

# Convexity in source separation: Models, geometry, and algorithms

Michael B. McCoy, Volkan Cevher, Quoc Tran Dinh,  
Afsaneh Asaei, and Luca Baldassarre

Source separation or *demixing* is the process of extracting multiple components entangled within a signal. Contemporary signal processing presents a host of difficult source separation problems, from interference cancellation to background subtraction, blind deconvolution, and even dictionary learning. Despite the recent progress in each of these applications, advances in high-throughput sensor technology place demixing algorithms under pressure to accommodate extremely high-dimensional signals, separate an ever larger number of sources, and cope with more sophisticated signal and mixing models. These difficulties are exacerbated by the need for real-time action in automated decision-making systems.

Recent advances in convex optimization provide a simple framework for efficiently solving numerous difficult demixing problems. This article provides an overview of the emerging field, explains the theory that governs the underlying procedures, and surveys algorithms that solve them efficiently. We aim to equip practitioners with a toolkit for constructing their own demixing algorithms that *work*, as well as concrete intuition for *why* they work.

## *Fundamentals of demixing*

The most basic model for mixed signals is a *superposition model*, where we observe a mixed signal  $\mathbf{z}_0 \in \mathbb{R}^d$  of the form

$$\mathbf{z}_0 = \mathbf{x}_0 + \mathbf{y}_0, \tag{1}$$

and we wish to determine the component signals  $\mathbf{x}_0$  and  $\mathbf{y}_0$ . This simple model appears in many guises. Sometimes, superimposed signals come from basic laws of nature. The amplitudes of electromagnetic waves, for example, sum together at a receiver, making the superposition model (1) common in wireless communications. Similarly, the additivity of sound waves makes superposition models natural in speech and audio processing.

Other times, a superposition provides a useful, if not literally true, model for more complicated nonlinear phenomena. Images, for example, can be modeled as the sum of constituent features—think of stars and galaxies that sum to create an image of a piece of the night sky [1]. In machine learning, superpositions can describe hidden structure [2], while in statistics, superpositions can model gross corruptions to data [3]. These models also appear in texture repair [4], graph clustering [5], and line-spectral estimation [6].

A conceptual understanding of demixing in all of these applications rests on two key ideas.

**Low-dimensional structures:** Natural signals in high dimensions often cluster around low-dimensional structures with few degrees of freedom relative to the ambient dimension [7]. Examples include bandlimited signals, array observations from seismic sources, and natural images. By identifying the convex functions that encourage these low-dimensional structures, we can derive convex programs that disentangle structured components from a signal.

The authors thank Joel A. Tropp for his helpful and detailed comments on this work. MBM is supported by ONR awards N00014-08-1-0883 and N00014-11-1002, AFOSR award FA9550-09-1-064. Work of VC, QTD, and LB is supported in part by the European Commission under Grant MIRG-268398, ERC Future Proof, SNF 200021-132548, SNF 200021-146750 and SNF CRSII2-147633. The work of AA is funded by SNF NCCR IM2.

**Incoherence:** Effective demixing requires more than just structure. To distinguish multiple elements in a signal, the components must look different from one another. We capture this idea by saying that two structured families of signal are *incoherent* if their constituents appear very different from each other. While demixing is impossible without incoherence, sufficient incoherence typically leads to provably correct demixing procedures.

The two notions of structure and incoherence above also appear at the core of recent developments in information extraction from incomplete data in compressive sensing and other linear inverse problems [8, 9]. The theory of demixing extends these ideas to a richer class of signal models, and it leads to a more coherent theory of convex methods in signal processing.

While this article primarily focuses on mixed signals drawn from the superposition model (1), recent extensions to *nonlinear* mixing models arise in blind deconvolution, source separation, and nonnegative matrix factorization [10, 11, 12]. We will see that the same techniques that let us demix superimposed signals reappear in nonlinear demixing problems.

### *The role of convexity*

Convex optimization provides a unifying theme for all of the demixing problems discussed above. This framework is based on the idea that many structured signals possess corresponding convex functions that encourage this structure [9]. By combining these functions in a sensible way, we can develop convex optimization procedures that demix a given observation. The geometry of these functions lets us understand when it is possible to demix a superimposed observation with incoherent components [13]. The resulting convex optimization procedures usually have both theoretical and practical guarantees of correctness and computational efficiency.

To illustrate these ideas, we consider a classical but surprisingly common demixing problem: separating impulsive signals from sinusoidal signals, called the *spikes and sines* model. This model appears in many applications, including star–galaxy separation in astronomy, interference cancellation in communications, inpainting and speech enhancement in signal processing [1, 14].

While individual applications feature additional structural assumptions on the signals, a simple low-dimensional signal model effectively captures the main idea present in all of these works: *sparsity*. A vector  $\mathbf{x}_0 \in \mathbb{R}^d$  is *sparse* if most of its entries are equal to zero. Similarly, a vector  $\mathbf{y}_0 \in \mathbb{R}^d$  is *sparse-in-frequency* if its discrete cosine transform (DCT)  $\mathbf{D}\mathbf{y}_0$  is sparse, where  $\mathbf{D} \in \mathbb{R}^{d \times d}$  is the matrix that encodes the DCT. Sparse vectors capture impulsive signals like pops in audio, while sparse-in-frequency vectors explain smooth objects like natural images. Clearly, such signals look different from one another. In fact, an arbitrary collection of spikes and sines is linearly independent or *incoherent* provided that the collection is not too big [14].

Is it possible to demix a superimposition  $\mathbf{z}_0 = \mathbf{x}_0 + \mathbf{y}_0$  of spikes and cosines into its constituents? One approach is to search for the *sparsest possible* constituents that generate the observation  $\mathbf{z}_0$ :

$$[\hat{\mathbf{x}}, \hat{\mathbf{y}}] := \arg \min_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^n} \{ \|\mathbf{x}\|_0 + \lambda \|\mathbf{D}\mathbf{y}\|_0 : \mathbf{z}_0 = \mathbf{x} + \mathbf{y} \}, \quad (2)$$

where the  $\ell_0$  “norm” measures the sparsity of its input, and  $\lambda > 0$  is a regularization parameter that trades the relative sparsity of solutions. Unfortunately, solving (2) involves an intractable computational problem. However, if we replace the  $\ell_0$  penalty with the convex  $\ell_1$ -norm, we arrive at a classical sparse approximation program [14]:

$$[\hat{\mathbf{x}}, \hat{\mathbf{y}}] := \arg \min_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^n} \{ \|\mathbf{x}\|_1 + \lambda \|\mathbf{D}\mathbf{y}\|_1 : \mathbf{z}_0 = \mathbf{x} + \mathbf{y} \}. \quad (3)$$

This key change to the combinatorial proposal (2) offers numerous benefits. First, the procedure (3) is a convex program, and a number of highly efficient algorithms are available for its solution.

