

---

# A random matrix analysis of online learning: coping with limited memory resources

---

Anonymous Authors<sup>1</sup>

## Abstract

This article introduces a random matrix framework for the analysis of online learning, a particularly relevant setting for a more sober processing of large amounts of data with limited memory and energy resources. Assuming data  $\mathbf{x}_1, \mathbf{x}_2, \dots$  arrives as a continuous flow and a small number  $L$  of them can be kept in the learning pipeline, one has only access to the diagonal elements of the Gram kernel matrix:  $[\mathbf{K}_L]_{i,j} = \frac{1}{p} \mathbf{x}_i^\top \mathbf{x}_j \mathbf{1}_{|i-j| < L}$ . Under a large dimensional data regime, we derive the limiting spectral distribution of the punctured kernel matrix  $\mathbf{K}_L$  and study its isolated eigenvalues and eigenvectors, which behave in an unfamiliar way. We detail how these results can be used to perform efficient online kernel spectral clustering and provide theoretical performance guarantees. Our findings are empirically confirmed on image classification tasks. Leveraging on optimality results of spectral methods for clustering, this work offers insights on efficient online clustering techniques for high-dimensional data.

## 1. Introduction

The ever-increasing amount of data coupled with the need for a more sober use of computational power puts online learning in the spotlight, as a way to deal with numerous and very large data with low memory resources. Be it because the volume of data is too high to be stored or because one is restricted to the sole use of a regular laptop, online learning appears as a handy and frugal way to process information. As data arrives in the learning pipeline, it is processed at a low computational cost before being discarded altogether, thus inducing a limited memory footprint.

Numerous works have proposed various algorithms to clus-

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

ter data streams in an unsupervised manner (see, e.g., (Ghesmoune et al., 2016) and references therein). Among standard methods are the construction of a graph (Fritzke, 1995) or a tree of clusters (Zhang et al., 1996) which is updated as new data arrives, or else, the formation of clusters using a distance function, as in  $k$ -means, (Aggarwal et al., 2003) or a density-based method (Ester et al., 1996). Such algorithms are often adaptations of existing offline algorithms, like OpticsStream (Tasoulis et al., 2007), StreamKM++ (Ackermann et al., 2012), online  $k$ -means (Liberty et al., 2015), etc. These techniques operate on the entire feature space and their performance deteriorate as the dimension of the data increases. Therefore, (Aggarwal et al., 2004) proposed to cluster data streams after a projection on a lower-dimensional space. Sketching methods (Keriven et al., 2017; Gribonval et al., 2021) are also convenient to perform large-scale learning on data streams with a limited memory budget; the idea being to summarize the dataset into a single vector computed in one pass over the data.

Adapted from the standard spectral clustering algorithm (von Luxburg, 2007), techniques like incremental spectral clustering (Ning et al., 2010; Dhanjal et al., 2014) have been proposed to handle evolving data. Yet, they become quite memory-demanding when the number of samples grow large. Better suited to streaming applications, the spectral clustering algorithm of (Yoo et al., 2016) constructs a spectral embedding of the stream in one pass by adapting ideas from matrix sketching (Liberty, 2012).

Spectral clustering has indeed remarkably good performances on high-dimensional data as it manages to greatly reduce the dimensionality by keeping just a few leading spectral components. It is therefore computationally less demanding than many other classical clustering algorithms. Moreover, it reaches the optimal phase transition threshold (i.e., it performs better than random guess as soon as theoretically possible) (Onatski et al., 2013) and achieves the optimal clustering error rate in the Gaussian mixture model (Löffler et al., 2020).

It is also of particular interest from a random matrix theory perspective. Following the works of (El Karoui, 2010; Cheng & Singer, 2012) on the spectrum of kernel random matrices, (Couillet & Benaych-Georges, 2016) propose an

analysis of kernel spectral clustering with numerous high-dimensional data. Then, (Mai & Couillet, 2017) demonstrate that many standard machine learning algorithms in fact suffer from being ill-used when dealing with such data. Besides, given some data matrix  $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n] \in \mathbb{R}^{p \times n}$ , (Couillet et al., 2021) show that it is possible to get huge reductions in computational and storage costs with almost no performance loss by puncturing the data, i.e., keeping only a few elements of  $\mathbf{X}$  and computing only a few elements of the Gram kernel matrix  $\mathbf{K} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$ . In addition, (Liao et al., 2020) demonstrate that, when carefully employed, sparsification and quantization of  $\mathbf{K}$  incur negligible performance loss, while providing a great computational gain.

In the light of these numerous benefits of spectral clustering when dealing with high-dimensional data, of the practicality of online learning to handle large data streams with limited memory, and of the promising path shown by random matrix theory towards resource-efficient learning with performance guarantees, the present work introduces an “online spectral learning” algorithm to which we attach a rigorous performance analysis using random matrix theory.

The algorithm goes as follows: supposing that, due to memory limitations, only a small number  $L$  of data points can be kept in the pipeline, the computation of the  $n \times n$  Gram kernel matrix is limited to the elements which are in a radius  $L$  around the diagonal of  $\mathbf{K}$ . This results in the following punctured kernel matrix model

$$\mathbf{K}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{T}$$

where  $\odot$  denotes the Hadamard product and  $\mathbf{T} \in \{0, 1\}^{n \times n}$  is a Toeplitz mask:  $\mathbf{T}_{i,j} = \mathbf{1}_{|i-j| < L}$ . A careful adaption of spectral clustering is then performed on  $\mathbf{K}_L$  to retrieve the class information.

In technical terms, the present analysis derives the spectral distribution of  $\mathbf{K}_L$  and analyzes the behavior of a few isolated eigenvalues (called *spikes*) which carry information (that is, indicators for the data classes) in their associated eigenvectors. Two new interesting behaviors are observed: unlike classical spectral clustering, due to the Toeplitz filter, the number of informative spikes can potentially grow very large even in the case of binary classification. In addition, the eigenvectors are strongly tainted (in a way “convolved”) by the eigenvectors of the Toeplitz mask, which then requires some careful post-processing for classification. Our results particularly shed light on how the learning performance is altered by the dimension of the data and the size of the pipeline, thus providing an analysis of the performance versus cost trade-off of online learning.

In a nutshell, our main contributions may be listed as follows

- we derive the limiting eigenvalue distribution of  $\mathbf{K}_L$  as  $n, p, L \rightarrow +\infty$  for data arising from a Gaussian mixture model:  $\mathbf{x}_i \sim \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \mathbf{I}_p)$ ;
- for centered data drawn from a two-class mixture  $\mathbf{x}_i \sim \mathcal{N}(\pm \boldsymbol{\mu}, \mathbf{I}_p)$ , we show that a phase transition phenomenon occurs: depending on the signal power  $\|\boldsymbol{\mu}\|$ , some eigenvalues of  $\mathbf{K}_L$  isolate and their eigenvectors carry information about the classes;
- we propose an algorithm to retrieve information from isolated eigenvectors, thus performing high-dimensional “online spectral clustering”;
- simulations of online spectral clustering on Fashion-MNIST and BigGAN-generated images confirm the predicted good behavior of the algorithm and support our theoretical findings.

**Proofs and simulations** All proofs are deferred to the appendix. Python codes to reproduce simulations are available as supplementary material.

## 2. Online learning model and problem setting

### 2.1. General framework

Let  $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be a collection of  $n$  data samples of dimension  $p$ . They are noisy observations of  $K$  unknown classes whose centroids are  $[\boldsymbol{\mu}_1 \ \dots \ \boldsymbol{\mu}_K] \equiv \mathbf{M} \in \mathbb{R}^{p \times K}$ . Also define the  $n \times K$  binary matrix  $\mathbf{J}$  such that  $\mathbf{J}_{i,j} = 1$  if  $\mathbf{x}_i$  belongs to class  $j$  and 0 otherwise.

We make the following assumptions.

**Assumption 2.1.** The rows of  $\mathbf{J}$  are independent realizations of a multinomial distribution with one trial and  $K$  outcomes, i.e., the class of  $\mathbf{x}_i$  does not depend on the class of  $\{\mathbf{x}_j\}_{j \neq i}$ .

**Assumption 2.2** (Non-triviality condition).  $\mathbf{M}$  is uniformly bounded in spectral norm as  $n, p \rightarrow +\infty$ .

**Assumption 2.3.** The random matrix  $\mathbf{X}$  can be decomposed into a deterministic signal matrix  $\mathbf{P} = \mathbf{M}\mathbf{J}^\top$  and a random standard Gaussian noise matrix  $\mathbf{Z}$  with independent entries<sup>1</sup>:  $\mathbf{X} = \mathbf{P} + \mathbf{Z}$ .

*Remark 2.4.* The non-triviality condition (assumption 2.2) places the work under scenarios of practical relevance, in the sense that the problem is asymptotically (as  $n, p, L \rightarrow +\infty$ ) neither too easy nor too hard. The clustering error rate is therefore *not* expected to vanish asymptotically.

In the considered online setting, only the  $L$  previously seen data points are kept in memory. Thus, the element

<sup>1</sup>The “interpolation trick” from (Lytova & Pastur, 2009) enables to interpolate the results to non-Gaussian noise, but we keep the Gaussian assumption for simplicity of exposition here.

$\mathbf{K}_{i,j} = \frac{1}{p} \mathbf{x}_i^\top \mathbf{x}_j$  of the Gram kernel matrix can be computed only for  $|i - j| < L$ . This is represented by the pointwise application of a Toeplitz mask  $\mathbf{T} = (\mathbf{1}_{|i-j| < L})_{1 \leq i,j \leq n}$  resulting in

$$\mathbf{K}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{T} \quad \text{with} \quad \mathbf{T} = \begin{bmatrix} 1 & \dots & 1 & & 0 \\ & \ddots & & \ddots & \\ & & \ddots & & \\ 1 & & & \ddots & 1 \\ & \ddots & & & \\ 0 & 1 & \dots & 1 & \end{bmatrix}.$$

As standard (offline) spectral clustering is “optimal”<sup>2</sup>, we argue that spectral clustering on  $\mathbf{K}_L$  ought to achieve good performance at least for not too small  $(2L - 1)/n$  ratios. Our technical goal is thus to first provide a description of the spectral behavior of  $\mathbf{K}_L$  as  $n, p$  and  $L$  are large. To this end, we place ourselves under the regime  $n, p, L \rightarrow +\infty$  with  $p/n \rightarrow c \in ]0, +\infty[$  and  $(2L - 1)/n \rightarrow \varepsilon \in ]0, +\infty[$ .

## 2.2. The circulant approximation

An important trick to derive our main result lies in the fact that the Toeplitz matrix  $\mathbf{T}$  can be approximated (Gray, 2006) to some extent by its circulant “version”  $\mathbf{C} = (\mathbf{1}_{|i-j| < L} + \mathbf{1}_{|i-j| > n-L})_{1 \leq i,j \leq n}$ . Denoting  $\{\tau_k\}_{0 \leq k < n}$  and  $\{\psi_k\}_{0 \leq k < n}$  their respective eigenvalues (which depend on  $n$  and  $L$ ), then for fixed  $L$  and any continuous function  $f : \mathbb{R} \rightarrow \mathbb{R}$ ,

$$\lim_{n \rightarrow +\infty} \frac{1}{n} \sum_{k=0}^{n-1} |f(\psi_k) - f(\tau_k)| = 0.$$

*Remark 2.5.* Keep in mind that, in our case,  $n$  and  $L$  grow together at the same rate. Therefore, approximating  $\mathbf{T}$  by  $\mathbf{C}$  is only reasonable if  $\varepsilon$  is sufficiently small.

The core advantage of  $\mathbf{C}$  is that, unlike  $\mathbf{T}$ , its eigendecomposition is well-known:

$$\mathbf{C} = \mathbf{F} \boldsymbol{\Psi} \mathbf{F}^*$$

where  $\mathbf{F}$  is the  $n \times n$  Fourier matrix ( $\mathbf{F}_{i,j} = \frac{1}{\sqrt{n}} e^{-2i\pi \frac{ij}{n}}$ ) and  $\boldsymbol{\Psi} = \text{diag}(\psi_k)_{0 \leq k < n}$  is the diagonal matrix of eigenvalues. The latter are a sampling of the Dirichlet kernel:

$$\psi_k = \nu_L \left( \frac{2k\pi}{n} \right) \quad \text{with} \quad \nu_L(x) = \frac{\sin((2L-1)\frac{x}{2})}{\sin(\frac{x}{2})}.$$

In Figure 1 are superimposed to the graph of  $\nu_L$  the eigenvalues of  $\mathbf{C}$  and<sup>3</sup>  $\mathbf{T}$ . The  $\tau_k$ ’s roughly follow the graph of  $\nu_L$ , as if they were noisy versions of the  $\psi_k$ ’s.

<sup>2</sup>In that it performs better than random guess as soon as theoretically possible (Onatski et al., 2013).

<sup>3</sup>Although there is a natural order for the eigenvalues of  $\mathbf{C}$  given by  $\psi_k = \nu_L(\frac{2k\pi}{n})$ , we use a small trick to get the corre-

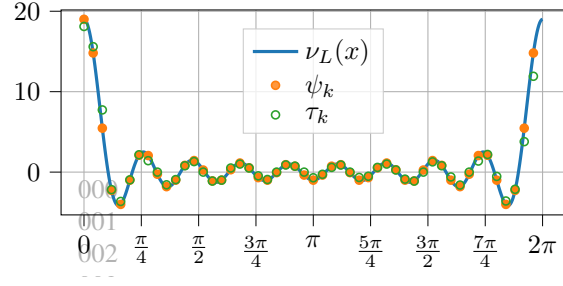


Figure 1. Graph of  $\nu_L$  on  $[0, 2\pi[$  (one period) with a plot of  $\psi_k = \nu_L(\frac{2k\pi}{n})$  and  $\tau_k$  for  $0 \leq k < n$  (the eigenvalues of  $\mathbf{C}$  and  $\mathbf{T}$  respectively). **Experimental setting:**  $n = 50, L = 10$ .

