

# Multireference Alignment is Easier with an Aperiodic Translation Distribution

Emmanuel Abbe<sup>1,2</sup>, Tamir Bendory<sup>1</sup>, William Leeb<sup>1</sup>, João Pereira<sup>1</sup>, Nir Sharon<sup>1</sup>, and Amit Singer<sup>1,3</sup>

<sup>1</sup>The Program in Applied and Computational Mathematics, Princeton University, Princeton, NJ, USA

<sup>2</sup>Electrical Engineering Department, Princeton University, Princeton, NJ, USA

<sup>3</sup>Department of Mathematics, Princeton University, Princeton, NJ, USA

**Abstract**—Multireference alignment refers to the problem of estimating a signal from its circularly translated copies in the presence of noise. Previous papers showed that if the translations are drawn from the uniform distribution, then the sample complexity of the problem scales as  $1/\text{SNR}^3$  in the low SNR regime. In this work, we show that the sample complexity for any aperiodic translation distribution scales as  $1/\text{SNR}^2$  in the low SNR regime. This rate is achieved by a simple spectral algorithm. We propose two additional algorithms based on non-convex optimization and expectation-maximization. We also draw a connection between the multireference alignment problem and the spiked covariance model.

**Index Terms**—multireference alignment, spectral algorithm, method of moments, spiked covariance model, non-convex optimization, expectation-maximization, cryo-EM

## I. INTRODUCTION

The problem of multireference alignment (MRA) arises in a variety of engineering and scientific applications, among them structural biology [1], [2], [3], [4], [5], radar [6], [7], robotics [8] and image processing [9], [10], [11]. In these applications, one aims to estimate a signal from its translated or rotated noisy copies. The problem also serves as a simplified model for more general problems like single-particle cryo-electron microscopy, in which a 3D volume is recovered from unknown 2D projections [12], [13].

In this paper, we focus on the one-dimensional (1D) discrete MRA problem on a circle. In this model, we acquire  $N$  measurements from the model

$$y = R_r x + \varepsilon, \quad (\text{I.1})$$

where  $\varepsilon$  is a vector of i.i.d. normal variables with zero mean and variance  $\sigma^2$ . The operator  $R_r$  translates a signal  $x \in \mathbb{R}^L$  circularly by  $r$  elements, namely,  $(R_r x)[i] = x[i-r]$ , where all indices should be considered as modulo  $L$ . The translation  $r$  is drawn from some unknown distribution  $\rho$  on  $[0, \dots, L-1]$ . Our goal is then to estimate the signal  $x$ , up to global cyclic translation, from

$$y_j = R_{r_j} x + \varepsilon_j, \quad j = 1, \dots, N. \quad (\text{I.2})$$

EA was partly supported by the Bell Labs Prize, the NSF CAREER Award CCF-1552131, ARO grant W911NF-16-1-0051, NSF Center for the Science of Information CCF-0939370, and the Google Faculty Research Award. TB, WL, JP, NS, and AS were partially supported by Award Number R01GM090200 from the NIGMS, the Simons Foundation Investigator Award and Simons Collaborations on Algorithms and Geometry, the Moore Foundation Data-Driven Discovery Investigator Award, and AFOSR FA9550-17-1-0291.

Note that while the translations  $r_j$  are unknown, their estimation is not the primary goal of the problem; the task is to estimate the signal  $x$ . The translations are frequently called *latent* or *hidden* variables. Figure I.1 illustrates the MRA problem in different noise levels.

Previous approaches for estimating  $x$  from (I.2) can be broadly classified into two main categories. The first approach is based on estimating the translations  $r_j$ , aligning all observations and averaging them to suppress the noise. However, alignment is impossible in low signal-to-noise ratios (SNRs) [14], defined here as  $\text{SNR} = \|x\|^2/\sigma^2$ . An alternative approach aims at estimating the signal  $x$  directly. Existing methods bypass the need to estimate the translations by employing expectation-maximization (EM) methods or by using features that are invariant under translation [15]. Section II is devoted to a detailed discussion on existing results and algorithms for MRA. In this paper, we take a different route by trying to estimate both the signal and the distribution of translations  $\rho$  simultaneously. When  $\rho$  is aperiodic, it turns out this is a simpler problem than ignoring  $\rho$  estimating  $x$  alone.

In [16], [17], it was proven that when  $\rho$  is the uniform distribution, then in the low SNR, or large  $\sigma$ , regime, the number of measurements  $N$  needs to scale like  $1/\text{SNR}^3$  in order to keep a constant estimation error for signals with non-vanishing Discrete Fourier Transform (DFT). In other words, the sample complexity of the problem scales as  $1/\text{SNR}^3$ . In this work, we show that any signal with non-vanishing DFT can be estimated at sample complexity scaling like  $1/\text{SNR}^2$  if the translation distribution  $\rho$  is aperiodic, meaning there is no  $0 \leq \ell \leq L-1$  where  $\rho[i+\ell] + \rho[i]$  for all  $i$ . This rate is optimal and can be provably achieved by a spectral algorithm based on the first two moments of the data. The main result of this paper is stated as follows:

**Main Result (informal):** Consider the model (I.2) and suppose that  $x \in \mathbb{R}^L$  has a non-vanishing DFT. If  $\rho$  is aperiodic, then  $x$  can be estimated, up to circular translation, from the first two moments of the data. As a consequence, the sample complexity grows like  $1/\text{SNR}^2$ . This sample complexity can be achieved by a spectral algorithm. Conversely, the sample complexity for any periodic distribution with periodicity smaller than  $L/2$  scales like  $1/\text{SNR}^3$ .

From a computational perspective, the proposed framework is based on a reliable estimation of the first two moments of the

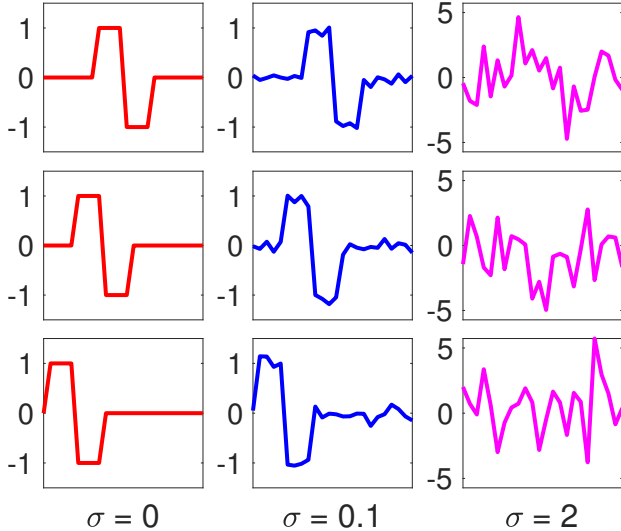


Fig. 1.1. The figures illustrates the MRA measurements according to (I.2). The left column presents three measurements with different translations in the absence of noise. In this case, because the solution is defined up to translation, each measurement is a solution. The middle and right columns show measurements with the same translations and low and high noise levels, respectively.

data. Hence, it requires only one pass over the measurements, low storage resources and is computationally efficient. To estimate the signal from the estimated moments, we propose, in addition to the aforementioned spectral algorithm, a non-convex least-squares (LS) algorithm. While the problem is non-convex, it empirically converges to the underlying signal, in the absence of noise, from a random initialization. We also suggest an expectation-maximization (EM) algorithm.

The outline of the paper is as follows. Section II provides a detailed discussion of existing results and algorithms for MRA. In Section III we show that if the distribution is aperiodic, then any signal with non-vanishing DFT can be estimated from its first and second moments, namely, the sample complexity is upper bounded proportionally to  $1/\text{SNR}^2$ . Section IV draws the connections between the MRA model and the well-studied spiked covariance model [18], [19], [20], [21], [22]. In Section V we prove that the sample complexity is also lower bounded proportionally to  $1/\text{SNR}^2$ . We also show that the sample complexity of any periodic distribution of translations with period less than  $L/2$  scales as  $1/\text{SNR}^3$ . This is an extension of the results of [16] which considered the uniform distribution case. Section VI discusses and analyzes alternative algorithmic methods based on LS and EM. Section VII examines the performance of the proposed algorithms by numerical simulations. Section VIII concludes the paper and proposes potential future extensions.

Throughout the paper we use the following notation. An estimator of a signal  $z \in \mathbb{R}^L$  is denoted by  $\hat{z}$ . We assume throughout that all signals are defined cyclically; that is, all indices should be considered modulo  $L$ . The indices will range from 0 to  $L - 1$ . The DFT of  $z$  is defined by  $(Fz)[k] = \sum_{i=0}^{L-1} z[i]e^{-2\pi i k i/L}$ , where  $\iota = \sqrt{-1}$ . We use  $C_z$  for a circulant matrix whose first column is  $z$ , namely,  $C_z[i, j] = z[i - j]$ . A diagonal matrix whose diagonal is  $z$

is denoted by  $D_z$ . We reserve  $\mathbb{E}, *$  and  $\odot$  for expectation, convolution and entry-wise product, respectively. The  $L$ -simplex is denoted by  $\Delta^L$ . That is to say,  $z \in \Delta^L$  implies that  $z[i] \geq 0$  for all  $i$  and  $\sum_{i=0}^{L-1} z[i] = 1$ . The first and the second moments of the data are denoted by  $\mu := \mathbb{E}\{y\}$  and  $M := \mathbb{E}\{yy^T\}$ , respectively.

Because the observations are invariant to a circular shift of  $x$  and  $\rho$ , we define the normalized squared error of an estimator  $\hat{x}$  of  $x$  by:

$$\min_{s \in \mathbb{Z}_L} \frac{\|R_s \hat{x} - x\|_2^2}{\|x\|_2^2}. \quad (\text{I.3})$$

That is, we look at the relative error of the best-aligned version of  $\hat{x}$  with  $x$ . This measure of error is invariant to translations of  $x$  and  $\rho$ .

## II. RELATED WORK

### A. Multireference alignment via synchronization

Given the translations  $r_j$ , the MRA problem (I.2) is easy. One trivial unbiased estimator of  $x$  is given by aligning all measurements and then averaging to suppress the noise, namely,

$$\hat{x} = \frac{1}{N} \sum_{j=1}^N R_{r_j}^{-1} y_j. \quad (\text{II.1})$$

The variance of this estimator is  $\sigma^2/N$  and therefore the number of measurements  $N$  needs to scale like  $\sigma^2$  to retain a constant estimation error. In other words, the sample complexity grows like  $1/\text{SNR}$ . One can replace (II.1) with other estimators, such as James-Stein shrinkage [23], [24], [25], but it will not change the asymptotic sample complexity. In practice, we do not have access to the underlying translations. However, if one can obtain a reliable estimation of the unknown translations  $\hat{r}_j$ , then one can estimate  $x$  by the sample mean as in (II.1) at sample complexity growing as  $1/\text{SNR}$ . This motivates the design of synchronization methods that aim to estimate the translations  $r_j$  from the data  $y_j$ .

A naïve approach for synchronization could be to fix one observation as a template, say  $y_1$ , and estimate the relative translation of each  $y_j$ , with respect to  $y_1$ , by the peak of their cross-correlation:

$$\hat{r}_j = \arg \max_{\ell} \sum_{i=0}^{L-1} y_1[i] y_j[i + \ell].$$

This approach may work in the high SNR regimes, but fails as the noise level increases (see for instance Figure I.1 in [15]). Many alternative synchronization methods were proposed in the literature. For instance, the angular synchronization method aims at aligning all pairwise observations simultaneously [26], [27], [28], [29], [30], [31]. Other methods propose to align through different semidefinite programs (SDPs) [32], [33], [34], [35]. However, alignment is impossible below a critical SNR threshold, no matter how many measurements are acquired. For instance, for the continuous counterpart of (I.1), it has been shown that the Crámer-Rao lower bound is proportional to  $\sigma^2$  and does not depend on  $N$ . This bound holds even if the sought signal is known [14].

### B. Multireference alignment in low SNR

This section reviews recent works on MRA in the low SNR regime, in which alignment is impossible. The key idea is to estimate the signal directly, without estimating the translations beforehand. As will be emphasized throughout, all these works did not consider the translation distribution  $\rho$ , or forced it to be uniform.

In [16], [17], it was shown that if the translations are uniformly distributed, namely,  $r \sim \text{Uniform}[0, 1, \dots, L-1]$ , then the number of measurements  $N$  needs to scale like  $1/\text{SNR}^3$  to retain a constant estimation error. The analysis of the uniform distribution is of particular interest since, no matter what  $\rho$  is, one can always enforce it to be uniform. This can be done simply by reshuffling all measurements by  $z_j = R_{r'_j} y_j$ , where  $r'_j$  are drawn from the uniform distribution. The new set of measurements  $z_j$  obeys the MRA model (I.2) with uniform translation distribution. However, as will be shown, this is in general a bad strategy, since the uniform distribution has a sample complexity scaling as  $1/\text{SNR}^3$ .

From the algorithmic point-of-view, a recent paper [15] proposes a method that completely overcomes the need to estimate the translations. The core idea is to estimate features of the underlying signal that are invariant under cyclic translation. Particularly, it was proposed to estimate the mean, power spectrum and bispectrum of the signal from the moments of the data. Since these invariant features are third degree polynomials in the signal, they can be estimated accurately at sample complexity growing like  $1/\text{SNR}^3$ . Using these invariant features, one can recover the signal to arbitrary accuracy as  $N \rightarrow \infty$  using a variety of algorithms [15]. Since this method requires only one pass over the data, it can be performed in a streaming mode, can be parallelized, requires low storage resources of  $\mathcal{O}(L^2)$  and has low computational load. The framework proposed in this paper is also based on estimating moments of the data and therefore enjoys the same advantages; however, since we only require second-order moments, we bring the sample complexity down to  $1/\text{SNR}^2$ .

Another approach for MRA is to apply an EM algorithm [36]. EM is an iterative algorithm that aims to find the marginal maximum-likelihood estimator. It is used ubiquitously in many statistical models, such as the Gaussian mixture model [37]. For the MRA model (I.1), this algorithm takes a simple form and consists of two steps at each iteration [15]. Given a current estimation  $x_{k-1}$ , the first step (called the E-step) computes a set of weights which can be understood as the translation distribution of each measurement  $y_j$ , if  $x_{k-1}$  was the underlying signal. These weights are computed by

$$w_k^{\ell,j} = C_k^j e^{-\frac{1}{2\sigma^2} \|R_\ell x_{k-1} - y_j\|_2^2},$$

where  $C_k^j$  is a normalization factor so that  $\sum_\ell w_k^{\ell,j} = 1$ . Then, the signal estimation is updated by marginalizing over the distributions and averaging (called the M-step):

$$x_k = \frac{1}{N} \sum_{j=1}^N \sum_{\ell=0}^{L-1} w_k^{\ell,j} R_\ell^{-1} y_j. \quad (\text{II.2})$$

The EM algorithm enjoys an excellent numerical performance; however its computational load and storage requirements are

heavy since it passes through all the data at each iteration. In Section VI-B, we modify the standard EM algorithm to take the distribution into account.

### III. PROVABLE ALGORITHM BASED ON THE FIRST TWO MOMENTS

In this section, we provide a spectral algorithm to estimate the signal, up to cyclic translation, from the first and second moments of the data, if the translation distribution is aperiodic. We prove that this algorithm estimates the signal exactly with high probability in the limit of SNR tending to 0 with a growing number of samples; we will describe the asymptotic model more precisely in Section III-C. Because the method relies on only second-order information, its sample complexity in this case only grows like  $1/\text{SNR}^2$ , compared to sample complexity growing as  $1/\text{SNR}^3$  if the translation distribution is periodic (with period smaller than  $L/2$ ; see Section III-D). As we prove in Section V,  $1/\text{SNR}^2$  is indeed the sample complexity for aperiodic distributions.

