

# Factor Analysis for Spectral Estimation

Joakim Andén

Program in Applied and Computational Mathematics  
Princeton University  
Princeton, NJ, 08544  
Email: janden@math.princeton.edu

Amit Singer

Department of Mathematics and  
Program in Applied and Computational Mathematics  
Princeton University  
Princeton, NJ, 08544

**Abstract**—Power spectrum estimation is an important tool in many applications, such as the whitening of noise. The popular multitaper method enjoys significant success, but fails for short signals with few samples. We propose a statistical model where a signal is given by a random linear combination of fixed, yet unknown, stochastic sources. Given multiple such signals, we estimate the subspace spanned by the power spectra of these fixed sources. Projecting individual power spectrum estimates onto this subspace increases estimation accuracy. We provide accuracy guarantees for this method and demonstrate it on simulated and experimental data from cryo-electron microscopy.<sup>1</sup>

## I. INTRODUCTION

Estimating the power spectrum of a stationary field arises in many applications [1]–[4]. For example, noise statistics play an important role in maximum likelihood estimation of signal parameters. Hence, its power spectrum must be estimated.

The most basic estimator of the power spectrum is the periodogram, which is asymptotically unbiased but has very high variance. To remedy this, spectral smoothness constraints or parametric models are often imposed [5], [6]. The multitaper estimator [6]–[8] is particularly popular, which multiplies the data with fixed tapers, computes their periodograms, and averages. This reduces variance at the expense of smoothing the estimated power spectrum, potentially increasing bias.

In this work, we consider the problem of estimating the power spectra of multiple independent signals. Given identically distributed signals, an ensemble average of their periodograms or multitaper estimates provides a low-variance power spectrum estimate. However, in certain applications, signals are not identically distributed, but combine a small number of identically distributed stochastic sources. For example, a non-stationary signal may be divided up into approximately stationary parts, each of which is a linear combination of some fixed random signals. Alternatively, noise signals may arise from measurements subject to differing experimental conditions which are described by combining various noise sources. This is the case for cryo-electron microscopy (cryo-EM) images, where large molecules are frozen in a thin layer of ice and then imaged by exposing them to an electron beam and recording the transmitted electrons [9], [10]. Variations

in experimental conditions result in different noise characteristics. To whiten the projection images, accurate estimates of their noise power spectrum are therefore necessary.

While traditional methods such as multitaper estimates are applied in this context, an estimator which takes into account the low-dimensional variability of the noise distributions could yield improved accuracy. We propose a factor analysis model in which the  $d$ -dimensional field  $\mathbf{X} : \mathbb{Z}^d \rightarrow \mathbb{R}$  is a linear combination of  $r$  fixed fields

$$\mathbf{X}[\vec{v}] := \sum_{l=1}^r a_l \mathbf{Z}_l[\vec{v}] \quad \vec{v} \in \mathbb{Z}^d, \quad (1)$$

where  $a_1, \dots, a_r$  are independent random variables and  $\mathbf{Z}_1, \dots, \mathbf{Z}_r$  are independent linear random fields. Typical values for  $d$  include  $d = 1$  for stationary processes and  $d = 2$  for random images.

It follows that the power spectrum of  $\mathbf{X}$  conditioned on  $a_1, \dots, a_r$  is a linear combination of the power spectra of  $\mathbf{Z}_1, \dots, \mathbf{Z}_r$ . That is, the conditional power spectra reside in a low-dimensional subspace and are therefore determined by a small number of factors. We introduce a method for estimating this subspace and these factors given multiple realizations of  $\mathbf{X}$ . The method allows for estimation of the number of factors, which is small for most applications. Linear projection of individual power spectrum estimates onto the factor subspace then yields improved accuracy.

Section II discusses the power spectrum estimation problem for a single signal and describes the basic periodogram and multitaper methods. The factor analysis model and associated power spectrum estimation method are introduced and analyzed in Section III. Finally, Section IV presents numerical results on simulations and experimental data taken from cryo-EM images. Software implementing our proposed method and reproducing the figures of this paper is available at <http://github.com/janden/fase/>.

## II. SINGLE POWER SPECTRUM ESTIMATION

Given a zero-mean stationary random field  $\mathbf{Y} : \mathbb{Z}^d \rightarrow \mathbb{R}$  with finite second moments, we define its autocovariance as

$$R_{\mathbf{Y}}[\vec{j}] := \mathbb{E}[\mathbf{Y}[\vec{v}]\mathbf{Y}[\vec{v} + \vec{j}]] \quad \vec{j} \in \mathbb{Z}^d. \quad (2)$$

Its Fourier transform is the power spectrum of  $\mathbf{Y}$ :

$$P_{\mathbf{Y}}(\vec{\xi}) := \sum_{\vec{v} \in \mathbb{Z}^d} R_{\mathbf{Y}}[\vec{v}] e^{-2\pi i \langle \vec{\xi}, \vec{v} \rangle} \quad \vec{\xi} \in [-1/2, 1/2]^d. \quad (3)$$

<sup>1</sup>The authors were partially supported by Award Number R01GM090200 from the NIGMS, FA9550-12-1-0317 from AFOSR, Simons Investigator Award and Simons Collaboration on Algorithms and Geometry from Simons Foundation, and the Moore Foundation Data-Driven Discovery Investigator Award.