## 3. Main results

Following standard methods in random matrix theory (Couillet & Liao, 2021), the large dimensional spectral behavior of  $\mathbf{K}_L$  is accessible through an analysis of the resolvent matrix

$$\mathbf{Q}(z) = (\mathbf{K}_L - z\mathbf{I}_n)^{-1}$$

defined for all  $z \in \mathbb{C} \setminus \text{Sp}(\mathbf{K}_L)$ , where  $\text{Sp}(\mathbf{K}_L)$  denotes the set of eigenvalues of  $\mathbf{K}_L$ . Notably, the Stieltjes transform of the empirical spectral measure  $\mu_n = \frac{1}{n} \sum_{\lambda \in \text{Sp}(\mathbf{K}_L)} \delta_\lambda$  of  $\mathbf{K}_L$  (from which the spectral measure itself can be recovered) is the normalized trace of its resolvent:

$$m_n(z) \equiv \int_{\mathbb{R}} \frac{\mu_n(dt)}{t - z} = \frac{1}{n} \text{tr} \mathbf{Q}(z).$$

The resolvent also encapsulates information about the eigenvectors of  $\mathbf{K}_L$ : given a closed positively-oriented complex contour  $\Gamma$  circling around an eigenvalue  $\lambda$  of  $\mathbf{K}_L$  and leaving all the other eigenvalues outside,  $-\frac{1}{2i\pi} \oint_{\Gamma} \mathbf{Q}(z) dz = \mathbf{u}\mathbf{u}^*$ , where  $\mathbf{u}$  is a unit eigenvector associated to  $\lambda$ .<sup>4</sup>

### 3.1. Large dimensional spectral behavior

Our main theorem provides a deterministic equivalent of the resolvent when the Toeplitz mask  $\mathbf{T}$  is approximated by its circulant version  $\mathbf{C}$ , i.e.,  $\tilde{\mathbf{Q}}(z) = (\tilde{\mathbf{K}}_L - z\mathbf{I}_n)^{-1}$  with  $\tilde{\mathbf{K}}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{C}$ . Namely, we find a deterministic matrix  $\bar{\mathbf{Q}}(z)$  such that, for any sequence of deterministic matrices  $\mathbf{A}_n \in \mathbb{R}^{n \times n}$  and vectors  $\mathbf{a}_n, \mathbf{b}_n \in \mathbb{R}^n$  of unit norm (spectral norm and Euclidean norm respectively),  $\frac{1}{n} \text{tr} \mathbf{A}_n (\tilde{\mathbf{Q}}(z) - \bar{\mathbf{Q}}(z)) \rightarrow 0$  and  $\mathbf{a}_n^\top (\tilde{\mathbf{Q}}(z) - \bar{\mathbf{Q}}(z)) \mathbf{b}_n \rightarrow 0$  almost surely as  $n, p, L \rightarrow +\infty$ . This will be simply denoted  $\tilde{\mathbf{Q}}(z) \leftrightarrow \bar{\mathbf{Q}}(z)$ .

sponding order for the eigenvalues of  $\mathbf{T}$ : after numerically computing them in descending order, we apply the same permutation that maps the eigenvalues of  $\mathbf{C}$  in descending order to  $(\psi_0, \dots, \psi_{n-1})$ . This yields the corresponding  $(\tau_0, \dots, \tau_{n-1})$ .

<sup>4</sup>This is only true if  $\lambda$  has multiplicity 1. In the general case, the integral equals the projection matrix on the eigenspace associated to  $\lambda$ .

**Theorem 3.1** (Deterministic equivalent of  $\tilde{\mathbf{Q}}$ ). *Let  $z \in \mathbb{C} \setminus \limsup_{n,p,L \rightarrow +\infty} \text{Sp}(\tilde{\mathbf{K}}_L)$  and define  $m(\cdot)$  as the unique Stieltjes transform solution to*

$$1 + zm(z) = \frac{1}{n} \sum_{k=0}^{n-1} \frac{m(z)\psi_k}{1 + \frac{m(z)}{p}\psi_k}. \quad (1)$$

*Under assumptions 2.1 – 2.3, if  $\left| \frac{2L-1}{p} m(z) \right| < 1$ , then*

$$\tilde{\mathbf{Q}}(z) \leftrightarrow \bar{\mathbf{Q}}(z) \equiv m(z) \left( \mathbf{I}_n + \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{F} \mathbf{\Lambda} \mathbf{F}^* \right)^{-1}$$

*where  $\mathbf{\Lambda} = m(z)\Psi \left( \mathbf{I}_n + \frac{m(z)}{p}\Psi \right)^{-1}$  is a diagonal matrix, thus  $\mathbf{F} \mathbf{\Lambda} \mathbf{F}^*$  is circulant.*

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

A first observation from theorem 3.1 is that  $\bar{\mathbf{Q}}(z)$  is the inverse of a perturbation of the identity which is *not* low rank. This strikingly differs from standard spiked random matrix models (Baik & Silverstein, 2006; Benaych-Georges & Nadakuditi, 2011) where a low-rank perturbation of the identity in the “population” matrix (here  $\mathbf{P}$ ) usually results in the presence of only a few isolated eigenvalues in the “sample” matrix (here  $\tilde{\mathbf{K}}_L$ ). This being said, here, in standard settings, most eigenvalues of  $\frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{F} \mathbf{\Lambda} \mathbf{F}^*$  are small enough for only a few number of corresponding isolated eigenvalues in the spectrum of  $\tilde{\mathbf{K}}_L$  to appear.

Yet, as a result of this full-rank perturbation of the identity property, it may be unclear whether  $m(\cdot)$  is the Stieltjes transform of the limiting spectral distribution of  $\tilde{\mathbf{K}}_L$  or not (this is important to ensure that isolated eigenvalues are truly informative). This issue is handled along the proof of theorem 3.1 (see proposition B.4) where it is shown that  $\frac{1}{n} \text{tr} \bar{\mathbf{Q}}(z) \rightarrow m(z)$  almost surely as  $n, p, L \rightarrow +\infty$ , thereby proving the (almost sure) convergence of the empirical spectral measure of  $\tilde{\mathbf{K}}_L$  to a measure  $\mu$  which is the inverse Stieltjes transform<sup>5</sup> of  $m$ .

*Remark 3.2* (Link with (Marčenko & Pastur, 1967)). In the particular case  $n = 2L - 1$ , the mask becomes  $\mathbf{C} = \mathbf{1}_n \mathbf{1}_n^\top$  and  $\tilde{\mathbf{K}}_L = \mathbf{K}$ . And, since  $\psi_0 = n$  and  $\psi_k = 0$  for  $1 \leq k < n$ , equation 1 becomes

$$zm^2(z) + (cz - c + 1)m(z) + c = 0$$

which is the canonical equation defining the Stieltjes transform of the Marčenko-Pastur distribution. The closer  $\varepsilon$  is to 1, the closer to the Marčenko-Pastur distribution is the limiting spectral distribution of  $\tilde{\mathbf{K}}_L$ .

<sup>5</sup>If  $\mu$  has a density  $d(x)$  at  $x \in \mathbb{R}$ , then  $d(x) = \lim_{y \downarrow 0} \frac{1}{\pi} \Im m(x + iy)$ .

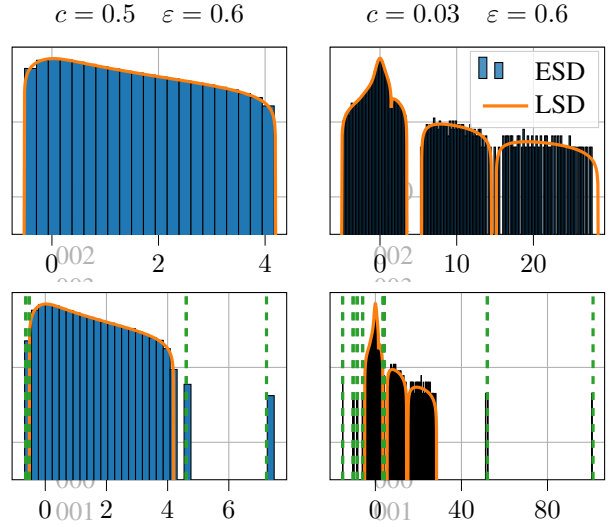


Figure 2. Empirical spectral distribution (ESD) and limiting spectral distribution (LSD) of  $\tilde{\mathbf{K}}_L$ . **The y-axis is in log scale.** **Top:** noise only,  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . **Bottom:** two-class mixture,  $\mathbf{x}_i \sim \mathcal{N}(\pm\boldsymbol{\mu}, \mathbf{I}_p)$  with  $\|\boldsymbol{\mu}\| = 2$ . Green dashed lines are the asymptotic positions of the spikes  $\tilde{\xi}_k$ . **Experimental setting:**  $n = 2500$ ,  $L = 750$  and  $p = 1250$  (left) or  $p = 75$  (right).

In practice, rather than computing  $m(z)$  directly from equation 1, it is easier to solve numerically the following fixed-point equation in  $\eta_0$

$$\eta_0 = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\psi_k^2}{p(1 - z - \eta_0) + \psi_k}$$

and deduce  $m(z) = \frac{1}{1 - z - \eta_0}$ .

Figure 2 displays, in log scale, the empirical spectral distribution of  $\tilde{\mathbf{K}}_L$  under two different settings with its limiting spectral distribution computed by inverting the Stieltjes transform given by theorem 3.1. Two kinds of data are presented: noise-only,  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ , (top row) and a two-class mixture  $\mathbf{x}_i \sim \mathcal{N}(\pm\boldsymbol{\mu}, \mathbf{I}_p)$  (bottom row). Notice how the shape of the distribution on the left column resembles the Marčenko-Pastur one (yet, some eigenvalues are negative here) while the second distribution has a completely different shape (there even are several bulks) for the same value of  $\varepsilon$ . This reveals that the parameter  $c$  also affects the closeness of the limiting spectral distribution to the Marčenko-Pastur one. Also note that, under the two-class mixture setting, more than one isolated eigenvalue pop out of the limiting support. It now remains to give a close look to their associated eigenvectors to understand how to exploit the latter in a spectral clustering perspective.



### 3.2. Phase transition and spike behavior

In this section, we thus focus back on our original classification objective. We consider two classes  $\mathcal{C}^\pm$  whose centroids are  $\pm\boldsymbol{\mu}$ , i.e.<sup>6</sup>,  $\mathbf{P} = \boldsymbol{\mu}\mathbf{j}^\top$  with  $\mathbf{j}_i = +1$  if  $\mathbf{x}_i \in \mathcal{C}^+$  and  $\mathbf{j}_i = -1$  if  $\mathbf{x}_i \in \mathcal{C}^-$ . This corresponds to a two-class mixture with globally centered data.

Because of the rank-one structure, using the relation  $\mathbf{M} \odot \mathbf{ab}^* = [\text{diag } \mathbf{a}] \mathbf{M} [\text{diag } \mathbf{b}]^*$ , the deterministic equivalent of the resolvent has a much simpler expression:

$$\bar{\mathbf{Q}}(z) = m(z) [\mathbf{D}_j \mathbf{F}] \left( \mathbf{I}_n + \frac{\|\boldsymbol{\mu}\|^2}{p} \boldsymbol{\Lambda} \right)^{-1} [\mathbf{D}_j \mathbf{F}]^*$$

where  $\mathbf{D}_j = \text{diag } \mathbf{j}$  is the diagonal matrix induced by vector  $\mathbf{j}$ . Now,  $\bar{\mathbf{Q}}(z)$  no longer involves a Hadamard product and we already have its eigendecomposition since  $\mathbf{I}_n + \frac{\|\boldsymbol{\mu}\|^2}{p} \boldsymbol{\Lambda}$  is diagonal and  $\mathbf{D}_j \mathbf{F}$  is unitary. Note that the columns of  $\mathbf{D}_j \mathbf{F}$  are simply the vectors of the Fourier basis with their sign switched at coordinates  $i$  such that  $\mathbf{x}_i \in \mathcal{C}^-$ .

With a deeper analysis of the resolvent  $\bar{\mathbf{Q}}(z)$ , the following theorem provides the position of the isolated eigenvalues and the shape of their associated eigenvectors.

**Theorem 3.3** (Phase transition, isolated eigenvalues and eigenvector alignments.). *Given an integer  $0 \leq k < n$ , let*

$$\bar{\xi}_k = \frac{\|\boldsymbol{\mu}\|^2 + 1}{p} \psi_k \left( 1 + \frac{p}{n} \sum_{l=0}^{n-1} \frac{\psi_l}{(\|\boldsymbol{\mu}\|^2 + 1) \psi_k - \psi_l} \right)$$

and

$$\bar{\zeta}_k = \frac{\|\boldsymbol{\mu}\|^2}{\|\boldsymbol{\mu}\|^2 + 1} \left( 1 - \frac{p}{n} \sum_{l=0}^{n-1} \left[ \frac{\psi_l}{(\|\boldsymbol{\mu}\|^2 + 1) \psi_k - \psi_l} \right]^2 \right).$$

The following propositions are equivalent.

1.  $\psi_k \neq 0$  and<sup>7</sup>  $\bar{\xi}_k \notin \text{supp } \mu$ .
2.  $\bar{\zeta}_k > 0$ .
3.  $\bar{\xi}_k$  is a singular point of  $\bar{\mathbf{Q}}(z)$ , i.e., the almost sure asymptotic position of an isolated eigenvalue of  $\tilde{\mathbf{K}}_L$ .

Then, in this case,  $\mathbf{U}_k = [\mathbf{u}_l]_{\substack{\psi_k = \psi_l \\ 0 \leq l < n}}$  is an isometric matrix gathering all the eigenvectors of  $\tilde{\mathbf{K}}_L$  whose associated eigenvalues converge to  $\bar{\xi}_k$  and

$$\mathbf{U}_k \mathbf{U}_k^* \leftrightarrow \bar{\zeta}_k [\mathbf{D}_j \mathbf{F}] \mathcal{D}_k [\mathbf{D}_j \mathbf{F}]^*$$

where  $\mathbf{D}_j = \text{diag } \mathbf{j}$  and  $\mathcal{D}_k = \text{diag} (\mathbf{1}_{\psi_k = \psi_l})_{0 \leq l < n}$ .

<sup>6</sup>Consistently with the previous setting,  $\mathbf{M} = [+ \boldsymbol{\mu} \quad - \boldsymbol{\mu}]$  and

$\mathbf{J}_{i \cdot} = [\mathbf{1}_{\mathbf{x}_i \in \mathcal{C}^+} \quad \mathbf{1}_{\mathbf{x}_i \in \mathcal{C}^-}]$ .

<sup>7</sup>Since  $\mu$  is the limiting spectral distribution of  $\tilde{\mathbf{K}}_L$ ,  $\text{supp } \mu = \limsup_{n,p,L \rightarrow +\infty} \text{Sp}(\tilde{\mathbf{K}}_L)$ .

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

To better understand this theorem, recall that, in theorem 3.1, we predicted the presence of a few isolated eigenvalues in the spectrum of  $\tilde{\mathbf{K}}_L$ . Theorem 3.3 details this assertion by specifying the number of spikes ( $\#\{\bar{\zeta}_k > 0\}$ ) and their position  $\bar{\xi}_k$ . The quantity  $\bar{\zeta}_k$  can really be seen as an ‘‘indicator of spike’’ as it tells whether an isolated eigenvalue exists and, if it does, the closer  $\bar{\zeta}_k$  is to 1, the better is the ‘‘quality’’ of the information carried in the corresponding eigenvector, i.e., the greater is the signal-to-noise ratio (see Figure 3).