#### A. Moments of $R_\ell x$

Before describing the algorithm, we will describe a few basic properties of the moments of the random vectors  $R_\ell x$ . We will first consider the first moment of the translated signal (without additive noise), which we denote  $\mu = \mathbb{E}[R_r x]$ , where the expectation is over the random  $r \sim \rho$ . This is equal to the convolution of  $x$  with  $\rho$ ; that is,

$$\mu = x * \rho = C_x \rho, \quad (\text{III.1})$$

where  $C_x$  is a circulant matrix with  $x$  as its first column. In this case, the convolution theorem implies

$$F\mu = Fx \odot F\rho, \quad (\text{III.2})$$

where  $\odot$  and  $F$  denote entry-wise product and Fourier transform, respectively. We can estimate the first moment from the noisy observations  $y_j = R_{r_j} x + \varepsilon_j$  by

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N y_i. \quad (\text{III.3})$$

Note that if  $L$  and  $\sigma$  are fixed, then  $\hat{\mu}$  is a consistent estimator of  $\mu$  as  $N \rightarrow \infty$ .

The second moment of  $R_r x$  is defined as

$$M = \mathbb{E} \left[ (R_r x)(R_r x)^T \right],$$

where the expectation is again taken over the random  $\ell \sim \rho$ . It can be verified that

$$M = C_x D_\rho C_x^T, \quad (\text{III.4})$$

where  $D_\rho$  is a diagonal matrix of  $\rho$ . The unbiased second moment of  $R_r x$  is then estimated from the observations  $y_j$  by:

$$\hat{M} = \frac{1}{N} \sum_{i=1}^N y_i y_i^T - \sigma^2 I. \quad (\text{III.5})$$

As with the first moment, when  $L$  and  $\sigma$  are fixed then  $\hat{M}$  is a consistent estimator of  $M$  as  $N \rightarrow \infty$ .

The following result shows that under a rather weak condition on  $\rho$ , there exists only one signal (up to translation) that agrees with the second moment data exactly. For this condition, recall that a number  $a$  is relatively prime to  $b$  if their largest common divisor is 1.

**Proposition III.1.** *Let  $\rho \in \Delta^L$  be a distribution with  $(F\rho)[k] \neq 0$  for some  $k$  that is relatively prime to  $L$ . If the DFT of  $x$  is non-vanishing, then it is uniquely determined (up to translation) from the first two moments  $\mu$  and  $M$ .*

*Proof.* See Appendix A.  $\square$

The power spectrum of the signal  $P_x[k] := |(Fx)[k]|^2$ , which is the Fourier transform of the signal's auto-correlation, plays an important role in the analysis. Recall that the  $j$ th entry of the auto-correlation of each measurement is  $a[j] = \sum_{i=0}^{L-1} y[i]y[i+j]$ , and therefore  $P_x$  can be estimated directly from  $M$  by averaging over the measurement's auto-correlations. Alternatively, it can be estimated in the Fourier domain directly [15]:

$$\hat{P}_x = \frac{1}{N} \sum_{j=1}^N P_{y_j} - \sigma^2 L \mathbf{1}, \quad (\text{III.6})$$

where  $\mathbf{1}$  is a vector of ones. This is an unbiased estimator, with each coordinate having variance  $\mathcal{O}(\sigma^4/N)$ .

### B. Moment inversion when $\rho$ has a unique entry

The key observation driving the algorithm we will describe is that when  $\rho$  has at least one distinct entry, and if  $x$  has non-zero DFT, then  $x$  can be recovered exactly from the first two moments  $\mu$  and  $M$  and the power spectrum  $P_x$ .

We first recall the factorization  $M = C_x D_\rho C_x^T$  from equation (III.4). The circulant matrix  $C_x$  is diagonalized by the Fourier matrix  $F$  as,

$$C_x = F^{-1} D_{Fx} F.$$

Consequently, conjugating  $M$  by the matrix  $F^{-1} D_{1/|P_x|^{1/2}} F$ , we obtain the matrix  $\tilde{M} = C_{\tilde{x}} D_\rho C_{\tilde{x}}^T$ , where  $\tilde{x}$  is the vector with the normalized Fourier transform

$$(F\tilde{x})[k] = \frac{(Fx)[k]}{|(Fx)[k]|}. \quad (\text{III.7})$$

Therefore, the matrix  $C_{\tilde{x}}$  is both circulant and real orthonormal, i.e.,  $C_{\tilde{x}}^{-1} = C_{\tilde{x}}^T$ . Consequently, the decomposition  $\tilde{M} = C_{\tilde{x}} D_\rho C_{\tilde{x}}^T$  is an eigendecomposition of  $\tilde{M}$ , and the eigenvectors are translations of  $\tilde{x}$ .

If  $\rho$  has at least one distinct entry, then the associated eigenvector  $v$  will be a translation of  $\tilde{x}$ , with arbitrary scaling; that is,  $v = \alpha \cdot R_r \tilde{x}$  for some number  $\alpha$ . Since the Fourier coefficients are still normalized, we multiply  $|P_x|^{1/2}$  and  $Fv$  coordinate-wise to get

$$\tilde{v} = \alpha \cdot F^{-1} \left( F(R_r \tilde{x}) \odot |P_x|^{1/2} \right) = \alpha \cdot R_r x.$$

Letting  $\text{Sum}(x)$  denote the sum of all elements in  $x$ , we have  $\alpha = \text{Sum}(\tilde{v})/\text{Sum}(x)$ . To uncover  $\alpha$ , note that the zeroth

Fourier coefficient of  $\mu = x * \rho$  is  $(F\mu)[0] = (Fx)[0] \cdot (F\rho)[0]$ . But since  $\rho$  is a probability vector,  $(F\rho)[0] = 1$ , and so  $\text{Sum}(\mu) = (F\mu)[0] = (Fx)[0] = \text{Sum}(x)$ . Consequently,  $\alpha = \text{Sum}(\tilde{v})/\text{Sum}(\mu)$ , and  $R_r x = \tilde{v}/\alpha$ .

Note that once we have determined  $x$ , we can also determine  $\rho$  from  $\mu = x * \rho$  by deconvolution; indeed, since  $\mu = C_x \rho$ , we have  $\rho = C_x^{-1} \mu$ .

---

### Algorithm 1 Exact recovery from the first two moments

---

Input: Moments  $\mu$  and  $M$ ; power spectrum  $P_x$ .

Output: The signal  $x$  and distribution  $\rho$ .

```

// Normalize  $Fx$ 
1.1:  $p \leftarrow (P_x)^{-1/2}$ 
1.2:  $\tilde{Q} \leftarrow F^{-1} D_p F$ 
1.3:  $\tilde{M} \leftarrow Q M Q^*$ 
// Extract eigenvector and rescale
2.1:  $v \leftarrow \text{UniqEig}(\tilde{M})$ 
2.2:  $\tilde{v} \leftarrow F^{-1} \left( (P_x)^{1/2} \odot Fv \right)$ 
2.3:  $x \leftarrow (\text{Sum}(\mu) / \text{Sum}(\tilde{v})) \tilde{v}$ 
2.4:  $\rho \leftarrow C_x^{-1} \mu$ 
2.5: return  $x$  and  $\rho$ 

```

---

We have proved the following result:

**Proposition III.2.** *Suppose  $x$  has non-vanishing DFT and  $\rho$  has at least one distinct entry. Let  $\mu = \mathbb{E}[R_r x]$  and  $M = \mathbb{E}[(R_r x)(R_r x)^T]$  be the first two moments, and  $P_x$  the power spectrum of  $x$ . Then Algorithm 1 returns the signal  $x$  and the distribution  $\rho$  exactly (up to translation).*

The following corollary states that Algorithm 1 is stable to perturbations of the moments and power spectrum:

**Corollary III.3.** *Suppose  $x$  has non-vanishing DFT and  $\rho$  has at least one distinct entry. Suppose that  $\hat{\mu}$ ,  $\hat{M}$ , and  $\hat{P}_x$  are estimates of, respectively,  $\mu$ ,  $M$  and  $P_x$ . Suppose that  $\|\hat{\mu} - \mu\| \leq \varepsilon$ ,  $\|\hat{M} - M\|_F \leq \varepsilon$ , and  $\|\hat{P}_x - P_x\| \leq \varepsilon$ , for sufficiently small  $\varepsilon > 0$ . Then running Algorithm 1 with input data  $\hat{\mu}$ ,  $\hat{M}$ , and  $\hat{P}_x$  returns an estimate  $\hat{x}$  of  $x$  with an error  $C\varepsilon$ , where  $C = C(x, \rho, L) > 0$  is a finite constant.*

*Proof.* This follows immediately from the variant of the Davis-Kahan theorem found in [38] (Theorem 2).  $\square$

### C. Estimating $x$ from finitely many samples in low SNR

Section III-B shows that Algorithm 1 recovers  $x$  exactly from the exact values of  $\mu$ ,  $M$  and  $P_x$ , so long as the DFT of  $x$  is non-vanishing and  $\rho$  has at least one distinct entry. We also showed that Algorithm 1 is stable under small perturbations of the moments and power spectrum, under the same conditions. In this section, we estimate the error when Algorithm 1 is applied to the unbiased estimators  $\hat{\mu}$ ,  $\hat{M}$  and  $\hat{P}_x$  from equations (III.3), (III.5) and (III.6), respectively.

We first observe that whenever  $\rho$  is aperiodic, we can modify the observations to assume that  $\rho$  in fact has all distinct entries. Indeed, we generate a new set of measurements  $z_j = R_{r_j} y_j$ , where  $r_j$  are drawn from a new, known distribution  $\theta$ . In this case, the translations are distributed according to  $\rho * \theta$ . The following lemma shows that by choosing  $\theta$  as a random

probability distribution on the simplex, we can ensure that all entries of  $\rho * \theta$  are distinct with probability 1. Note that if the DFT of  $\theta$  is non-vanishing (which holds with probability 1 for random  $\theta$ ), then one can recover fully  $\rho$  from  $\rho * \theta$ .

**Lemma III.4.** *Let  $\rho$  be an aperiodic vector on the simplex and let  $\theta$  be a random probability density function on the simplex. Then, all entries of  $\rho * \theta$  are distinct with probability 1.*

*Proof.* See Appendix B.  $\square$

Using this lemma, we will assume from now on that all entries of  $\rho$  are distinct.

We study the low-SNR regime when  $\sigma \rightarrow \infty$ . Concretely, we suppose we have a sequence of noise levels  $\sigma_1 \leq \sigma_2 \leq \dots$ , such that  $\sigma_n \rightarrow \infty$ . We also suppose we have a sequence of sample sizes  $N_1 \leq N_2 \leq \dots$ . For each  $n$ , we draw observations  $y_1, \dots, y_{N_n}$  at noise level  $\sigma_n$ .

The following theorem shows that if  $N_n$  grows like  $\sigma_n^4$ , the MSE of the estimator  $\hat{x}$  can be controlled.

**Theorem III.5.** *Let  $\hat{\mu}_n$ ,  $\hat{M}_n$  and  $\hat{P}_{x,n}$  denote the sample moments and power spectrum from  $y_1, \dots, y_{N_n}$ , as defined by equations (III.3), (III.5) and (III.6). Let  $\hat{x}_n$  denote the estimate of  $x$  from Algorithm 1 applied to  $\hat{\mu}_n, \hat{M}_n$  and  $\hat{P}_{x,n}$ . Then for all sufficiently small  $t > 0$ :*

$$\mathbb{P} \left[ \min_{s \in \mathbb{Z}_L} \|R_s \hat{x}_n - x\|^2 \geq t \right] \leq C_1 \exp \left\{ -C_2 \frac{N_n}{\sigma_n^4} t \right\}, \quad (\text{III.8})$$

where  $C_1 = C_1(x, \rho, L)$  and  $C_2 = C_2(x, \rho, L)$  are finite, positive constants.

*Proof.* This follows from the fact that the residuals  $\hat{\mu}_n - \mu$ ,  $\hat{M}_n - M$  and  $\hat{P}_{x,n} - P_x$  are subexponential; the Bernstein-type inequality for subexponential random variables found in [39]; and Corollary III.3.  $\square$

From Theorem III.5, we see that if  $N_n \geq c_n \sigma_n^4$  for some sequence  $c_n \rightarrow \infty$ , then the error of  $\hat{x}_n$  converges to 0 in probability as  $n \rightarrow \infty$ . Furthermore, if  $N_n \geq K \log(n) \sigma_n^4$  for a sufficiently large constant  $K$ , then the error of  $\hat{x}_n$  converges to 0 almost surely as  $n \rightarrow \infty$ .

Algorithm 2 describes the entire pipeline for estimating  $x$  from the noisy measurements  $y_j = R_{r_j} x + \varepsilon_j$ , including randomly shifting the observations, estimating the moments, and using Algorithm 1 to estimate  $x$  from the estimated moments.

#### D. Non-uniqueness for periodic $\rho$

We have shown that the first and the second moments suffice to determine the signal if the distribution is aperiodic. In this section we provide a complementary result, showing that if the distribution is periodic, then having the first two moments is not enough to uniquely determine the signal. In particular, given a distribution  $\rho$  with period  $\ell$ , a signal  $x_2$  (with non-vanishing DFT) has the same first two moments as  $x_1$  if it satisfies:

$$(Fx_2)[k] = \begin{cases} (Fx_1)[k], & k = t\frac{L}{\ell}, \quad t = 0, \dots, \ell - 1, \\ -(Fx_1)[k], & \text{otherwise.} \end{cases} \quad (\text{III.9})$$

---

#### Algorithm 2 Estimating $x$ and $\rho$ from noisy data

---

Input:  $y_j, j = 1, \dots, N$  of (I.2) and noise variance  $\sigma^2$

Output: An estimated signal  $\hat{x}$  and estimated distribution  $\hat{\rho}$

// **Reshuffling observations (optional)**

1.1: draw a random distribution  $\theta \in \Delta^L$

1.2: for each  $j = 1, \dots, N$ :  $y_j \leftarrow R_{r_j} y_j$  for  $r_j \sim \theta$

// **Moment estimation**

2.1:  $\hat{\mu} \leftarrow \frac{1}{N} \sum_{j=1}^N y_j$

2.2:  $\hat{M} \leftarrow \frac{1}{N} \sum_{j=1}^N y_j y_j^T - \sigma^2 I$

2.3:  $\hat{P}_x \leftarrow \frac{1}{N} \sum_{j=1}^N |F y_j|^2 - \sigma^2 L I$

// **Eigendecomposition and normalization**

3.1: call Algorithm 1 with  $\hat{\mu}$ ,  $\hat{M}$  and  $\hat{P}_x$ .

3.2: return  $\hat{x}$  and  $\hat{\rho}$

---

This construction is demonstrated in Figure III.1.

**Proposition III.6.** *Let  $\ell < L/2$  be a divisor of  $L > 1$ . Suppose that  $\rho$  is periodic, with period  $\ell$ . Then, for a given real signal  $x_1$  with non-vanishing DFT, there exists a different real signal  $x_2$  (which is not a translation of  $x_1$ ) such that both signals have the same first and second moments. Therefore, if the distribution is periodic, then any signal with non-vanishing DFT is not uniquely determined from its first and second moments.*

*Proof.* See Appendix C.  $\square$

In Section V we establish this result from an information-theoretic perspective by showing that the sample complexity for periodic distribution grows like  $1/\text{SNR}^3$ .