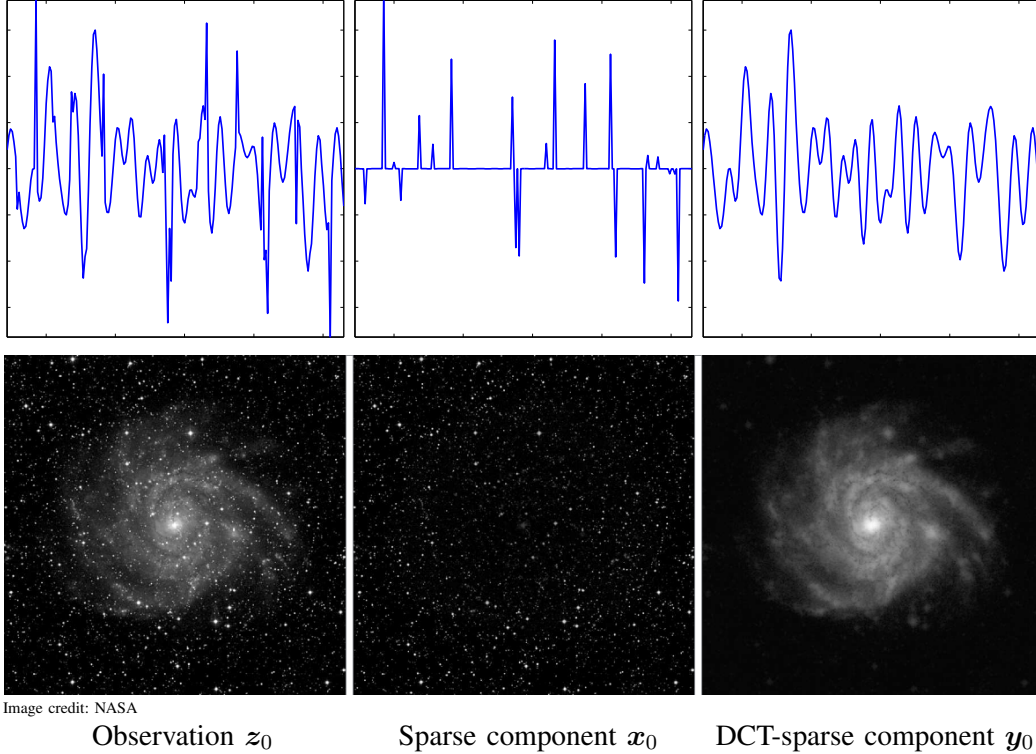


Fig. 1: [Top] A perfect separation of spikes from sinusoids from their additive mixture with (3). The original signal (left) is perfectly separated into its sparse component (center) and its DCT-sparse component (right) [Bottom] Star-galaxy separation using (3) on a real astronomical image. The original (left) is separated into a starfield (center) corresponding to a nearly sparse component and a galaxy (right) corresponding to a nearly DCT-sparse component.

Second, this procedure admits provable guarantees of correctness and noise-stability under incoherence. Finally, the demixing procedure (3) often performs admirably in practice.

Figure 1 illustrates the performance of (3) on both a synthetic signal drawn from the spikes-and-sines model above, as well as on a real astronomical image. The resulting performance for the basic model is quite appealing even for real data that mildly violates the modeling assumptions. Last but not least, this strong baseline performance can be obtained in fractions of seconds with simple and efficient algorithms.

### Outline

The combination of efficient algorithms, rigorous theory, and impressive real-world performance are a hallmark of the convex demixing paradigm described in this article. Below, we provide a unified treatment of demixing problems using convex geometry and optimization starting with Section I. Section II describes some emerging connections between statistics and geometry that characterizes the success and the failure of convex demixing. Section III describes scalable algorithms for practical demixing. Sections IV and V trace the recent frontier in source separation. We not only ground the new theory on compelling signal processing applications but also point out how we can tackle *nonlinear* demixing problems.

## I. DEMIXING MADE EASY

This section provides a recipe to generate a convex program that accepts a mixed signal  $z_0 = x_0 + y_0$  and returns a set of demixed components. The approach requires two ingredients.

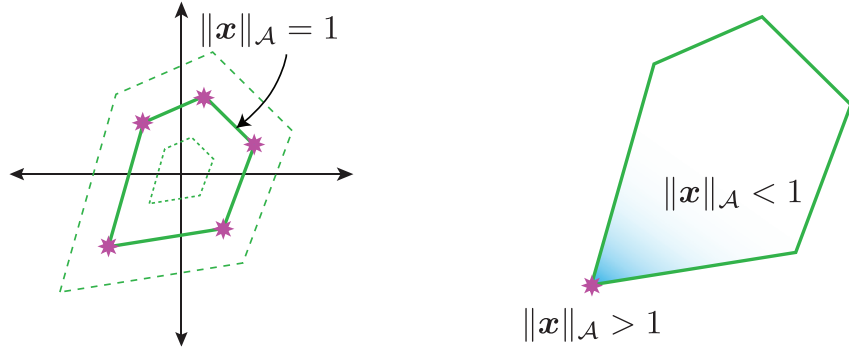


Fig. 2: [Left] An atomic set  $\mathcal{A}$  consisting of five atoms (stars). The “unit ball” of the atomic gauge  $\|\cdot\|_{\mathcal{A}}$  is the closed convex hull of  $\mathcal{A}$  (heavy line). Other level sets (dashed lines) of the gauge are dilations of the unit ball. [Right] At an atom (star), the unit ball of  $\|\cdot\|_{\mathcal{A}}$  tends to have sharp corners. Most perturbations away from this atom increase the value of  $\|\cdot\|_{\mathcal{A}}$ , so the atomic gauge often penalizes complex signals that are comprised of a large number of atoms.

First, we must identify convex functions that promote the structure we expect in  $\mathbf{x}_0$  and  $\mathbf{y}_0$ . Second, we combine these functions together into a convex objective. This simple and versatile approach easily extends to multiple signal components and undersampled observations.

#### Structure-inducing convex functions

We say that a signal has structure when it has fewer degrees of freedom than the ambient space. Familiar examples of structured objects include sparse vectors, sign vectors, and low-rank matrices. It turns out that each of these structured families have an associated convex function, called an atomic gauge, adapted to their specific features [9].

The general principle is simple. Given a set of atoms  $\mathcal{A} \subset \mathbb{R}^d$ , we say that a signal  $\mathbf{x} \in \mathbb{R}^d$  is *atomic* if it is formed by a sum of a small number of scaled atoms. For example, sparse vectors are atomic relative to the set of standard basis vectors because every sparse vector is the sum of just a few standard basis vectors. For a more sophisticated example, recall that the singular value decomposition implies that low-rank matrices are the sum of a few rank-one matrices. Hence, low-rank matrices are atomic relative to the set  $\mathcal{A}$  of all rank-one matrices.

We can define a function that measures the inherent complexity of signals relative to a given set  $\mathcal{A}$ . One natural measure is the *fewest* number of scaled atoms required to write a signal using atoms from  $\mathcal{A}$ , but unfortunately, computing this quantity can be computationally intractable. Instead, we define the *atomic gauge*  $\|\mathbf{x}\|_{\mathcal{A}}$  of a signal  $\mathbf{x} \in \mathbb{R}^d$  by

$$\|\mathbf{x}\|_{\mathcal{A}} := \inf \{ \lambda > 0 : \mathbf{x} \in \lambda \cdot \text{conv}(\mathcal{A}) \},$$

where  $\text{conv}(\mathcal{A})$  is the convex hull of  $\mathcal{A}$ . In other words, the level sets of the atomic gauge are the scaled versions of the convex hull of all the atoms  $\mathcal{A}$  (Figure 2 [Left]).

