mathcal{M}$. Then

1. Every representation $\Gamma = (\alpha_i)_{i=1}^P \in \rho(x)$ satisfies $\text{supp } \Gamma \in \mathcal{C}_{\mathcal{G}}(P)$, and therefore

$$\Sigma_{\mathcal{M}} \subseteq \mathcal{R}_{\mathcal{G}}(P). \quad (30)$$

2. The model \mathcal{M} can be characterized “intrinsically” by the dictionary as follows:

$$\mathcal{M} = \bigcup_{\mathcal{G} \in \mathcal{R}_{\text{eg}}(P)} \ker A_{\mathcal{G}}. \quad (31)$$

Proof. Let $\text{supp } \Gamma = (S_1, \dots, S_P)$ with $S_i = \text{supp } \alpha_i$ if $\alpha_i \neq \mathbf{0}$, and $S_i = \emptyset$ if $\alpha_i = \mathbf{0}$. Then by Proposition 14, we must have that

$$\text{rank} [S_B D_{S_i} - S_T D_{S_{i+1}}] < |S_i| + |S_{i+1}| \leq 2s.$$

Furthermore, since $\Gamma \in \rho(x)$ we must have that D_{S_i} is full rank for each $i = 1, \dots, P$. Thus $(S_i, S_{i+1}) \in \mathcal{G}_{D,s}$, and so $\text{supp } \Gamma \in \mathcal{R}_{\text{eg}}(P)$. Since by assumption $\text{supp } \Gamma \in \Sigma_{\mathcal{M}}$, this proves (30).

To show (31), notice that if $\text{supp } \Gamma \text{ amma} \in \mathcal{R}_{\text{eg}}(P)$, then for every $x \in \ker A_{\text{supp } \Gamma}$, we have $R_i x = P_{S_i} R_i x$, i.e., $R_i x = D \alpha_i$ for some α_i with $\text{supp } \alpha_i \subseteq S_i$. Clearly in this case $|\text{supp } \alpha_i| \leq s$ and therefore $x \in \mathcal{M}$. The other direction of (31) follows immediately from the definitions. \square

Definition 20. The dictionary D is called “ (s, P) -good” if

$$|\mathcal{R}_{\text{eg}}(P)| > 0.$$

Theorem 10. The set of “ (s, P) -good” dictionaries has measure zero in the space of all $n \times m$ matrices.

Proof. Every low-rank condition defines a finite number of algebraic equations on the entries of D (given by the vanishing of all the $2s \times 2s$ minors of $[S_B D_{S_i} \ S_T D_{S_i}]$). Since the number of possible graphs is finite (given fixed n, m and s), the resulting solution set is a finite union of semi-algebraic sets of low dimension and hence has measure zero. \square

Constructing “Good” Dictionaries

The above considerations suggest that the good dictionaries are hard to come by; here we provide an example of an explicit construction.

We start by defining an abstract graph \mathcal{G} with some desirable properties, and subsequently look for a nontrivial realization D of the graph, so that in addition $\mathcal{R}_{\text{eg}} \neq \emptyset$.

In this context, we would want \mathcal{G} to contain *sufficiently many different long cycles*, which would correspond to long signals and a rich resulting model \mathcal{M} . In contrast with the models from Subsection 4.2 (where all the graphs consist of a single cycle), one therefore should allow for some branching mechanism. An example of a possible \mathcal{G} is given in Figure 13 on page 46. Notice that due to the

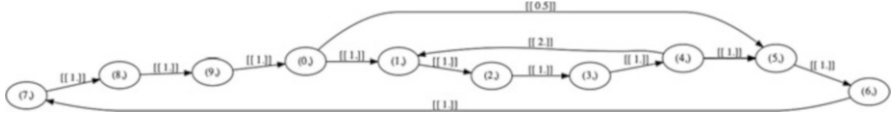


Fig. 13 A possible dependency graph \mathcal{G} with $m = 10$. In this example, $|\mathcal{C}_{\mathcal{G}}(70)| = 37614$.

structure of \mathcal{G} , there are many possible paths in $\mathcal{C}_{\mathcal{G}}(P)$. In fact, a direct search algorithm yields $|\mathcal{C}_{\mathcal{G}}(70)| = 37614$.