The power spectrum  $P_{\mathbf{Y}}$  is symmetric and non-negative.

While  $\mathbf{Y}$  is defined over  $\mathbb{Z}^d$ , we are only given samples over some finite domain  $M_N^d$ , where  $M_N = \{-\lfloor N/2 \rfloor + 1, \dots, \lfloor N/2 \rfloor\}$ . Let us denote this restriction of  $\mathbf{Y}$  to  $M_N^d$  by  $\mathbf{Y}^{(N)}$ . Our task is then to estimate  $P_{\mathbf{Y}}$  given  $\mathbf{Y}^{(N)}$ .

### A. Periodogram

The most basic estimator for the power spectrum of a stationary field is the periodogram, given by

$$\widehat{P}_{\mathbf{Y}^{(N)}}[\vec{k}] = \frac{1}{N^d} \left| \sum_{\vec{i} \in M_N^d} \mathbf{Y}^{(N)}[\vec{i}] e^{-2\pi i \langle \vec{k}, \vec{i} \rangle / N} \right|^2 \quad \vec{k} \in M_N^d. \quad (4)$$

Since  $\mathbf{Y}^{(N)}$  is real,  $\widehat{P}_{\mathbf{Y}^{(N)}}$  is symmetric and only needs to be computed on half of  $M_N^d$ . Let us therefore choose a subdomain  $M_{N,+}^d$  of  $M_N^d$  such that  $M_{N,+}^d \cup -M_{N,+}^d = M_N^d \pmod{N}$  and  $M_{N,+}^d \cap -M_{N,+}^d = \{\vec{0}\}$ .

To study the properties of the periodogram, we first define a linear stationary field.

**Definition 1.** A stationary field  $\mathbf{Y} : \mathbb{Z}^d \rightarrow \mathbb{R}$  is linear if

$$\mathbf{Y}[\vec{i}] = \sum_{\vec{j} \in \mathbb{Z}^d} W[\vec{i}] \psi[\vec{i} - \vec{j}], \quad (5)$$

where  $W : \mathbb{Z}^d \rightarrow \mathbb{R}$  is a stationary, zero-mean, white noise field such that  $\mathbb{E}|W[\vec{i}]|^2 = 1$ ,  $\mathbb{E}|W[\vec{i}]|^4 < \infty$  and  $\psi : \mathbb{Z}^d \rightarrow \mathbb{R}$  is a kernel satisfying  $\sum_{\vec{i} \in \mathbb{Z}^d} |\psi[\vec{i}]| \|\vec{i}\|_1^{1/2} < \infty$ .

For such a field, the periodogram is an asymptotically unbiased estimator for the power spectrum:

**Lemma 1.** Let  $\mathbf{Y} : \mathbb{Z}^d \rightarrow \mathbb{R}$  be a linear field and let  $\mathbf{Y}^{(N)}$  be its restriction  $M_N^d$ . Then we have, for  $\vec{k} \in M_{N,+}^d$ ,

$$\mathbb{E} \widehat{P}_{\mathbf{Y}^{(N)}}[\vec{k}] = P_{\mathbf{Y}}(\vec{k}/N) + O(N^{-1/2}). \quad (6)$$

The proof is obtained as a generalization of the case  $d = 1$  found in Proposition 10.3.1 of [5]. Recasting the periodogram as the Fourier transform of sums of pairs  $\mathbf{Y}^{(N)}[\vec{i}] \mathbf{Y}^{(N)}[\vec{i} + \vec{j}]$ , the bias is due to not weighting these sums by the number of available pairs. Reweighting eliminates the bias, but at a significant increase in variance, increasing the mean squared error of the estimator. For this reason, the biased periodogram estimator is often preferred over its unbiased variant [6].

Despite its low bias for large  $N$ , the periodogram is not suitable for power spectrum estimation due its high variance:

**Lemma 2.** Let  $\mathbf{Y}$  and  $\mathbf{Y}^{(N)}$  be defined as in Lemma 1 and

$$\delta_N[\vec{k}_1, \vec{k}_2] = \begin{cases} 2 & \vec{k}_1 = \vec{k}_2 \text{ and } \vec{k}_1 \in \frac{N}{2} \mathbb{Z}^d \\ 1 & \vec{k}_1 = \vec{k}_2 \text{ and } \vec{k}_1 \notin \frac{N}{2} \mathbb{Z}^d \\ 0 & \vec{k}_1 \neq \vec{k}_2 \end{cases}. \quad (7)$$