Another difference with classical spiked random matrix models is that each asymptotic spike  $\bar{\xi}_k$ , which has the same multiplicity as the population spike  $\psi_k$ , is rarely simple<sup>8</sup>. However, for finite values of  $n, p, L$ , the corresponding eigenvalues of  $\tilde{\mathbf{K}}_L$  are not necessarily degenerate (with probability one, they are not), but they have the same limit<sup>9</sup>.

One also notices from theorem 3.3 that the number of isolated eigenvalues could potentially grow very large as  $\|\boldsymbol{\mu}\|$  increases. Indeed, the value of  $\|\boldsymbol{\mu}\|$  at which  $\bar{\zeta}_k$  changes sign (i.e., when one or more eigenvalues isolate from the bulk around  $\bar{\xi}_k$  during the *phase transition*) is given by

$$1 - \frac{p}{n} \sum_{l=0}^{n-1} \left[ \frac{\psi_l}{(\|\boldsymbol{\mu}\|^2 + 1) \psi_k - \psi_l} \right]^2 = 0.$$

Therefore, potentially any eigenvalue could leave the bulk, but this is prevented by the non-triviality condition (assumption 2.2):  $\|\boldsymbol{\mu}\| = \mathcal{O}_{n,p,L \rightarrow +\infty}(1)$ . Moreover, since most  $\psi_k$ 's are small (see Figure 1), the corresponding  $\bar{\xi}_k$ 's fall into the bulk and there are only a few spikes visible in practice. Yet, it is common to see negative isolated eigenvalues (see Figure 2). Indeed, since  $\psi_k$  can be negative, there can be spikes on *both sides* of the spectrum.

When positive, the quantity  $\bar{\zeta}_k$  is the asymptotic alignment between the empirical eigenvector  $\mathbf{u}_k$  and the corresponding information vector  $\mathbf{v}_k = [\mathbf{D}_j \mathbf{F}]_{\cdot,k} = \mathbf{F}_{\cdot,k} \odot \mathbf{j}$ , i.e.,

$$|\mathbf{u}_k^* \mathbf{v}_k|^2 \xrightarrow[n,p,L \rightarrow +\infty]{\text{a.s.}} \bar{\zeta}_k.$$

Thus,  $\bar{\zeta}_k$  measures the quality of the empirical eigenvector  $\mathbf{u}_k$ . Said differently,  $\mathbf{u}_k$  is a noisy version of a vector  $\mathbf{v}_k = \mathbf{F}_{\cdot,k} \odot \mathbf{j}$  and the noise level is indicated by  $0 \leq 1 - \bar{\zeta}_k < 1$ .

Figure 3 displays the value of  $\bar{\zeta}_k^+ = \max(\bar{\zeta}_k, 0)$  as a function of  $\|\boldsymbol{\mu}\|$  for the setting corresponding to the bottom right part of Figure 2. The empirical alignment of the dominant eigenvector  $\mathbf{u}_0$  with  $\mathbf{v}_0 = \frac{1}{\sqrt{n}} \mathbf{j}$  fits perfectly with the curve

<sup>8</sup> $\psi_0$ , and  $\psi_{n/2}$  when  $n$  is even, are the only simple eigenvalues of  $\mathbf{C}$ .

<sup>9</sup>In this case,  $\bar{\xi}_k = \bar{\xi}_l$  for all  $l$  such that  $\psi_k = \psi_l$ .

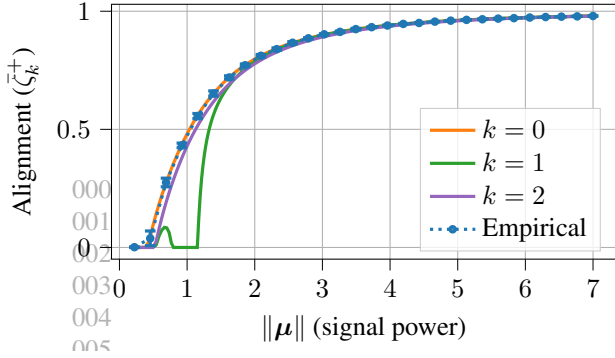


Figure 3. Asymptotic alignment  $\bar{\zeta}_k^+$  versus  $\|\mu\|$  for three values of  $k$ . The empirical alignment is computed as the mean of  $|\mathbf{u}_0^* \mathbf{v}_0|^2$  on 10 realizations (error bars indicate the standard deviation). **Experimental setting:**  $n = 2500$ ,  $p = 75$ ,  $L = 750$ .

of  $\bar{\zeta}_0^+$  predicted by theorem 3.3. Moreover, notice the interesting fact that  $\bar{\zeta}_1$  has several phase transitions: as  $\|\mu\|$  grows, it appears once, then disappears and appears once again! This is due to the limiting spectral distribution having several bulks under this setting (see Figure 2). The first time this spike appears, it is located between two bulks. It then goes through the rightmost bulk (so it is no longer an isolated eigenvalue thus  $\bar{\zeta}_1 \leq 0$ ), and finally goes out on the right edge of the distribution.

This last result may sound awkward and possibly testify of the suboptimality of our approach (when the signal-to-noise ratio increases, the information attached to some eigenvectors vanishes). This conclusion is not so immediate though, as the classification information is still contained within other eigenvectors which, as  $\|\mu\|$  increases, do carry increasingly clearer information.

### 3.3. Discussion on the circulant approximation

The approximation of the Toeplitz mask  $\mathbf{T}$  by the circulant mask  $\mathbf{C}$  used in the previous theorems 3.1 and 3.3 can be seen as a way to remove undesired edge effects, whose size is governed by  $L$ .<sup>10</sup> If  $L$  is chosen small compared to  $n$ , edge effects are expected to be negligible and the previous results can plausibly be extended to the original setting.

To adapt the previous results from  $\mathbf{C}$  to  $\mathbf{T}$ , one only needs to change the eigenvalues and eigenvectors, i.e., replace  $\psi_k$  by  $\tau_k$ , the eigenvalues of  $\mathbf{T}$ , and replace  $\mathbf{F}$  by  $[\mathbf{g}_0 \ \dots \ \mathbf{g}_{n-1}] \equiv \mathbf{G}$ , an eigenbasis of  $\mathbf{T}$ .

Very precise predictions on the original model can be made with these simple changes. Comparisons between these and reality are provided in appendix D.

<sup>10</sup>Remove the first and last  $L - 1$  rows and columns of  $\mathbf{C}$  and  $\mathbf{T}$  and we are left with the same two Toeplitz matrices.

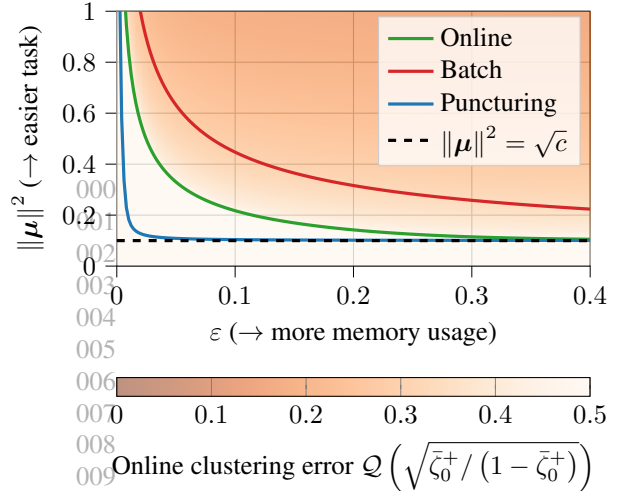


Figure 4. Phase transition position ( $\|\mu\|^2$ ) of the dominant eigenvector of the kernel matrix against the sparsity parameter ( $\epsilon$ ) with  $n/p = 100$ . Classification with a given method is only possible above the corresponding curve. The black dashed line is the optimal phase transition (with all the data available). **Green:** online kernel spectral clustering (circulant mask). **Red:** regular kernel spectral clustering with  $L = \frac{n\epsilon+1}{2}$  points. **Blue:** punctured (off-line) kernel spectral clustering (Bernoulli mask).

## 4. Online spectral clustering of large data

The previous results find direct applications to the online clustering of high-dimensional data.

### 4.1. Performance vs. cost trade-off in online learning

The phase transition position provided by theorem 3.3 lets us determine under which setting classification is possible or not. Consider the dominant eigenvector  $\mathbf{u}_0$ . If  $\bar{\zeta}_0 \leq 0$  then no eigenvalue isolates from the bulk and classification cannot be performed. After the phase transition,  $\bar{\zeta}_0 > 0$  and the closer it is to 1, the closer  $\mathbf{u}_0$  is to  $\mathbf{v}_0 = \frac{1}{\sqrt{n}}\mathbf{j}$ . The fluctuations of the entries of  $\mathbf{u}_0$  happen to be asymptotically Gaussian and pairwise independent (Kadavankandy & Couillet, 2019) with, for equal-size classes, mean  $\pm\sqrt{\bar{\zeta}_0/n}$  and variance  $(1 - \bar{\zeta}_0)/n$ . Thus, the asymptotic clustering error is given by  $\mathcal{Q}\left(\sqrt{\bar{\zeta}_0^+ / (1 - \bar{\zeta}_0^+)}\right)$ , where  $\mathcal{Q}$  is the Gaussian tail function:  $\mathcal{Q}(x) = \frac{1}{\sqrt{2\pi}} \int_x^{+\infty} e^{-t^2/2} dt$ .

Figure 4 shows the phase transition position  $\|\mu\|^2$  as a function of  $\epsilon = \frac{2L-1}{n}$ , with the asymptotic clustering error of online kernel spectral clustering, when  $n/p = 100$ . For comparison, the phase transition curves of the following two methods are also represented:

- *Batch clustering*, i.e., standard  $L \times L$  kernel spectral clustering with the  $L$  data points available in memory.

- *Punctured kernel spectral clustering* (Zarrouk et al., 2020; Couillet et al., 2021), i.e., *offline* clustering performed with a sparsified kernel matrix  $\mathbf{K}_\varepsilon = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{B}$ , where  $\mathbf{B}_{i,j} = \mathbf{B}_{j,i} \sim \text{Bern}(\varepsilon)$  and  $\mathbf{B}_{i,i} = 1$ .<sup>11</sup>

As  $\varepsilon$  grows, the phase transition position of online spectral clustering rapidly reaches the optimal threshold  $\|\boldsymbol{\mu}\|^2 = \sqrt{c}$  under which no information can be recovered (regardless of the method used and the data available). Moreover, the clustering error decreases to 0 as  $\|\boldsymbol{\mu}\|$  increases (the signal is more powerful). A good compromise between memory usage and performance appears to be  $0.1 \lesssim \varepsilon \lesssim 0.2$ , i.e.,  $\frac{n}{20} \lesssim L \lesssim \frac{n}{10}$  as it keeps  $L$  (i.e., the memory usage) small while not impairing much the performance.

Our method performs better (i.e., the phase transition occurs earlier) than the naive setting performing standard clustering on batches of  $L$  points available in memory. It is also able to classify the  $n$  previous points (and not only the  $L$  previous ones) at any time, *although the corresponding data points have left memory*. It is instructive to see that, under the same sparsity level of the kernel matrix, the puncturing method performs better. Yet, this requires the access to  $n$  data points to compute  $\mathbf{K}_\varepsilon$ , which is not possible in online learning.

## 4.2. Online clustering algorithm

Before diving into the simulations, we detail a clustering algorithm based on our previous results. We here use the banded version of the kernel matrix:  $\mathbf{K}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{T}$  (the circulant mask was only useful for theoretical considerations) and recall the notation of the eigenbasis of  $\mathbf{T}$ :  $[\mathbf{g}_0 \ \dots \ \mathbf{g}_{n-1}] \equiv \mathbf{G}$ .

We consider a data stream of length  $T$  (possibly infinite). At each time step, a new vector  $\mathbf{x}_t$  arrives and the kernel matrix is updated:  $[\mathbf{K}_L^{(t)}]_{i,j} = \frac{1}{p} \mathbf{x}_{t-n+i}^\top \mathbf{x}_{t-n+j} \mathbf{1}_{|i-j| < L}$ .

*Remark 4.1* (Choice of  $n$  and eigenvector localization). It is important to emphasize that  $n$  is *not* the length of the data stream (given by the newly-introduced parameter  $T \geq n$ ). As  $\mathbf{K}_L$  has size  $n \times n$ , one can “only” classify the last  $n$  points of the stream, even when discarded from the length- $L$  memory (older points are no longer classified though).

The parameter  $n$  is left for the user to choose, accounting for  $L$  and our previous considerations on the performance (Figure 4) and memory limitations:  $\mathcal{O}(Lp + Ln)$  space is needed to store the data *and* the kernel matrix. Moreover, as the graph associated to  $\mathbf{K}_L$  becomes sparser ( $n \gg L$ ), its eigenvectors tend to localize (Hata & Nakao, 2017), making classification more challenging.

<sup>11</sup>This of course is not doable with a memory bank of size  $L$  but the comparison is interesting as the number of entries in  $\mathbf{K}_\varepsilon$  and  $\mathbf{K}_L$  is the same.

As per standard kernel spectral clustering, we use the dominant eigenvectors of  $\mathbf{K}_L^{(t)}$  to estimate the classes. The last  $n$  points of the stream are classified at each time step so each point is classified  $n$  times. Then, the final class estimate can be chosen by a majority vote. However, because of the particular shape of the eigenvectors caused by the Toeplitz mask<sup>12</sup> (see Figure 6), standard clustering algorithms such as  $k$ -means perform poorly on such spectral embeddings. Therefore, we propose a new way to cluster the data.

*Remark 4.2.* The eigenvectors of  $\mathbf{K}_L^{(t)}$  can be quickly computed at a low cost with a warm start of the power iteration algorithm from the previously computed eigenvectors of  $\mathbf{K}_L^{(t-1)}$ .

In a binary setting with globally centered data, classification can be performed using only the dominant eigenvector  $\mathbf{u}_0^{(t)}$  of  $\mathbf{K}_L^{(t)}$ . Relying on the alignment of  $\mathbf{u}_0^{(t)}$  with  $\mathbf{v}_0^{(t)} = \mathbf{g}_0 \odot \mathbf{j}^{(t)}$  (theorem 3.3) and the fact that the coordinates of  $\mathbf{g}_0$  have constant sign, the class of  $\mathbf{x}_{t-n+i}$  can be estimated from the sign of  $[\mathbf{u}_0^{(t)}]_i$ . This online clustering procedure is summarized in Algorithm 1.