The uniform distribution is merely a special case of periodic distributions with minimal period  $\ell = 1$ . When  $\ell > 1$ , one can interpret the periodicity as having a uniform distribution over the different cosets of  $\mathbb{Z}_L$  with respect to the subgroup generated by a translation in  $\ell$  coordinates. These cosets are exactly the analogue of the sparsity pattern of  $F\rho$  attained by jumps of  $L/\ell$ . This also explains why uniformity is the only pathological case for a prime  $L$ . Therefore, if one can choose how to sample the signal, a prime number of samples should be considered.

As it turns out, there is one special case where the first two moments are enough to determine  $x$  uniquely up to cyclic translation, even when  $\rho$  is periodic. This special case occurs when  $L$  is even and  $\rho$  is  $L/2$ -periodic. This result is formulated in the following claim:

**Claim III.7.** *Suppose that  $x$  has non-vanishing DFT,  $L$  is even and  $\rho$  is  $L/2$ -periodic. Then,  $x$  is uniquely determined from its first two moments, up to global translation.*

*Proof.* See Appendix D.  $\square$

## IV. CONNECTION WITH THE SPIKED COVARIANCE MODEL

In this section, we point out a connection between the spectral algorithm presented in Section III and the spiked covariance model well-known in statistics [18], [19], [20], [21], [22]. Though somewhat informal, this analysis will provide insight into how the complexity of recovering  $x$  depends on

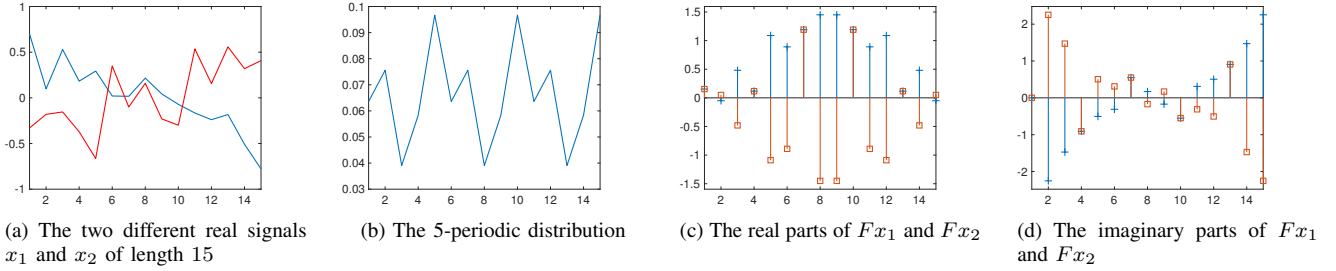


Fig. III.1. This example demonstrates the constrict of (III.9) and Proposition III.6. The figures present two different real signals of length 15 and a 5-periodic distribution. The Fourier transforms of the signals obey (III.9). The two signals have the same first two moments under the periodic distribution.

the dimension  $L$  when the distribution  $\rho$  has a fixed support size.

In the spiked model, we observe a matrix

$$Y = X + G \in \mathbb{R}^{L \times N} \quad (\text{IV.1})$$

where  $X$  is a rank  $r$  matrix and  $G = (g_{ij})$  with  $g_{ij} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ . This model is typically studied in the *high-dimensional regime*, in which  $L$  grows proportionally to  $N$ ; that is,  $L = L(N)$  and  $L/N \rightarrow \gamma > 0$  as  $N \rightarrow \infty$ . In this setting, there is a precise understanding of the limiting behavior of the data matrix  $Y$  and the low-rank matrix  $X = [X_1, \dots, X_N]$ .

In [21] (see also [19]), it is shown that when the low-rank matrix  $X$  is random (for instance, its columns may be drawn from a suitable low-rank, mean-zero distribution), then the limiting cosine  $c$  of the angles between the top eigenvector of  $XX^T$  and the top eigenvector of  $YY^T$  is given by the formula:

$$c^2 = \begin{cases} \frac{1 - \sigma^4 \gamma / \lambda^2}{1 + \sigma^2 \gamma / \lambda} & \text{if } \lambda > \sigma^2 \sqrt{\gamma}, \\ 0 & \text{otherwise,} \end{cases} \quad (\text{IV.2})$$

where  $\lambda$  is the top eigenvalue of  $XX^T/N$ .

The key phenomenon is the phase transition at

$$\lambda_{critical} = \sigma^2 \sqrt{\gamma}. \quad (\text{IV.3})$$

It is only when  $\lambda$  is greater than this critical value that we are guaranteed a non-trivial correlation between the top eigenvector of the observed matrix  $YY^T/N$  and the top eigenvector of  $XX^T/N$ .

We can view the observation model in the one-dimensional MRA model (I.1) as a special instance of the spiked model, by taking the  $i$ th column of  $X$  to be  $X_i = R_{r_i} x$ . As  $N \rightarrow \infty$ , we can write

$$\frac{1}{N} XX^T = C_x D_\rho C_x^T. \quad (\text{IV.4})$$

Consequently, under the assumption that the DFT of  $x$  does not vanish, the rank of  $X$  is the size of the support of  $\rho$ . When the support size of  $\rho$  is fixed at  $r$ , the MRA problem is an instance of the spiked model.

Let us assume that the  $|(Fx)[k]| = 1$  for all  $k$ . This can be done by estimating the power spectrum first and then normalizing all Fourier coefficients. In this case,  $C_x$  is an orthogonal matrix. In other words,  $x \perp R_\ell x$  for every  $\ell \neq 0$ ; consequently, the  $R_\ell x$  are precisely the top

$r$  eigenvectors of  $XX^T/N$ , with corresponding eigenvalues  $\|x\|^2 \rho[\ell]$ . Then, (IV.2) tells us exactly how well we expect the spectral algorithm to perform in recovering  $x$ ; indeed, the theory predicts a non-zero angle between  $x$  and the top eigenvector of  $YY^T/N$  whenever:

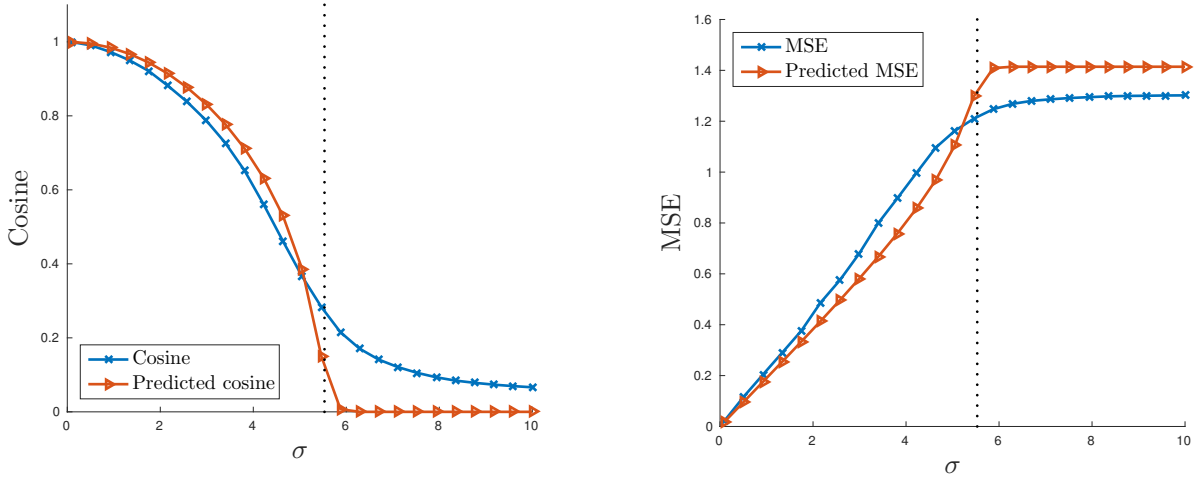
$$N \geq \frac{L\sigma^4}{\|x\|^4 (\max \rho)^2} = \frac{L}{(\max \rho)^2 \text{SNR}^2}. \quad (\text{IV.5})$$

Below this threshold, the output will be essentially random. We see that if the distribution is well-localized, then  $\max \rho = \Omega(1)$  (with respect to the growing value of  $L$ ) and then the sample complexity grows like  $\frac{L}{\text{SNR}^2}$ . On the other hand, if the distribution is almost uniform, then  $\max \rho = \mathcal{O}(1/L)$  as  $L \rightarrow \infty$ , and thus the sample complexity will be  $\frac{L^3}{\text{SNR}^2}$ .

To illustrate the relationship between the spiked model and MRA, we ran the following experiment. We generated a signal  $x \in \mathbb{R}^{400}$  with i.i.d. normal entries and normalized it so that  $\|x\|_2 = 10$ . For noise levels  $\sigma$  between 0.1 and 10, we drew  $N$  samples of  $x$  with noise at level  $\sigma$ , where  $N$  is chosen at 100 plus the critical threshold given by (IV.5) for  $\sigma = \lambda^{1/2} \gamma^{-1/4} = 5.5313$  according to (IV.3). For  $\sigma$  large enough,  $N$  will not be large enough for the spectral method to produce an estimate better than random. The distribution of translations  $\rho$  was taken to be  $\rho[i] \propto i^2$ , for  $i = 1, \dots, 5$ , and zero elsewhere. Each experiment was repeated 200 times. The plots in Figure IV.1 display the average values over these 200 runs.

For each draw, we compute the top eigenvalue of the clean data matrix (IV.4), denoted by  $\lambda$ , and the associated eigenvector, which is a translated copy of  $x$ . We also compute the top eigenvector of the data matrix  $YY^T/N$ . The angle between the two eigenvectors is predicted by (IV.2). In Figure IV.1(a), we plot the predicted cosine against the true cosine. Clearly, we never attain the predicted value of zero in finite samples, but we see a precipitous decline when the noise level  $\sigma$  exceeds its threshold value (the vertical dashed line).

We also measure the relative mean squared error defined by equation (I.3), where  $\hat{x}$  is the top eigenvector multiplied by  $\|x\|$ . In Figure IV.1(b), we plot this error as a function of  $\sigma$ . For reference, we also plot the ordinary error predicted by the spiked model (as derived from the predicted cosine between the vectors), without minimizing over shifts. Of course, minimizing over shifts will decrease the error; however, we still see the same qualitative behavior predicted from the spiked model, namely an increase in error as  $\sigma$  grows, until the critical threshold of  $\sigma$  is reached, after which the error plateaus.



(a) Empirical cosines between the top eigenvectors of the matrices  $\frac{1}{N}XX^T$  and  $\frac{1}{N}YY^T$  as a function of the noise level compared to asymptotic cosines predicted by spiked model; see (IV.2).

(b) Empirical MSE compared to the asymptotic MSE predicted by spiked model. The MSE is defined in (I.3).

Fig. IV.1. Experiments related to the connection between the spike model and the MRA problem as discussed in Section IV. The dashed line is the predicted threshold value of  $\sigma = \lambda^{1/2}\gamma^{-1/4} = 5.5313$ .

## V. INFORMATION THEORETIC LIMIT

In this section, we provide lower bounds for the mean square error (MSE) of an estimator of the signal in terms of the SNR and the number of observations  $N$ . In particular, we show that the MSE is bounded below by order of  $1/\text{SNR}^2$  under mild conditions on the signal. As described in Section III, Algorithm 2 achieves this sample complexity. In addition, if the distribution is periodic then the MSE is lower bounded by  $1/\text{SNR}^3$ . The framework proposed in [15] and described in Section II achieves this sample complexity for any distribution.

Recall that we can estimate the signal only up to cyclic translation. We define the best alignment of  $\hat{x}$  with  $x$  by

$$\phi_x(\hat{x}) = \underset{z \in \{R_s \hat{x}\}_{s \in \mathcal{Z}_L}}{\text{argmin}} \|z - x\|. \quad (\text{V.1})$$

Accordingly, we define the relative MSE to be

$$\text{relative MSE} = \frac{1}{\|x\|^2} \mathbb{E} \left[ \|\phi_x(\hat{x}) - x\|^2 \right]. \quad (\text{V.2})$$

The expectation is taken over  $\hat{x}$ , which is a function of the observations, which are themselves random. In the bounds we present, we assume that the estimator is asymptotically unbiased, i.e.,  $\mathbb{E}[\phi_x(\hat{x})] \rightarrow x$  as  $N \rightarrow \infty$ . While this assumption is unnecessary to obtain meaningful results, it simplifies the analysis by removing the dependence of the estimator on the expectation, as shown by Theorem V.4.

We now present the main results of this section as follows:

**Theorem V.1.** *Assume that  $x$  is not a constant vector. If  $\hat{x}$  is an asymptotically unbiased estimator of  $x$ , then*

$$\text{MSE} \geq \frac{1}{8N} \frac{1}{\text{SNR}^2} - \mathcal{O} \left( \frac{1}{N \text{SNR}^{1.5}} \right). \quad (\text{V.3})$$

Moreover, if  $\rho$  is periodic, with a period  $\ell < \frac{L}{2}$ , then

$$\text{MSE} \geq \frac{1}{54N} \frac{L-2\ell}{2\ell} \frac{1}{\text{SNR}^3} - \mathcal{O} \left( \frac{1}{N \text{SNR}^{2.5}} \right). \quad (\text{V.4})$$

Note that previous work [16] derived the sample complexity for the uniform distribution of translations. Theorem V.1 extends it to any distribution. In addition, we extend [16] by estimating the constant that multiplies  $\sigma^6$ .

In the rest of this section, we develop the main tools required to prove Theorem V.1. Specifically, we start by introducing an auxiliary notation and definitions. Then, in Section V-B we use an adaptation of the Chapman-Robbins lower bound [40], which is a generalization of the Cramér-Rao bound [41], to derive a lower bound on the MSE in terms of the  $\chi^2$  divergence. Finally, in Section V-C, we express the  $\chi^2$  divergence in terms of the Taylor expansion of the posterior probability density and generalized auto-correlations. This characterization connects the generalized notion of auto-correlations of the signal (which also depends on the distribution) with the sample complexity (as a function of the SNR). This point is analogous to the Boolean case, treated in [42]. The main proof of Theorem V.1 is given in Appendix H.

### A. Notations and Definitions

Similarly to Section IV, let  $Y \in \mathbb{R}^{L \times N}$  be the collection of all measurements as columns in a matrix. We also denote by  $f_{x,\rho}$  the probability density of the posterior distribution of  $Y$ ,

$$f_{x,\rho}(Y) = \prod_{i=1}^N f_{x,\rho}(y_i), \quad (\text{V.5})$$

and the expectation of a function  $g$  of the measurements under the measure  $f_{x,\rho}$  by

$$\mathbb{E}_{x,\rho} [g(Y)] := \int_{\mathbb{R}^{L \times N}} g(Y) f_{x,\rho}(Y) dY.$$

For ease of notation, hereinafter we write  $\mathbb{E} [g(Y)]$  when the signal and distribution are implicit. We also recall the following bias-variance trade-off of the MSE:

$$\text{MSE} = \frac{\text{tr}(\text{Cov}[\phi_x(\hat{x})])}{\|x\|^2} + \frac{\|\mathbb{E}[\phi_x(\hat{x})] - x\|^2}{\|x\|^2}. \quad (\text{V.6})$$

with

$$\text{Cov}[\phi_x(\hat{x})] = \mathbb{E} \left[ \phi_x(\hat{x})\phi_x(\hat{x})^T \right] - \mathbb{E}[\phi_x(\hat{x})]\mathbb{E}[\phi_x(\hat{x})]^T. \quad (\text{V.7})$$

We conclude this part with two definitions. First, we define the generalized notion of auto-correlations of a signal  $x$ .