Every edge in \mathcal{G} corresponds to a conditions of the form (29) imposed on the entries of D . As discussed in Theorem 10, this in turn translates to a set of algebraic equations. So the natural idea would be to write out the large system of such equations and look for a solution over the field \mathbb{R} by well-known algorithms in numerical algebraic geometry [5]. However, this approach is highly impractical because these algorithms have (single or double) exponential running time. We consequently propose a simplified, more direct approach to the problem.

In detail, we replace the low-rank conditions (29) with more explicit and restrictive ones below.

Assumptions(*) For each $(S_i, S_j) \in \mathcal{G}$ we have $|S_i| = |S_j| = k$. We require that $\text{span } S_B D_{S_i} = \text{span } S_T D_{S_j} = \Lambda_{i,j}$ with $\dim \Lambda_{i,j} = k$. Thus there exists a nonsingular transfer matrix $C_{i,j} \in \mathbb{R}^{k \times k}$ such that

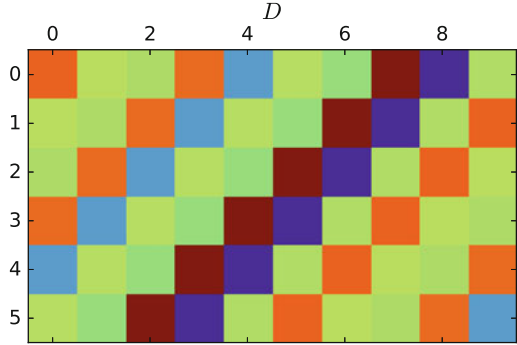
$$S_B D_{S_i} = C_{i,j} S_T D_{S_j}. \quad (32)$$

In other words, every column in $S_B D_{S_i}$ must be a specific linear combination of the columns in $S_T D_{S_j}$. This is much more restrictive than the low-rank condition, but on the other hand, given the matrix $C_{i,j}$, it defines a set of linear constraints on D . To summarize, the final algorithm is presented in Algorithm 5. In general, nothing guarantees that for a particular choice of \mathcal{G} and the transfer matrices, there is a nontrivial solution D ; however, in practice we do find such solutions. For example, taking the graph from Figure 13 on page 46 and augmenting it with the matrices $C_{i,j}$ (scalars in this case), we obtain a solution over \mathbb{R}^6 which is shown in Figure 14 on page 47. Notice that while the resulting dictionary has a Hankel-type structure similar to what we have seen previously, the additional dependencies between the atoms produce a rich signal space structure, as we shall demonstrate in the following section.

Algorithm 5 Finding a realization D of the graph \mathcal{G}

1. Input: a graph \mathcal{G} satisfying the **Assumptions(*)** above, and the dimension n of the realization space \mathbb{R}^n .
 2. Augment the edges of \mathcal{G} with arbitrary nonsingular transfer matrices $C_{i,j}$.
 3. Construct the system of linear equations given by (32).
 4. Find a nonzero D solving the system above over \mathbb{R}^n .
-

Fig. 14 A realization $D \in \mathbb{R}^{6 \times 10}$ of \mathcal{G} from Figure 13 on page 46.



Generating Signals

Now suppose the graph \mathcal{G} is known (or can be easily constructed). Then this gives a simple procedure to generate signals from \mathcal{M} , presented in Algorithm 6.

Algorithm 6 Constructing a signal from \mathcal{M} via \mathcal{G}

1. Construct a path $\mathcal{S} \in \mathcal{C}_{\mathcal{G}}(P)$.
 2. Construct the matrix $A_{\mathcal{S}}$.
 3. Find a nonzero vector in $\ker A_{\mathcal{S}}$.
-

Let us demonstrate this on the example in Figure 13 on page 46 and Figure 14 on page 47. Not all paths in $\mathcal{C}_{\mathcal{G}}$ are realizable, but it turns out that in this example we have $|\mathcal{R}_{\mathcal{G}}(70)| = 17160$. Three different signals and their supports \mathcal{S} are shown in Figure 15 on page 48. As can be seen from these examples, the resulting model \mathcal{M} is indeed much richer than the signature-type construction from Subsection 4.2.