---

### Algorithm 1 Online kernel spectral clustering (binary)

---

**Output:** class estimators  $\{\hat{\mathcal{C}}_t^+, \hat{\mathcal{C}}_t^-\}_{n \leq t \leq T}$ .

**for**  $t = 1$  to  $T$  **do**

Get a new point  $\mathbf{x}_t$  into the pipeline.

Compute  $\mathbf{x}_t^* \mathbf{x}_{t-l}$  for  $l = 0$  to  $L - 1$ .

Update  $\mathbf{K}_L^{(t-1)}$  into  $\mathbf{K}_L^{(t)}$ .

**if**  $t \geq n$  **then**

$\mathbf{u}_0^{(t)} \leftarrow \text{PowerIteration}(\mathbf{K}_L^{(t)}, \mathbf{u}_0^{(t-1)})$ .

$\hat{\mathcal{C}}_t^\pm \leftarrow \{\mathbf{x}_{t-n+i} \mid [\mathbf{u}_0^{(t)}]_i \gtrless 0\}$ .

**end if**

**end for**

---

In appendix E, we also propose a (more complex) online spectral clustering algorithm capable of handling  $K$ -class mixtures and test it on Fashion-MNIST images.

Note that this algorithm can easily be adapted to a setting where more than one vector  $\mathbf{x}_t$  arrives at each time step (and this quantity does not need to be constant in time).

## 4.3. Simulations on real-world images

We illustrate our findings with two applications on image classification tasks. We first apply Algorithm 1 on globally centered and scaled VGG-features (Simonyan & Zisserman, 2015) of randomly BigGAN-generated images (Brock et al., 2019) of tabby cats and collie dogs (see Figure 5). The vectors thus generated have dimension  $p = 4096$

<sup>12</sup>The dominant eigenvector of  $\mathbf{T}$ , for example, is not constant, contrary to the first Fourier mode with the circulant mask.



Figure 5. Examples of BigGAN-generated images: collie (top) and tabby (bottom).

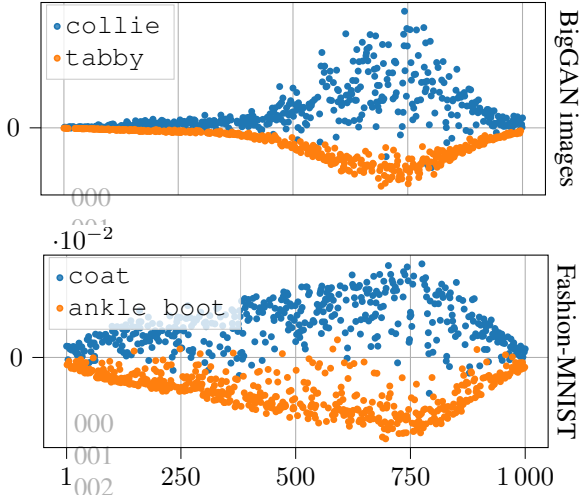


Figure 6. Dominant eigenvector of  $\mathbf{K}_L^{(t)}$  with BigGAN-generated images (top) and Fashion-MNIST images (bottom). **Experimental setting:**  $T = 20\,000$ ,  $n = 1\,000$ ,  $p = 4\,096$ ,  $L = 100$  (BigGAN images) and  $T = 14\,000$ ,  $n = 1\,000$ ,  $p = 784$ ,  $L = 100$  (Fashion-MNIST).

and simulate a stream of length  $T = 20\,000$  with evenly likely cats and dogs. In addition, our algorithm is applied to a stream made of  $T = 14\,000$  centered raw-images from the Fashion-MNIST dataset (Xiao et al., 2017). Their dimension is  $p = 784$  and we want to discriminate `coat` versus `ankle boot` in an online fashion. In both cases, we choose  $n = 1\,000$  and  $L = 100$ .

Figure 6 shows the shape of the dominant eigenvector  $\mathbf{u}_0^{(t)}$  at a given time during the execution of the algorithm. We clearly see a separation between the classes. For both settings, Figure 7 depicts the mean clustering error at  $t_0 + \Delta t$  of a data point seen at  $t_0$ , as well as the overall classification error obtained after a majority vote (green dashed line), to be compared with the classification error obtained with a  $T \times T$  offline kernel spectral clustering<sup>13</sup> (black dotted line). The mean classification error remains constant with  $\Delta t$ , thus showing that our algorithm does not lose any discriminative

<sup>13</sup>For which optimality results are known.

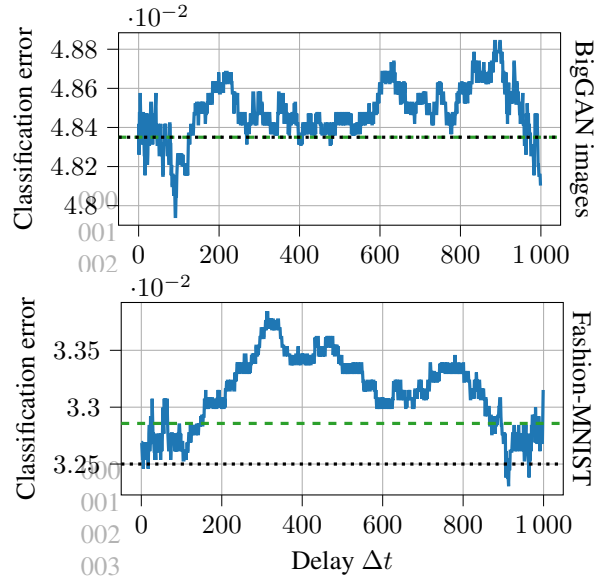


Figure 7. Classification error against delay  $\Delta t$  on BigGAN-generated images (top) and Fashion-MNIST images (bottom). This is the mean classification error at time  $t_0 + \Delta t$  of a point arrived at  $t_0$ . The green dashed line indicate the overall classification error when the class is chosen by a majority vote. The black dotted line is the classification error obtained with a  $T \times T$  offline kernel spectral clustering. **Experimental setting:** as in Figure 6.

power between  $t_0$  and  $t_0 + n - 1$ . The classification performances of our algorithm are very close to those of the full spectral clustering but require much less memory resources:  $\mathcal{O}(Lp + Ln)$  against  $\mathcal{O}(Tp + T^2)$  space for the storage of the data and the kernel matrix.

## 5. Concluding remarks

Leveraging tools from random matrix theory, the article shows that, under limited memory resources, near-optimal performances on high-dimensional data can be achieved using an online kernel spectral clustering algorithm. By means of a thorough asymptotic analysis, we specify the optimal performances achievable when learning on a data stream, which we exploit to propose a novel efficient clustering algorithm adapted to memory-limited systems.

The article not only introduces a novel algorithm for online classification, but also paves the path towards the question of large-dimensional learning in data streaming with theoretical guarantees. Still, here we miss an information-theoretic result of optimality for the proposed approach (which exists in the unbanded case), a key direction we currently investigate.



## References

- 440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494
- Ackermann, M. R., Märtens, M., Raupach, C., Swierkot, K., Lammersen, C., and Sohler, C. StreamKM++: A clustering algorithm for data streams. *ACM Journal of Experimental Algorithmics*, 17:2.4:2.1–2.4:2.30, May 2012. ISSN 1084-6654. doi: 10.1145/2133803.2184450. URL <https://doi.org/10.1145/2133803.2184450>.
- Aggarwal, C. C., Yu, P. S., Han, J., and Wang, J. - A Framework for Clustering Evolving Data Streams. In Freytag, J.-C., Lockemann, P., Abiteboul, S., Carey, M., Selinger, P., and Heuer, A. (eds.), *Proceedings 2003 VLDB Conference*, pp. 81–92. Morgan Kaufmann, San Francisco, January 2003. ISBN 978-0-12-722442-8. doi: 10.1016/B978-012722442-8/50016-1. URL <https://www.sciencedirect.com/science/article/pii/B9780127224428500161>.
- Aggarwal, C. C., Han, J., Wang, J., and Yu, P. S. - A Framework for Projected Clustering of High Dimensional Data Streams. In Nascimento, M. A., Özsu, M. T., Kossman, D., Miller, R. J., Blakeley, J. A., and Schiefer, B. (eds.), *Proceedings 2004 VLDB Conference*, pp. 852–863. Morgan Kaufmann, St Louis, January 2004. ISBN 978-0-12-088469-8. doi: 10.1016/B978-012088469-8.50075-9. URL <https://www.sciencedirect.com/science/article/pii/B9780120884698500759>.
- Bai, Z. and Silverstein, J. W. *Spectral analysis of large dimensional random matrices*, volume 20. Springer, 2010.
- Baik, J. and Silverstein, J. W. Eigenvalues of large sample covariance matrices of spiked population models. *Journal of Multivariate Analysis*, 97(6):1382–1408, July 2006. ISSN 0047-259X. doi: 10.1016/j.jmva.2005.08.003. URL <https://www.sciencedirect.com/science/article/pii/S0047259X0500134X>.
- Benaych-Georges, F. and Nadakuditi, R. R. The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices. *Advances in Mathematics*, 227(1):494–521, May 2011. ISSN 0001-8708. doi: 10.1016/j.aim.2011.02.007. URL <https://www.sciencedirect.com/science/article/pii/S0001870811000570>.
- Brock, A., Donahue, J., and Simonyan, K. Large Scale GAN Training for High Fidelity Natural Image Synthesis. *arXiv:1809.11096 [cs, stat]*, February 2019. URL <http://arxiv.org/abs/1809.11096>. arXiv: 1809.11096.
- Cheng, X. and Singer, A. The Spectrum of Random Inner-product Kernel Matrices. *arXiv:1202.3155 [math]*, March 2012. URL <http://arxiv.org/abs/1202.3155>. arXiv: 1202.3155.
- Couillet, R. and Benaych-Georges, F. Kernel spectral clustering of large dimensional data. *Electronic Journal of Statistics*, 10(1):1393–1454, 2016. doi: 10.1214/16-EJS1144. URL <https://hal.archives-ouvertes.fr/hal-01215343>. Publisher: Shaker Heights, OH : Institute of Mathematical Statistics.
- Couillet, R. and Liao, Z. *Random Matrix Methods for Machine Learning: When Theory meets Applications*. 2021.
- Couillet, R., Chatelain, F., and Le Bihan, N. Two-way kernel matrix puncturing: towards resource-efficient PCA and spectral clustering. *arXiv:2102.12293 [cs, stat]*, May 2021. URL <http://arxiv.org/abs/2102.12293>. arXiv: 2102.12293.
- Dhanjal, C., Gaudel, R., and Cléménçon, S. Efficient Eigenupdating for Spectral Graph Clustering. *arXiv:1301.1318 [stat]*, January 2014. URL <http://arxiv.org/abs/1301.1318>. arXiv: 1301.1318.
- El Karoui, N. The spectrum of kernel random matrices. *The Annals of Statistics*, 38(1):1–50, February 2010. ISSN 0090-5364, 2168-8966. doi: 10.1214/08-AOS648. URL <https://projecteuclid.org/journals/annals-of-statistics/volume-38/issue-1/The-spectrum-of-kernel-random-matrices/10.1214/08-AOS648.full>. Publisher: Institute of Mathematical Statistics.
- Ester, M., Krieger, H., Sander, J., and Xu, X. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *KDD*, 1996.
- Fritzke, B. A Growing Neural Gas Network Learns Topologies. In *Advances in Neural Information Processing Systems*, volume 7. MIT Press, 1995. URL <https://papers.nips.cc/paper/1994/hash/d56b9fc4b0f1be8871f5e1c40c0067e7-Abstract.html>.
- Ghesmoune, M., Lebbah, M., and Azzag, H. State-of-the-art on clustering data streams. *Big Data Analytics*, 1(1): 13, December 2016. ISSN 2058-6345. doi: 10.1186/s41044-016-0011-3. URL <https://doi.org/10.1186/s41044-016-0011-3>.
- Gray, R. M. Toeplitz and Circulant Matrices: A Review. *Foundations and Trends® in Communications and Information Theory*, 2(3):155–239, January 2006. ISSN 1567-2190, 1567-2328. doi: 10.1561/0100000006.