**Definition V.2.** The  $d$ -autocorrelation of the pair  $(x, \rho)$  is a tensor of order  $d$ , defined by

$$A_{x,\rho}^d := \mathbb{E}_r \left[ (R_r x)^{\otimes d} \right],$$

where  $r \sim \rho$ .

This generalized notion of auto-correlation also appears in previous works [16], [42]. This notion is a generalization of the first two moments defined in (III.1) and (III.4), respectively, according to

$$A_{x,\rho}^1 = \mu \quad \text{and} \quad A_{x,\rho}^2 = M - \sigma^2 I.$$

Our last definition is of the  $\chi^2$  divergence between two distributions which, in a sense, measures the difference between two probability distributions.

**Definition V.3.** The  $\chi^2$  divergence between the distributions  $f_{\tilde{x},\tilde{\rho}}$  and  $f_{x,\rho}$  is defined by

$$\chi_N^2(f_{\tilde{x},\tilde{\rho}}||f_{x,\rho}) := \mathbb{E}_{x,\rho} \left[ \left( \frac{f_{\tilde{x},\tilde{\rho}}(Y)}{f_{x,\rho}(Y)} - 1 \right)^2 \right].$$

Due to equation (V.5), the relation between the  $\chi^2$  divergence for  $N$  observations and only one observation is given by

$$\chi_N^2(f_{\tilde{x},\tilde{\rho}}||f_{x,\rho}) = (\chi_1^2(f_{\tilde{x},\tilde{\rho}}||f_{x,\rho}) + 1)^N - 1.$$

To ease notation, hereinafter we use  $\chi^2$  for  $\chi_1^2$ .

### B. Chapman-Robbins lower bound for an orbit

The classical Chapman-Robbins gives a lower bound on an error metric of the form  $\mathbb{E}[\|\hat{x} - x\|^2]$ , i.e., it does not take into consideration a translation-invariant error metric as appears naturally in the MRA problem. Hence, we modified the Chapman-Robbins bound to accommodate error of the form (V.2). We point out that  $\text{Cov}[\phi_x(\hat{x})]$  is related to the MSE by (V.6).

**Theorem V.4** (Chapman-Robbins). *For any  $\tilde{x} \in \mathbb{R}^L$  and  $\tilde{\rho} \in \Delta^L$ , we have*

$$\text{Cov}[\phi_x(\hat{x})] \succeq \frac{zz^T}{\chi_N^2(f_{\tilde{x},\tilde{\rho}}||f_{x,\rho})},$$

where  $z = \mathbb{E}_{\tilde{x},\tilde{\rho}}[\phi_x(\hat{x})] - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})]$ .

*Proof.* See Appendix E.  $\square$

### C. Fisher information and auto-correlations

Instead of considering the posterior probability density of  $Y$ , we will consider  $\tilde{Y} = Y/\sigma$ . This change of variable does not change the  $\chi^2$  divergence. We then have

$$\tilde{Y}_j = \gamma R_{r_j} x + G_j, \quad (\text{V.8})$$

where  $\gamma = 1/\sigma$ ,  $r_j \sim \rho$  and  $G_j \sim \mathcal{N}(0, I)$ . We can now take the Taylor expansion of the probability density around  $\gamma = 0$ , that is,

$$f_{x,\rho}(y; \gamma) = f_G(y) \sum_{j=0}^{\infty} \alpha_{x,\rho}^j(y) \frac{\gamma^j}{j!}, \quad (\text{V.9})$$

where  $f_G(y) = f_{x,\rho}(y; 0)$  is the probability density of  $G_j$  (since when  $\gamma = 0$ ,  $Y_j = G_j$ ) and  $\alpha_{x,\rho}^0(y) = 1$  since

$$\alpha_{x,\rho}^j(y) := \frac{1}{f_G(y)} \frac{\partial^j f_{x,\rho}(y; 0)}{\partial \gamma^j}. \quad (\text{V.10})$$

We now use the Taylor expansion (V.9) to give an expression of the  $\chi^2$  divergence.

**Lemma V.5.** *The divergence  $\chi^2(f_{\tilde{x},\tilde{\rho}}||f_{x,\rho})$  can be expressed in terms of auto-correlations as:*

$$\begin{aligned} \chi^2(f_{\tilde{x},\tilde{\rho}}||f_{x,\rho}) &= \frac{\sigma^{-2d}}{(d!)^2} \mathbb{E}_G \left[ \left( \alpha_{\tilde{x},\tilde{\rho}}^d(G) - \alpha_{x,\rho}^d(G) \right)^2 \right] + \mathcal{O}(\sigma^{-2d-1}), \\ &= \frac{\sigma^{-2d}}{d!} \|A_{\tilde{x},\tilde{\rho}}^d - A_{x,\rho}^d\|^2 + \mathcal{O}(\sigma^{-2d-1}), \end{aligned} \quad (\text{V.11})$$

$$= \frac{\sigma^{-2d}}{d!} \|A_{\tilde{x},\tilde{\rho}}^d - A_{x,\rho}^d\|^2 + \mathcal{O}(\sigma^{-2d-1}), \quad (\text{V.12})$$

where  $d = \inf \left\{ n : \|A_{\tilde{x},\tilde{\rho}}^n - A_{x,\rho}^n\|^2 > 0 \right\}$ .

*Proof.* See Appendix F  $\square$

Equation (V.11) is not specific to MRA: one can always obtain this expression as long we are considering the low SNR regime and the observations are independent of the signal in the limit of SNR tending to 0. The particularization to MRA happens in (V.12), due to (V.8) and (V.10).

Using Theorem V.4 together with Lemma V.5 we can obtain lower bounds in terms of the auto-correlations. For example, to obtain (V.3), one could provide  $\tilde{x}$  and  $\tilde{\rho}$  which have  $A_{\tilde{x},\tilde{\rho}}^1 = A_{x,\rho}^1$  and bound  $\|A_{\tilde{x},\tilde{\rho}}^2 - A_{x,\rho}^2\|$  from above. Moreover, to obtain (V.4) when  $\rho$  is periodic one could provide  $\tilde{x}$  and  $\tilde{\rho}$  which have  $A_{\tilde{x},\tilde{\rho}}^d = A_{x,\rho}^d$  for  $d = 1, 2$ , similarly to Proposition III.6, and bound  $\|A_{\tilde{x},\tilde{\rho}}^3 - A_{x,\rho}^3\|$  from above.

Nevertheless, to prove Theorem V.1, we use intermediate results which explore the limit  $(\tilde{x}, \tilde{\rho}) \rightarrow (x, \rho)$ , and provide tighter bounds. An informal explanation for the tighter bounds is that it is easier to confuse  $x$  and  $\rho$  with signals and distributions that are closer, rather than farther. However, since considering the limit introduces some technical details, we analyze it in Appendix G.

## VI. ADDITIONAL ALGORITHMS

While the spectral algorithm (Algorithm 2) is asymptotically optimal as  $\sigma, N \rightarrow \infty$  and for signals with non-vanishing DFT, it may not perform well in small sample size or low DFT values. Therefore, in this section, we present two additional algorithms based on non-convex least-squares minimization and a modification of the EM algorithm presented in Section II that takes the distribution into account. In Appendix J, we also describe and analyze a convex relaxation approach based on semidefinite programming.



### A. Non-convex least-squares minimization

The following method aims to find a signal in  $\mathbb{R}^L$  and a distribution in  $\Delta^L$  that fit the observed data as well as possible in the LS sense. We formulate the problem as a smooth, non-convex, optimization problem with the constraint that the distribution lies on a simplex. Given estimators  $\hat{M}$  and  $\hat{\mu}$  of the first two moments  $M = \mathbb{E}[(R_\ell x)(R_\ell x)^T]$  and  $\mu = \mathbb{E}[R_\ell x]$ , the problem reads

$$\min_{\hat{x} \in \mathbb{R}^L, \hat{\rho} \in \Delta^L} \|\hat{M} - C_{\hat{x}} D_{\hat{\rho}} C_{\hat{x}}^T\|_F^2 + \lambda \|\hat{\mu} - C_{\hat{x}} \hat{\rho}\|_2^2, \quad (\text{VI.1})$$

where  $\lambda > 0$  is a predefined parameter. It can be verified that, by omitting signal-dependent terms, the variance of the elements of the first moment estimator is proportional to  $\sigma^2$ . It can be also shown that the variance of the elements of the second moment is proportional to  $3L\sigma^4$  and  $L\sigma^2$  in the low and the high SNR regimes, respectively (again, by omitting signal-dependent terms). Therefore, we set  $\lambda = \frac{1}{L(1+3\sigma^2)}$  in our implementation.

As will be shown empirically in Section VII, given the exact first two moments of the data, the non-convex problem (VI.1) converges to the sought signal from a random initialization.

### B. An expectation-maximization algorithm for estimating $x$ and $\rho$ simultaneously

In Section II, we reviewed the EM algorithm for MRA from [15], which is invariant to the distribution of translations. In this section, we modify the algorithm to take the distribution into account.

The forward model of the MRA model (I.1) reads:

$$p(y, r|x, \rho) = \prod_{j=1}^N \frac{1}{(2\pi\sigma^2)^{L/2}} e^{-\frac{1}{2\sigma^2} \|R_{r_j} x - y_j\|^2} \rho[r_j].$$

The log-likelihood function is then given, up to a constant, by

$$\log L(x, \rho|y, r) = \sum_{j=1}^N \left\{ \log \rho[r_j] - \frac{1}{2\sigma^2} \|R_{r_j} x - y_j\|^2 \right\}.$$

The goal of the EM algorithm is to compute the maximum marginal likelihood  $L(x, \rho|y) = \sum_r L(x, \rho|y, r)$ . The algorithm proceeds as follows. Start with some initial guesses  $x_0$  and  $\rho_0$  for the signal and distribution. Given  $x_k$  and  $\rho_k$ , the next guess is given as follows:

$$(x_{k+1}, \rho_{k+1}) = \arg \max_{x, \rho} Q(x, \rho|x_k, \rho_k),$$

where

$$Q(x, \rho|x_k, \rho_k) := \mathbb{E}_{r|y, x_k, \rho_k} [\log L(x, \rho|y, r)]$$

Therefore, (by omitting the constant term)

$$\begin{aligned} Q(x, \rho|x_k, \rho_k) &= \sum_{j=1}^N \mathbb{E}_{r|y, x_k, \rho_k} \left[ \log \rho_k[r_j] - \frac{1}{2\sigma^2} \|R_{r_j} x_k - y_j\|^2 \right] \\ &= \sum_{j=1}^N \sum_{\ell=0}^{L-1} w_k^{\ell, j} \left\{ \log \rho_k[\ell] - \frac{1}{2\sigma^2} \|R_\ell x_k - y_j\|^2 \right\}, \end{aligned}$$

where the values  $w_k^{\ell, j}$  are defined by the formula:

$$w_k^{\ell, j} = \mathbb{P}[r_j = \ell|y, x_k, \rho_k] = C_k^j e^{-\frac{1}{2\sigma^2} \|R_r x_k - y_j\|^2} \rho_k[\ell],$$

where  $C_k^j$  is a normalization factor so that  $\sum_{\ell} w_k^{\ell, j} = 1$ .

To maximize  $Q$  over  $x$  and  $\rho$  is simple, since the first term depends only on  $\rho$  and the second term depends only on  $x$ . Specifically, it is easy to see that the maximum over  $x$  is given by a weighted average of the translated observations:

$$x_{k+1} = \frac{1}{N} \sum_{j=1}^N \sum_{\ell=0}^{L-1} w_k^{\ell, j} R_r^{-1} y_j. \quad (\text{VI.2})$$

This step is almost identical (up to the values of the weights) to the standard EM update step (II.2).

The maximizing value of  $\rho$  also has a closed formula. First, observe we can write:

$$\rho_{k+1} = \arg \max_{\rho \in \Delta^L} \sum_{\ell=0}^{L-1} W_k[\ell] \log(\rho_k[\ell])$$

where  $W_k[\ell] = \sum_{j=1}^N w_k^{\ell, j}$ . To maximize a positive weighted combination of logarithms over the simplex, we use the following lemma:

**Lemma VI.1.** *If  $w[\ell] > 0$  are positive weights, then the maximizer of  $\sum_{\ell} w[\ell] \log(q[\ell])$  over all  $q \in \Delta^L$  is*

$$q^*[\ell] = w[\ell] / \sum_{\ell'} w[\ell'].$$

*Proof.* See Appendix I.  $\square$

From this lemma, the maximizing  $\rho$  is given by the formula:

$$\rho_{k+1}[\ell] = \frac{W_k[\ell]}{\sum_{\ell'=0}^{L-1} W_k[\ell']}. \quad (\text{VI.3})$$

To conclude, the modified EM updates the signal and the distribution estimations by (VI.2) and (VI.3), respectively. However, compared to the methods which are based on moments estimation like Algorithm 2 or the LS, it passes through the data at each iteration. Therefore, for large sample size, its computational cost may be substantially heavier.

## VII. NUMERICAL EXPERIMENTS

In this section, we present numerical results for the algorithms described in Section VI and Algorithm 2. To measure the accuracy of an estimator  $\hat{x}$ , we define the recovery relative error as

$$\text{relative error} = \min_{\ell \in \mathbb{Z}_L} \frac{\|R_r \hat{x} - x\|_2}{\|x\|_2}.$$

The code of this section, including Matlab implementations and examples, is available online<sup>1</sup>.

<sup>1</sup><https://github.com/nirsharon/aperiodicMRA>

### A. Influence of the number of samples

In the first example, we use a Haar-like signal of length  $L = 24$ , depicted in Figure VII.1(a). Next, we generate its noisy, translated copies according to the MRA model (I.2), with noise variance of  $\sigma = 1$ . One example of a data sample corrupted with such noise is illustrated in Figure VII.1(b).

We use the least squares (LS) method of Section VI-A to estimate the signal. This process is repeated three times for different number of samples,  $N = 10^3$ ,  $N = 10^5$ , and  $N = 10^7$ . The estimates are presented in Figure VII.1(c)–VII.1(e). As expected, the quality of the estimation improves significantly as  $N$  grows.

### B. Comparison of EM algorithms

In [15], it is shown that in most cases, an expectation-maximization (EM) method as described in Section II-B, achieves the smallest estimation error compared to the competitor algorithms. The EM algorithm described in that paper is invariant to the distribution  $\rho$ . In particular, it treats the data as if it were drawn from the uniform distribution, which requires sample complexity that grows like  $1/\text{SNR}^3$  rather than  $1/\text{SNR}^2$ . By contrast, the EM algorithm we propose in Section VI-B estimates the distribution  $\rho$ , in addition to the signal  $x$ .

To demonstrate the importance of including the distribution into the model of the estimator, we consider a family of distributions

$$\rho[t] \propto \exp(-t^2/s^2), \quad s > 0. \quad (\text{VII.1})$$

The parameter  $s$  controls the concentration of  $\rho$ , or alternatively its uniformity; the larger  $s$  is, the more uniform  $\rho$  is. In general, we expect our algorithms to provide better estimations when  $\rho$  is more concentrated; see Section IV.

We compared the standard EM with the EM algorithm described in Section VI-B (Adapted EM). The experiments were conducted as follows. We fixed a random signal of length  $L = 25$  with i.i.d. normal entries and unit norm, and a series of distributions of the form (VII.1) with the parameter  $s$  varying between 3 and 9. Then, for each distribution we generated  $N = 2,000$  samples drawn with a fixed level of noise  $\sigma = 1$ . We repeated the test independently 20 times and averaged the errors. In Figure VII.2, we plot the relative errors of the methods as a function of the uniformity parameter  $s$ . As expected, the standard EM is invariant to  $s$ . On the other hand, the adapted version of the EM exploits the varying distribution and performs better under more concentrated distributions. As the distribution becomes more uniform, the two methods exhibit similar error rates.