An interesting question arises: given $\mathcal{S} \in \mathcal{C}_{\mathcal{G}}(P)$, can we say something about $\dim \ker A_{\mathcal{S}}$? In particular, when is it strictly positive (i.e., when $\mathcal{S} \in \mathcal{R}_{\mathcal{G}}(P)$?) While in general the question seems to be difficult, in some special cases this number can be estimated using only the properties of the local connections (S_i, S_{i+1}) , by essentially counting the additional “degrees of freedom” when moving from patch i to patch $i + 1$. To this effect, we prove two results.

Proposition 15. *For every $\mathcal{S} \in \mathcal{R}_{\mathcal{G}}(P)$, we have*

$$\dim \ker A_{\mathcal{S}} = \dim \ker M_*^{(\mathcal{S})}.$$

Proof. Notice that

$$\ker A_{\mathcal{S}} = \left\{ D_G^{(\mathcal{S})} \Gamma_{\mathcal{S}}, M_*^{(\mathcal{S})} \Gamma_{\mathcal{S}} = 0 \right\} = \text{im} \left(D_G^{(\mathcal{S})} \Big|_{\ker M_*^{(\mathcal{S})}} \right),$$

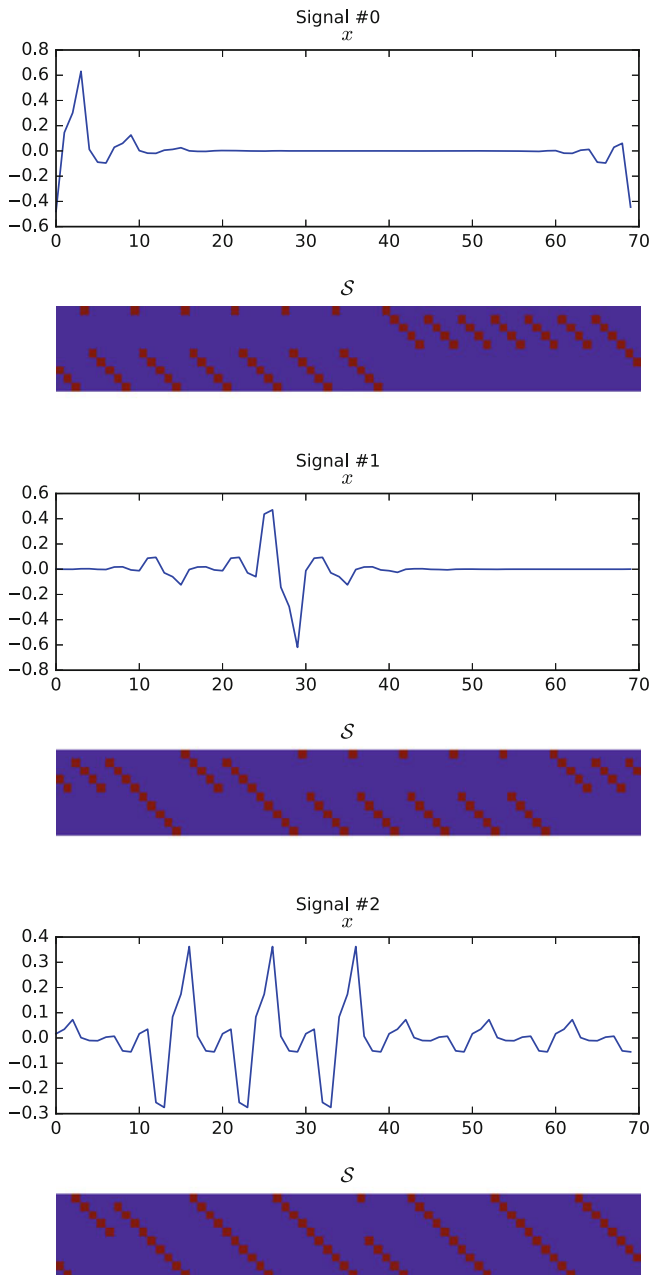


Fig. 15 Examples of signals from \mathcal{M} and the corresponding supports \mathcal{S} .