We then have, for all  $\vec{k}_1, \vec{k}_2 \in M_{N,+}^d$ ,

$$\begin{aligned} \text{Cov} \left[ \widehat{P}_{\mathbf{Y}^{(N)}}[\vec{k}_1], \widehat{P}_{\mathbf{Y}^{(N)}}[\vec{k}_2] \right] \\ = \delta_N[\vec{k}_1, \vec{k}_2] P_{\mathbf{Y}}(\vec{k}_1/N)^2 + O(N^{-1/2}). \end{aligned} \quad (8)$$

The proof again follows the one-dimensional case, Proposition 10.3.2 in [5]. It relies on the fact that as  $N \rightarrow \infty$ ,  $\widehat{P}_{\mathbf{Y}^{(N)}}[\vec{k}]$  converges to a  $\chi^2$ -distributed random variable. Asymptotically, the standard deviation of  $\widehat{P}_{\mathbf{Y}^{(N)}}[\vec{k}]$  is therefore approximately proportional to its expectation  $P_{\mathbf{Y}}(\vec{k}/N)$ .

### B. The Multitaper Method

The multitaper method for spectral estimation was introduced by D. J. Thomson as an alternative to the periodogram with reduced variance at the cost of increased bias [7]. It relies on multiplying the sample  $\mathbf{Y}^{(N)}$  by a number of data tapers, computing their periodograms, and averaging.

Given a target spectral resolution  $\frac{1}{2N} < W < \frac{1}{2}$ , the data tapers for  $d = 1$  are given by  $K = \lfloor 2NW \rfloor$  discrete prolate spheroidal sequences  $v_{N,W}^{(l)}[i]$  for  $i \in M_N$ , where  $l \in 1, \dots, K$ . Abusing notation slightly, we denote their tensor products by

$$v_{N,W}^{(\vec{l})}[\vec{i}] = \prod_{p=1}^d v_{N,W}^{(l_p)}[i_p] \quad \vec{i} = (i_1, \dots, i_d) \in M_N^d, \quad (9)$$

where  $\vec{l} = (l_1, \dots, l_d) \in \{1, \dots, K\}^d$  [11]. The associated multitaper estimator is then given by

$$\widehat{P}_{\mathbf{Y}^{(N)}}^{(\text{MT})}[\vec{k}] = \frac{1}{K^d} \sum_{\vec{l} \in \{1, \dots, K\}^d} \widehat{P}_{v_{N,W}^{(\vec{l})} \cdot \mathbf{Y}^{(N)}}[\vec{k}], \quad (10)$$

where  $(v_{N,W}^{(\vec{l})} \cdot \mathbf{Y}^{(N)})[\vec{i}] = v_{N,W}^{(\vec{l})}[\vec{i}] \mathbf{Y}^{(N)}[\vec{i}]$  for all  $\vec{i} \in M_N^d$ .

The resolution parameter  $W$  specifies the bias-variance trade-off of the estimator. As  $W$  is increased, variance is reduced, but the estimate is smoothed, potentially increasing the bias of the estimator.

## III. POWER SPECTRUM FACTOR ANALYSIS

In certain applications, we do not have just one signal whose power spectrum we would like to estimate, but a set of signals. To estimate their power spectra, we can apply the methods described in the previous section to each signal individually. If in addition the signals are identically distributed, we can reduce the variance in the estimate by an ensemble average.

On the other hand, if the signals are not identically distributed, an ensemble average will destroy the variability in the power spectra. A different approach is applicable if the signals can be written as linear combinations of a small set of fixed random fields. In this case, factor analysis allows us to identify the subspace spanned by the power spectra of these fixed fields, allowing for increased accuracy by projecting multitaper estimates onto this space.

### A. Motivation: Cryo-Electron Microscopy

Single-particle cryo-EM images are typically very noisy, with signal-to-noise ratios typically below 1/10 [9], [10]. Many reconstruction algorithms assume that this noise is close to white, but this is rarely the case. As a result, a whitening step is typically performed prior to reconstruction, in which the noise power spectrum is estimated and used to compute a whitening filter that is then applied. Power spectrum estimation

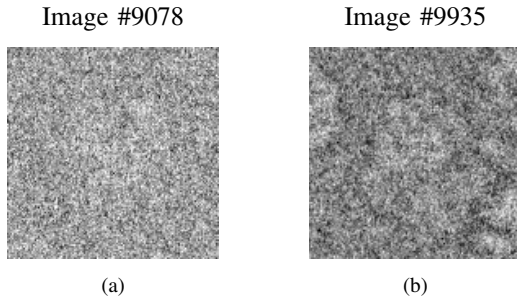


Fig. 1. Two sample cryo-electron microscopy images with noise power spectra (a) close to white noise and (b) mostly concentrated in the lower frequencies. Both images are taken from the same experimental dataset [14].

is typically done by extracting image patches containing mostly noise and averaging their periodograms [12]. Other methods estimate the noise power spectrum as part of the reconstruction algorithm [13].