### C. Comparison of the different methods

This paper presents three main approaches for solving the MRA: the spectral method described in Algorithm 2, the least squares (LS) optimization of Section VI-A, and the adapted expectation-maximization (adapted EM) of Section VI-B. In this comparison, we examined the estimation error of these three methods with different noise levels. In detail, we use

a random signal of length  $L = 25$  with i.i.d. normal entries with unit norm, and a random distribution. We fix the number of samples to be  $N = 10,000$ . Then, we increase the level of noise  $\sigma$  from 0.001 to 1. In Figure VII.3 we plot the average error. As can be seen, the LS and the adapted EM are more robust to noise than the spectral method. In addition, the gap between these two methods becomes small as the SNR decreases.

### D. Numerical error rates for the LS method

The LS method uses the first two moments as its input data. Since the error in approximating the second moment matrix scales like  $O(\sigma^2/\sqrt{N})$ , we expect the error of the solution to grow at no less than the same rate. In Figure VII.4 we plot the average error over 50 experiments as a function of  $\sigma$ . When  $\sigma > 1$ , the curve in a log scale is a line with slope close to 2, which is the expected rate. However, when  $\sigma < 1$ , the curve is a line with slope close to 1; namely, the error behaves approximately like  $O(\sigma/\sqrt{N})$ , rather than  $O(\sigma^2/\sqrt{N})$ . The moderate slope for high SNR suggests that in this regime the recovery problem is easier; for example, we know that alignment is possible in high SNR, as described in Section II-A.

## VIII. DISCUSSION

In this paper, we have shown that the sample complexity for MRA with an aperiodic distribution of translations grows like  $1/\text{SNR}^2$ . This sample complexity can be achieved by a simple spectral algorithm. We also examined empirically the LS and EM algorithms. Additionally, we extended previous works by showing that the sample complexity for any periodic distribution scales as  $1/\text{SNR}^3$ .

We drew connections between the MRA problem and the spiked covariance model. This connection implies that the sample complexity is inversely proportional to the square of the maximal value of the distribution. Therefore, the more uniform the distribution is, the higher the sample complexity of the problem.

One of the motivations for considering the MRA model arises from the imaging technique called single particle cryo-electron microscopy (cryo-EM), allowing to visualize molecules at near-atomic resolution [12], [13]. In cryo-EM, noisy two-dimensional tomographic projections of the three-dimensional underlying molecule, taken at unknown viewing direction, are collected. The distribution of viewing directions in cryo-EM is typically non-uniform, as many molecules exhibit some preferred orientation [43].

The MRA model (I.1) can be thought of as a simplified model for the cryo-EM problem, where cyclic translations replace actions of elements of the group  $SO(3)$ . The tomographic projection does not appear in (I.1). Our technique for MRA, based on the low-order moments of the data, is similar to the framework proposed by Zvi Kam in [44] for cryo-EM. In particular, Kam suggested a method to estimate a molecule directly from the statistics of the projections, rather than estimating the viewing directions. Our work is one step towards understanding the sample complexity of Kam's method in particular, and the cryo-EM problem in general.

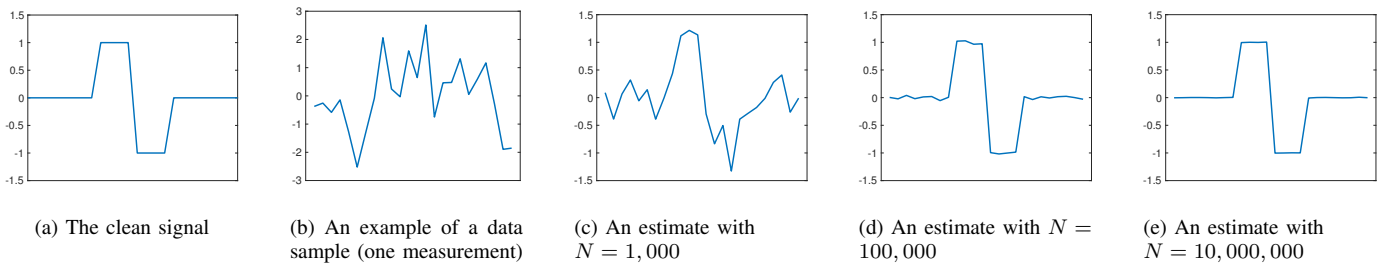


Fig. VII.1. An example of the estimation quality of a Haar-like signal with different number of samples ( $N$ ), using the LS method. In these tests,  $\sigma = 1$ .

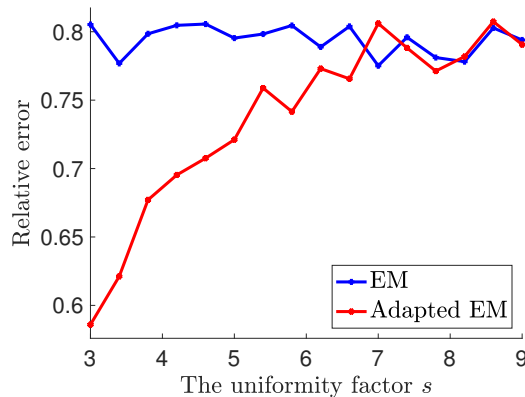


Fig. VII.2. EM comparison: the standard EM (EM) described in Section II versus the modified EM that includes distribution estimation (adapted EM) described in Section VI-B. The algorithms were compared with different distributions of the form VII.1 as a function of the parameter  $s$ .

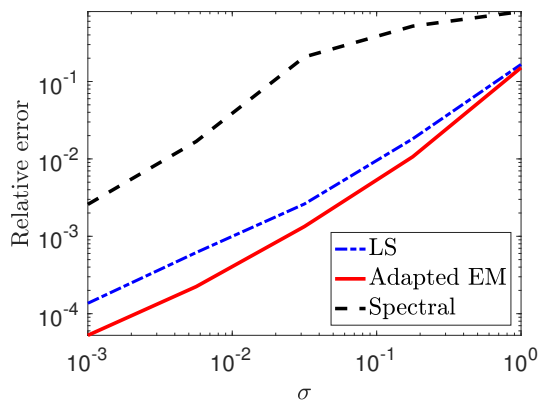


Fig. VII.3. A comparison of three methods: least squares (LS), adapted EM, and spectral method, under varying level of noise.

#### ACKNOWLEDGMENTS

We would like to thank Afonso Bandeira, Nicolas Boumal, Joseph Kileel, Roy Lederman and Zhizhen Zhao for many insightful discussions.

#### REFERENCES

- [1] R. Diamond, "On the multiple simultaneous superposition of molecular structures by rigid body transformations," *Protein Science*, vol. 1, no. 10, pp. 1279–1287, 1992.
- [2] D. L. Theobald and P. A. Steindel, "Optimal simultaneous superpositioning of multiple structures with missing data," *Bioinformatics*, vol. 28, no. 15, pp. 1972–1979, 2012.

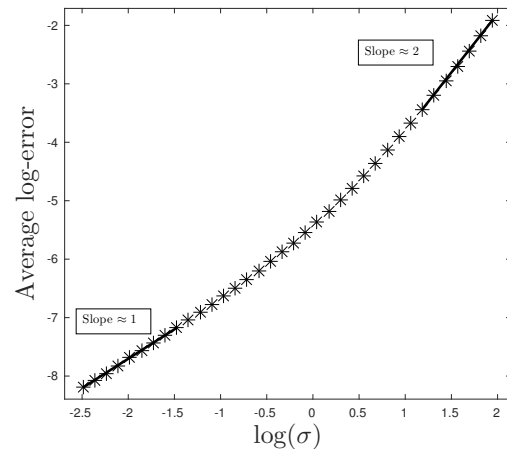


Fig. VII.4. Log-log plot of the error of the LS method versus  $\sigma$

- [3] W. Park, C. R. Midgett, D. R. Madden, and G. S. Chirikjian, "A stochastic kinematic model of class averaging in single-particle electron microscopy," *The International journal of robotics research*, vol. 30, no. 6, pp. 730–754, 2011.
- [4] W. Park and G. S. Chirikjian, "An assembly automation approach to alignment of noncircular projections in electron microscopy," *IEEE Transactions on Automation Science and Engineering*, vol. 11, no. 3, pp. 668–679, 2014.
- [5] S. H. Scheres, M. Valle, R. Nuñez, C. O. Sorzano, R. Marabini, G. T. Herman, and J.-M. Carazo, "Maximum-likelihood multi-reference refinement for electron microscopy images," *Journal of molecular biology*, vol. 348, no. 1, pp. 139–149, 2005.
- [6] J. P. Zwart, R. van der Heiden, S. Gelsema, and F. Groen, "Fast translation invariant classification of HRR range profiles in a zero phase representation," *IEE Proceedings-Radar, Sonar and Navigation*, vol. 150, no. 6, pp. 411–418, 2003.
- [7] R. Gil-Pita, M. Rosa-Zurera, P. Jarabo-Amores, and F. López-Ferreras, "Using multilayer perceptrons to align high range resolution radar signals," in *International Conference on Artificial Neural Networks*, pp. 911–916, Springer, 2005.
- [8] D. M. Rosen, L. Carlone, A. S. Bandeira, and J. J. Leonard, "A certifiably correct algorithm for synchronization over the special euclidean group," *arXiv preprint arXiv:1611.00128*, 2016.
- [9] I. L. Dryden and K. V. Mardia, *Statistical shape analysis*, vol. 4. J. Wiley Chichester, 1998.
- [10] H. Foroosh, J. B. Zerubia, and M. Berthod, "Extension of phase correlation to subpixel registration," *IEEE transactions on image processing*, vol. 11, no. 3, pp. 188–200, 2002.
- [11] D. Robinson, S. Farsiu, and P. Milanfar, "Optimal registration of aliased images using variable projection with applications to super-resolution," *The Computer Journal*, vol. 52, no. 1, pp. 31–42, 2009.
- [12] A. Bartesaghi, A. Merk, S. Banerjee, D. Matthies, X. Wu, J. L. Milne, and S. Subramaniam, "2.2 Å resolution cryo-EM structure of  $\beta$ -galactosidase in complex with a cell-permeant inhibitor," *Science*, vol. 348, no. 6239, pp. 1147–1151, 2015.
- [13] D. Sirohi, Z. Chen, L. Sun, T. Klose, T. C. Pierson, M. G. Rossmann, and R. J. Kuhn, "The 3.8 Å resolution cryo-EM structure of Zika virus," *Science*, vol. 352, no. 6284, pp. 467–470, 2016.

- [14] C. Aguerrebere, M. Delbracio, A. Bartesaghi, and G. Sapiro, "Fundamental limits in multi-image alignment," *IEEE Transactions on Signal Processing*, vol. 64, no. 21, pp. 5707–5722, 2016.
- [15] T. Bendory, N. Boumal, C. Ma, Z. Zhao, and A. Singer, "Bispectrum inversion with application to multireference alignment," *arXiv preprint arXiv:1705.00641*, 2017.
- [16] A. Bandeira, P. Rigollet, and J. Weed, "Optimal rates of estimation for multi-reference alignment," *arXiv preprint arXiv:1702.08546*, 2017.
- [17] A. Perry, J. Weed, A. Bandeira, P. Rigollet, and A. Singer, "The sample complexity of multi-reference alignment," *arXiv preprint at arXiv:1707.00943*, 2017.
- [18] I. M. Johnstone, "On the distribution of the largest eigenvalue in principal components analysis," *Annals of Statistics*, vol. 29, no. 2, pp. 295–327, 2001.
- [19] D. Paul, "Asymptotics of sample eigenstructure for a large dimensional spiked covariance model," *Statistica Sinica*, vol. 17, no. 4, pp. 1617–1642, 2007.
- [20] M. Gavish and D. L. Donoho, "Optimal shrinkage of singular values," *IEEE Transactions on Information Theory*, vol. 63, no. 4, pp. 2137–2152, 2017.
- [21] F. Benaych-Georges and R. R. Nadakuditi, "The singular values and vectors of low rank perturbations of large rectangular random matrices," *Journal of Multivariate Analysis*, vol. 111, pp. 120–135, 2012.
- [22] E. Dobriban, W. Leeb, and A. Singer, "Optimal prediction in the linearly transformed spiked model," *arXiv preprint arXiv:1709.03393*, 2017.
- [23] W. James and C. Stein, "Estimation with quadratic loss," in *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, pp. 361–379, 1961.
- [24] B. Efron and C. Morris, "Stein's estimation rule and its competitors—an empirical Bayes approach," *Journal of the American Statistical Association*, vol. 68, no. 341, pp. 117–130, 1973.
- [25] B. Efron and C. Morris, "Data analysis using Stein's estimator and its generalizations," *Journal of the American Statistical Association*, vol. 70, no. 350, pp. 311–319, 1975.
- [26] A. Singer, "Angular synchronization by eigenvectors and semidefinite programming," *Applied and computational harmonic analysis*, vol. 30, no. 1, pp. 20–36, 2011.
- [27] N. Boumal, "Nonconvex phase synchronization," *SIAM Journal on Optimization*, vol. 26, no. 4, pp. 2355–2377, 2016.
- [28] A. Perry, A. S. Wein, A. S. Bandeira, and A. Moitra, "Message-passing algorithms for synchronization problems over compact groups," *arXiv preprint arXiv:1610.04583*, 2016.
- [29] Y. Chen and E. Candes, "The projected power method: An efficient algorithm for joint alignment from pairwise differences," *arXiv preprint arXiv:1609.05820*, 2016.
- [30] A. S. Bandeira, N. Boumal, and A. Singer, "Tightness of the maximum likelihood semidefinite relaxation for angular synchronization," *Mathematical Programming*, vol. 163, no. 1, pp. 145–167, 2017.
- [31] Y. Zhong and N. Boumal, "Near-optimal bounds for phase synchronization," *arXiv preprint arXiv:1703.06605*, 2017.
- [32] A. S. Bandeira, Y. Chen, and A. Singer, "Non-unique games over compact groups and orientation estimation in cryo-em," *arXiv preprint arXiv:1505.03840*, 2015.
- [33] A. S. Bandeira, M. Charikar, A. Singer, and A. Zhu, "Multireference alignment using semidefinite programming," in *Proceedings of the 5th conference on Innovations in theoretical computer science*, pp. 459–470, ACM, 2014.
- [34] Y. Chen, L. Guibas, and Q. Huang, "Near-optimal joint object matching via convex relaxation," in *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pp. 100–108, 2014.
- [35] A. S. Bandeira, N. Boumal, and V. Voroninski, "On the low-rank approach for semidefinite programs arising in synchronization and community detection," in *Conference on Learning Theory*, pp. 361–382, 2016.
- [36] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the em algorithm," *Journal of the royal statistical society. Series B (methodological)*, pp. 1–38, 1977.
- [37] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. Springer, 2nd ed., 2009.
- [38] Y. Yu and R. J. Samworth, "A useful variant of the Davis-Kahan theorem for statisticians," *Biometrika*, vol. 102, no. 2, pp. 315–323, 2015.
- [39] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," *arXiv preprint arXiv:1011.3027*, 2010.
- [40] D. G. Chapman and H. Robbins, "Minimum variance estimation without regularity assumptions," *Ann. Math. Statist.*, vol. 22, pp. 581–586, 12 1951.
- [41] H. Cramér, *Mathematical Methods of Statistics (PMS-9)*, vol. 9. Princeton university press, 2016.
- [42] E. Abbe, J. a. M. Pereira, and A. Singer, "Sample complexity of the Boolean multireference alignment problem," in *International Symposium in Information Theory (ISIT)*, IEEE, 2017.
- [43] M. Radermacher, T. Wagenknecht, A. Verschoor, and J. Frank, "Three-dimensional reconstruction from a single-exposure, random conical tilt series applied to the 50S ribosomal subunit of Escherichia coli," *Journal of Microscopy*, vol. 146, no. 2, pp. 113–136, 1987.
- [44] Z. Kam, "The reconstruction of structure from electron micrographs of randomly oriented particles," *Journal of Theoretical Biology*, vol. 82, no. 1, pp. 15–39, 1980.
- [45] M. X. Goemans and D. P. Williamson, "Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming," *Journal of the ACM (JACM)*, vol. 42, no. 6, pp. 1115–1145, 1995.
- [46] M. Grant, S. Boyd, and Y. Ye, "CVX: Matlab software for disciplined convex programming," 2008.