and therefore $\dim \ker A_{\mathcal{S}} \leq \dim \ker M_*^{(\mathcal{S})}$. Furthermore, the map $D_G^{(\mathcal{S})}|_{\ker M_*^{(\mathcal{S})}}$ is injective, because if $D_G^{(\mathcal{S})}\Gamma_{\mathcal{S}} = 0$ and $M_*^{(\mathcal{S})}\Gamma_{\mathcal{S}} = 0$, we must have that $D_{S_i}\alpha_i|_{S_i} = 0$ and, since D_{S_i} has full rank, also $\alpha_i = 0$. The conclusion follows. \square

Proposition 16. *Assume that the model satisfies **Assumptions**(*) above. Then for every $\mathcal{S} \in \mathcal{R}_{\mathcal{G}}(P)$*

$$\dim \ker A_{\mathcal{S}} \leq k.$$

Proof. The idea is to construct a spanning set for $\ker M_*^{(\mathcal{S})}$ and invoke Proposition 15. Let us relabel the nodes along \mathcal{S} to be $1, 2, \dots, P$. Starting from an arbitrary α_1 with support $|S_1| = k$, we use (32) to obtain, for $i = 1, 2, \dots, P - 1$, a formula for the next portion of the global representation vector Γ

$$\alpha_{i+1} = C_{i,i+1}^{-1}\alpha_i. \quad (33)$$

This gives a set Δ consisting of overall k linearly independent vectors Γ with support $\Gamma_i = \mathcal{S}$. It may happen that equation (33) is not satisfied for $i = P$. However, every Γ with $\text{supp } \Gamma = \mathcal{S}$ and $M_*^{(\mathcal{S})}\Gamma = 0$ must belong to $\text{span } \Delta$, and therefore

$$\dim \ker M_*^{(\mathcal{S})} \leq \dim \text{span } \Delta = k.$$

\square

We believe that Proposition 16 can be extended to more general graphs, not necessarily satisfying **Assumptions**(*). In particular, the following estimate appears to hold for a general model \mathcal{M} and $\mathcal{S} \in \mathcal{R}_{\mathcal{G}}(P)$:

$$\dim \ker A_{\mathcal{S}} \leq |S_1| + \sum_i (|S_{i+1}| - \text{rank} [S_B D_{S_i} \ S_T D_{S_{i+1}}]).$$

We leave the rigorous proof of this result to a future work.

Further Remarks

While the model presented in this section is the hardest to analyze theoretically, even in the restricted case of **Assumptions**(*) (when does a nontrivial realization of a given \mathcal{G} exist? How does the answer depend on n ? When $\mathcal{R}_{\mathcal{G}}(P) \neq \emptyset$? etc?), we hope that this construction will be most useful in applications such as denoising of natural signals.