- 495 URL [https://www.nowpublishers.com/](https://www.nowpublishers.com/article/Details/CIT-006)  
 496 [article/Details/CIT-006](https://www.nowpublishers.com/article/Details/CIT-006). Publisher: Now  
 497 Publishers, Inc.
- 498 Gribonval, R., Chatalic, A., Keriven, N., Schellekens, V.,  
 499 Jacques, L., and Schniter, P. Sketching Data Sets for  
 500 Large-Scale Learning: Keeping only what you need. *IEEE Signal Processing Magazine*, 38(5):12–36, Septem-  
 501 ber 2021. doi: 10.1109/MSP.2021.3092574. URL  
 502 <https://hal.inria.fr/hal-03350599>. Pub-  
 503 lisher: Institute of Electrical and Electronics Engineers.  
 504
- 506 Hata, S. and Nakao, H. Localization of Laplacian  
 507 eigenvectors on random networks. *Scientific*  
 508 *Reports*, 7(1):1121, April 2017. ISSN 2045-  
 509 2322. doi: 10.1038/s41598-017-01010-0. URL  
 510 [https://www.nature.com/articles/](https://www.nature.com/articles/s41598-017-01010-0)  
 511 [s41598-017-01010-0](https://www.nature.com/articles/s41598-017-01010-0). Bandiera\_abtest: a  
 512 Cc\_license\_type: cc\_by Cg\_type: Nature Research  
 513 Journals Number: 1 Primary\_atype: Research Publisher:  
 514 Nature Publishing Group Subject\_term: Complex  
 515 networks;Nonlinear phenomena Subject\_term\_id:  
 516 complex-networks;nonlinear-phenomena.  
 517
- 518 Kadavankandy, A. and Couillet, R. Asymptotic Gaus-  
 519 sian Fluctuations of Spectral Clustering Eigenvectors.  
 520 In *2019 IEEE 8th International Workshop on Computa-*  
 521 *tional Advances in Multi-Sensor Adaptive Process-*  
 522 *ing (CAMSAP)*, pp. 694–698, December 2019. doi:  
 523 10.1109/CAMSAP45676.2019.9022474.
- 524 Keriven, N., Bourrier, A., Gribonval, R., and Pérez, P.  
 525 Sketching for Large-Scale Learning of Mixture Mod-  
 526 els. *arXiv:1606.02838 [cs, stat]*, May 2017. URL  
 527 <http://arxiv.org/abs/1606.02838>. arXiv:  
 528 1606.02838.  
 529
- 530 Liao, Z., Couillet, R., and Mahoney, M. W. Sparse Quan-  
 531 tized Spectral Clustering. *arXiv:2010.01376 [cs, math,*  
 532 *stat]*, October 2020. URL [http://arxiv.org/](http://arxiv.org/abs/2010.01376)  
 533 [abs/2010.01376](http://arxiv.org/abs/2010.01376). arXiv: 2010.01376.
- 534 Liberty, E. Simple and Deterministic Matrix Sketch-  
 535 ing. *arXiv:1206.0594 [cs]*, July 2012. URL [http://](http://arxiv.org/abs/1206.0594)  
 536 [arxiv.org/abs/1206.0594](http://arxiv.org/abs/1206.0594). arXiv: 1206.0594.  
 537
- 538 Liberty, E., Sriharsha, R., and Sviridenko, M. An Algorithm  
 539 for Online K-Means Clustering. In *2016 Proceedings of*  
 540 *the Meeting on Algorithm Engineering and Experiments*  
 541 *(ALENEX)*, Proceedings, pp. 81–89. Society for Industrial  
 542 and Applied Mathematics, December 2015. doi: 10.1137/  
 543 1.9781611974317.7. URL [https://epubs.siam.](https://epubs.siam.org/doi/abs/10.1137/1.9781611974317.7)  
 544 [org/doi/abs/10.1137/1.9781611974317.7](https://epubs.siam.org/doi/abs/10.1137/1.9781611974317.7).
- 545 Lytova, A. and Pastur, L. Central limit theorem for linear  
 546 eigenvalue statistics of random matrices with independent  
 547 entries. *The Annals of Probability*, 37(5):1778–1840,  
 548 2009. Publisher: Institute of Mathematical Statistics.  
 549
- Löffler, M., Zhang, A. Y., and Zhou, H. H. Optimality  
 of Spectral Clustering in the Gaussian Mixture Model.  
*arXiv:1911.00538 [cs, math, stat]*, August 2020. URL  
<http://arxiv.org/abs/1911.00538>. arXiv:  
 1911.00538.
- Mai, X. and Couillet, R. A random matrix analysis and  
 improvement of semi-supervised learning for large di-  
 mensional data. *arXiv:1711.03404 [cs, stat]*, Novem-  
 ber 2017. URL [http://arxiv.org/abs/1711.](http://arxiv.org/abs/1711.03404)  
 03404. arXiv: 1711.03404.
- Marčenko, V. A. and Pastur, L. A. Distribution of eigen-  
 values for some sets of random matrices. *Mathematics*  
*of the USSR-Sbornik*, 1(4):457, 1967. Publisher: IOP  
 Publishing.
- Ning, H., Xu, W., Chi, Y., Gong, Y., and Huang,  
 T. S. Incremental spectral clustering by effi-  
 ciently updating the eigen-system. *Pattern Recog-*  
*nition*, 43(1):113–127, January 2010. ISSN  
 0031-3203. doi: 10.1016/j.patcog.2009.06.001.  
 URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S0031320309002209)  
[science/article/pii/S0031320309002209](https://www.sciencedirect.com/science/article/pii/S0031320309002209).
- Onatski, A., Moreira, M. J., and Hallin, M. Asymptotic  
 power of sphericity tests for high-dimensional data. *The*  
*Annals of Statistics*, 41(3):1204–1231, June 2013. ISSN  
 0090-5364, 2168-8966. doi: 10.1214/13-AOS1100.  
 URL [https://projecteuclid.org/](https://projecteuclid.org/journals/annals-of-statistics/volume-41/issue-3/Asymptotic-power-of-sphericity-tests-for-high-dim-10.1214/13-AOS1100.full)  
[journals/annals-of-statistics/](https://projecteuclid.org/journals/annals-of-statistics/volume-41/issue-3/Asymptotic-power-of-sphericity-tests-for-high-dim-10.1214/13-AOS1100.full)  
[volume-41/issue-3/](https://projecteuclid.org/journals/annals-of-statistics/volume-41/issue-3/Asymptotic-power-of-sphericity-tests-for-high-dim-10.1214/13-AOS1100.full)  
[Asymptotic-power-of-sphericity-tests-for-high-dim-](https://projecteuclid.org/journals/annals-of-statistics/volume-41/issue-3/Asymptotic-power-of-sphericity-tests-for-high-dim-10.1214/13-AOS1100.full)  
[10.1214/13-AOS1100.full](https://projecteuclid.org/journals/annals-of-statistics/volume-41/issue-3/Asymptotic-power-of-sphericity-tests-for-high-dim-10.1214/13-AOS1100.full). Publisher: Institute  
 of Mathematical Statistics.
- Simonyan, K. and Zisserman, A. Very Deep Convo-  
 lutional Networks for Large-Scale Image Recognition.  
*arXiv:1409.1556 [cs]*, April 2015. URL [http://](http://arxiv.org/abs/1409.1556)  
[arxiv.org/abs/1409.1556](http://arxiv.org/abs/1409.1556). arXiv: 1409.1556.
- Stein, C. M. Estimation of the Mean of a Multivari-  
 ate Normal Distribution. *The Annals of Statistics*,  
 9(6):1135–1151, November 1981. ISSN 0090-  
 5364, 2168-8966. doi: 10.1214/aos/1176345632.  
 URL [https://projecteuclid.org/](https://projecteuclid.org/journals/annals-of-statistics/volume-9/issue-6/Estimation-of-the-Mean-of-a-Multivariate-Normal-D-10.1214/aos/1176345632.full)  
[journals/annals-of-statistics/](https://projecteuclid.org/journals/annals-of-statistics/volume-9/issue-6/Estimation-of-the-Mean-of-a-Multivariate-Normal-D-10.1214/aos/1176345632.full)  
[volume-9/issue-6/](https://projecteuclid.org/journals/annals-of-statistics/volume-9/issue-6/Estimation-of-the-Mean-of-a-Multivariate-Normal-D-10.1214/aos/1176345632.full)  
[Estimation-of-the-Mean-of-a-Multivariate-Normal-D-](https://projecteuclid.org/journals/annals-of-statistics/volume-9/issue-6/Estimation-of-the-Mean-of-a-Multivariate-Normal-D-10.1214/aos/1176345632.full)  
[10.1214/aos/1176345632.full](https://projecteuclid.org/journals/annals-of-statistics/volume-9/issue-6/Estimation-of-the-Mean-of-a-Multivariate-Normal-D-10.1214/aos/1176345632.full). Publisher:  
 Institute of Mathematical Statistics.
- Tasoulis, D. K., Ross, G., and Adams, N. M. Visualising  
 the Cluster Structure of Data Streams. In R. Berthold,  
 M., Shawe-Taylor, J., and Lavrač, N. (eds.), *Advances*  
*in Intelligent Data Analysis VII*, Lecture Notes in Com-  
 puter Science, pp. 81–92, Berlin, Heidelberg, 2007.

550 Springer. ISBN 978-3-540-74825-0. doi: 10.1007/  
551 978-3-540-74825-0\_8.

552  
553 Trench, W. F. Some Spectral Properties of Hermitian  
554 Toeplitz Matrices. *SIAM Journal on Matrix Analysis*  
555 *and Applications*, 15(3):938–942, July 1994. ISSN  
556 0895-4798. doi: 10.1137/S0895479892239007. URL  
557 [https://epubs.siam.org/doi/10.1137/  
558 S0895479892239007](https://epubs.siam.org/doi/10.1137/S0895479892239007). Publisher: Society for  
559 Industrial and Applied Mathematics.

560 von Luxburg, U. A tutorial on spectral cluster-  
561 ing. *Statistics and Computing*, 17(4):395–416, De-  
562 cember 2007. ISSN 1573-1375. doi: 10.1007/  
563 s11222-007-9033-z. URL [https://doi.org/10.  
564 1007/s11222-007-9033-z](https://doi.org/10.1007/s11222-007-9033-z).

565  
566 Xiao, H., Rasul, K., and Vollgraf, R. Fashion-MNIST: a  
567 Novel Image Dataset for Benchmarking Machine Learn-  
568 ing Algorithms. *arXiv:1708.07747 [cs, stat]*, Septem-  
569 ber 2017. URL [http://arxiv.org/abs/1708.  
570 07747](http://arxiv.org/abs/1708.07747). arXiv: 1708.07747.

571  
572 Yoo, S., Huang, H., and Kasiviswanathan, S. P. Streaming  
573 spectral clustering. In *2016 IEEE 32nd International*  
574 *Conference on Data Engineering (ICDE)*, pp. 637–648,  
575 May 2016. doi: 10.1109/ICDE.2016.7498277.

576 Zarrouk, T., Couillet, R., Chatelain, F., and Le Bihan,  
577 N. Performance-Complexity Trade-Off in Large Dimen-  
578 sional Statistics. In *2020 IEEE 30th International Work-*  
579 *shop on Machine Learning for Signal Processing (MLSP)*,  
580 pp. 1–6, September 2020. doi: 10.1109/MLSP49062.  
581 2020.9231568. ISSN: 1551-2541.

582  
583 Zhang, T., Ramakrishnan, R., and Livny, M. BIRCH: an  
584 efficient data clustering method for very large databases.  
585 In *Proceedings of the 1996 ACM SIGMOD international*  
586 *conference on Management of data, SIGMOD '96*, pp.  
587 103–114, New York, NY, USA, June 1996. Association  
588 for Computing Machinery. ISBN 978-0-89791-794-0.  
589 doi: 10.1145/233269.233324. URL [https://doi.  
590 org/10.1145/233269.233324](https://doi.org/10.1145/233269.233324).

591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604

**A. Useful identities**

Let  $\mathbf{P} + \mathbf{Z} = \mathbf{X} \in \mathbb{R}^{p \times n}$  where  $\mathbf{P}$  is a deterministic matrix and  $\mathbf{Z}$  is a random matrix with independent entries  $Z_{i,j} \sim \mathcal{N}(0, 1)$ .

Let  $\mathbf{Q} = \left( \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{R} - z \mathbf{I}_n \right)^{-1}$  where  $\mathbf{R} \in \mathbb{R}^{n \times n}$  is symmetric with bounded entries and  $z \in \mathbb{C} \setminus \text{Sp} \left( \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{R} \right)$ .

**Proposition A.1.**

$$\frac{\partial \mathbf{Q}_{k,l}}{\partial Z_{i,j}} = -\frac{1}{p} \left( [\mathbf{X} \mathbf{D}_{\mathbf{R},j} \mathbf{Q}]_{i,k} \mathbf{Q}_{j,l} + \mathbf{Q}_{k,j} [\mathbf{X} \mathbf{D}_{\mathbf{R},j} \mathbf{Q}]_{i,l} \right)$$

*Proof.*

$$\begin{aligned} \frac{\partial \mathbf{Q}_{k,l}}{\partial Z_{i,j}} &= \left[ \frac{\partial \mathbf{Q}}{\partial Z_{i,j}} \right]_{k,l} \\ &= -\frac{1}{p} \left[ \mathbf{Q} \frac{\partial (\mathbf{X}^\top \mathbf{X} \odot \mathbf{R})}{\partial Z_{i,j}} \mathbf{Q} \right]_{k,l} \\ &= -\frac{1}{p} \sum_{r,s=1}^n \mathbf{Q}_{k,r} \mathbf{Q}_{s,l} \frac{\partial}{\partial Z_{i,j}} \sum_{t=1}^n \mathbf{X}_{t,r} \mathbf{X}_{t,s} \mathbf{R}_{r,s} \\ &= -\frac{1}{p} \sum_{r,s,t=1}^n \mathbf{Q}_{k,r} \mathbf{Q}_{s,l} \mathbf{R}_{r,s} \left[ \frac{\partial \mathbf{X}_{t,r}}{\partial Z_{i,j}} \mathbf{X}_{t,s} + \mathbf{X}_{t,r} \frac{\partial \mathbf{X}_{t,s}}{\partial Z_{i,j}} \right] \\ &= -\frac{1}{p} \left( \sum_{s=1}^n \mathbf{Q}_{k,j} \mathbf{Q}_{s,l} \mathbf{R}_{j,s} \mathbf{X}_{i,s} + \sum_{r=1}^n \mathbf{Q}_{k,r} \mathbf{Q}_{j,l} \mathbf{R}_{r,j} \mathbf{X}_{i,r} \right) \\ \frac{\partial \mathbf{Q}_{k,l}}{\partial Z_{i,j}} &= -\frac{1}{p} \left( \mathbf{Q}_{k,j} [\mathbf{X} \mathbf{D}_{\mathbf{R},j} \mathbf{Q}]_{i,l} + \mathbf{Q}_{j,l} [\mathbf{X} \mathbf{D}_{\mathbf{R},j} \mathbf{Q}]_{i,k} \right) \quad \text{since } \mathbf{Q}^\top = \mathbf{Q}. \end{aligned}$$

□

**Lemma A.2 ((Stein, 1981)).** Let  $Z \sim \mathcal{N}(0, 1)$  and  $f : \mathbb{R} \rightarrow \mathbb{R}$  be a differentiable function such that  $\mathbb{E}[|f'(Z)|] < +\infty$  and  $f(z) = o_{z \rightarrow \pm\infty}(e^{z^2})$ . Then,

$$\mathbb{E}[Zf(Z)] = \mathbb{E}[f'(Z)].$$

*Proof.* Using integration by parts,

$$\mathbb{E}[Zf(Z)] = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} z f(z) e^{-\frac{z^2}{2}} dz = \frac{1}{\sqrt{2\pi}} \left[ -f(z) e^{-\frac{z^2}{2}} \right]_{-\infty}^{+\infty} + \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} f'(z) e^{-\frac{z^2}{2}} dz.$$

In the right-hand side, the first term vanishes since  $f(z) = o_{z \rightarrow \pm\infty}(e^{z^2})$  and the second term equals  $\mathbb{E}[f'(Z)]$ . □

**Proposition A.3.** Let  $\mathbf{A} \in \mathbb{R}^{n \times n}$  be a matrix with bounded entries.

1.

$$\mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} = \mathbb{E}[\mathbf{A}_{i,i} \mathbf{Q}_{i,j}] - \mathbb{E} \left[ \frac{1}{p} \text{tr} \left( \mathbf{Q} \left( \frac{(\mathbf{P} + \mathbf{Z})^\top \mathbf{Z}}{p} \odot [\mathbf{R}_{\cdot,i} \mathbf{A}_{i,\cdot}] \right) \right) \mathbf{Q}_{i,j} \right] + o_{n,p \rightarrow +\infty} \left( \frac{1}{p} \right)$$

2.

$$\mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{P}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} = -\mathbb{E} \left[ \frac{1}{p} \text{tr} \left( \mathbf{Q} \left( \frac{(\mathbf{P} + \mathbf{Z})^\top \mathbf{P}}{p} \odot [\mathbf{R}_{\cdot,i} \mathbf{A}_{i,\cdot}] \right) \right) \mathbf{Q}_{i,j} \right] + o_{n,p \rightarrow +\infty} \left( \frac{1}{p} \right)$$

3.

$$\mathbb{E} \left[ \left( \frac{\mathbf{P}^\top \mathbf{Z}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} = -\mathbb{E} \left[ \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \mathcal{D}_{\mathbf{A}, \mathbf{R}}^{(i)} \mathbf{Q} \right]_{i,j} + o_{n,p \rightarrow +\infty} \left( \frac{1}{p} \right)$$

605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659



where  $\mathcal{D}_{\mathbf{A},\mathbf{R}}^{(i)}$  is a diagonal matrix such that  $[\mathcal{D}_{\mathbf{A},\mathbf{R}}^{(i)}]_{k,k} = \frac{1}{p} \sum_{l=1}^n \mathbf{Q}_{l,l} \mathbf{A}_{i,l} \mathbf{R}_{l,k} = \frac{1}{p} \text{tr} \mathbf{D}_{\mathbf{R},\cdot,k} \mathbf{Q} \mathbf{D}_{\mathbf{A}_{i,\cdot}}$ .