## APPENDIX

### A. Proof of Proposition III.1

In the proof, we aim to show that under the stated conditions, there exists only one pair  $x \in \mathbb{R}^L$  (up to global translation) and  $\rho \in \Delta_L$  that satisfies the data constraints.

Suppose  $C_x D_\alpha C_x^T = C_y D_\beta C_y^T$  for vectors  $x, y, \alpha, \beta \in \mathbb{R}^L$ . Then, applying the Fourier transform to each side, we get:

$$D_{Fx} C_{F\alpha} D_{Fx}^* = D_{Fy} C_{F\beta} D_{Fy}^*.$$

Now suppose that  $Fx$  never vanishes. Then

$$C_{F\alpha} = D_{Fy/Fx} C_{F\beta} D_{Fy/Fx}^* = C_{F\beta} \odot r r^*,$$

where  $r = Fy/Fx$ . This equation implies

$$(F\alpha)[k - \ell] = (F\beta)[k - \ell] r[k] \overline{r[\ell]}, \quad \forall k, \ell. \quad (\text{A.1})$$

Now, since  $\alpha$  and  $\beta$  are probability distributions, they satisfy  $(F\alpha)[0] = (F\beta)[0] \neq 0$ . Then when  $k = \ell$ , equation (A.1) becomes  $1 = |r[k]|^2$ , so that  $r[k]$  has unit modulus entries.

Suppose now that  $(F\alpha)[M] \neq 0$  for some  $M$  that is relatively prime to  $L$ . From (A.1),  $(F\beta)[M] \neq 0$  too. Let  $\tilde{\omega} = (F\alpha)[M]/(F\beta)[M]$ . Then  $r[\ell + M] = \tilde{\omega} \cdot r[\ell]$  for every  $\ell$ . Since  $M$  is relatively prime to  $L$ , there is some integer  $K$  so that  $K \cdot M = 1 \pmod L$ . Therefore,  $r[\ell + 1] = \tilde{\omega}^K \cdot r[\ell]$ .

Let  $\omega = \tilde{\omega}^K$  so that  $r[\ell + 1] = \omega \cdot r[\ell]$ . We can then get  $r[\ell] = \omega^\ell \cdot r[0]$ . In particular,  $r[0] = r[L] = \omega^L \cdot r[0]$ , and so  $\omega^L = 1$ ; that is,  $\omega$  is an  $L$ th root of unity. But also, by definition of  $r[\ell] = (Fy)[\ell]/(Fx)[\ell]$ , so we have

$$(Fy)[\ell] = \omega^\ell \cdot r[0] \cdot (Fx)[\ell].$$

Thus,  $y$  is just a shifted and rescaled version of  $x$ , which is the desired result.

### B. Proof of Lemma III.4

For any  $0 \leq i \leq L - 1$ , we can write

$$(\rho * \theta)[i] = e_i^T C_\rho \theta,$$

with  $e_i$  the unit vector with one in its  $i$ th entry. Consequently, equality of two distinct entries  $i$  and  $j$  implies

$$(e_i - e_j)^T C_\rho \theta = 0. \quad (\text{B.1})$$

However, for a random choice of  $\theta$ , if (B.1) holds with non-zero probability, then

$$(e_i - e_j)^T C_\rho = 0,$$

or,

$$C_\rho^T e_i = C_\rho^T e_j.$$

The latter implies that  $\rho$  shifted by  $i$  equals  $\rho$  shifted by  $j$ , i.e.,  $\rho[k - i] = \rho[k - j]$ , or

$$\rho[k] = \rho[k + i - j], \quad \forall k.$$

Therefore,  $\rho$  is periodic.

### C. Proof of Proposition III.6

We prove the proposition by construction. Let  $x_1$  be a signal with a non-vanishing DFT  $Fx_1$ . Define a second signal  $x_2$  so that

$$(Fx_2)[k] = \begin{cases} (Fx_1)[k], & k = t\frac{L}{\ell}, \quad t = 0, \dots, \ell - 1, \\ -(Fx_1)[k], & \text{otherwise.} \end{cases} \quad (\text{C.1})$$

Clearly, as  $L > 1$ ,  $x_1 \neq x_2$ . In addition, since  $x_1$  is real, the construction ensures that  $x_2$  is real as well.

The  $\ell$  periodicity of  $\rho$  means a sparsity pattern for  $F\rho$ . Particularly,  $F\rho$  is zero everywhere besides

$$(F\rho)[kL/\ell] \neq 0 \iff kL/\ell \text{ is integer}, \quad (\text{C.2})$$

for  $k = 0, \dots, \ell - 1$ . It is easy to verify that

$$(Fx_1)[k](F\rho)[k] = (Fx_2)[k](F\rho)[k], \quad k = 0, \dots, L - 1.$$

Therefore,  $x_1$  and  $x_2$  share the same first moment.

For the second moments, we will show the equality

$$C_{x_1} D_\rho C_{x_1}^T = C_{x_2} D_\rho C_{x_2}^T.$$

Applying the Fourier matrix, due to the realness of  $\rho$ , the latter is equivalent to

$$D_{Fx_1} C_{F\rho} D_{Fx_1} = D_{Fx_2} C_{F\rho} D_{Fx_2},$$

or, for all  $i, j = 0, \dots, L - 1$ ,

$$(Fx_1)[i](F\rho)[i - j](Fx_1)[j] = (Fx_2)[i](F\rho)[i - j](Fx_2)[j].$$

By the sparsity pattern of (C.2), this equality should hold only in:

$$(Fx_1)[i](Fx_1)[i + tL/\ell] = (Fx_2)[i](Fx_2)[i + tL/\ell],$$

for all  $t = 0, \dots, \ell$  and  $i = 0, \dots, L - 1$ . By the construction (C.1), this equation holds true.

### D. Proof of Claim III.7

Throughout the proof, we assume that each period has no repeated values. This property is guaranteed by reshuffling the measurements with random  $\theta \in \Delta_L$ ; see Lemma III.4. Additionally, we assume without loss of generality that  $|Fx[k]| = 1$  for all  $k$ .

Observe that both  $x$  and  $R_{L/2}x$  are eigenvectors of  $\hat{M} = C_x D_\rho C_x^T$  (we assume exact knowledge of the moments) with the same eigenvalue. Also,  $x$  and  $R_{L/2}x$  are orthogonal as columns in the orthogonal matrix  $C_x$ . Then, if  $u$  is an eigenvector, we can write for some scalars  $\alpha, \beta \in \mathbb{R}^L$ :

$$u = \alpha x + \beta R_{L/2}x,$$

and therefore,

$$R_{L/2}u = \alpha R_{L/2}x + \beta x,$$

as  $R_{L/2} = R_{L/2}^{-1}$ . Then, one can verify that the inner product of  $u$  and  $R_{L/2}u$  is  $2\alpha\beta\|x\|^2$ . Since the signals are orthogonal, their inner product is zero. This means that  $\alpha$  or  $\beta$  must be zero. This in turn implies that  $u$  was either  $x$  or  $R_{L/2}x$  in the first place. So  $x$  is the unique eigenvector of  $\hat{M}$  that is orthogonal to its translation by  $L/2$ . This completes the proof.

### E. Proof of Theorem V.4

The proof mimics the original proof by Chapman and Robbins. Recalling equation (V.7) and the definition of positive semidefinite matrices, the statement is equivalent to

$$\begin{aligned} & \mathbb{E}_{x,\rho}[(w^T(\phi_x(\hat{x}) - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})]))^2] \\ & \geq \frac{(w^T(\mathbb{E}_{\tilde{x},\tilde{\rho}}[\phi_x(\hat{x})] - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})]))^2}{\chi_N^2(f_{\tilde{x},\tilde{\rho}}\|f_{x,\rho})}, \end{aligned} \quad (\text{E.1})$$

for all  $w, \tilde{x} \in \mathbb{R}^L$  and  $\tilde{\rho} \in \Delta^L$ . Let

$$Z = \frac{f_{\tilde{x},\tilde{\rho}}(Y)}{f_{x,\rho}(Y)}.$$

This random variable has the properties  $\mathbb{E}_{x,\rho}[g(Y)Z] = \mathbb{E}_{\tilde{x},\tilde{\rho}}[g(Y)]$  and  $\mathbb{E}_{x,\rho}[(Z - 1)^2] = \chi_N^2(f_{\tilde{x},\tilde{\rho}}\|f_{x,\rho})$ . Thus

$$\begin{aligned} & w^T(\mathbb{E}_{\tilde{x},\tilde{\rho}}[\phi_x(\hat{x})] - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})]) \\ & = \mathbb{E}_{\tilde{x},\tilde{\rho}}[w^T \phi_x(\hat{x})] - \mathbb{E}_{x,\rho}[w^T \phi_x(\hat{x})] \\ & = \mathbb{E}_{x,\rho}[w^T \phi_x(\hat{x})(Z - 1)] \\ & = \mathbb{E}_{x,\rho}[w^T(\phi_x(\hat{x}) - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})])(Z - 1)], \end{aligned}$$

and by Cauchy-Schwarz

$$\begin{aligned} & (w^T(\mathbb{E}_{\tilde{x},\tilde{\rho}}[\phi_x(\hat{x})] - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})]))^2 \\ & \leq \mathbb{E}_{x,\rho}[(w^T(\phi_x(\hat{x}) - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})]))^2] \chi_N^2(f_{\tilde{x},\tilde{\rho}}\|f_{x,\rho}). \end{aligned}$$

### F. Proof of Proposition V.5

Equation (V.11) follows from some algebraic manipulations.

$$\begin{aligned}
& \chi^2(f_{\tilde{x},\tilde{\rho}}||f_{x,\rho}) \\
&= \int_{\mathbb{R}^L} \left( \frac{f_{\tilde{x},\tilde{\rho}}(y;\gamma)}{f_{x,\rho}(y;\gamma)} - 1 \right)^2 f_{x,\rho}(y;\gamma) dy \\
&= \int_{\mathbb{R}^L} \frac{\left( \sum_{i=0}^{\infty} (\alpha_{\tilde{x},\tilde{\rho}}^i(y) - \alpha_{x,\rho}^i(y)) \frac{\gamma^i}{i!} \right)^2}{\sum_{i=0}^{\infty} \alpha_{x,\rho}^i(y) \frac{\gamma^i}{i!}} f_Z(y) dy \\
&= \int_{\mathbb{R}^L} \frac{\left( \sum_{i=d}^{\infty} (\alpha_{\tilde{x},\tilde{\rho}}^i(y) - \alpha_{x,\rho}^i(y)) \frac{\gamma^i}{i!} \right)^2}{1 + \sum_{i=1}^{\infty} \alpha_{x,\rho}^i(y) \frac{\gamma^i}{i!}} f_Z(y) dy \\
&= \frac{\gamma^{2d}}{(d!)^2} \int_{\mathbb{R}^L} \left( \alpha_{\tilde{x},\tilde{\rho}}^d(y) - \alpha_{x,\rho}^d(y) \right)^2 f_Z(y) dy + \mathcal{O}(\gamma^{2d+1}) \\
&= \frac{\gamma^{2d}}{(d!)^2} \mathbb{E} \left[ \left( \alpha_{\tilde{x},\tilde{\rho}}^d(Z) - \alpha_{x,\rho}^d(Z) \right)^2 \right] + \mathcal{O}(\gamma^{2d+1}),
\end{aligned}$$

where the third equation follows from  $\alpha_{\tilde{x},\tilde{\rho}}^n(z) = \alpha_{x,\rho}^n(z)$  almost surely for all  $n < d$ , by the definition of  $d$ . (V.11) now follows from  $\gamma = 1/\sigma$ .

We now prove (V.12). We can write  $f_{x,\rho}$  explicitly by

$$\begin{aligned}
f_{x,\rho}(y;\gamma) &= \frac{1}{\sqrt{2\pi}^L} \sum_{\ell=0}^{L-1} \rho[\ell] \exp\left(-\frac{\|y - \gamma R_\ell x\|^2}{2}\right) \\
&= \mathbb{E}_S[f_G(y - \gamma R_S x)]. \tag{F.1}
\end{aligned}$$

We show that

$$\mathbb{E}_G \left[ \alpha_{\tilde{x},\tilde{\rho}}^d(G) \alpha_{x,\rho}^d(G) \right] = d! \left\langle A_{\tilde{x},\tilde{\rho}}^d, A_{x,\rho}^d \right\rangle,$$

implies (V.12). Let  $S$  and  $\tilde{S}$  be two independent random variables such that  $S \sim \rho$  and  $\tilde{S} \sim \rho$ . We have

$$\begin{aligned}
\left\langle A_{\tilde{x},\tilde{\rho}}^d, A_{x,\rho}^d \right\rangle &= \left\langle \mathbb{E}_{\tilde{S}}(R_{\tilde{S}\tilde{x}}^{\otimes d}), \mathbb{E}_S[(R_S x)^{\otimes d}] \right\rangle \\
&= \mathbb{E}_{\tilde{S},S} \left[ \left\langle (R_{\tilde{S}\tilde{x}})^{\otimes d}, (R_S x)^{\otimes d} \right\rangle \right] \\
&= \mathbb{E}_{\tilde{S},S} \left[ \left\langle R_{\tilde{S}\tilde{x}}, R_S x \right\rangle^d \right].
\end{aligned}$$

On the other hand, by equations (V.10) and (F.1):

$$\begin{aligned}
& \mathbb{E}_G \left[ \alpha_{\tilde{x},\tilde{\rho}}^d(G) \alpha_{x,\rho}^d(G) \right] \\
&= \mathbb{E}_G \left[ \frac{\partial^d}{\partial \tilde{\gamma}^d} \left( \frac{f_{\tilde{x},\tilde{\rho}}(G;\tilde{\gamma})}{f_G(G)} \right) \frac{\partial^d}{\partial \gamma^d} \left( \frac{f_{x,\rho}(G;\gamma)}{f_G(G)} \right) \right]_{\tilde{\gamma}=0, \gamma=0} \\
&= \frac{\partial^{2d}}{\partial \tilde{\gamma}^d \partial \gamma^d} \mathbb{E}_G \left[ \frac{f_{\tilde{x},\tilde{\rho}}(G;\tilde{\gamma})}{f_G(G)} \frac{f_{x,\rho}(G;\gamma)}{f_G(G)} \right]_{\tilde{\gamma},\gamma=0} \\
&= \frac{\partial^{2d}}{\partial \tilde{\gamma}^d \partial \gamma^d} \mathbb{E}_{G,\tilde{S},S} \left[ \frac{f_G(G - \tilde{\gamma} R_{\tilde{S}\tilde{x}})}{f_G(G)} \frac{f_G(G - \gamma R_S x)}{f_G(G)} \right]_{\tilde{\gamma},\gamma=0}.
\end{aligned}$$

We have

$$\begin{aligned}
& \mathbb{E}_G \left[ \frac{f_G(G - \tilde{\gamma} R_{\tilde{S}\tilde{x}})}{f_G(G)} \frac{f_G(G - \gamma R_S x)}{f_G(G)} \right] \\
&= \frac{1}{\sqrt{2\pi}^L} \int_{\mathbb{R}^L} \exp\left(-\frac{\|z - \tilde{\gamma} R_{\tilde{S}\tilde{x}}\|^2 + \|z - \gamma R_S x\|^2 - \|z\|^2}{2}\right) dz \\
&= \frac{1}{\sqrt{2\pi}^L} \int_{\mathbb{R}^L} \exp\left(-\frac{\|z - \tilde{\gamma} R_{\tilde{S}\tilde{x}} - \gamma R_S x\|^2}{2} + \gamma \tilde{\gamma} \langle R_{\tilde{S}\tilde{x}}, R_S x \rangle\right) dz \\
&= \exp\left(\gamma \tilde{\gamma} \langle R_{\tilde{S}\tilde{x}}, R_S x \rangle\right).
\end{aligned}$$

The proof of (V.12) now follows

$$\begin{aligned}
& \frac{\partial^{2d}}{\partial \tilde{\gamma}^d \partial \gamma^d} \mathbb{E}_{\tilde{S},S} \left[ \exp\left(\gamma \tilde{\gamma} \langle R_{\tilde{S}\tilde{x}}, R_S x \rangle\right) \right]_{\tilde{\gamma},\gamma=0} \\
&= d! \mathbb{E}_{\tilde{S},S} \left[ \langle R_{\tilde{S}\tilde{x}}, R_S x \rangle^d \right].
\end{aligned}$$

### G. Analog results for derivatives

This section provides analog results to the ones presented in section V, but involving the limit  $(\tilde{x}, \tilde{\rho}) \rightarrow (x, \rho)$ . Instead of the  $\chi^2$  divergence and the autocorrelations, these results use the Fisher information matrix and directional derivatives of the auto correlations, respectively.