References

1. M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G.S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, X. Zheng, TensorFlow: large-scale machine learning on heterogeneous systems (2015). <http://tensorflow.org/>. Software available from tensorflow.org
2. R. Aceska, J.L. Bouchot, S. Li, Local sparsity and recovery of fusion frames structured signals. preprint (2015). <http://www.mathc.rwth-aachen.de/~bouchot/files/pubs/FusionCSfinal.pdf>
3. M. Aharon, M. Elad, Sparse and redundant modeling of image content using an image-signature-dictionary. *SIAM J. Imag. Sci.* **1**(3), 228–247 (2008)
4. U. Ayaz, S. Dirksen, H. Rauhut, Uniform recovery of fusion frame structured sparse signals. *Appl. Comput. Harmon. Anal.* **41**(2), 341–361 (2016). <https://doi.org/10.1016/j.acha.2016.03.006>. <http://www.sciencedirect.com/science/article/pii/S1063520316000294>
5. S. Basu, R. Pollack, M.F. Roy, *Algorithms in Real Algebraic Geometry*. Algorithms and Computation in Mathematics, 2nd edn., vol. 10 (Springer, Berlin, 2006)
6. T. Blumensath, M. Davies, Sparse and shift-invariant representations of music. *IEEE Trans. Audio Speech Lang. Process.* **14**(1), 50–57 (2006). http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=1561263
7. T. Blumensath, M.E. Davies, Sampling theorems for signals from the union of finite-dimensional linear subspaces. *IEEE Trans. Inf. Theory* **55**(4), 1872–1882 (2009). http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4802322
8. P. Boufounos, G. Kutyniok, H. Rauhut, Sparse recovery from combined fusion frame measurements. *IEEE Trans. Inf. Theory* **57**(6), 3864–3876 (2011). <https://doi.org/10.1109/TIT.2011.2143890>
9. S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* **3**(1), 1–122 (2011). <http://dx.doi.org/10.1561/22000000016>
10. H. Bristow, A. Eriksson, S. Lucey, Fast convolutional sparse coding. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2013), pp. 391–398
11. A.M. Bruckstein, D.L. Donoho, M. Elad, From sparse solutions of systems of equations to sparse modeling of signals and images. *SIAM Rev.* **51**(1), 34–81 (2009). <http://epubs.siam.org/doi/abs/10.1137/060657704>
12. E.J. Candes, Modern statistical estimation via oracle inequalities. *Acta Numer.* **15**, 257–325 (2006). http://journals.cambridge.org/abstract_S0962492906230010
13. S. Chen, S.A. Billings, W. Luo, Orthogonal least squares methods and their application to non-linear system identification. *Int. J. Control.* **50**(5), 1873–1896 (1989)
14. W. Dong, L. Zhang, G. Shi, X. Li, Nonlocally centralized sparse representation for image restoration. *IEEE Trans. Image Process.* **22**(4), 1620–1630 (2013)
15. D.L. Donoho, M. Elad, Optimally sparse representation in general (nonorthogonal) dictionaries via ℓ_1 minimization. *Proc. Natl. Acad. Sci.* **100**(5), 2197–2202 (2003). [doi:10.1073/pnas.0437847100](https://doi.org/10.1073/pnas.0437847100). <http://www.pnas.org/content/100/5/2197>
16. C. Ekanadham, D. Tranchina, E.P. Simoncelli, A unified framework and method for automatic neural spike identification. *J. Neurosci. Methods* **222**, 47–55 (2014). [doi:10.1016/j.jneumeth.2013.10.001](https://doi.org/10.1016/j.jneumeth.2013.10.001). <http://www.sciencedirect.com/science/article/pii/S0165027013003415>
17. M. Elad, *Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing* (Springer, New York, 2010)