*Proof.* We start with the first equation.

$$\begin{aligned} \mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} &= \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} [\mathbf{Z}_{s,i} \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \mathbf{Q}_{r,j}] \\ &= \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \frac{\partial (\mathbf{Z}_{s,r} \mathbf{Q}_{r,j})}{\partial \mathbf{Z}_{s,i}} \mathbf{A}_{i,r} \right] \quad \text{using Stein's lemma} \\ &= \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \frac{\partial \mathbf{Z}_{s,r}}{\partial \mathbf{Z}_{s,i}} \mathbf{A}_{i,r} \mathbf{Q}_{r,j} + \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \frac{\partial \mathbf{Q}_{r,j}}{\partial \mathbf{Z}_{s,i}} \right] \\ \mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} &= \mathbb{E} [\mathbf{A}_{i,i} \mathbf{Q}_{i,j}] + \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \frac{\partial \mathbf{Q}_{r,j}}{\partial \mathbf{Z}_{s,i}} \right]. \end{aligned}$$

From proposition A.1, we know that

$$\frac{\partial \mathbf{Q}_{r,j}}{\partial \mathbf{Z}_{s,i}} = -\frac{1}{p} \left( [\mathbf{X} \mathbf{D}_{\mathbf{R},\cdot,i} \mathbf{Q}]_{s,r} \mathbf{Q}_{i,j} + \mathbf{Q}_{r,i} [\mathbf{X} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q}]_{s,j} \right)$$

therefore,

$$\frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \frac{\partial \mathbf{Q}_{r,j}}{\partial \mathbf{Z}_{s,i}} \right] = -\frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \frac{1}{p} \mathbf{Z}_{s,r} \mathbf{A}_{i,r} [\mathbf{X} \mathbf{D}_{\mathbf{R},\cdot,i} \mathbf{Q}]_{s,r} \mathbf{Q}_{i,j} \right] - \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \frac{1}{p} \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \mathbf{Q}_{r,i} [\mathbf{X} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q}]_{s,j} \right].$$

Moreover,

$$\begin{aligned} \sum_{r=1}^n \sum_{s=1}^p \mathbf{Z}_{s,r} \mathbf{A}_{i,r} [\mathbf{X} \mathbf{D}_{\mathbf{R},\cdot,i} \mathbf{Q}]_{s,r} &= \sum_{r=1}^n \sum_{s=1}^p \sum_{t=1}^n \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \mathbf{X}_{s,t} \mathbf{R}_{t,i} \mathbf{Q}_{t,r} \\ &= \sum_{r=1}^n \sum_{t=1}^n \mathbf{Q}_{t,r} [\mathbf{X}^\top \mathbf{Z}]_{t,r} \mathbf{R}_{t,i} \mathbf{A}_{i,r} \\ \sum_{r=1}^n \sum_{s=1}^p \mathbf{Z}_{s,r} \mathbf{A}_{i,r} [\mathbf{X} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q}]_{s,r} &= \text{tr} \left( \mathbf{Q} (\mathbf{X}^\top \mathbf{Z} \odot [\mathbf{R}_{\cdot,i} \mathbf{A}_{i,\cdot}]) \right) \end{aligned}$$

and

$$\begin{aligned} \sum_{r=1}^n \sum_{s=1}^p \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \mathbf{Q}_{r,i} [\mathbf{X} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q}]_{s,j} &= \sum_{r=1}^n \mathbf{Q}_{i,r} \mathbf{A}_{i,r} [\mathbf{Z}^\top \mathbf{X} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q}]_{r,j} \\ \sum_{r=1}^n \sum_{s=1}^p \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \mathbf{Q}_{r,i} [\mathbf{X} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q}]_{s,j} &= [\mathbf{Q} \mathbf{D}_{\mathbf{A}_{i,\cdot}} \mathbf{Z}^\top \mathbf{X} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q}]_{i,j}. \end{aligned}$$

So we finally have

$$\mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} = \mathbb{E} [\mathbf{A}_{i,i} \mathbf{Q}_{i,j}] - \frac{1}{p} \mathbb{E} \left[ \text{tr} \left( \mathbf{Q} \left( \frac{\mathbf{X}^\top \mathbf{Z}}{p} \odot [\mathbf{R}_{\cdot,i} \mathbf{A}_{i,\cdot}] \right) \right) \mathbf{Q}_{i,j} \right] - \underbrace{\frac{1}{p} \mathbb{E} \left[ \mathbf{Q} \mathbf{D}_{\mathbf{A}_{i,\cdot}} \frac{\mathbf{Z}^\top \mathbf{X}}{p} \mathbf{D}_{\mathbf{R},i,\cdot} \mathbf{Q} \right]_{i,j}}_{= \mathcal{O}_{n,p \rightarrow +\infty} \left( \frac{1}{p} \right)}.$$

The second equation can be shown in the same way.

$$\begin{aligned}
 \mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{P}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} &= \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} [\mathbf{Z}_{s,i} \mathbf{P}_{s,r} \mathbf{A}_{i,r} \mathbf{Q}_{r,j}] \\
 &= \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \mathbf{P}_{s,r} \mathbf{A}_{i,r} \frac{\partial \mathbf{Q}_{r,j}}{\partial \mathbf{Z}_{s,i}} \right] \quad \text{using Stein's lemma} \\
 \mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{P}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} &= -\frac{1}{p} \mathbb{E} \left[ \text{tr} \left( \mathbf{Q} \left( \frac{\mathbf{X}^\top \mathbf{P}}{p} \odot [\mathbf{R}_{\cdot,i} \mathbf{A}_{i,\cdot}] \right) \right) \mathbf{Q}_{i,j} \right] - \frac{1}{p} \mathbb{E} \left[ \underbrace{\mathbf{Q} \mathbf{D}_{\mathbf{A}_{i,\cdot}} \frac{\mathbf{P}^\top \mathbf{X}}{p} \mathbf{D}_{\mathbf{R}_{i,\cdot}} \mathbf{Q}}_{=\mathcal{O}_{n,p \rightarrow +\infty} \left( \frac{1}{p} \right)} \right]_{i,j}.
 \end{aligned}$$

We are left to show the third equation.

$$\begin{aligned}
 \mathbb{E} \left[ \left( \frac{\mathbf{P}^\top \mathbf{Z}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} &= \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} [\mathbf{P}_{s,i} \mathbf{Z}_{s,r} \mathbf{A}_{i,r} \mathbf{Q}_{r,j}] \\
 &= \frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \mathbf{P}_{s,i} \mathbf{A}_{i,r} \frac{\partial \mathbf{Q}_{r,j}}{\partial \mathbf{Z}_{s,r}} \right] \quad \text{using Stein's lemma} \\
 &= -\frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \mathbb{E} \left[ \frac{1}{p} \mathbf{P}_{s,i} \mathbf{A}_{i,r} \left( [\mathbf{X} \mathbf{D}_{\mathbf{R}_{\cdot,r}} \mathbf{Q}]_{s,r} \mathbf{Q}_{r,j} + \mathbf{Q}_{r,r} [\mathbf{X} \mathbf{D}_{\mathbf{R}_{\cdot,r}} \mathbf{Q}]_{s,j} \right) \right] \\
 &= -\frac{1}{p} \sum_{r=1}^n \sum_{s=1}^p \sum_{t=1}^n \mathbb{E} \left[ \frac{1}{p} \mathbf{P}_{s,i} \mathbf{A}_{i,r} (\mathbf{X}_{s,t} \mathbf{R}_{t,r} \mathbf{Q}_{t,r} \mathbf{Q}_{r,j} + \mathbf{Q}_{r,r} \mathbf{X}_{s,t} \mathbf{R}_{r,t} \mathbf{Q}_{t,j}) \right] \\
 &= -\frac{1}{p} \sum_{r=1}^n \sum_{t=1}^n \mathbb{E} \left[ \frac{1}{p} [\mathbf{P}^\top \mathbf{X}]_{i,t} \mathbf{A}_{i,r} (\mathbf{R}_{t,r} \mathbf{Q}_{t,r} \mathbf{Q}_{r,j} + \mathbf{Q}_{r,r} \mathbf{R}_{r,t} \mathbf{Q}_{t,j}) \right] \\
 &= -\frac{1}{p} \sum_{t=1}^n \mathbb{E} \left[ [\mathbf{P}^\top \mathbf{X}]_{i,t} \left( \frac{1}{p} [(\mathbf{Q} \odot \mathbf{R}) \mathbf{D}_{\mathbf{A}_{i,\cdot}} \mathbf{Q}]_{t,j} + [\mathcal{D}_{\mathbf{A},\mathbf{R}}^{(i)}]_{t,t} \mathbf{Q}_{t,j} \right) \right] \\
 \mathbb{E} \left[ \left( \frac{\mathbf{P}^\top \mathbf{Z}}{p} \odot \mathbf{A} \right) \mathbf{Q} \right]_{i,j} &= -\mathbb{E} \left[ \frac{\mathbf{P}^\top \mathbf{X}}{p} \mathcal{D}_{\mathbf{A},\mathbf{R}}^{(i)} \mathbf{Q} \right]_{i,j} - \frac{1}{p} \mathbb{E} \left[ \underbrace{\frac{\mathbf{P}^\top \mathbf{X}}{p} (\mathbf{Q} \odot \mathbf{R}) \mathbf{D}_{\mathbf{A}_{i,\cdot}} \mathbf{Q}}_{=\mathcal{O}_{n,p \rightarrow +\infty} \left( \frac{1}{p} \right)} \right]_{i,j}.
 \end{aligned}$$

□

## B. Proof of theorem 3.1

The study the spectral behavior of  $\tilde{\mathbf{K}}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{C}$  is made through its resolvent

$$\mathbf{Q} = \left( \tilde{\mathbf{K}}_L - z \mathbf{I}_n \right)^{-1}$$

where we have dropped the dependence in  $z$  to ease notations.

### B.1. Analysis of the model with noise only: $\mathbf{X} = \mathbf{Z}$

In order to find the limiting spectral distribution of  $\tilde{\mathbf{K}}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{C}$ , we first consider the simpler model with noise only, i.e.,

$$\mathbf{X} = \mathbf{Z}, \quad \tilde{\mathbf{K}}_L = \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C}, \quad \mathbf{Q} = \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} - z \mathbf{I}_n \right)^{-1}.$$

## B.1.1. A FIRST EQUIVALENT OF THE RESOLVENT

Let us first consider the following expression of the resolvent

$$\mathbf{Q} = -\frac{1}{z}\mathbf{I}_n + \frac{1}{z} \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} \right) \mathbf{Q}$$

which is a rewriting of  $\mathbf{Q}^{-1}\mathbf{Q} = \mathbf{I}_n$ .

In order to find a deterministic equivalent of  $\mathbf{Q}$ , we study its expected value

$$\mathbb{E}[\mathbf{Q}_{i,j}] = -\frac{1}{z}\delta_{i,j} + \frac{1}{z}\mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} \right) \mathbf{Q} \right]_{i,j}.$$

Taking  $\mathbf{A} = \mathbf{C}$  in the first equation of proposition A.3, we have

$$z\mathbb{E}[\mathbf{Q}_{i,j}] = -\delta_{i,j} + \mathbf{C}_{i,i}\mathbb{E}[\mathbf{Q}_{i,j}] - \mathbb{E}[\eta_{i,i}\mathbf{Q}_{i,j}] + \mathcal{O}_{n,p,L \rightarrow +\infty} \left( \frac{1}{p} \right)$$

with  $\eta \in \mathbb{C}^{n \times n}$  such that

$$\eta_{r,s} = \frac{1}{p} \text{tr} \left( \mathbf{Q} \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot [\mathbf{C}_{\cdot,r}\mathbf{C}_{s,\cdot}] \right) \right).$$

Thus, denoting  $\mathbf{D}_\eta = \eta \odot \mathbf{I}_n$ , we have the matrix equivalence  $z\mathbf{Q} \leftrightarrow -\mathbf{I}_n + \mathbf{Q} - \mathbf{D}_\eta\mathbf{Q}$  from which we deduce that the resolvent is equivalent to a diagonal matrix:

$$\mathbf{Q} \leftrightarrow (\mathbf{I}_n - z\mathbf{I}_n - \mathbf{D}_\eta)^{-1}.$$

 B.1.2. ANALYSIS OF THE MATRIX  $\eta$ 

Now taking  $\mathbf{A} = \mathbf{C}_{\cdot,r}\mathbf{C}_{s,\cdot}$  in the first equation of proposition A.3, we have

$$\begin{aligned} \mathbb{E}[\eta_{r,s}] &= \frac{1}{p} \sum_{t=1}^n \mathbf{C}_{t,r}\mathbf{C}_{s,t}\mathbb{E}[\mathbf{Q}_{t,t}] - \frac{1}{p} \sum_{t=1}^n \mathbb{E} \left[ \frac{1}{p} \text{tr} \left( \mathbf{Q} \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot [\mathbf{C}_{\cdot,t}\mathbf{C}_{s,\cdot}] \right) \right) \mathbf{C}_{t,r}\mathbf{Q}_{t,t} \right] + \mathcal{O}_{n,p,L \rightarrow +\infty} \left( \frac{1}{p} \right) \\ &= \frac{1}{p} \sum_{t=1}^n \mathbf{C}_{t,r} (\mathbf{C}_{s,t}\mathbb{E}[\mathbf{Q}_{t,t}] - \mathbb{E}[\eta_{t,s}\mathbf{Q}_{t,t}]) + \mathcal{O}_{n,p,L \rightarrow +\infty} \left( \frac{1}{p} \right) \end{aligned}$$

and, using the previous matrix equivalent of  $\mathbf{Q}$ , we can write  $\eta \leftrightarrow \bar{\eta}$  where  $\bar{\eta}$  is a deterministic matrix such that

$$\bar{\eta}_{r,s} = \frac{1}{p} \sum_{t=1}^n \mathbf{C}_{t,r} \frac{\mathbf{C}_{s,t} - \bar{\eta}_{t,s}}{1 - z - \bar{\eta}_{t,t}}.$$

Therefore,  $\bar{\eta}$  has a circulant structure. Indeed, for all  $d \in \mathbb{Z}$ ,

$$\bar{\eta}_{r+d,s+d} = \frac{1}{p} \sum_{t=r-L+1}^{r+L-1} \frac{\mathbf{C}_{s+d,t+d} - \bar{\eta}_{t+d,s+d}}{1 - z - \bar{\eta}_{t+d,t+d}} \quad d \in \mathbb{Z}.$$

where we write  $\bar{\eta}_{i,j}$  for any  $i, j \in \mathbb{Z}$  to represent  $\bar{\eta}_{(i \bmod n), (j \bmod n)}$ .