By construction, atomic gauges are “pointy” at atomic vectors. This property means that most deviations away from the atoms result in a rapid increase in the value of the gauge, so that the function tends to penalize deviations away from simple signals (Figure 2 [Right]). The pointy geometry plays an important role in the theoretical understanding of demixing, as we will see in Section II.

A number of common structured families and their associated gauge functions appear in Table I. More sophisticated examples include gauges for probability measures, cut matrices, and low-rank tensors. We caution, however, that not every atomic gauge is easy to compute, and so we must take care in order to develop *tractable* forms of atomic gauges [9, 16]. Surprisingly, it is sometimes

TABLE I: Example signal structures and their atomic gauges [9, 15]. *The top two rows correspond to vectors while the bottom three refer to matrices. The vector norms extend to matrix norms by treating  $m \times n$  matrices as length- $mn$  vectors. The expression  $\|\mathbf{x}\|_2$  denotes the Euclidean norm of the vector  $\mathbf{x}$ , while  $\sigma_i(\mathbf{X})$  returns the  $i$ th singular value of the matrix  $\mathbf{X}$ .*

Structure	Atomic set	Atomic gauge $\ \cdot\ _{\mathcal{A}}$
Sparse vector	Signed basis vectors $\{\pm \mathbf{e}_i\}$	$\ell_1$ norm $\ \mathbf{x}\ _{\ell_1} = \sum_i  x_i $
Binary sign vector	Sign vectors $\{\pm 1\}^d$	$\ell_\infty$ norm $\ \mathbf{x}\ _{\ell_\infty} = \max_i  x_i $
Low-rank matrix	Rank-1 matrices $\{\mathbf{u}\mathbf{v}^t : \ \mathbf{u}\mathbf{v}^t\ _F = 1\}$	Schatten 1-norm $\ \mathbf{X}\ _{S_1} = \sum_i \sigma_i(\mathbf{X})$
Orthogonal matrix	Orthogonal matrices $\{\mathbf{O} : \mathbf{O}\mathbf{O}^t = \mathbf{I}\}$	Schatten $\infty$ -norm $\ \mathbf{X}\ _{S_\infty} = \sigma_1(\mathbf{X})$
Row-sparse matrix	Matrices w/one nonzero row $\{\mathbf{e}_i \mathbf{v}^t : \ \mathbf{v}\ _2 = 1\}$	Row- $\ell_1$ norm $\ \mathbf{X}\ _{\ell_1/\ell_2}$

easier to compute the value of atomic gauges than it is to compute the (possibly nonunique) decomposition of a vector into its atoms [12]. We will return to the discussion of tractable gauges when we discuss numerical schemes further in Section III.

#### The basic demixing program

Suppose that we know the signal components  $\mathbf{x}_0$  and  $\mathbf{y}_0$  are atomic with respect to the known atomic sets  $\mathcal{A}_x$  and  $\mathcal{A}_y$ . In this section, we describe how to use the atomic gauge functions  $\|\cdot\|_{\mathcal{A}_x}$  and  $\|\cdot\|_{\mathcal{A}_y}$  defined above to help us demix the components  $\mathbf{x}_0$  and  $\mathbf{y}_0$  from the observation  $\mathbf{z}_0$ .

Our intuition developed above indicates that the values  $\|\mathbf{x}_0\|_{\mathcal{A}_x}$  and  $\|\mathbf{y}_0\|_{\mathcal{A}_y}$  are relatively small because the vectors  $\mathbf{x}_0$  and  $\mathbf{y}_0$  are atomic with respect to the atomic sets  $\mathcal{A}_x$  and  $\mathcal{A}_y$ . This suggests that we search for constituents that generate the observation *and* have small atomic gauges. That is, we determine the demixed constituents  $\hat{\mathbf{x}}, \hat{\mathbf{y}}$  by solving

$$[\hat{\mathbf{x}}, \hat{\mathbf{y}}] =: \arg \min_{\mathbf{x}, \mathbf{y} \in \mathbb{R}^d} \{ \|\mathbf{x}\|_{\mathcal{A}_x} + \lambda \|\mathbf{y}\|_{\mathcal{A}_y} : \mathbf{x} + \mathbf{y} = \mathbf{z}_0 \}. \quad (4)$$

The parameter  $\lambda > 0$  negotiates a tradeoff between the relative importance of the atomic gauges, and the constraint  $\mathbf{x} + \mathbf{y} = \mathbf{z}_0$  ensures that our estimates  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{y}}$  satisfy the observation model (1). The hope, of course, is that  $\hat{\mathbf{x}} = \mathbf{x}_0$  and  $\hat{\mathbf{y}} = \mathbf{y}_0$ , so that the demixing program (4) actually identifies the true components in the observation  $\mathbf{z}_0$ .

The demixing program (4) is closely related to linear inverse problems and compressive sampling (CS) [8, 9]. Indeed, the summation map  $(\mathbf{x}, \mathbf{y}) \mapsto \mathbf{x} + \mathbf{y}$  is a linear operator, so demixing amounts to inverting an underdetermined linear system using structural assumptions. The main conceptual difference between demixing and standard CS is that demixing treats the components  $\mathbf{x}_0$  and  $\mathbf{y}_0$  as unrelated structures. Also, unlike conventional CS, demixing does not require exact knowledge of the atomic decomposition, but only the value of the gauge.

The only link between the structures that appears in our recipe comes through the choice of tuning parameter  $\lambda$  in (4), which makes these convex demixing procedures easily adaptable to new problems. In general, determining an optimal value of  $\lambda$  may involve fine tuning or cross-validation, which can be quite computationally demanding in practice. Some theoretical guidance on the explicit choices regularization appears, for example, in [2, 3, 17].

### Extensions

There are many extensions of the linear superposition model (1). In some applications, we are confronted with a signal that is only partially observed—*compressive* demixing. In others, we might consider an observation with additive noise, for instance, or a signal with more than two components. The same ingredients that we introduced above can be used to demix signals from these more elaborate models.

For example, if we only see  $z_0 = \Phi(x_0 + y_0)$ , a linear mapping of the superposition, then we simply update the consistency constraint in the usual demixing program (4) and solve instead

$$[\hat{x}, \hat{y}] =: \arg \min_{x, y \in \mathbb{R}^d} \{ \|x\|_{\mathcal{A}_x} + \lambda \|y\|_{\mathcal{A}_y} : \Phi(x + y) = z_0 \}. \quad (5)$$

Some applications for this undersampled demixing model appear in image alignment [18], robust statistics [5], and graph clustering [19].