18. M. Elad, M. Aharon, Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Trans. Image Process.* **15**(12), 3736–3745 (2006)
19. Y.C. Eldar, M. Mishali, Block sparsity and sampling over a union of subspaces, in *2009 16th International Conference on Digital Signal Processing* (IEEE, New York, 2009), pp. 1–8. http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5201211
20. Y.C. Eldar, M. Mishali, Robust recovery of signals from a structured union of subspaces. *IEEE Trans. Inf. Theory* **55**(11), 5302–5316 (2009)
21. Finite Frames - Theory and Applications. <http://www.springer.com/birkhauser/mathematics/book/978-0-8176-8372-6>
22. S. Foucart, H. Rauhut, *A Mathematical Introduction to Compressive Sensing* (Springer, New York, 2013). <http://link.springer.com/content/pdf/10.1007/978-0-8176-4948-7.pdf>
23. D. Gabay, B. Mercier, A dual algorithm for the solution of nonlinear variational problems via finite element approximation. *Comput. Math. Appl.* **2**(1), 17–40 (1976)
24. R. Glowinski, On alternating direction methods of multipliers: a historical perspective, in *Modeling, Simulation and Optimization for Science and Technology* (Springer, Dordrecht, 2014), pp. 59–82
25. R. Grosse, R. Raina, H. Kwong, A.Y. Ng, Shift-invariance sparse coding for audio classification (2012). arXiv preprint arXiv: 1206.5241
26. R. Grosse, R. Raina, H. Kwong, A.Y. Ng, Shift-invariance sparse coding for audio classification. arXiv: 1206.5241 [cs, stat] (2012). <http://arxiv.org/abs/1206.5241>. arXiv: 1206.5241
27. S. Gu, W. Zuo, Q. Xie, D. Meng, X. Feng, L. Zhang, Convolutional sparse coding for image super-resolution, in *Proceedings of the IEEE International Conference on Computer Vision* (2015), pp. 1823–1831
28. F. Heide, W. Heidrich, G. Wetzstein, Fast and flexible convolutional sparse coding, in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (IEEE, New York, 2015), pp. 5135–5143
29. J. Huang, T. Zhang, D. Metaxas, Learning with structured sparsity. *J. Mach. Learn. Res.* **12**, 3371–3412 (2011)
30. J. Huang, T. Zhang, et al., The benefit of group sparsity. *Ann. Stat.* **38**(4), 1978–2004 (2010)
31. K. Kavukcuoglu, P. Sermanet, Y.L. Boureau, K. Gregor, M. Mathieu, Y.L. Cun, Learning convolutional feature hierarchies for visual recognition, in *Advances in Neural Information Processing Systems*, ed. by J.D. Lafferty, C.K.I. Williams, J. Shawe-Taylor, R.S. Zemel, A. Culotta, vol. 23 (Curran Associates, Red Hook, 2010), pp. 1090–1098. <http://papers.nips.cc/paper/4133-learning-convolutional-feature-hierarchies-for-visual-recognition.pdf>
32. A. Kyrillidis, L. Baldassarre, M.E. Halabi, Q. Tran-Dinh, V. Cevher, Structured sparsity: discrete and convex approaches, in *Compressed Sensing and Its Applications*. Applied and Numerical Harmonic Analysis, ed. by H. Boche, R. Calderbank, G. Kutyniok, J. Vybíral (Springer, Cham, 2015), pp. 341–387. http://link.springer.com/chapter/10.1007/978-3-319-16042-9_12. https://doi.org/10.1007/978-3-319-16042-9_12
33. P.L. Lions, B. Mercier, Splitting algorithms for the sum of two nonlinear operators. *SIAM J. Numer. Anal.* **16**(6), 964–979 (1979)
34. M.A. Little, N.S. Jones, Generalized methods and solvers for noise removal from piecewise constant signals. II. New methods. *Proc. R. Soc. Lond. A: Math. Phys. Eng. Sci.* [doi:https://doi.org/10.1098/rspa.2010.0674](https://doi.org/10.1098/rspa.2010.0674). <http://rspa.royalsocietypublishing.org/content/early/2011/06/07/rspa.2010.0674>
35. Y.M. Lu, M.N. Do, A theory for sampling signals from a union of subspaces. *IEEE Trans. Signal Process.* **56**, 2334–2345 (2007)
36. J. Mairal, G. Sapiro, M. Elad, Learning multiscale sparse representations for image and video restoration. *Multiscale Model. Simul.* **7**(1), 214–241 (2008)
37. J. Mairal, F. Bach, J. Ponce, G. Sapiro, A. Zisserman, Non-local sparse models for image restoration. in *2009 IEEE 12th International Conference on Computer Vision* (IEEE, New York, 2009), pp. 2272–2279