However, the noise distribution typically varies between projection images. Indeed, different projection images are obtained under different experimental conditions (ice thickness, electron dose, microscope defocus, electron scattering properties of the molecules, and so on) which affect the characteristics of the noise [15]–[17]. The variation in the noise signal can be modeled as a random linear combination of fixed noise sources due to background electrons, electrons scattered by the ice, and those inelastically scattered by the molecule. Experimental images with different noise power spectra are shown in Figure 1. As we will see in the following, the power spectra of these noise signals can be accurately estimated using the proposed factor analysis.

### B. Problem Setup

Let us denote the zero-mean stationary random field of interest by  $\mathbf{X} : \mathbb{Z}^d \rightarrow \mathbb{R}$ . We model it in (1) as a linear combination of  $r$  fixed independent linear fields  $\mathbf{Z}_1, \dots, \mathbf{Z}_r$ , with independent random coefficients  $a_1, \dots, a_r$ .

Instead of a single field  $\mathbf{X}$ , we are concerned with a set of  $n$  independent copies  $\mathbf{X}_1, \dots, \mathbf{X}_n$  of  $\mathbf{X}$  defined by

$$\mathbf{X}_s[\vec{v}] := \sum_{l=1}^r a_{s,l} \mathbf{Z}_{s,l}[\vec{v}], \quad (11)$$

for  $s = 1, \dots, n$ . Here  $\mathbf{Z}_{1,l}, \dots, \mathbf{Z}_{n,l}$  and  $a_{1,l}, \dots, a_{n,l}$  are independent and identically distributed copies of  $\mathbf{Z}_l$  and  $a_l$ , respectively, for  $l = 1, \dots, r$ . The given data consists of restrictions  $\mathbf{X}_1^{(N)}, \dots, \mathbf{X}_n^{(N)}$  of  $\mathbf{X}_1, \dots, \mathbf{X}_n$  to  $M_N^d$ .

The coefficient vectors  $\vec{a}_1, \dots, \vec{a}_n$  are unknown, but fixed for each signal. As a result, the power spectrum of interest for  $\mathbf{X}_s$  is the conditional power spectrum  $P_{\mathbf{X}_s|\vec{a}_s}$ , where we have replaced the expectation  $\mathbb{E}$  in (2) by the conditional expectation  $\mathbb{E}_{\mathbf{Z}_{s,1}, \dots, \mathbf{Z}_{s,r}|\vec{a}_s}$ . Our problem is therefore estimating  $P_{\mathbf{X}_1|\vec{a}_1}, \dots, P_{\mathbf{X}_n|\vec{a}_n}$  given  $\mathbf{X}_1^{(N)}, \dots, \mathbf{X}_n^{(N)}$ .

One approach is to use the multitaper method described in Section II-B applied to each signal  $\mathbf{X}_1^{(N)}, \dots, \mathbf{X}_n^{(N)}$ . However,

we can obtain improved estimates by exploiting the fact that, for  $s = 1, \dots, n$ ,

$$P_{\mathbf{X}_s|\vec{a}_s}(\vec{\xi}) = \sum_{l=1}^r a_{s,l}^2 P_{\mathbf{Z}_l}(\vec{\xi}). \quad (12)$$

That is, the conditional power spectra are contained in a subspace of dimension at most  $r$  spanned by  $P_{\mathbf{Z}_1}, \dots, P_{\mathbf{Z}_r}$ . These power spectra can thus be described using a small set (at most  $r$ ) of common factors. We can exploit this to increase the accuracy of a power spectrum estimate by projecting it onto this low-dimensional subspace. In the following, we propose a method for estimating this subspace.

### C. Power Spectrum Covariance Estimation

The non-trivial eigenvectors of the covariance matrix  $\text{Cov}[P_{\mathbf{X}_s|\vec{a}_s}]$  span the subspace containing the conditional power spectra. To simplify our expressions, we define  $P_{\mathbf{X}|\vec{a}}$  in the same way as  $P_{\mathbf{X}_s|\vec{a}_s}$  for  $s = 1, \dots, n$  and note that these are all identically distributed. We consider estimating the covariance matrix  $\text{Cov}[P_{\mathbf{X}|\vec{a}}]$  using the following theorem:

**Theorem 1.** *Let  $\mathbf{X}$  and  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be as in (1) and (11), respectively. Furthermore, let  $\mathbf{X}_1^{(N)}, \dots, \mathbf{X}_n^{(N)}$  be the restrictions of  $\mathbf{X}_1, \dots, \mathbf{X}_n$  to  $M_N^d$ . In this case, the quantities*

$$\mu_n[\vec{k}] = \frac{1}{n} \sum_{s=1}^n \hat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}] \quad (13)$$

and

$$C_n[\vec{k}_1, \vec{k}_2] = \frac{1}{n} \sum_{s=1}^n \hat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}_1] \hat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}_2], \quad (14)$$

satisfy

$$\mu_n[\vec{k}] = P_{\mathbf{X}}(\vec{k}/N) + \epsilon_1[\vec{k}] \quad (15)$$

and

$$\begin{aligned} & \frac{1}{1 + \delta[\vec{k}_1, \vec{k}_2]} C_n[\vec{k}_1, \vec{k}_2] - \mu_n[\vec{k}_1] \mu_n[\vec{k}_2] \\ &= \text{Cov} \left[ P_{\mathbf{X}|\vec{a}}(\vec{k}_1/N), P_{\mathbf{X}|\vec{a}}(\vec{k}_2/N) \right] + \epsilon_2[\vec{k}_1, \vec{k}_2], \end{aligned} \quad (16)$$

where  $\vec{k}, \vec{k}_1, \vec{k}_2 \in M_{N,+}^d$ . The expected magnitudes  $\mathbb{E}|\epsilon_1[\vec{k}]|$  and  $\mathbb{E}|\epsilon_2[\vec{k}_1, \vec{k}_2]|$  are each bounded by  $O(n^{-1/2} + N^{-1/2})$ .

*Proof.* First, we condition Lemma 1 on  $\vec{a}$  to obtain

$$\mathbb{E}_{\mathbf{Z}_1, \dots, \mathbf{Z}_r|\vec{a}} \hat{P}_{\mathbf{X}^{(N)}}[\vec{k}] = P_{\mathbf{X}|\vec{a}}(\vec{k}/N) + O(N^{-1/2}), \quad (17)$$

for  $\vec{k} \in M_{N,+}^d$ . Taking expectation with respect to  $\vec{a}$  then gives

$$\mathbb{E} \hat{P}_{\mathbf{X}^{(N)}}[\vec{k}] = P_{\mathbf{X}}(\vec{k}/N) + O(N^{-1/2}). \quad (18)$$

Averaging  $\hat{P}_{\mathbf{X}_1^{(N)}}[\vec{k}], \dots, \hat{P}_{\mathbf{X}_n^{(N)}}[\vec{k}]$  then gives us,

$$\mu_n[\vec{k}] = \frac{1}{n} \sum_{s=1}^n \hat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}] = \mathbb{E} \hat{P}_{\mathbf{X}^{(N)}}[\vec{k}] + \epsilon_1[\vec{k}], \quad (19)$$

according to the central limit theorem, where  $\mathbb{E}|\epsilon_1[\vec{k}]| = O(n^{-1/2})$ . Plugging (18) into (19) then proves (15).

Conditioning Lemma 2 with respect to  $\vec{a}$  gives

$$\begin{aligned} & \text{Cov}_{\mathbf{Z}_1, \dots, \mathbf{Z}_r | \vec{a}} \left[ \widehat{P}_{\mathbf{X}^{(N)}}[\vec{k}_1], \widehat{P}_{\mathbf{X}^{(N)}}[\vec{k}_2] \right] \\ &= \delta[\vec{k}_1, \vec{k}_2] P_{\mathbf{X} | \vec{a}}(\vec{k}_1/N)^2 + O(N^{-1/2}). \end{aligned} \quad (20)$$

Consequently

$$\begin{aligned} & \mathbb{E}_{\mathbf{Z}_1, \dots, \mathbf{Z}_r | \vec{a}} \left[ \widehat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}_1] \widehat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}_2] \right] \\ &= (1 + \delta[\vec{k}_1, \vec{k}_2]) P_{\mathbf{X} | \vec{a}}(\vec{k}_1/N) P_{\mathbf{X} | \vec{a}}(\vec{k}_2/N) + O(N^{-1/2}). \end{aligned}$$

Taking the expectation with respect to  $\vec{a}$  now gives

$$\begin{aligned} & \mathbb{E} \left[ \widehat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}_1] \widehat{P}_{\mathbf{X}_s^{(N)}}[\vec{k}_2] \right] \\ &= (1 + \delta[\vec{k}_1, \vec{k}_2]) \mathbb{E} \left[ P_{\mathbf{X} | \vec{a}}(\vec{k}_1/N) P_{\mathbf{X} | \vec{a}}(\vec{k}_2/N) \right] + O(N^{-1/2}). \end{aligned}$$

We now have

$$\begin{aligned} & \frac{1}{1 + \delta[\vec{k}_1, \vec{k}_2]} \mathbb{E} \left[ \widehat{P}_{\mathbf{X}^{(N)}}[\vec{k}_1] \widehat{P}_{\mathbf{X}^{(N)}}[\vec{k}_2] \right] - P_{\mathbf{X}}(\vec{k}_1/N) P_{\mathbf{X}}(\vec{k}_2/N) \\ &= \text{Cov} \left[ P_{\mathbf{X} | \vec{a}}(\vec{k}_1/N), P_{\mathbf{X} | \vec{a}}(\vec{k}_2/N) \right] + O(N^{-1/2}). \end{aligned} \quad (21)$$