Another straightforward extension involves demixing more than two signals. For example, if we observe  $z_0 = x_0 + y_0 + w_0$ , the sum of three structured components, we can determine the components by solving

$$[\hat{x}, \hat{y}, \hat{w}] := \arg \min_{x, y, w \in \mathbb{R}^d} \{ \|x\|_{\mathcal{A}_x} + \lambda_1 \|y\|_{\mathcal{A}_y} + \lambda_2 \|w\|_{\mathcal{A}_w} : x + y + w = z_0 \}, \quad (6)$$

where  $\mathcal{A}_w$  is an atomic set tuned to  $w_0$ , and as before, the parameters  $\lambda_i > 0$  trade off the relative importance of the regularizers. This model appears, for example, in image processing applications where multiple basis representations, such as curvelets, ridgelets, shearlets, etc., explain different morphological components [1]. Further modifications along the lines above extend the demixing framework to a massive number of problems relevant to modern signal processing.

## II. GEOMETRY OF DEMIXING

A critical question we can ask about a demixing program is “When does it work?” Answers to this question can be found by studying the underlying geometry of convex demixing programs. Surprisingly, we can characterize the success and failure of convex demixing *precisely* by leveraging a basic randomized model for incoherence. Indeed, the geometric viewpoint reveals a tight characterization of the success and failure of demixing in terms of geometric parameters that act as the “degrees-of-freedom” of the mixed signal. The consequences for demixing are intuitive: demixing succeeds if and only if the dimensionality of the observation exceeds the total degrees-of-freedom in the signal.

### *Descent cones and the statistical dimension*

Our study of demixing begins with a basic object that encodes the local geometry of a convex function. The *descent cone*  $\mathcal{D}(\mathcal{A}, x)$  at a point  $x$  with respect to an atomic set  $\mathcal{A} \subset \mathbb{R}^d$  consists of the directions where the gauge function  $\|\cdot\|_{\mathcal{A}}$  does not increase near  $x$ . Mathematically, the descent cone is given by

$$\mathcal{D}(\mathcal{A}, x) := \{h : \|x + \tau h\|_{\mathcal{A}} \leq \|x\|_{\mathcal{A}} \text{ for some } \tau > 0\}.$$

The descent cone encodes detailed information about the *local* behavior of the atomic gauge  $\|\cdot\|_{\mathcal{A}}$  near  $x$ . Since local optimality implies global optimality in convex optimization, we can characterize when demixing succeeds in terms of a configuration of descent cones. See Figure 3 for a precise description of this optimality condition.

In order to understand when the geometric optimality condition is likely to hold, we need a measure for the “size” of cones. The most apparent measure of size is perhaps the solid angle, which quantifies the amount of space occupied by a cone. The solid angle, however, proves

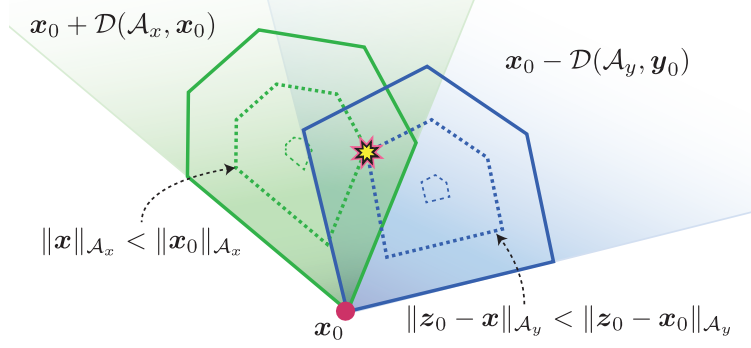


Fig. 3: Geometric characterization of demixing. When the descent cones  $\mathcal{D}(\mathcal{A}_x, \mathbf{x}_0)$  and  $\mathcal{D}(\mathcal{A}_y, \mathbf{y}_0)$  share a line, then there is an optimal point  $\hat{\mathbf{x}}$  (star) for the demixing program (4) not equal to  $\mathbf{x}_0$ . Conversely, demixing can succeed for some value of  $\lambda > 0$  if the two descent cones touch only at the origin. In other words, demixing can succeed if and only if  $\mathcal{D}(\mathcal{A}_x, \mathbf{x}_0) \cap -\mathcal{D}(\mathcal{A}_y, \mathbf{y}_0) = \{\mathbf{0}\}$  [13].

inadequate for describing the intersection of cones even in the simple case of linear subspaces. Indeed, linear subspaces are cones that take up no space at all, but when their dimensions are large enough, any two subspaces will always intersect along a line. Imagine trying to arrange two flat sheets of paper so that they only touch at their centers: impossible!

It turns out that we find a much more informative statistic for demixing when we measure the proportion of space *near* a cone, rather than the proportion of space *inside* the cone.

*Definition 1:* Let  $C \subset \mathbb{R}^d$  be a closed convex cone, and denote by  $\mathbf{II}_C(\mathbf{x}) := \arg \min_{\mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\|$  the closest point in  $C$  to  $\mathbf{x}$ . We define the *statistical dimension*  $\delta(C)$  of a convex cone  $C \subset \mathbb{R}^d$  by

$$\delta(C) := \mathbb{E} \|\mathbf{II}_C(\mathbf{g})\|_2^2, \quad (7)$$

where  $\mathbf{g} \sim \text{NORMAL}(\mathbf{0}, \mathbf{I})$  is a standard Gaussian random variable and the letter  $\mathbb{E}$  denotes the expected value.

The statistical dimension gets its name because it extends many properties of the usual dimension of linear subspaces to convex cones [20], and it is closely related to the Gaussian width used in [9]. Our interest here, however, comes from the interpretation of the statistical dimension as a “size” of a cone. A large statistical dimension  $\delta(C) \approx d$  means that  $\|\mathbf{II}_C(\mathbf{x})\|_2^2$  is large for most  $\mathbf{x} \in \mathbb{R}^d$ —that is, most points lie near the cone. On the other hand, a small statistical dimension implies that most points lie far from  $C$ . We will see below that the statistical dimension of descent cones provides the key parameter for understanding the success and failure of demixing procedures. Of course, a parameter is only useful if we can compute it. Fortunately, the statistical dimension of descent cones is often easy to compute or approximate. Several ready-made statistical dimension formulas and a step-by-step recipe for accurately deriving new formulas appear in [20]. Some useful approximate statistical dimension calculations can also be found in the works [9, 17]. As an added bonus, recent work indicates that statistical dimension calculations are closely related to the problem of finding optimal regularization parameters [17, Thm. 2].