 B.1.3. FROM  $\bar{\eta}$  TO THE LIMITING SPECTRAL DISTRIBUTION

Since  $\bar{\eta}$  is circulant, it has a constant diagonal:  $\bar{\eta}_{k,k} = \eta_0$ . Then, we can recognize a matrix product in the expression of  $\bar{\eta}_{r,s}$ :

$$\bar{\eta}_{r,s} = \frac{1}{p} \frac{1}{1 - z - \eta_0} \sum_{t=1}^n \mathbf{C}_{t,r} (\mathbf{C}_{s,t} - \bar{\eta}_{t,s}) = \frac{1}{p} \frac{1}{1 - z - \eta_0} [\mathbf{C}(\mathbf{C} - \bar{\eta})]_{r,s}$$

thus,  $\bar{\eta} = (p(1-z-\eta_0)\mathbf{I}_n + \mathbf{C})^{-1} \mathbf{C}^2$  and, since  $\eta_0 = \frac{1}{n} \text{tr } \bar{\eta}$ ,

$$\eta_0 = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\psi_k^2}{p(1-z-\eta_0) + \psi_k}.$$

Recalling that  $\mathbf{Q} \leftrightarrow (\mathbf{I}_n - z\mathbf{I}_n - \mathbf{D}_\eta)^{-1}$ , we can state the following theorem.

**Theorem B.1** (Deterministic equivalent of  $\mathbf{Q}$  when  $\mathbf{X} = \mathbf{Z}$ ). *Let  $z \in \mathbb{C} \setminus \limsup_{n,p,L \rightarrow +\infty} \text{Sp}(\tilde{\mathbf{K}}_L)$ . Then,*

$$\mathbf{Q} \leftrightarrow m(z)\mathbf{I}_n \quad \text{with} \quad m(z) = \frac{1}{1-z-\eta_0}$$

and  $\eta_0$  is solution to the fixed-point equation

$$\eta_0 = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\psi_k^2}{p(1-z-\eta_0) + \psi_k}.$$

$m$  is the Stieljes transform of the limiting spectral distribution of  $\tilde{\mathbf{K}}_L = \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C}$ .

*Remark B.2.* Notice that  $m(z) \neq 0$ . Indeed, for a given  $z \in \mathbb{C} \setminus \limsup_{n,p,L \rightarrow +\infty} \text{Sp}(\tilde{\mathbf{K}}_L)$ ,  $(1-z-\eta_0)m(z) = 1$  and the fixed-point equation prevent  $\eta_0$  from going to  $\pm\infty$ .

**Proposition B.3** (Fixed-point equation for  $m(z)$ ). *Under the setting of theorem B.1,  $m(z)$  is also solution to a fixed-point equation:*

$$1 + zm(z) = \frac{1}{n} \sum_{k=0}^{n-1} \frac{m(z)\psi_k}{1 + \frac{m(z)}{p}\psi_k}.$$

*Proof.* A rewriting of  $\bar{\eta} = (p(1-z-\eta_0)\mathbf{I}_n + \mathbf{C})^{-1} \mathbf{C}^2$  yields another interesting formula:

$$\begin{aligned} \bar{\eta} &= \left( \frac{p}{m(z)}\mathbf{I}_n + \mathbf{C} \right)^{-1} \left( \frac{p}{m(z)}\mathbf{I}_n + \mathbf{C} - \frac{p}{m(z)}\mathbf{I}_n \right) \mathbf{C} \\ \bar{\eta} &= \mathbf{C} - \left( \mathbf{I}_n + \frac{m(z)}{p}\mathbf{C} \right)^{-1} \mathbf{C} \end{aligned}$$

therefore,

$$\eta_0 = \underbrace{\frac{1}{n} \text{tr } \mathbf{C}}_{=1} - \frac{1}{n} \sum_{k=0}^{n-1} \frac{\psi_k}{1 + \frac{m(z)}{p}\psi_k}$$

and, since  $1 - \eta_0 = \frac{1}{m(z)} + z$ , we get the result.  $\square$

## B.2. Analysis of the full model: $\mathbf{X} = \mathbf{P} + \mathbf{Z}$

So far, we have been able to find a deterministic equivalent of the resolvent under the setting where  $\mathbf{X} = \mathbf{Z}$ , i.e., when the observations are composed of noise only.

Now, we consider the setting where the observations are composed of a signal corrupted with additive noise:

$$\mathbf{X} = \mathbf{P} + \mathbf{Z}, \quad \tilde{\mathbf{K}}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{C}, \quad \mathbf{Q} = \left( \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1}.$$

Let us first prove that the limiting spectral distribution is unchanged.

**Proposition B.4.**

$$\left| \frac{1}{n} \text{tr} \left( \frac{(\mathbf{P} + \mathbf{Z})^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1} - \frac{1}{n} \text{tr} \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1} \right| \xrightarrow[n,p,L \rightarrow +\infty]{a.s.} 0.$$



Proof.

$$\frac{(\mathbf{P} + \mathbf{Z})^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} = \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} + \mathbf{A} \odot \mathbf{C}$$

with  $\mathbf{A} = \frac{\mathbf{Z}^\top \mathbf{P}}{p} + \frac{\mathbf{P}^\top \mathbf{Z}}{p} + \frac{\mathbf{P}^\top \mathbf{P}}{p}$ . Notice that,  $\frac{\mathbf{Z}^\top \mathbf{P}}{p}$ ,  $\frac{\mathbf{P}^\top \mathbf{Z}}{p}$  and  $\frac{\mathbf{P}^\top \mathbf{P}}{p}$  are uniformly bounded in spectral norm (from the non-triviality condition) and their rank is at most  $K$ . Thus  $\mathbf{A}$  is also uniformly bounded in spectral norm and has rank at most  $3K$ .

Let  $\mathbf{Q}_A = \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} + \mathbf{A} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1}$  and  $\mathbf{Q}_0 = \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1}$ . Using singular-value inequalities which can be found in theorems A.12 and A.14 of (Bai & Silverstein, 2010),

$$\begin{aligned} \left| \frac{1}{n} \operatorname{tr} \mathbf{Q}_A - \frac{1}{n} \operatorname{tr} \mathbf{Q}_0 \right| &= \frac{1}{n} |\operatorname{tr} \mathbf{Q}_0 (\mathbf{A} \odot \mathbf{C}) \mathbf{Q}_A| \\ &\leq \frac{1}{n} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \sum_{i=1}^n s_i(\mathbf{A} \odot \mathbf{C}) \quad \text{from theorems A.12 and A.14} \\ &\leq \frac{1}{n} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \sqrt{n \sum_{i=1}^n s_i^2(\mathbf{A} \odot \mathbf{C})} \quad \text{from Jensen's inequality} \\ &\leq \frac{1}{n} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \sqrt{n \sum_{i=1}^n s_i^2(\mathbf{A})} \quad \text{since } \|\mathbf{A} \odot \mathbf{C}\|_F \leq \|\mathbf{A}\|_F \end{aligned}$$

$$\left| \frac{1}{n} \operatorname{tr} \mathbf{Q}_A - \frac{1}{n} \operatorname{tr} \mathbf{Q}_0 \right| \leq \sqrt{\frac{3K}{n}} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \|\mathbf{A}\| = \mathcal{O}_{n,p,L \rightarrow +\infty} \left( \frac{1}{\sqrt{n}} \right) \quad \text{since } \mathbf{A} \text{ has rank at most } 3K.$$

□

Since  $\left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1} \leftrightarrow m(z)\mathbf{I}_n$  according to theorem B.1, proposition B.4 justifies that the limiting spectral distribution is unchanged by the presence of signal.

As previously, we consider a rewriting of  $\mathbf{Q}^{-1}\mathbf{Q} = \mathbf{I}_n$ ,

$$\mathbf{Q} = -\frac{1}{z}\mathbf{I}_n + \frac{1}{z} \left( \frac{(\mathbf{P} + \mathbf{Z})^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} \right) \mathbf{Q}$$

and we study the expected value of  $\mathbf{Q}_{i,j}$ ,

$$\mathbb{E} [\mathbf{Q}_{i,j}] = -\frac{1}{z}\delta_{i,j} + \frac{1}{z} \mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} \right) \mathbf{Q} \right]_{i,j} + \frac{1}{z} \mathbb{E} \left[ \left( \frac{\mathbf{Z}^\top \mathbf{P}}{p} \odot \mathbf{C} \right) \mathbf{Q} \right]_{i,j} + \frac{1}{z} \mathbb{E} \left[ \left( \frac{\mathbf{P}^\top \mathbf{Z}}{p} \odot \mathbf{C} \right) \mathbf{Q} \right]_{i,j} + \frac{1}{z} \mathbb{E} \left[ \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{C} \right) \mathbf{Q} \right]_{i,j}.$$

$\mathbf{P}^\top \mathbf{P}$  is deterministic so there is no work to do on the last term of the sum. Expanding the other terms yields (see proposition A.3)

$$z \mathbb{E} [\mathbf{Q}_{i,j}] = -\delta_{i,j} + \mathbb{E} [\mathbf{C}_{i,i} \mathbf{Q}_{i,j}] - \mathbb{E} [\kappa_{i,i} \mathbf{Q}_{i,j}] - \mathbb{E} \left[ \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \mathcal{D}_{\mathbf{C},\mathbf{C}}^{(i)} \mathbf{Q} \right]_{i,j} + \mathbb{E} \left[ \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{C} \right) \mathbf{Q} \right]_{i,j} + \mathcal{O}_{n,p,L \rightarrow +\infty} \left( \frac{1}{p} \right)$$

with  $\kappa \in \mathbb{C}^{n \times n}$  such that

$$\kappa_{r,s} = \frac{1}{p} \operatorname{tr} \left( \mathbf{Q} \left( \frac{(\mathbf{P} + \mathbf{Z})^\top (\mathbf{P} + \mathbf{Z})}{p} \odot [\mathbf{C}_{\cdot,r} \mathbf{C}_{s,\cdot}] \right) \right)$$

and  $\mathcal{D}_{\mathbf{C},\mathbf{C}}^{(i)}$  is a diagonal matrix such that  $[\mathcal{D}_{\mathbf{C},\mathbf{C}}^{(i)}]_{k,k} = \frac{1}{p} \sum_{l=1}^n \mathbf{Q}_{l,l} \mathbf{C}_{i,l} \mathbf{C}_{l,k}$ .

**Proposition B.5.**

$$\kappa \leftrightarrow \bar{\eta}.$$

*Proof.* Similarly to the proof of proposition B.4, we consider a matrix  $\mathbf{A}$  uniformly bounded in spectral norm whose rank is at most  $K$ , representing  $\frac{\mathbf{Z}^\top \mathbf{P}}{p}$ ,  $\frac{\mathbf{P}^\top \mathbf{Z}}{p}$  or  $\frac{\mathbf{P}^\top \mathbf{P}}{p}$  and we make use of singular-value inequalities.

$$\begin{aligned} \frac{1}{p} |\operatorname{tr}(\mathbf{Q}(\mathbf{A} \odot [\mathbf{C}_{\cdot,r} \mathbf{C}_{s,\cdot}]))| &= \frac{1}{p} |\operatorname{tr}(\mathbf{Q} \mathbf{D}_{\mathbf{C}_{\cdot,r}} \mathbf{A} \mathbf{D}_{\mathbf{C}_{s,\cdot}})| \\ &\leq \frac{1}{p} \sum_{i=1}^n s_i(\mathbf{Q} \mathbf{D}_{\mathbf{C}_{\cdot,r}} \mathbf{A} \mathbf{D}_{\mathbf{C}_{s,\cdot}}) \\ &\leq \frac{1}{p} \sum_{i=1}^n s_i(\mathbf{Q}) s_i(\mathbf{D}_{\mathbf{C}_{\cdot,r}} \mathbf{A} \mathbf{D}_{\mathbf{C}_{s,\cdot}}) \\ &\leq \frac{1}{p} \|\mathbf{Q}\| \sum_{i=1}^n s_i(\mathbf{D}_{\mathbf{C}_{\cdot,r}} \mathbf{A} \mathbf{D}_{\mathbf{C}_{s,\cdot}}) \\ &\leq \frac{K}{p} \|\mathbf{Q}\| \|\mathbf{D}_{\mathbf{C}_{\cdot,r}} \mathbf{A} \mathbf{D}_{\mathbf{C}_{s,\cdot}}\| \quad \text{since } \mathbf{A} \text{ has rank at most } K \\ \frac{1}{p} |\operatorname{tr}(\mathbf{Q}(\mathbf{A} \odot [\mathbf{C}_{\cdot,r} \mathbf{C}_{s,\cdot}]))| &\leq \frac{K}{p} \|\mathbf{Q}\| \|\mathbf{A}\| = \frac{\mathcal{O}}{n,p,L \rightarrow +\infty} \left(\frac{1}{p}\right) \quad \text{since } \|\mathbf{D}_{\mathbf{C}_{\cdot,r}}\| = \|\mathbf{D}_{\mathbf{C}_{s,\cdot}}\| = 1. \end{aligned}$$

Hence,

$$\kappa_{r,s} = \frac{1}{p} \operatorname{tr} \left( \underbrace{\mathbf{Q} \left( \frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot [\mathbf{C}_{\cdot,r} \mathbf{C}_{s,\cdot}] \right)}_{=\eta_{r,s}} \right) + \frac{\mathcal{O}}{n,p,L \rightarrow +\infty} \left(\frac{1}{p}\right)$$

□

So far, we have

$$z\mathbf{Q} \leftrightarrow -\mathbf{I}_n + \mathbf{Q} - \eta_0 \mathbf{Q} + \left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} \right) \mathbf{Q} \quad \text{i.e.} \quad \mathbf{Q} \leftrightarrow \left( \frac{1}{m(z)} \mathbf{I}_n + \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} \right)^{-1}.$$

The analysis of  $\left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} \right) \mathbf{Q}$  is summarized in the following proposition.