Replacing the left-hand side expectations with sums using the central limit theorem yields (16).  $\square$

The theorem allows us to estimate the covariance matrix of the conditional power spectra. Traditionally, this would be done by forming  $C_n[\vec{k}_1, \vec{k}_2] - \mu_n[\vec{k}_1] \mu_n[\vec{k}_2]$ . However, due to the  $\chi^2$  distribution of the periodogram coordinates, the variances of  $C_n$  are larger than they should be, so a correction is needed. That correction takes the form of  $(1 + \delta[\vec{k}_1, \vec{k}_2])^{-1}$ .

Let us denote this covariance matrix estimate by

$$\Sigma_n[\vec{k}_1, \vec{k}_2] = \frac{1}{1 + \delta[\vec{k}_1, \vec{k}_2]} C_n[\vec{k}_1, \vec{k}_2] - \mu_n[\vec{k}_1] \mu_n[\vec{k}_2], \quad (22)$$

for  $\vec{k}_1, \vec{k}_2 \in M_{N,+}^d$ . Calculating its top eigenvectors allows us to estimate the low-dimensional subspace containing the conditional power spectrum  $P_{\mathbf{X} | \vec{a}}$ . Specifically, these eigenvectors define a linear space that approximates the subspace spanned by  $P_{\mathbf{Z}_1}, \dots, P_{\mathbf{Z}_r}$ .

#### D. Linear Projection Estimators

We now use the information gained in the previous section to obtain a better estimate of the conditional power spectrum  $P_{\mathbf{X} | \vec{a}}$  from the multitaper estimate  $\widehat{P}_{\mathbf{X}^{(N)}}^{(\text{MT})}$ . It satisfies

$$\widehat{P}_{\mathbf{X}^{(N)}}^{(\text{MT})} = P_{\mathbf{X} | \vec{a}} + R^{(\text{MT})}, \quad (23)$$

where the size of the residual  $R^{(\text{MT})}$  depends on  $N$  and the regularity of  $P_{\mathbf{X} | \vec{a}}$  with respect to the target multitaper resolution  $W$  [6]. We know that  $P_{\mathbf{X} | \vec{a}}$  is in the subspace spanned by the non-trivial eigenvectors of  $\text{Cov} [P_{\mathbf{X} | \vec{a}}]$ . We can use this to reduce the variance of the residual term.

While the exact covariance is not known to us, we have an estimator  $\Sigma_n$ . Let  $v_1, \dots, v_r$  be the leading eigenvectors of  $\Sigma_n$ . We note that  $r$  does not need to be known beforehand, but can be estimated from the eigenspectrum of  $\Sigma_n$ . Here, we set  $r$  to be the number of dominant eigenvalues of  $\Sigma_n$  by locating the ‘‘knee’’ of the spectrum, where there is a significant drop

in amplitude from the  $r$ th eigenvalue to the  $(r + 1)$ th. In addition, since we expect  $\mu_n$  to be contained in the span of  $v_1, \dots, v_r$ , we can also check that the projection of  $\mu_n$  to that span preserves most of its energy.

Projecting  $\widehat{P}_{\mathbf{X}^{(N)}}^{(\text{MT})}$  onto the span of  $v_1, \dots, v_r$  yields

$$\widehat{P}_{\mathbf{X}^{(N)}}^{(\text{MT}, \text{lin})}[\vec{k}] = \sum_{l=1}^r v_l[\vec{k}] \langle v_l, \widehat{P}_{\mathbf{X}^{(N)}}^{(\text{MT})} \rangle. \quad (24)$$

This significantly reduces the estimation error  $R$  since it will typically be almost orthogonal to  $v_l$  and thus  $\langle v_l, R \rangle \approx 0$  for  $l = 1, \dots, r$ . The orthogonality follows from the fact that  $R$  is approximately isotropically distributed (its coordinates are close to independent random variables), so its inner product with a small set of fixed vectors is small with high probability. Note that  $\widehat{P}^{(\text{MT}, \text{lin})}$  is not guaranteed to be non-negative. To remedy this, negative components can be set to zero. However, this does not greatly affect the error of the estimate.

## IV. NUMERICAL RESULTS

To evaluate our method, we test it on noisy images simulated using (1) for  $r = 2$ . The first noise source  $\mathbf{Z}_1$  has its energy concentrated in the low frequencies, with a power spectrum of  $P_{\mathbf{Z}_1}(\xi) = 2\text{rect}(4\|\xi\|)$ , where  $\text{rect}$  is the rectangular function with support  $[-1/2, 1/2]$ . The second source  $\mathbf{Z}_2$  has a power spectrum of  $P_{\mathbf{Z}_2}(\xi) = 1/(1 + 4\|\xi\|)$ . These are combined with normally distributed coefficients  $a_1, a_2 \sim N(0, 1)$  in (1) to yield  $\mathbf{X}$ . Figure 2(a) shows two sample images of size  $N = 32$ .