#### *Phase transitions in convex demixing*

The true power of the statistical dimension comes from its ability to predict *phase transitions* in demixing programs. By phase transition, we mean the peculiar behavior where demixing programs switch from near-certain failure to near-certain success within a narrow range of model parameters. While the optimality condition from Figure 3 characterizes the success and failure of demixing, but it is often difficult to certify directly. To understand how demixing operates in

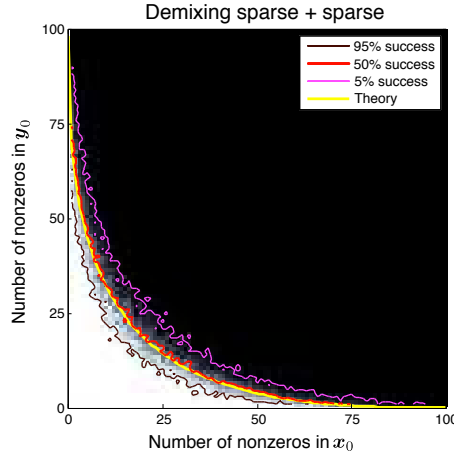


Fig. 4: Phase transitions in demixing. *Phase transition diagram for demixing two sparse signals using  $\ell_1$  minimization [13, 20]. This experiment replaces the DCT matrix  $\mathbf{D}$  in (3) with a random rotation  $\mathbf{Q}$ . The colormap shows the transition from pure success (white) to complete failure (black). The 95%, 50%, and 5% empirical success contours (tortuous curves) appear above the theoretical phase transition curve (yellow) where  $\Delta = 1$ . See [13] for experimental details.*

typical situations, we need an incoherence model. One proposal to model incoherence assumes that the structured signals are oriented generically relative to one another. This is achieved, for example, by assuming that the structured components are drawn structured relative to a rotated atomic set  $\mathbf{Q}\mathcal{A}$ , where  $\mathbf{Q} \in \mathbb{R}^{d \times d}$  is a random orthogonal matrix [13]. Surprisingly, this basic randomized model of incoherence leads to a rich theory with precise guarantees and predict typical behaviors well, and complements other phase transition characterizations in linear inverse problems [21, 22]. Many works propose alternative incoherence models applicable to specific cases, including [3, 9], but these specific choices do not possess known phase transitions. Under the random model of [13], however, a very general theory is available.

*Theorem 1 ([20]):* Suppose that the atomic set of  $\mathbf{x}_0$  is randomly rotated, i.e., that  $\mathcal{A}_x = \mathbf{Q}\tilde{\mathcal{A}}_x$  for some random rotation  $\mathbf{Q}$  and some fixed atomic set  $\tilde{\mathcal{A}}_x$ . Fix a probability tolerance  $\eta \in (0, 1)$ , and define the normalized total statistical dimension  $\Delta := d^{-1}(\delta(\mathcal{D}(\tilde{\mathcal{A}}_x, \mathbf{x}_0)) + \delta(\mathcal{D}(\mathcal{A}_y, \mathbf{y}_0)))$ . Then there is a scalar  $C > 0$  that depends only on  $\eta$  such that

$$\Delta \leq 1 - C/\sqrt{d} \implies \text{demixing can succeed with probability } \geq 1 - \eta$$

$$\Delta \geq 1 + C/\sqrt{d} \implies \text{demixing always fails with probability } \geq 1 - \eta.$$

By “demixing can succeed,” we mean that there exists a regularization parameter  $\lambda > 0$  so that  $(\mathbf{x}_0, \mathbf{y}_0)$  is an optimal point of (4). “Demixing always fails” means that  $(\mathbf{x}_0, \mathbf{y}_0)$  is not an optimal point of (4) fails for *any* parameter  $\lambda > 0$ .

Theorem 1 indicates that demixing exhibits a *phase transition* as the total statistical dimension increases beyond the ambient dimension. Indeed, if the total statistical dimension is slightly less than the ambient dimension, we can be confident that demixing will succeed, but if the total statistical dimension is slightly larger than the ambient dimension, then demixing is hopeless. See Figure 4 for an example of the accuracy of this theory for the MCA model from the introduction when the DCT matrix  $\mathbf{D}$  is replaced with a random rotation  $\mathbf{Q}$ . The agreement between the empirical 50% success line and the curve where  $\Delta = 1$  is remarkable.

This theory extends analogously to the compressive and multiple demixing models (5) and (6). Under a similar incoherence model as above, compressive and multiple demixing are likely to



succeed if and only if the total statistical dimension is slightly less than the number of (possibly compressed) measurements [23, Thm. A]. This fact lets us interpret the statistical dimension  $\delta(\mathcal{D}(\mathcal{A}, \mathbf{x}_0))$  as the degrees-of-freedom of the signal  $\mathbf{x}_0$  with respect to the atomic set  $\mathcal{A}$ . The message is clear: Incoherent demixing can succeed if and only if the total dimension of the observation exceeds the total degrees-of-freedom of the constituent signals.

### III. PRACTICAL DEMIXING ALGORITHMS

*In theory*, many demixing problem instances of the form (4) admit efficient numerical solutions. Indeed, if we can transform these problems into standard linear, cone, or semidefinite formulations, we can apply black-box interior point methods to obtain high-accuracy solutions in polynomial time [24]. *In practice*, however, the computational burden of interior point methods makes these methods impracticable as the dimension  $d$  of the problem grows. Fortunately, a simple and effective iterative algorithm for computing approximate solutions to the demixing program (4) and its extensions can be implemented with just a few lines of high-level code.

#### *Splitting the work*

The simplest and most popular method for iteratively solving demixing programs goes by the name *alternating direction method of multipliers* (ADMM). The key object in this algorithm is the *augmented Lagrangian* function  $L_\rho$  defined by

$$L_\rho(\mathbf{x}, \mathbf{y}, \mathbf{w}) := \|\mathbf{x}\|_{\mathcal{A}_x} + \lambda \|\mathbf{y}\|_{\mathcal{A}_y} + \langle \mathbf{w}, \mathbf{x} + \mathbf{y} - \mathbf{z}_0 \rangle + \frac{1}{2\rho} \|\mathbf{x} + \mathbf{y} - \mathbf{z}_0\|^2,$$

where  $\langle \cdot, \cdot \rangle$  denotes the usual inner product between two vectors and  $\rho > 0$  is a parameter that can be tuned to the problem. Starting with arbitrary points  $\mathbf{x}^1, \mathbf{y}^1, \mathbf{w}^1 \in \mathbb{R}^d$ , the ADMM method generates a sequence of points iteratively as

$$\begin{cases} \mathbf{x}^{k+1} &= \arg \min_{\mathbf{x} \in \mathbb{R}^d} L_\rho(\mathbf{x}, \mathbf{y}^k, \mathbf{w}^k) \\ \mathbf{y}^{k+1} &= \arg \min_{\mathbf{y} \in \mathbb{R}^d} L_\rho(\mathbf{x}^{k+1}, \mathbf{y}, \mathbf{w}^k) \\ \mathbf{w}^{k+1} &= \mathbf{w}^k + (\mathbf{x}^{k+1} + \mathbf{y}^{k+1} - \mathbf{z}_0) / \rho. \end{cases} \quad (8)$$

In other words, the  $\mathbf{x}$ - and  $\mathbf{y}$ -updates iteratively minimize the Lagrangian over just *one* parameter, leaving all others fixed. The alternating minimization of  $L_\rho$  gives the method its name. Despite the simple updates, the sequence  $(\mathbf{x}^k, \mathbf{y}^k)$  of iterates generated in this manner converges to the minimizers  $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$  of the demixing program (4) under fairly general conditions [25].