38. J. Mairal, F. Bach, J. Ponce, Sparse modeling for image and vision processing. *Found. Trends Comput. Graph. Vis.* **8**(2–3), 85–283 (2014). <https://doi.org/10.1561/06000000058>. <http://www.nowpublishers.com/article/Details/CGV-058>
39. Maplesoft, a division of Waterloo Maple Inc. <http://www.maplesoft.com>
40. V. Papyan, M. Elad, Multi-scale patch-based image restoration. *IEEE Trans. Image Process.* **25**(1), 249–261 (2016). <https://doi.org/10.1109/TIP.2015.2499698>
41. V. Papyan, Y. Romano, M. Elad, Convolutional neural networks analyzed via convolutional sparse coding. *J. Mach. Learn. Res.* **18**(83), 1–52 (2017)
42. V. Papyan, J. Sulam, M. Elad, Working locally thinking globally: theoretical guarantees for convolutional sparse coding. *IEEE Trans. Signal Process.* **65**(21), 5687–5701 (2017)
43. Y.C. Pati, R. Rezaifar, P. Krishnaprasad, Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition, in *Asilomar Conference on Signals, Systems and Computers* (IEEE, New York, 1993), pp. 40–44
44. R. Quiroga, Spike sorting. *Scholarpedia* **2**(12), 3583 (2007). <https://doi.org/10.4249/scholarpedia.3583>
45. Y. Romano, M. Elad, Boosting of image denoising algorithms. *SIAM J. Imag. Sci.* **8**(2), 1187–1219 (2015)
46. Y. Romano, M. Elad, Patch-disagreement as a way to improve K-SVD denoising, in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (IEEE, New York, 2015), pp. 1280–1284
47. Y. Romano, M. Protter, M. Elad, Single image interpolation via adaptive nonlocal sparsity-based modeling. *IEEE Trans. Image Process.* **23**(7), 3085–3098 (2014)
48. L.I. Rudin, S. Osher, E. Fatemi, Nonlinear total variation based noise removal algorithms. *Physica D* **60**(1), 259–268 (1992). <http://www.sciencedirect.com/science/article/pii/016727899290242F>
49. C. Rusu, B. Dumitrescu, S. Tsafaris, Explicit shift-invariant dictionary learning. *IEEE Signal Process. Lett.* **21**, 6–9 (2014). http://www.schur.pub.ro/IdeI2011/Articole/SPL_2014_shifts.pdf
50. E. Smith, M.S. Lewicki, Efficient coding of time-relative structure using spikes. *Neural Comput.* **17**(1), 19–45 (2005). <http://dl.acm.org/citation.cfm?id=1119614>
51. A.M. Snijders, N. Nowak, R. Seagraves, S. Blackwood, N. Brown, J. Conroy, G. Hamilton, A.K. Hindle, B. Huey, K. Kimura, S. Law, K. Myambo, J. Palmer, B. Ylstra, J.P. Yue, J.W. Gray, A.N. Jain, D. Pinkel, D.G. Albertson, Assembly of microarrays for genome-wide measurement of DNA copy number. *Nat. Genet.* **29**(3), 263–264 (2001). <https://doi.org/10.1038/ng754>. <https://www.nature.com/ng/journal/v29/n3/full/ng754.html>
52. J. Sulam, M. Elad, Expected patch log likelihood with a sparse prior, in *International Workshop on Energy Minimization Methods in Computer Vision and Pattern Recognition* (Springer, New York, 2015), pp. 99–111
53. J. Sulam, B. Ophir, M. Elad, Image denoising through multi-scale learnt dictionaries, in *2014 IEEE International Conference on Image Processing (ICIP)* (IEEE, New York, 2014), pp. 808–812
54. J.J. Thiagarajan, K.N. Ramamurthy, A. Spanias, Shift-invariant sparse representation of images using learned dictionaries, in *IEEE Workshop on Machine Learning for Signal Processing, 2008, MLSP 2008* (2008), pp. 145–150 <https://doi.org/10.1109/MLSP.2008.4685470>
55. J.A. Tropp, A.C. Gilbert, M.J. Strauss, Algorithms for simultaneous sparse approximation. Part i: greedy pursuit. *Signal Process.* **86**(3), 572–588 (2006)
56. J. Yang, J. Wright, T.S. Huang, Y. Ma, Image super-resolution via sparse representation. *IEEE Trans. Image Process.* **19**(11), 2861–2873 (2010)
57. G. Yu, G. Sapiro, S. Mallat, Solving inverse problems with piecewise linear estimators: From gaussian mixture models to structured sparsity. *IEEE Trans. Image Process.* **21**(5), 2481–2499 (2012). <https://doi.org/10.1109/TIP.2011.2176743>
58. M.D. Zeiler, D. Krishnan, G.W. Taylor, R. Fergus, Deconvolutional networks, in *2010 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, (IEEE, New York, 2010), pp. 2528–2535. http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5539957

59. M. Zeiler, G. Taylor, R. Fergus, Adaptive deconvolutional networks for mid and high level feature learning, in *2011 IEEE International Conference on Computer Vision (ICCV)*, pp. 2018–2025 (2011). doi:10.1109/ICCV.2011.6126474
60. D. Zoran, Y. Weiss, From learning models of natural image patches to whole image restoration, in *2011 IEEE International Conference on Computer Vision (ICCV)* (IEEE, New York, 2011), pp. 479–486. http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=6126278