We calculate noise images of different sizes  $N$  to study the effect of this parameter on the accuracy of the estimation. Figure 2(b) shows the top eigenvalues of  $\Sigma_n$  for a dataset of size  $n = 1024$  with image size  $N = 32$ . There is a good separation of the top two eigenvalues from the bulk of the spectrum, with a spectral gap of  $\lambda_2/\lambda_3$  of about 3.4.

Figure 2(c,d) shows the mean absolute error (MAE) of  $\widehat{P}_{\mathbf{X}^{(N)}}^{(\text{MT})}, \dots, \widehat{P}_{\mathbf{X}_n^{(N)}}^{(\text{MT})}$  and  $\widehat{P}_{\mathbf{X}_1^{(N)}}^{(\text{MT}, \text{lin})}, \dots, \widehat{P}_{\mathbf{X}_n^{(N)}}^{(\text{MT}, \text{lin})}$  with respect to the conditional power spectra  $P_{\mathbf{X}_1 | \vec{a}_1}, \dots, P_{\mathbf{X}_n | \vec{a}_n}$  as functions of  $n$  for image sizes  $N = 32$  and  $N = 64$ , respectively. The unprojected multitaper estimator was computed with  $W = 1/16$  while the others used  $W = 1/64$ . For comparison, we plot the baseline estimator of assigning  $\mu_n$  to each power spectrum and the oracle projection of the multitaper estimate, where the top  $r$  eigenvectors of the population covariance  $\text{Cov}[P_{\mathbf{X} | \vec{a}}]$  are used in the definition of  $\widehat{P}^{(\text{MT}, \text{lin})}$ . For small values of  $n$ , the covariance estimation fails, so the projected estimator  $\widehat{P}^{(\text{MT}, \text{lin})}$  performs worse compared to the unprojected variant  $\widehat{P}^{(\text{MT})}$ . As  $n$  increases, however, the linear projection improves results, eventually converging to the error obtained using the subspace obtained from the population covariance.

We also evaluate our algorithm using two experimental datasets from cryo-EM, one containing  $n = 10000$  images of size  $N = 130$  depicting a 70S ribosome [14] and another containing  $n = 105247$  images with  $N = 360$  depicting an 80S ribosome [18]. The top eigenvalues of  $\Sigma_n$  are shown for the two datasets in Figure 3(a,b). As in the simulation, we see that a few eigenvalues dominate, with the estimation making

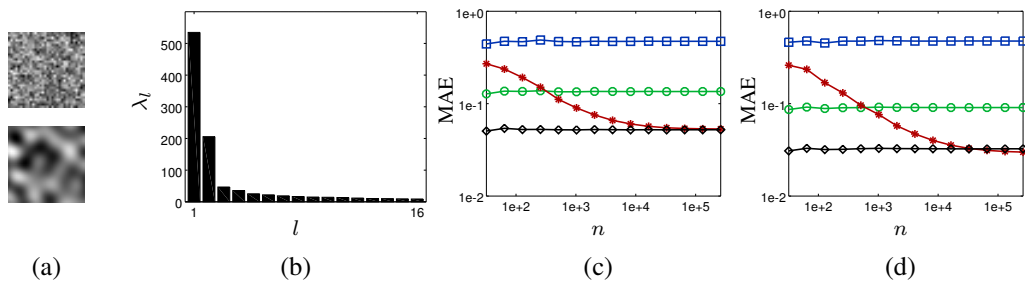


Fig. 2. (a) Two sample noise images for  $N = 32$  and  $r = 2$ . (b) The top 16 eigenvalues of  $\Sigma_n$  for  $n = 1024$  images of size  $N = 32$ . The estimation error of the conditional power spectra as a function of  $n$  for images of size (c)  $N = 32$  and (d)  $N = 128$ . We have the ensemble average multitaper estimator (blue square), the unprojected multitaper (green circle), the projected multitaper (red star), and the oracle projected multitaper (black diamond).

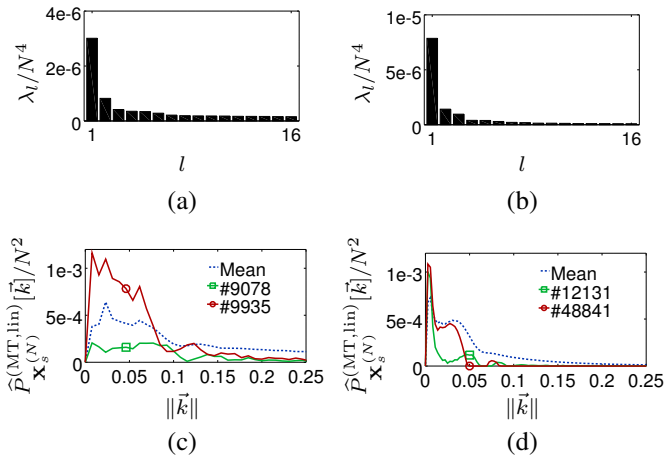


Fig. 3. Results experimental datasets. The top 16 eigenvalues of  $\Sigma_n$  estimated from a dataset depicting a (a) 70S ribosome from *E. Coli* where  $N = 130$  and  $n = 10000$  [14] and (b) an 80S ribosome from *P. Falciparum* where  $N = 360$  and  $n = 105247$  [18]. Sample conditional power spectrum estimates from the (c) 70S ribosome and (d) 80S ribosome datasets.