The key to the efficiency of ADMM comes from the fact that the updates are often easy to compute. By completing the square, the  $\mathbf{x}$ - and  $\mathbf{y}$ -updates above amount to evaluating *proximal operators* of the form

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{x}\|_{\mathcal{A}_x} + \frac{1}{2\rho} \|\mathbf{u}^k - \mathbf{x}\|^2 \quad \text{and} \quad \mathbf{y}^{k+1} = \arg \min_{\mathbf{y} \in \mathbb{R}^d} \lambda \|\mathbf{y}\|_{\mathcal{A}_y} + \frac{1}{2\rho} \|\mathbf{v}^k - \mathbf{y}\|^2, \quad (9)$$

where  $\mathbf{u}^k := \mathbf{z}_0 - \mathbf{y}^k - \rho \mathbf{w}^k$  and  $\mathbf{v}^k := \mathbf{z}_0 - \mathbf{x}^{k+1} - \rho \mathbf{w}^k$ . When solutions to the proximal minimizations (9) are simple to compute, each iteration of ADMM is highly efficient.

Fortunately, proximal operators are easy to compute for many atomic gauges. For example, when the atomic gauge is the  $\ell_1$  norm, the proximal operator corresponds to soft-thresholding by  $\rho$ :

$$\arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{x}\|_{\ell_1} + \frac{1}{2\rho} \|\mathbf{u} - \mathbf{x}\|^2 = \text{soft}(\mathbf{u}, \rho) = \begin{cases} u_i - \rho, & u_i > \rho, \\ 0, & |u_i| \leq \rho, \\ u_i + \rho, & u_i < -\rho. \end{cases}$$

If we replace the  $\ell_1$  norm above with the Schatten-1 norm, then the corresponding proximal operator amounts to soft thresholding the singular values. Numerous other explicit examples of proximal operations appear in [25, Sec. 2.6].

Not all atomic gauges, however, have efficient proximal operations. Even sets with finite number of atoms do not necessarily lead to more efficient proximal maps than sets with an infinite number of atoms. For instance, when the atomic set consists of rank-one matrices with unit Frobenius norm, we have an infinite set of atoms and yet the proximal map can be efficiently obtained via singular value thresholding. On the other hand, when the atomic set consists of rank-one matrices with binary  $\pm 1$  entries, we have a finite set of atoms and yet the best-known algorithm for computing the proximal map requires an intractable amount of computation.

There is some hope, however, even for difficult gauges. Recent algebraic techniques for approximating atomic gauges provide computable proximal operators in a relatively efficient manner, which opens the door to additional demixing algorithms for richer signal structures [9, 16].

### Extensions

While the ADMM method is the prime candidate for solving problem (4), it is not usually the best method for the extensions (5) or (6). In the first case, if  $\Phi$  is a general linear operator, it creates a major computational bottleneck since we need an additional loop to solve the subproblems within the ADMM algorithm. In the latter case, ADMM even loses convergence guarantees [26].

One possible way to handle both problems (5) and (6) is to use decomposition methods. Roughly speaking, these methods decompose problems (5) or (6) into smaller components and then solve the convex subproblem corresponding to each term simultaneously. For example, we can use the decomposition method from [27]:

$$\begin{cases} \mathbf{v}^k &= \mathbf{w}^k + \rho(\Phi(\mathbf{x}^k + \mathbf{y}^k) - \mathbf{z}_0) \\ \mathbf{x}^{k+1} &= \arg \min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{x}\|_{\mathcal{A}_x} + \langle \mathbf{v}^k, \Phi \mathbf{x} \rangle + \frac{1}{2\rho} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \\ \mathbf{y}^{k+1} &= \arg \min_{\mathbf{y} \in \mathbb{R}^d} \lambda \|\mathbf{y}\|_{\mathcal{A}_y} + \langle \mathbf{v}^k, \Phi \mathbf{y} \rangle + \frac{1}{2\rho} \|\mathbf{y} - \mathbf{y}^k\|_2^2 \\ \mathbf{w}^{k+1} &= \mathbf{w}^k + \rho(\Phi(\mathbf{x}^{k+1} + \mathbf{y}^{k+1}) - \mathbf{z}_0). \end{cases} \quad (10)$$

When the parameter  $\rho$  is chosen appropriately, the generated sequence  $\{(\mathbf{x}^k, \mathbf{y}^k)\}$  in (10) converges to the solution of (5). Since the second and the third lines of (10) are independent, it is even possible to solve them in parallel. This scheme easily extends to demixing three or more signals (6).

Another practical method appears in [28]. In essence, this approach combines a dual formulation, Nesterov's smoothing technique, and the fast gradient method [24]. This technique works both for problems (5) and (6), and it possesses a rigorous  $\mathcal{O}(1/k)$  convergence rate.

## IV. EXAMPLES

The ideas above apply to a large number of examples. Here, we highlight some recent applications of convex demixing in signal processing. The first example, texture inpainting, uses a low-rank and sparse decomposition to discover and repair axis-aligned texture in images. The second example explores an application of demixing to direction-of-arrival estimation, where we demix a source covariance from a noise covariance to improve beamforming.

### Texture inpainting

Many natural and man-made images include highly regular textures. These repeated patterns, when aligned with the image frame, tend to have very low rank. Of course, rarely does a natural image consist solely of a texture. Often, though, a background texture is *sparsely* occluded by

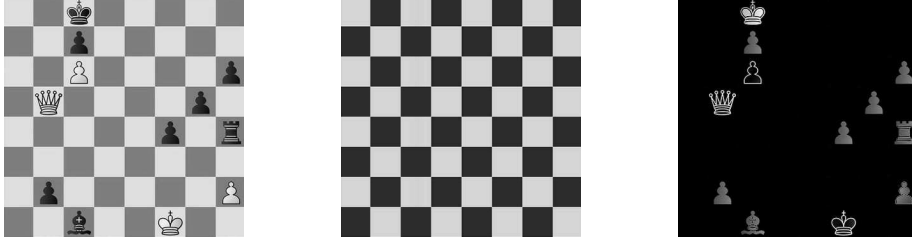


Fig. 5: Texture inpainting (White to move, checkmate in 2). *The rank-sparsity decomposition (11) perfectly separates the chessboard from the pieces. (Left) Original image. (Center) Low-rank component. (Right) Sparse component.*

a untextured component. By modeling the occlusion as an additive error, we can use convex demixing to solve for the underlying texture and extract the occlusion [4].

In this model, we treat the observed digital image  $\mathbf{Z}_0 \in \mathbb{R}^{m \times n}$  as a matrix formed by the sum  $\mathbf{Z}_0 = \mathbf{X}_0 + \mathbf{Y}_0$ , where the textured component  $\mathbf{X}_0$  has low rank and  $\mathbf{Y}_0$  is a sparse corruption or occlusion. The natural demixing program in this setting is the rank-sparsity decomposition [2, 3]:

$$[\hat{\mathbf{X}}, \hat{\mathbf{Y}}] = \arg \min_{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times n}} \|\mathbf{X}\|_{S_1} + \lambda \|\mathbf{Y}\|_1 \quad \text{subject to} \quad \mathbf{X} + \mathbf{Y} = \mathbf{Z}_0, \quad (11)$$

This unsupervised texture-repair method exhibits state-of-the-art performance, exceeding even the quality of a supervised procedure built in to Adobe Photoshop® on some images [4]. When applied, for example, to an image of a chessboard, the method flawlessly recovers the checkerboard from the pieces (Figure 5).