For the rest of the section, let  $v = (z, \theta) \in \mathbb{R}^{2L}$  such that  $\mathbf{1}^T \theta = 0$  and  $\theta[i] \geq 0$  whenever  $\rho[i] = 0$ . The Fisher information matrix is the  $2L \times 2L$  matrix defined by

$$\Gamma_{x,\rho}^N := \text{Cov}[\nabla \log f_{x,\rho}].$$

The Fisher information matrix is also the Hessian of the  $\chi^2$  divergence, i.e.,

$$\lim_{h \rightarrow 0} \frac{\chi_N^2(f_{x+hv, \rho+h\theta} || f_{x,\rho})}{h^2} = v^T \Gamma_{x,\rho}^N v. \tag{G.1}$$

The Fisher information matrix of  $N$  observations is related to the one observation version by

$$\Gamma_{x,\rho}^N = N \Gamma_{x,\rho}. \tag{G.2}$$

We define the Jacobian  $J_{x,\rho}$  as the  $L \times 2L$  matrix such that

$$J_{x,\rho} v = \lim_{h \rightarrow 0} \frac{\mathbb{E}_{x+hv, \rho+h\theta}[\phi_x(\hat{x})] - \mathbb{E}_{x,\rho}[\phi_x(\hat{x})]}{h}. \tag{G.3}$$

We also define the directional derivative of  $A_{x,\rho}^d$  along  $v$  as the  $d$ -dimensional tensor

$$\nabla_v A_{x,\rho}^d := \lim_{h \rightarrow 0} \frac{A_{x+hv, \rho+h\theta}^d - A_{x,\rho}^d}{h}.$$

The next corollary is an analog of the Cramér-Rao bound for estimation of an orbit in MRA.

**Corollary G.1.** *For any  $v = (z, \theta) \in \mathbb{R}^{2L}$ , such that  $\mathbf{1}^T \theta = 0$  and  $\theta[i] \geq 0$ , whenever  $\rho[i] = 0$ , we have*

$$\text{Cov}[\phi_x(\hat{x})] \succeq \frac{J_{x,\rho} v v^T J_{x,\rho}^T}{N v^T \Gamma_{x,\rho} v}.$$

*Proof.* If  $\theta$  is under the hypothesis of the theorem, then there exists  $h_0 > 0$  such that for all  $0 \leq h \leq h_0$ ,  $\rho + h\theta \in \Delta^L$ .

Letting  $(\tilde{x}, \tilde{\rho}) = hv + (x, \rho)$  in Theorem V.4 we obtain for any  $w \in \mathbb{R}^L$

$$\begin{aligned} & w^T \text{Cov}[\phi_x(\hat{x})]w \\ & \geq \lim_{h \rightarrow 0} \frac{(w^T (\mathbb{E}_{x+h\tilde{z}, \rho+h\tilde{\theta}}[\phi_x(\hat{x})] - \mathbb{E}_{x, \rho}[\phi_x(\hat{x})]))^2}{\chi_N^2(f_{x+h\tilde{z}, \rho+h\tilde{\theta}} \| f_{x, \rho})} \\ & = \frac{(w^T J_{x, \rho} v)^2}{N v^T \Gamma_{x, \rho} v}, \end{aligned}$$

by equations (G.1), (G.2) and (G.3), and the corollary follows.  $\square$

We now use (V.9) to give an expression of the Fisher information.

**Lemma G.2.** For any  $v = (z, \theta) \in \mathbb{R}^{2L}$ ,

$$v^T \Gamma_{x, \rho} v = \frac{\sigma^{-2d}}{(d!)^2} \mathbb{E}_G \left[ \left( v^T \nabla \alpha_{x, \rho}^d(G) \right)^2 \right] + \mathcal{O}(\sigma^{-2d-1}), \quad (\text{G.4})$$

$$\text{where } d = \inf \left\{ n : \mathbb{E}_G \left[ \left( v^T \nabla \alpha_{x, \rho}^n(G) \right)^2 \right] > 0 \right\}.$$

*Proof.* In this case we cannot just take the limit  $(\tilde{x}, \tilde{\rho}) \rightarrow (x, \rho)$ , since  $\mathcal{O}(\sigma^{-2d-1})$  might blow up. Instead we proceed by doing similar algebraic manipulations.

$$\begin{aligned} v^T \Gamma_{x, \rho} v &= v^T \text{Cov}[\nabla \log f_{x, \rho}]v \\ &= \mathbb{E}_{x, \rho} \left[ \left( \frac{v^T \nabla f_{x, \rho}(y; \gamma)}{f_{x, \rho}(y; \gamma)} \right)^2 \right] \\ &= \int_{\mathbb{R}^L} \left( \frac{v^T \nabla f_{x, \rho}(y; \gamma)}{f_{x, \rho}(y; \gamma)} \right)^2 f_{x, \rho}(y; \gamma) dy \\ &= \int_{\mathbb{R}^L} \frac{\left( \sum_{i=0}^{\infty} v^T \nabla \alpha_{x, \rho}^i(y) \frac{\gamma^i}{i!} \right)^2}{\sum_{i=0}^{\infty} \alpha_{x, \rho}^i(y) \frac{\gamma^i}{i!}} f_G(y) dy \\ &= \int_{\mathbb{R}^L} \frac{\left( \sum_{i=d}^{\infty} v^T \nabla \alpha_{x, \rho}^i(y) \frac{\gamma^i}{i!} \right)^2}{1 + \sum_{i=i}^{\infty} \alpha_{x, \rho}^i(y) \frac{\gamma^i}{i!}} f_G(y) dy \\ &= \frac{\gamma^{2d}}{(d!)^2} \int_{\mathbb{R}^L} \left( v^T \nabla \alpha_{x, \rho}^d(y) \right)^2 f_G(y) dy + \mathcal{O}(\gamma^{2d+1}) \\ &= \frac{\gamma^{2d}}{(d!)^2} \mathbb{E}_G \left[ \left( v^T \nabla \alpha_{x, \rho}^d(G) \right)^2 \right] + \mathcal{O}(\gamma^{2d+1}), \end{aligned}$$

where the second and fifth lines follow from

$$\mathbb{E}_{x, \rho} \left[ \frac{\nabla f_{x, \rho}(y; \gamma)}{f_{x, \rho}(y; \gamma)} \right] = 0 \quad \text{and} \quad v^T \nabla \alpha_{x, \rho}^n(z) = 0,$$

almost surely for  $n < d$ , respectively. The result follows since  $\gamma = 1/\sigma$ .  $\square$

We now use expression (V.8) to calculate (G.4).

**Lemma G.3.**

$$\mathbb{E}_G \left[ \left( v^T \nabla \alpha_{x, \rho}^d(G) \right)^2 \right] = d! \|\nabla_v A_{x, \rho}^d\|^2.$$

*Proof.* We let  $(\tilde{x}, \tilde{\rho}) = (x, \rho) + hv$  in (V.12) and take the limit  $h \rightarrow 0$  to get

$$\begin{aligned} & \mathbb{E} \left[ \left( v^T \nabla \alpha_{x, \rho}^d(G) \right)^2 \right] \\ &= \lim_{h \rightarrow 0} \frac{\mathbb{E} \left[ \left( \alpha_{x+h\tilde{z}, \rho+h\tilde{\theta}}^d(G) - \alpha_{x, \rho}^d(G) \right)^2 \right]}{h^2} \\ &= d! \lim_{h \rightarrow 0} \frac{\|A_{x+h\tilde{z}, \rho+h\tilde{\theta}}^d - A_{x, \rho}^d\|^2}{h^2} \\ &= d! \|\nabla_v A_{x, \rho}^d\|^2. \end{aligned}$$

$\square$

Next, from Corollary G.1 and Lemmas V.5 and G.3 we obtain

**Lemma G.4.** For any  $v = (z, \theta) \in \mathbb{R}^{2L}$ , such that  $\mathbf{1}^T \theta = 0$  and  $\theta[i] \geq 0$  whenever  $\rho[i] = 0$ , we have

$$\text{Cov}[\phi_x(\hat{x})] \succeq \frac{1}{d!} \frac{\sigma^{2d}}{N} \frac{J_{x, \rho} v v^T J_{x, \rho}^T}{\|\nabla_v A_{x, \rho}^d\|^2} - \mathcal{O}(\sigma^{2d-1}),$$

where  $d = \inf \left\{ n : \|\nabla_v A_{x, \rho}^n\|^2 > 0 \right\}$ .

#### H. Proof of Theorem V.1

Before proving Theorem V.1, we need the following lemma.

**Lemma H.1.** The entries with index  $\mathbf{k} = (k_1, k_2, \dots, k_d) \in \mathbb{Z}_L^d$  of  $A_{x, \rho}^d$  and  $\nabla_v A_{x, \rho}^d$  can be explicitly written as

$$A_{x, \rho}^d[\mathbf{k}] := \sum_{\ell=0}^L \rho[\ell] \prod_{i=1}^d x[k_i - \ell], \quad (\text{H.1})$$

and

$$(\nabla_v A_{x, \rho}^d)[\mathbf{k}] = \sum_{\ell=0}^{L-1} \left( \rho[\ell] \sum_{i=1}^d \frac{z[k_i - \ell]}{x[k_i - \ell]} + \theta[\ell] \right) \prod_{i=1}^d x[k_i - \ell], \quad (\text{H.2})$$

where we use the notation  $x[k_i - \ell]/x[k_i - \ell] = 1$  when  $x[k_i - \ell] = 0$ . Moreover, denote the  $d$ -dimensional Fourier Transform by  $F_d$ . We have

$$\|A_{\tilde{x}, \tilde{\rho}}^d - A_{x, \rho}^d\|^2 = \frac{1}{L^d} \|F_d A_{\tilde{x}, \tilde{\rho}}^d - F_d A_{x, \rho}^d\|^2, \quad (\text{H.3})$$

and

$$\|\nabla_v A_{x, \rho}^d\|^2 = \frac{1}{L^d} \|F_d \nabla_v A_{x, \rho}^d\|^2. \quad (\text{H.4})$$

Also, for any  $\mathbf{a} = (a_1, a_2, \dots, a_d) \in \mathbb{Z}_L^d$  we have

$$F_d A_{x, \rho}^d[\mathbf{a}] = F \rho \left[ \sum_{j=1}^d a_j \right] \prod_{j=1}^d F x[a_j], \quad (\text{H.5})$$

and

$$F_d (\nabla_v A_{x, \rho}^d)[\mathbf{a}] = \quad (\text{H.6})$$

$$\left( \sum_{j=1}^d \frac{F z[a_j]}{F x[a_j]} + \frac{F \theta \left[ \sum_{j=1}^d a_j \right]}{F \rho \left[ \sum_{j=1}^d a_j \right]} \right) F_d A_{x, \rho}^d[\mathbf{a}], \quad (\text{H.7})$$

again using the notation  $Fx[a_j]/Fx[a_j] = 1$  even if  $Fx[a_j] = 0$ .

*Proof.* We first prove equation (H.1). By equation (V.2), we have

$$\begin{aligned} A_{x,\rho}^d[\mathbf{k}] &= \mathbb{E}_S \left[ \prod_{i=1}^d (R_S x)[k_i] \right] \\ &= \mathbb{E}_S \left[ \prod_{i=1}^d x[k_i - S] \right] \\ &= \sum_{\ell=0}^L \rho[\ell] \prod_{i=1}^d x[k_i - \ell]. \end{aligned}$$

Equation (H.2) follows from the formula of the derivative of the product:

$$\begin{aligned} &\left( \rho[\ell] \prod_{i=1}^d x[k_i - \ell] \right)' \\ &= \left( \rho'[\ell] + \rho[\ell] \sum_{i=1}^d \frac{x'[k_i - \ell]}{x[k_i - \ell]} \right) \prod_{i=1}^d x[k_i - \ell]. \end{aligned}$$

We finally prove (H.5), the proof of (H.6) is analogous.

$$\begin{aligned} F_d A_{x,\rho}^d[\mathbf{a}] &= \sum_{\mathbf{k} \in \mathbb{Z}_L^d} A_{x,\rho}^d[\mathbf{k}] \exp\left(-\frac{2\pi\iota}{L} \langle \mathbf{k}, \mathbf{a} \rangle\right) \\ &= \sum_{\mathbf{k} \in \mathbb{Z}_L^d} \sum_{\ell=0}^{L-1} \rho[\ell] \prod_{j=1}^d x[k_j - \ell] \exp\left(-\frac{2\pi\iota}{L} k_j a_j\right) \\ &= \sum_{\mathbf{k} \in \mathbb{Z}_L^d} \sum_{\ell=0}^{L-1} \rho[\ell] \prod_{j=1}^d x[k_j] \exp\left(-\frac{2\pi\iota}{L} a_j (k_j + \ell)\right) \\ &= \sum_{\mathbf{k} \in \mathbb{Z}_L^d} \sum_{\ell=0}^{L-1} \rho[\ell] \prod_{j=1}^d x[k_j] \exp\left(-\frac{2\pi\iota}{L} (k_j a_j + \ell a_j)\right) \\ &= \sum_{\ell=0}^{L-1} \rho[\ell] \exp\left(-\frac{2\pi\iota}{L} \left(\ell \sum_{j=1}^d a_j\right)\right) \prod_{j=1}^d Fx[a_j] \\ &= F\rho \left[ \sum_{j=1}^d a_j \right] \prod_{j=1}^d Fx[a_j]. \end{aligned}$$