up a separate bulk distribution. For the first dataset, a value of  $r = 2$  seems appropriate, while for the second dataset, we take  $r = 3$ . This number of components is in line with most noise models for cryo-EM, which typically include two or three noise components [15]–[17].

Two sample conditional power spectra estimates for each dataset are shown in Figure 3(c,d). The projection images corresponding to the two power spectra in 3(c) are the ones shown previously in Figure 1. We see that the estimated power spectra are quite reasonable, with the lower-frequency noise image corresponding to the power spectrum concentrated in the low frequencies while the image with “whiter” noise has a flatter estimated power spectrum.

## V. CONCLUSION

We have introduced a factor analysis model for noise fields. These arise when the nature of the noise varies between measurements such as in cryo-EM, where the noise is affected by various experimental factors. To estimate the individual power

spectra of each noise image we introduce a new estimation method which relies on approximating the covariance of these spectra. The method is shown to provide accurate results in both simulated and experimental datasets.

## ACKNOWLEDGMENT

The authors thank Fred Sigworth for discussions about cryo-EM noise models, which inspired this work. They also thank Frederik Simons for his helpful comments.

## REFERENCES

- [1] P. Stoica and R. L. Moses, *Introduction to Spectral Analysis*. Prentice Hall, 1997, vol. 1.
- [2] G. Bond *et al.*, “A pervasive millennial-scale cycle in North Atlantic Holocene and glacial climates,” *Science*, 278(5341): 1257–1266 (1997).
- [3] C. Harig and F. J. Simons, “Mapping Greenland’s mass loss in space and time,” *Proc. Natl. Acad. Sci.*, 109(49): 19934–19937 (2012).
- [4] A. Delorme and S. Makeig, “EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis,” *J. Neurosci. Meth.*, 134(1): 9–21 (2004).
- [5] P. J. Brockwell and R. A. Davis, *Time Series: Theory and Methods*, 2nd ed. Springer-Verlag New York, 1991.
- [6] D. B. Percival and A. T. Walden, *Spectral Analysis for Physical Applications*. Cambridge University Press, 1993.
- [7] D. J. Thomson, “Spectrum estimation and harmonic analysis,” *Proceedings of the IEEE*, 70(9): 1055–1096 (1982).
- [8] L. D. Abreu and J. L. Romero, “MSE estimates for multitaper spectral estimates and off-grid compressive sensing,” *arXiv:1703.08190*, 2017.
- [9] J. Frank, *Three-dimensional electron microscopy of macromolecular assemblies*. Academic Press, 1996.
- [10] E. Callaway, “The revolution will not be crystallized: a new method sweeps through structural biology,” *Nature*, 525(7568): 172–174 (2015).
- [11] A. Hanssen, “Multidimensional multitaper spectral estimation,” *Signal Processing*, 58(3): 327–332 (1997).
- [12] T. Bhamre, T. Zhang *et al.*, “Denosing and covariance estimation of single particle cryo-EM images,” *J. Struct. Biol.*, 195(1): 72–81 (2016).
- [13] S. H. W. Scheres, “RELION: Implementation of a Bayesian approach to Cryo-EM structure determination,” *J. Struct. Biol.*, 180(3): 519–520 (2012).
- [14] H. Liao and J. Frank, “Classification by bootstrapping in single particle methods,” in *Proc. ISBI*. IEEE, 169–172 (2010).
- [15] W. T. Baxter, R. A. Grassucci, H. Gao, and J. Frank, “Determination of signal-to-noise ratios and spectral snrs in cryo-em low-dose imaging of molecules,” *J. Struct. Biol.*, 166(2): 126–132 (2009).
- [16] P. A. Penczek, “Chapter two: Image restoration in cryo-electron microscopy,” *Methods Enzymol.*, 482: 35–72 (2010).
- [17] M. Vulović, R. B. Ravelli *et al.*, “Image formation modeling in cryo-electron microscopy,” *J. Struct. Biol.*, 183(1): 19–32 (2013).
- [18] W. Wong, X. Bai *et al.*, “Cryo-EM structure of the plasmodium falciparum 80s ribosome bound to the anti-protozoan drug emetine,” *Elife*, 3: e03080 (2014).