#### *Direction-of-arrival estimation*

We describe a convex demixing program for direction-of-arrival (DOA) estimation. In DOA estimation, we use an array of  $n$  sensors to determine the bearing of multiple sources in wireless communications [11]. When the sources are independent, the joint covariance matrix  $\mathbf{Z}_0$  of all of the signals takes the form  $\mathbf{Z}_0 = \mathbf{A}_0 \mathbf{A}_0^t + \mathbf{Y}_0$  in *expectation*, where the column space of the  $n \times r$  matrix  $\mathbf{A}_0$  encodes the bearing information from  $r$  sources, and  $\mathbf{Y}_0$  is the covariance matrix of the noise at the sensors.

When the number of sources  $r$  is much smaller than the number of sensors  $n$ , the matrix  $\mathbf{X}_0 := \mathbf{A}_0 \mathbf{A}_0^t$  is positive semidefinite and has low rank. Moreover, when the sensor noise is uncorrelated, the matrix  $\mathbf{Y}_0$  is diagonal. Using the atomic gauge recipe from above, we can demix  $\mathbf{X}_0$  and  $\mathbf{Y}_0$  from the empirical covariance matrix  $\mathbf{Z}_0$  by setting

$$[\hat{\mathbf{X}}, \hat{\mathbf{Y}}, \hat{\mathbf{E}}] = \arg \min_{\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times n}} \|\mathbf{X}\|_{S_1^+} + \|\mathbf{Y}\|_{\text{diag}} + \lambda \|\mathbf{E}\|_{\text{Fro}}^2 \quad \text{subject to} \quad \mathbf{X} + \mathbf{Y} + \mathbf{E} = \mathbf{Z}_0, \quad (12)$$

where  $\mathbf{E}$  absorbs the deviations in the expectation model due to the finite sample size. Here,  $\|\cdot\|_{S_1^+}$  is the atomic gauge generated by positive semidefinite rank-one matrices, which is equal to the trace for positive semidefinite matrices, but returns  $+\infty$  when its argument has a negative eigenvalue. Similarly, the gauge  $\|\cdot\|_{\text{diag}}$  is the atomic gauge generated by the set of all diagonal matrices, and so it is equal to zero on diagonal matrices but  $+\infty$  otherwise. The norm  $\|\cdot\|_{\text{Fro}}$  is the usual Frobenius norm on a matrix. The results of [11] relate the success of a similar problem to the geometric problem of ellipsoid fitting, and show that under some incoherence conditions convex demixing succeeds.

In DOA estimation, the source covariance matrix plays a key role in estimating the source directions [29]. For instance, the multiple signal classification (MUSIC) algorithm exploits the nullspace of the source covariance matrix to localize the sources. In the presence of white additive

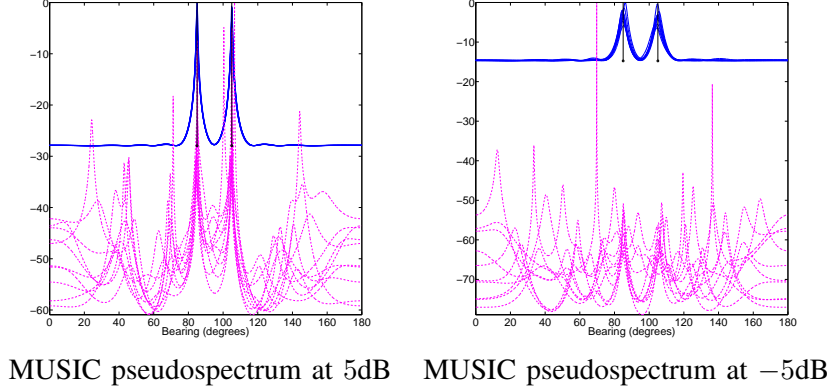


Fig. 6: Enhancing DOA estimation. The MUSIC pseudospectrum based on the demixed estimate  $\widehat{\mathbf{X}}$  (solid blue lines) from (12) is significantly more informative for the source bearings than the MUSIC pseudospectrum based on the raw covariance  $\mathbf{Z}_0$  (dashed magenta lines).

Gaussian noise, the empirical covariance estimate becomes corrupted, deteriorating the bearing estimates generated by MUSIC.

Figure 6 shows how the demixing procedure (12) can significantly boost the performance of MUSIC under additive noise. In this experiment, we generate an array data for  $r = 2$  sources and  $n = 10$  sensors with signal-to-noise ratios of 5dB and -5dB. We simulate the data and compute the empirical covariance matrix  $\mathbf{Z}_0$ . Then we estimate the source covariance  $\mathbf{X}_0$  using the demixed output  $\widehat{\mathbf{X}}$  of (12). We compare the performance of MUSIC with given the raw empirical covariance  $\mathbf{Z}_0$  and the demixed estimate  $\widehat{\mathbf{X}}$ .

At 5dB SNR, about one-third of the DOA estimates of the MUSIC algorithm with  $\mathbf{Z}_0$  are more than three degrees off of the true bearings. At -5dB, MUSIC's performance on the raw covariance is even worse: 90% of the estimated bearings are off by three degrees or more. In contrast, the MUSIC algorithm using the demixed estimate  $\widehat{\mathbf{X}}$  provides consistently accurate bearing estimates.

## V. HORIZONS: NONLINEAR SEPARATION

We conclude our demixing tutorial with some promising directions for the future. In many applications, the constituent signals are tangled together in a *nonlinear* fashion [10, 12]. While this situation would seem to rule out the linear superposition model considered above, we can leverage the same convex optimization tools to obtain demixing guarantees and often return to a linear model using a technique called *semidefinite relaxation*.

We describe the basic idea behind this maneuver with a concrete application: *blind deconvolution*. Convolved signals appear frequently in communications due, for example, to multipath channel effects. When the channel is known, removing the channel effects is a difficult but well-understood linear inverse problem. With blind deconvolution, however, we see only the convolved signal  $\mathbf{z}_0 = \mathbf{x}_0 * \mathbf{y}_0$  from which we must determine both the channel  $\mathbf{x}_0 \in \mathbb{R}^m$  and the source  $\mathbf{y}_0 \in \mathbb{R}^d$ .

While the convolution  $\mathbf{x}_0 * \mathbf{y}_0$  involves nonlinear interactions between  $\mathbf{x}_0$  and  $\mathbf{y}_0$ , the convolution is in fact *linear* in the matrix formed by the outer product  $\mathbf{x}_0 \mathbf{y}_0^t$ . In other words, there is a linear operator  $\mathcal{C}: \mathbb{R}^{m \times d} \rightarrow \mathbb{R}^{m+d}$  such that

$$\mathbf{z}_0 = \mathcal{C}(\mathbf{X}_0) \quad \text{where} \quad \mathbf{X}_0 := \mathbf{x}_0 \mathbf{y}_0^t.$$

The matrix  $\mathbf{X}_0$  has rank one by definition, so it is natural use the Schatten 1-norm to search for

low-rank matrices that generate the observed signal:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathbb{R}^{m \times d}} \|\mathbf{X}\|_{S_1} \quad \text{subject to} \quad \mathbf{z}_0 = \mathcal{C}(\mathbf{X}).$$

This is the basic idea behind the convex approach to blind deconvolution of [10].