□

We are now ready to prove Theorem V.1, starting by (V.3). Since  $\hat{x}$  is asymptotically unbiased,  $\mathbb{E}_{x,\rho}[\phi_x(\hat{x})] \rightarrow x$  and  $J_{x,\rho} \rightarrow [L \ 0_{L \times L}]$  as  $N \rightarrow \infty$ . By (V.6) and Corollary G.1

we have

$$\begin{aligned} &\lim_{N \rightarrow \infty} N \cdot \text{MSE} \\ &\geq \lim_{N \rightarrow \infty} \frac{N \text{tr}(\text{Cov}[\phi_x(\hat{x})])}{\|x\|^2} \quad (\text{H.8}) \\ &\geq \lim_{N \rightarrow \infty} \frac{\sigma^{2d} \|J_{x,\rho} v\|^2}{d! \|x\|^2} \frac{1}{\|\nabla_v A_{x,\rho}^d\|^2} - \mathcal{O}(\sigma^{2d-1}) \\ &= \frac{\sigma^{2d} \|z\|^2}{d! \|x\|^2} \frac{1}{\|\nabla_v A_{x,\rho}^d\|^2} - \mathcal{O}(\sigma^{2d-1}). \quad (\text{H.9}) \end{aligned}$$

We will choose  $z = x - \frac{1^T x}{L} \mathbf{1}$ , and  $\theta = \frac{1}{L} \mathbf{1} - \rho$ . This choice of  $\theta$  is under the theorem assumptions, since  $1^T \theta = 0$  and  $\theta[i] = \frac{1}{L} \geq 0$  whenever  $\rho[i] = 0$ . By the linearity of the Fourier transform, this definition is equivalent to  $Fz = Fx - Fx[0]\delta_0$  and  $F\theta = F\rho[0]\delta_0 - F\rho = \delta_0 - F\rho$ . Since the  $d$ -dimensional Fourier Transform is unitary, we can write using Lemma H.1

$$\|\nabla_v A_{x,\rho}^d\|^2 = \frac{1}{L^d} \sum_{\mathbf{a} \in \mathbb{Z}_L^d} |F_d \nabla_v A_{x,\rho}^d[\mathbf{a}]|^2. \quad (\text{H.10})$$

For  $d = 1, 2$  we have

$$F_1 \nabla_v A_{x,\rho}^1[a] = F\rho[a] Fz[a] + F\theta[a] Fx[a],$$

and

$$\begin{aligned} F_2 \nabla_v A_{x,\rho}^2[a_1, a_2] &= F\rho[a_1 + a_2] Fz[a_1] Fx[a_2] \\ &\quad + F\rho[a_1 + a_2] Fx[a_1] Fz[a_2] \\ &\quad + F\theta[a_1 + a_2] Fx[a_1] Fx[a_2]. \quad (\text{H.11}) \end{aligned}$$

Now by our choice of  $z$  and  $\theta$  we have  $F\rho[a] Fz[a] = -F\theta[a] Fx[a]$  for all  $a \in \mathbb{Z}_L$ , so  $\|\nabla_v A_{x,\rho}^1\| = 0$ . On the other hand, by some algebra manipulation of (H.11) we obtain

$$\begin{aligned} &\|\nabla_v A_{x,\rho}^2\|^2 \\ &= \frac{1}{L^2} \left( 3\|Fz\|_4^4 + \sum_{\mathbf{a} \in \mathbb{Z}_L^2} |F\rho[a_1 + a_2] Fz[a_1] Fz[a_2]|^2 \right) \\ &\leq \frac{4}{L^2} \|Fz\|^4 \\ &\leq 4\|z\|^2 \|x\|^2, \end{aligned}$$

where we used  $|F\rho[a_1 + a_2]| \leq 1$  and  $\|Fz\|_4 \leq \|Fz\| \leq \|Fx\|$ , and the result follows.

We now proceed to prove (V.4). Suppose that  $\rho$  is periodic with period  $\ell < \frac{L}{2}$ , and let  $b = \frac{L}{\ell}$ , so that  $b > 2$ . Then  $F\rho[k] = 0$  if  $b$  does not divide  $k$ . For a positive integer  $i \leq \lceil \frac{b-2}{2} \rceil$ , define  $z_i \in \mathbb{R}^L$  such that

$$Fz_i[k] = \begin{cases} Fx[k]\iota & \text{if } b|k - i, \\ -Fx[k]\iota & \text{if } b|k + i, \\ 0 & \text{otherwise,} \end{cases}$$

where  $b|k$  means that  $b$  divides  $k$ . Assume  $z_i \neq 0$ , let  $\theta_i = 0_L$  and  $v_i = (z_i, \theta_i)$ . Since  $\hat{x}$  is asymptotically unbiased and



$\{z_i\}_{1 \leq i \leq \lceil \frac{b-2}{2} \rceil}$  is a set of orthogonal vectors, we have by (H.8) and Corollary G.1:

$$\begin{aligned} & \lim_{N \rightarrow \infty} N \cdot \text{MSE} \\ & \geq \lim_{N \rightarrow \infty} \frac{N \text{tr}(\text{Cov}[\phi_x(\hat{x})])}{\|x\|^2} \\ & \geq \lim_{N \rightarrow \infty} \frac{1}{\|x\|^2} \sum_{i=1}^{\lceil \frac{b-2}{2} \rceil} \frac{N z_i^T \text{Cov}[\phi_x(\hat{x})] z_i}{\|z_i\|^2} \\ & \geq \frac{1}{\|x\|^2} \sum_{i=1}^{\lceil \frac{b-2}{2} \rceil} \frac{\sigma^{2d_i}}{d_i!} \frac{\|z_i\|^2}{\|\nabla_{v_i} A_{x,\rho}^{d_i}\|^2} - \mathcal{O}(\sigma^{2d_i-1}), \end{aligned}$$

where  $d_i = \inf \{n : \|\nabla_{v_i} A_{x,\rho}^n\|^2 > 0\}$ . Recalling equation (H.10), we have now for  $d = 1, 2, 3$ , since  $\theta_i = 0$ ,

$$F_1 \nabla_{v_i} A_{x,\rho}^1[a] = F\rho[a] Fz_i[a], \quad (\text{H.12})$$

$$\begin{aligned} F_2 \nabla_{v_i} A_{x,\rho}^2[a_1, a_2] &= F\rho[a_1 + a_2] (Fz_i[a_1] Fx[a_2] \\ &\quad + Fx[a_1] Fz_i[a_2]), \quad (\text{H.13}) \end{aligned}$$

and

$$\begin{aligned} F_3 \nabla_{v_i} A_{x,\rho}^3[a_1, a_2, a_3] &= F\rho[a_1 + a_2 + a_3] \\ &\quad (Fz_i[a_1] Fx[a_2] Fx[a_3] \\ &\quad + Fx[a_1] Fz_i[a_2] Fx[a_3] \\ &\quad + Fx[a_1] Fx[a_2] Fz_i[a_3]). \quad (\text{H.14}) \end{aligned}$$

Since  $F\rho[a] \neq 0 \Rightarrow b|a \Rightarrow Fz_i[a] = 0$ , (H.12) = 0  $\forall a \in \mathbb{Z}_L$ . Also  $F\rho[a_1 + a_2] \neq 0$  implies  $b|a_1 + a_2$ . Let  $\tilde{a}_j = \text{mod}(a_j, b)$  for  $j = 1$  and  $2$ . Since  $b|a_1 + a_2$ ,  $\tilde{a}_1 + \tilde{a}_2 = b$ , so assume with out loss of generality that  $\tilde{a}_1 \leq \frac{b}{2}$ . If  $\tilde{a}_1 \neq i$ , then  $Fz_i[a_1] = Fz_i[a_2] = 0$ . On the other hand, if  $\tilde{a}_1 = i$ , then

$$\begin{aligned} & Fz_i[a_1] Fx[a_2] + Fx[a_1] Fz_i[a_2] \\ &= \iota Fx[a_1] Fx[a_2] - \iota Fx[a_1] Fx[a_2] \\ &= 0, \end{aligned}$$

so (H.13) = 0  $\forall \mathbf{a} \in \mathbb{Z}_L^2$ . Finally since  $|F\rho[\cdot]| \leq 1$  we have

$$\begin{aligned} \|\nabla_{v_i} A_{x,\rho}^3\|^2 &\leq \frac{9}{L^3} \sum_{\mathbf{a} \in \mathbb{Z}_L^3} |Fz_i[a_1] Fx[a_2] Fx[a_3]|^2 \\ &= 9 \|z_i\|^2 \|x\|^4, \end{aligned}$$

and the result follows. Finally, if  $z_i = 0$ , we can alternatively choose

$$F\tilde{z}_i[k] = \begin{cases} \iota & \text{if } b|k - i, \\ -\iota & \text{if } b|k + i, \\ 0 & \text{otherwise.} \end{cases}$$

We still have (H.12) = 0  $\forall a \in \mathbb{Z}_L$  and (H.13) = 0 for all  $\mathbf{a} \in \mathbb{Z}_L^2$  except if  $\tilde{a}_1 = i$ . But  $z_i = 0$  implies  $Fx[a] = 0$  if  $\text{mod}(a, b) = \pm i$ , so (H.13) = 0 also if  $\tilde{a}_1 = i$ .

## I. Proof of Lemma VI.1

It is easy to check that the condition  $q[\ell] > 0$  is automatically enforced whenever  $w[\ell] > 0$  (otherwise the objective is  $-\infty$ ). So the simplex constraint is equivalent to  $\sum_{\ell=0}^{L-1} q[\ell] = 1$ . The Lagrangian for this problem is the function:

$$(q, \nu) = \sum_{\ell=0}^{L-1} w[\ell] \log(q[\ell]) + \nu \left( 1 - \sum_{\ell=0}^{L-1} q[\ell] \right),$$

and the KKT conditions imply  $q^*[\ell] = \frac{w[\ell]}{\nu^*}$ . Since  $q$  is on the simplex, we conclude that  $\nu^* = \sum_{\ell'=0}^{L-1} w[\ell']$ .

## J. Convex relaxation with semidefinite program

In this section, we propose an additional algorithm for non-uniform MRA based on a semidefinite programming (SDP) relaxation.

Since the power spectrum of the signal can be estimated from the data at sample complexity scaling as  $1/\text{SNR}^2$  according to (III.6), we assume in this section, without loss of generality, that  $|Fx[k]| = 1$  for all  $k$ . Note, that as in Algorithm 2, the normalization is done on the second moment matrix, not the individual observations, in order to retain the noise statistics.

The SDP relaxation is based on considering the second moment matrix in the Fourier domain, namely,

$$M' = F \left( M - \sigma^2 LI \right) F^{-1} = D_{Fx} C_{F\rho}^T D_{Fx}^*. \quad (\text{J.1})$$

The last expression can be also written as

$$M' = C_{F\rho}^T \odot (FxFx^*),$$

or

$$M' \odot \bar{X} = C_{F\rho}, \quad (\text{J.2})$$

where  $X = (Fx)(Fx)^*$ . and  $\bar{\rho} := F^{-1}(\bar{F}\rho)$ .

The formulation of (J.2) suggests to pose the recovery problem as,

$$\begin{aligned} & \min_{\bar{\rho}, \tilde{X}} \left\| \hat{M} \odot \bar{X} - C_{F\rho} \right\|_F^2 \\ & \text{subject to} \quad \text{diag}(\tilde{X}) = 1, \quad \text{rank}(\tilde{X}) = 1, \\ & \quad \tilde{X}[1, 0] = 1, \quad \tilde{X} \succeq 0, \quad \bar{\rho}[0] = 1, \\ & \quad \bar{\rho}[k] = \overline{\bar{\rho}[-k]}, \quad \forall k. \end{aligned} \quad (\text{J.3})$$

The constraint  $\tilde{X}[1, 0] = 1$  follows the assumption that  $(Fx)[0] = (Fx)[1] = 1$ . While we can easily estimate  $(Fx)[0]$  and therefore fix it, the assumption of fixed  $(Fx)[1] = 1$  is more delicate. Recall that the solution for the MRA problem is always up to cyclic translation. In the Fourier domain, it means that the first entry of the Fourier transform of the signal is determined up to an arbitrary modulation by  $e^{2\pi i \ell / L}$  for some  $\ell \in \mathbb{Z}$ . If  $L \rightarrow \infty$ , this allows us to fix this coefficient arbitrarily.

Similarly to the well-known SDP relaxation of the Max-Cut problem [45], the non-convex problem (J.3) can be relaxed to a convex program by omitting the rank constraint as follows,

$$\begin{aligned} \min_{\tilde{\rho}, \tilde{X}} \quad & \left\| \hat{M}' \odot \overline{\tilde{X}} - C_{F\tilde{\rho}} \right\|_F^2 \\ \text{subject to} \quad & \text{diag}(\tilde{X}) = 1, \quad \tilde{X}[1, 0] = 1, \\ & \tilde{X} \succeq 0, \quad \tilde{\rho}[0] = 1, \quad \tilde{\rho}[k] = \overline{\tilde{\rho}[-k]}, \forall k. \end{aligned} \quad (\text{J.4})$$

This relaxation is convex and can be solved in polynomial time using off-the-shelf software, such as CVX [46].

The SDP relaxation (J.4) recovers the Fourier phases of the signal and the distribution exactly for  $N \rightarrow \infty$  and fixed noise level, since in this regime we can estimate the first two moments arbitrarily well.

**Theorem J.1.** *Assume that  $|Fx[k]| = 1$  for all  $k$  and that  $F\rho$  is non-vanishing. In addition, assume that  $(Fx)[0] = (Fx)[1] = 1$ . Then, if  $N \rightarrow \infty$  and  $\sigma$  is fixed, the solution of (J.4) is given by  $\tilde{X} = (Fx)(Fx)^*$  and  $\tilde{\rho} = F\hat{\rho}$ .*

*Proof.* Since  $\sigma$  is fixed and  $N \rightarrow \infty$ , one can estimate  $M'$  as in (J.1) exactly. Then, since (J.4) admits at least one solution (the underlying signal and distribution), the objective is zero at the solution and we get the relation:

$$C_{\tilde{\rho}} = M' \odot \overline{\tilde{X}} = C_{F\tilde{\rho}} \odot (Fx Fx^*) \odot \overline{\tilde{X}}, \quad (\text{J.5})$$

where we use  $\hat{\rho} := F^{-1}(\overline{F\tilde{\rho}})$ . Let  $u = \tilde{\rho}/F\hat{\rho}$ . Since  $\tilde{X} \succeq 0$  we conclude that  $C_u \succeq 0$  and hence  $Fu \geq 0$  (the Fourier transform of  $u$  is non-negative). By the constraints of (J.4), we also have  $u[0] = 1$ . By examining the  $(1, 0)$ th entry of (J.5), we also conclude that

$$(Fx)[1] \overline{(Fx)[0]} (F\hat{\rho})[1] \overline{\tilde{X}[1, 0]} = \tilde{\rho}[1] \Rightarrow u[1] = \overline{\tilde{X}[1, 0]} = 1,$$

where the last equality holds because of the constraints of (J.4).

Until now, we have shown that the vector  $u$  satisfies  $u[0] = u[1] = 1$ , it is conjugate-symmetric and its Fourier transform is non-negative. Therefore, by Lemma IV.2 of [15], we conclude that  $u[n] = 1$  for all  $n$ , or  $\tilde{\rho} = F\hat{\rho}$ . Next, we substitute  $\tilde{\rho} = F\hat{\rho}$  in (J.5) and get

$$1 = (Fx Fx^*) \odot \overline{\tilde{X}},$$

where the equality holds entry-wise. Since all entries of  $\hat{x}$  are normalized, we conclude that  $\tilde{X} = (Fx)(Fx)^*$ . This concludes the proof.  $\square$