The implications of the non-linear demixing example above are far reaching. There are large classes of signal and mixing models that support efficient, provable, and stable demixing. Viewing different demixing problems within a common framework of convex optimization, we can leverage decades of research in various diverse disciplines from applied mathematics to signal processing, and from theoretical computer science to statistics. We expect that the diversity of convex demixing models and geometric tools will also inspire the development of new kinds of scalable optimization algorithms that handle non-conventional cost functions along with atomic gauges [30].

#### REFERENCES

- [1] J.-L. Starck, F. Murtagh, and J. M. Fadili, *Sparse image and signal processing*. Cambridge: Cambridge University Press, 2010, wavelets, curvelets, morphological diversity.
- [2] V. Chandrasekaran, S. Sanghavi, P. A. Parrilo, and A. S. Willsky, “Rank-sparsity incoherence for matrix decomposition,” *SIAM J. Optim.*, vol. 21, no. 2, pp. 572–596, 2011.
- [3] E. J. Candès, X. Li, Y. Ma, and J. Wright, “Robust principal component analysis?” *J. Assoc. Comput. Mach.*, vol. 58, no. 3, pp. 1–37, May 2011. [Online]. Available: <http://arxiv.org/pdf/0912.3599>
- [4] X. Liang, X. Ren, Z. Zhang, and Y. Ma, “Repairing sparse low-rank texture,” in *Computer Vision–ECCV 2012*. Springer, 2012, pp. 482–495.
- [5] Y. Chen, A. Jalali, S. Sanghavi, and C. Caramanis, “Low-rank matrix recovery from errors and erasures,” *IEEE Trans. Inform. Theory.*, 2013, to appear.
- [6] B. N. Bhaskar, G. Tang, and B. Recht, “Atomic norm denoising with applications to line spectral estimation,” *preprint*, 2013. [Online]. Available: <http://arxiv.org/abs/1204.0562>
- [7] R. G. Baraniuk, V. Cevher, and M. B. Wakin, “Low-dimensional models for dimensionality reduction and signal recovery: A geometric perspective,” *Proc. IEEE*, vol. 98, no. 6, pp. 959–971, 2010.
- [8] E. J. Candès and M. B. Wakin, “An introduction to compressive sampling,” *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 21–30, 2008.
- [9] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, “The convex geometry of linear inverse problems,” *Found. Comput. Math.*, vol. 12, no. 6, pp. 805–849, 2012.
- [10] A. Ahmed, B. Recht, and J. Romberg, “Blind deconvolution using convex programming,” *arXiv preprint arXiv:1211.5608*, 2012.
- [11] J. Saunderson, V. Chandrasekaran, P. A. Parrilo, and A. S. Willsky, “Diagonal and low-rank matrix decompositions, correlation matrices, and ellipsoid fitting,” *SIAM J. Matrix Anal. Appl.*, vol. 33, no. 4, pp. 1395–1416, 2012.
- [12] V. Bittorf, C. Ré, B. Recht, and J. A. Tropp, “Factoring nonnegative matrices with linear programs,” in *Advances in Neural Information Processing Systems 25 (NIPS)*, December 2012, pp. 1223–1231.
- [13] M. B. McCoy and J. A. Tropp, “Sharp recovery bounds for convex deconvolution, with applications,” *preprint*, 2012, [arXiv:1205.1580v1](http://arxiv.org/abs/1205.1580v1).
- [14] D. L. Donoho and X. Huo, “Uncertainty principles and ideal atomic decomposition,” *IEEE Trans. Inform. Theory*, vol. 47, no. 7, pp. 2845–2862, Aug. 2001.
- [15] S. S. Chen, D. L. Donoho, and M. A. Saunders, “Atomic decomposition by basis pursuit,” *SIAM J. Sci. Comput.*, vol. 20, no. 1, pp. 33–61, 1998.

- [16] F. Bach, “Structured sparsity-inducing norms through submodular functions,” *Advances in Neural Information Processing Systems*, pp. 118–126, 2010.
- [17] R. Foygel and L. Mackey, “Corrupted sensing: Novel guarantees for separating structured signals,” *preprint*, May 2013. [Online]. Available: <http://arxiv.org/abs/1305.2524>
- [18] Y. Peng, A. Ganesh, J. Wright, W. Xu, and Y. Ma, “RASL: Robust alignment by sparse and low-rank decomposition for linearly correlated images,” *IEEE Trans. Pattern Anal.*, vol. 34, no. 11, pp. 2233–2246, 2012.
- [19] Y. Chen, A. Jalali, S. Sanghavi, and C. Caramanis, “Clustering partially observed graphs via convex optimization,” in *International Symposium on Information Theory (ISIT)*, 2011.
- [20] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp, “Living on the edge: A geometric theory of phase transitions in convex optimization,” *preprint*, March 2013, arXiv:1303.6672.
- [21] D. L. Donoho and J. Tanner, “Precise undersampling theorems,” *Proc. IEEE*, vol. 98, no. 6, pp. 913–924, Jun. 2010.
- [22] M. Bayati, M. Lelarge, and A. Montanari, “Universality in polytope phase transitions and message passing algorithms,” *preprint*, July 2012, arXiv:1207.7321.
- [23] M. B. McCoy and J. A. Tropp, “The achievable performance of convex demixing,” *preprint*, 2013, arXiv:1309.7478.
- [24] Y. Nesterov, *Introductory lectures on convex optimization: a basic course*, ser. Applied Optimization. Kluwer Academic Publishers, 2004, vol. 87.
- [25] P. L. Combettes and V. R. Wajs, “Signal recovery by proximal forward-backward splitting,” *Multiscale Model. Simul.*, vol. 4, pp. 1168–1200, 2005.
- [26] C. Chen, B. S. He, Y. Ye, and X. Yuan, “The direct extension of admm for multi-block convex minimization problems is not necessarily convergent,” *Optimization Online*, 2013.
- [27] G. Chen and M. Teboulle, “A proximal-based decomposition method for convex minimization problems,” *Math. Program.*, vol. 64, pp. 81–101, 1994.
- [28] I. Necoara and J. Suykens, “Applications of a smoothing technique to decomposition in convex optimization,” *IEEE Trans. Automatic control*, vol. 53, no. 11, pp. 2674–2679, 2008.
- [29] H. L. V. Trees, *Optimum Array Processing: Part IV of Detection, Estimation, and Modulation Theory*. John Wiley and Sons, Inc., 2002, vol. Print ISBN: 9780471093909.
- [30] Q. T. Dinh, A. Kyrillidis, and V. Cevher, “Composite self-concordant minimization,” Lab. Inform. Infer. Sys. (LIONS), EPFL, Switzerland, Tech. Report, January 2013.