**Proposition B.6.**

$$\left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} \right) \mathbf{Q} \leftrightarrow \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{F} \Psi \left( \mathbf{I}_n + \frac{m(z)}{p} \Psi \right)^{-1} \mathbf{F}^* \right) \mathbf{Q}.$$

*Proof.* From assumption 2.1, all diagonal entries of  $\mathbf{Q}$  are statistically equivalent. Thus, we can have a simple matrix equivalent of  $\mathcal{D}_{\mathbf{C}^t, \mathbf{C}}^{(i)}$  for all integer  $t \geq 1$ :

$$\left[ \mathcal{D}_{\mathbf{C}^t, \mathbf{C}}^{(i)} \right]_{k,k} = \frac{1}{p} \sum_{l=1}^n \mathbf{Q}_{l,l} [\mathbf{C}^t]_{i,l} \mathbf{C}_{l,k} \leftrightarrow \frac{1}{p} \frac{\operatorname{tr} \mathbf{Q}}{n} \sum_{l=1}^n [\mathbf{C}^t]_{i,l} \mathbf{C}_{l,k} \leftrightarrow \frac{m(z)}{p} [\mathbf{C}^{t+1}]_{i,k}$$

where the last equivalence is justified by proposition B.4.

Now, we can notice the following recurrence relation

$$\begin{aligned} \left[ \left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C}^t \right) \mathbf{Q} \right]_{i,j} &\leftrightarrow - \left[ \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \mathcal{D}_{\mathbf{C}^t, \mathbf{C}}^{(i)} \mathbf{Q} \right]_{i,j} + \left[ \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{C}^t \right) \mathbf{Q} \right]_{i,j} \\ &\leftrightarrow - \frac{m(z)}{p} \sum_{k=1}^n \left[ \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \right]_{i,k} [\mathbf{C}^{t+1}]_{i,k} \mathbf{Q}_{k,j} + \left[ \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{C}^t \right) \mathbf{Q} \right]_{i,j} \\ \left[ \left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C}^t \right) \mathbf{Q} \right]_{i,j} &\leftrightarrow - \frac{m(z)}{p} \left[ \left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C}^{t+1} \right) \mathbf{Q} \right]_{i,j} + \left[ \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{C}^t \right) \mathbf{Q} \right]_{i,j}. \end{aligned}$$

In particular, for all integer  $T \geq 1$ ,

$$\left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} \right) \mathbf{Q} \leftrightarrow \left( -\frac{m(z)}{p} \right)^T \left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C}^{T+1} \right) \mathbf{Q} + \sum_{t=0}^{T-1} \left( -\frac{m(z)}{p} \right)^t \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{C}^{t+1} \right) \mathbf{Q}.$$

We know that  $\|\mathbf{C}\| = (2L - 1)$  and, using the fact that the spectral norm of a pointwise product (as well as a regular matrix product) can be bounded by the product of the spectral norms (see theorem A.19 of (Bai & Silverstein, 2010)), we have

$$\left\| \left( -\frac{m(z)}{p} \right)^T \left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C}^{T+1} \right) \mathbf{Q} \right\| \leq \left| \frac{m(z)}{p} \right|^T \left\| \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \right\| |2L - 1|^{T+1} \|\mathbf{Q}\|.$$

Thus, if  $\left| \frac{2L-1}{p} m(z) \right| < 1$ ,

$$\left( \frac{\mathbf{P}^\top (\mathbf{P} + \mathbf{Z})}{p} \odot \mathbf{C} \right) \mathbf{Q} \leftrightarrow \left( \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \left[ \sum_{t=0}^{+\infty} \left( -\frac{m(z)}{p} \right)^t \mathbf{C}^{t+1} \right] \right) \mathbf{Q}.$$

And, since  $\mathbf{C} = \mathbf{F}\Psi\mathbf{F}^*$ ,

$$\sum_{t=0}^{+\infty} \left( -\frac{m(z)}{p} \right)^t \mathbf{C}^{t+1} = \mathbf{F}\Psi \left( \mathbf{I}_n + \frac{m(z)}{p} \Psi \right)^{-1} \mathbf{F}^*$$

which completes the proof.  $\square$

We can now state the following theorem.

**Theorem B.7** (Deterministic equivalent of  $\mathbf{Q}$  when  $\mathbf{X} = \mathbf{P} + \mathbf{Z}$ ). *Let  $z \in \mathbb{C} \setminus \limsup_{n,p,L \rightarrow +\infty} \text{Sp}(\tilde{\mathbf{K}}_L)$ . If  $\left| \frac{2L-1}{p} m(z) \right| < 1$ , then*

$$\mathbf{Q} \leftrightarrow m(z) \left( \mathbf{I}_n + \frac{\mathbf{P}^\top \mathbf{P}}{p} \odot \mathbf{F}\Lambda\mathbf{F}^* \right)^{-1} \quad \text{with} \quad \Lambda = m(z) \Psi \left( \mathbf{I}_n + \frac{m(z)}{p} \Psi \right)^{-1}.$$

*Remark B.8.* This is coherent with the result of theorem B.1 when  $\mathbf{P} = \mathbf{0}$ .

*Remark B.9.* From proposition B.3, we see that  $1 + zm(z) = \frac{1}{n} \text{tr} \Lambda$ .

### C. Proof of theorem 3.3

In this section, we use the following notation:

$$\text{Sp}_\infty(\tilde{\mathbf{K}}_L) = \limsup_{n,p,L \rightarrow +\infty} \text{Sp}(\tilde{\mathbf{K}}_L)$$

#### C.1. Spikes

Here,  $\mathbf{P} = \mu \mathbf{j}^\top$ . Let us state a much more tractable expression of the deterministic equivalent of the resolvent.

**Theorem C.1** (Deterministic equivalent of  $\mathbf{Q}$  when  $\mathbf{P} = \mu \mathbf{j}^\top$ ). *Let  $z \in \mathbb{C} \setminus \text{Sp}_\infty(\tilde{\mathbf{K}}_L)$ . If  $\left| \frac{2L-1}{p} m(z) \right| < 1$ , then*

$$\mathbf{Q} \leftrightarrow \bar{\mathbf{Q}} = m(z) [\mathbf{D}_j \mathbf{F}] \left( \mathbf{I}_n + \frac{\|\mu\|^2}{p} \Lambda \right)^{-1} [\mathbf{D}_j \mathbf{F}]^*$$

where  $\mathbf{D}_j$  is the diagonal matrix induced by vector  $\mathbf{j}$ .

1045 *Proof.* From theorem B.7,

$$\begin{aligned}
 1046 \quad \mathbf{Q} &\leftrightarrow m(z) \left( \mathbf{I}_n + \frac{\|\boldsymbol{\mu}\|^2}{p} \mathbf{j}\mathbf{j}^\top \odot \mathbf{F}\boldsymbol{\Lambda}\mathbf{F}^* \right)^{-1} \\
 1047 \quad & \\
 1048 \quad &\leftrightarrow m(z) \left( \mathbf{I}_n + \frac{\|\boldsymbol{\mu}\|^2}{p} \mathbf{D}_j \mathbf{F} \boldsymbol{\Lambda} \mathbf{F}^* \mathbf{D}_j \right)^{-1} \\
 1049 \quad & \\
 1050 \quad & \\
 1051 \quad & \\
 1052 \quad & \\
 1053 \quad \mathbf{Q} &\leftrightarrow m(z) [\mathbf{D}_j \mathbf{F}] \left( \mathbf{I}_n + \frac{\|\boldsymbol{\mu}\|^2}{p} \boldsymbol{\Lambda} \right)^{-1} [\mathbf{D}_j \mathbf{F}]^* \quad \text{since } \mathbf{D}_j \mathbf{F} \text{ is a unitary matrix.} \\
 1054 \quad & \\
 1055 \quad & \\
 1056 \quad & \square
 \end{aligned}$$

1058 The sought-after *spikes* which encapsulate the information about our data are the singular points of the resolvent. Therefore, their asymptotical position is given by the solution in  $z$  to

$$1059 \quad 1 + \frac{\|\boldsymbol{\mu}\|^2}{p} \frac{m(z)\psi_k}{1 + \frac{m(z)}{p}\psi_k} = 0 \quad 0 \leq k < n.$$

1060 Since  $\tilde{\mathbf{K}}_L$  is symmetric, all solutions are real. Moreover, there cannot be any spike inside  $\text{Sp}_\infty(\tilde{\mathbf{K}}_L)$  (the eigenvalue must be isolated). Therefore, we are only interested in solutions outside  $\text{Sp}_\infty(\tilde{\mathbf{K}}_L)$ .

1061 If  $\psi_k = 0$ , there is no solution, whereas if  $\psi_k \neq 0$ ,

$$1062 \quad 1 + \frac{\|\boldsymbol{\mu}\|^2}{p} \frac{m(z)\psi_k}{1 + \frac{m(z)}{p}\psi_k} = 0 \iff m(z) = \frac{-1}{\frac{\|\boldsymbol{\mu}\|^2+1}{p}\psi_k}$$

1063 and, supposing  $z \in \mathbb{C} \setminus \text{Sp}_\infty(\tilde{\mathbf{K}}_L)$ , we have, from proposition B.3,

$$1064 \quad z = \frac{\|\boldsymbol{\mu}\|^2+1}{p}\psi_k + \frac{1}{n} \sum_{l=0}^{n-1} \frac{\psi_l}{1 - \frac{\psi_l}{(\|\boldsymbol{\mu}\|^2+1)\psi_k}}.$$

1065 **Proposition C.2** (Singular points of  $\bar{\mathbf{Q}}$ ). *Let*

$$1066 \quad \bar{\xi}_k = \frac{\|\boldsymbol{\mu}\|^2+1}{p}\psi_k \left( 1 + \frac{p}{n} \sum_{l=0}^{n-1} \frac{\psi_l}{(\|\boldsymbol{\mu}\|^2+1)\psi_k - \psi_l} \right) \quad 0 \leq k < n.$$

1067 *The set of singular points of  $\bar{\mathbf{Q}}$  is  $\{\bar{\xi}_k \mid 0 \leq k < n, \psi_k \neq 0\} \cap (\mathbb{C} \setminus \text{Sp}_\infty(\tilde{\mathbf{K}}_L))$ .*

## 1068 C.2. Alignments and phase transition

1069 Let us denote  $\{(\xi_k, \mathbf{u}_k)\}_{0 \leq k < n}$  the pairs eigenvalue-eigenvector of  $\tilde{\mathbf{K}}_L$ . From the definition of the resolvent, we know that

$$1070 \quad \mathbf{Q} = \sum_{l=0}^{n-1} \frac{\mathbf{u}_l \mathbf{u}_l^*}{\xi_l - z}.$$

1071 Therefore, with Cauchy's integral formula and a positively oriented closed contour  $\Gamma_k$  circling around  $\xi_k$  and leaving the other eigenvalues outside, we can have access to the quantity

$$1072 \quad \sum_{\substack{0 \leq l \leq n-1 \\ \xi_l = \xi_k}} \mathbf{u}_l \mathbf{u}_l^* = -\frac{1}{2i\pi} \oint_{\Gamma_k} \mathbf{Q}(z) dz$$

1073



which is simply  $\mathbf{u}_k \mathbf{u}_k^*$  when the associated eigenvalue has multiplicity one. Then, we can calculate the alignment of any vector  $\mathbf{v} \in \mathbb{C}^n$  with the eigenspace associated to  $\xi_k$ :

$$\sum_{\substack{0 \leq l \leq n-1 \\ \xi_l = \xi_k}} |\mathbf{v}^* \mathbf{u}_l|^2 = -\frac{1}{2i\pi} \oint_{\Gamma_k} \mathbf{v}^* \mathbf{Q}(z) \mathbf{v} dz.$$

Using the deterministic equivalent of  $\mathbf{Q}$ , we have the following result.

**Proposition C.3** (Spike alignments). *For  $0 \leq k < n$  such that  $\bar{\xi}_k$  is a singular point of  $\bar{\mathbf{Q}}$ , let  $\Gamma_k$  be a positively oriented closed contour circling around  $\bar{\xi}_k$  and leaving all the  $\xi_l \neq \bar{\xi}_k$  outside.*

$$-\frac{1}{2i\pi} \oint_{\Gamma_k} \bar{\mathbf{Q}}(z) dz = \bar{\zeta}_k [\mathbf{D}_j \mathbf{F}] \mathcal{D}_k [\mathbf{D}_j \mathbf{F}]^*$$

where

$$\bar{\zeta}_k = \frac{\|\boldsymbol{\mu}\|^2}{\|\boldsymbol{\mu}\|^2 + 1} \left( 1 - \frac{p}{n} \sum_{l=0}^{n-1} \left[ \frac{\psi_l}{(\|\boldsymbol{\mu}\|^2 + 1) \psi_k - \psi_l} \right]^2 \right) \quad \text{and} \quad \mathcal{D}_k = \text{diag}(\mathbf{1}_{\psi_k = \psi_l})_{0 \leq l < n}.$$

*Proof.* By residue calculus,

$$-\frac{1}{2i\pi} \oint_{\Gamma_k} \bar{\mathbf{Q}}(z) dz = -[\mathbf{D}_j \mathbf{F}] \left[ \lim_{z \rightarrow \bar{\xi}_k} (z - \bar{\xi}_k) m(z) \left( \mathbf{I}_n + \frac{\|\boldsymbol{\mu}\|^2}{p} \Lambda(z) \right)^{-1} \right] [\mathbf{D}_j \mathbf{F}]^*.$$

Let  $0 \leq l < n$ . If  $\psi_l \neq \psi_k$ , then

$$\lim_{z \rightarrow \bar{\xi}_k} \frac{(z - \bar{\xi}_k) m(z)}{1 + \frac{\|\boldsymbol{\mu}\|^2}{p} \frac{m(z) \psi_l}{1 + \frac{m(z)}{p} \psi_l}} = 0$$

whereas if  $\psi_l = \psi_k$ , L'Hôpital's rule yields

$$\begin{aligned} \lim_{z \rightarrow \bar{\xi}_k} \frac{(z - \bar{\xi}_k) m(z)}{1 + \frac{\|\boldsymbol{\mu}\|^2}{p} \frac{m(z) \psi_l}{1 + \frac{m(z)}{p} \psi_l}} &= \frac{m(\bar{\xi}_k)}{\frac{d}{dz} \left[ 1 + \frac{\|\boldsymbol{\mu}\|^2}{p} \frac{m(z) \psi_l}{1 + \frac{m(z)}{p} \psi_l} \right]_{z=\bar{\xi}_k}} \\ &= \frac{m(\bar{\xi}_k) \left( 1 + \frac{m(\bar{\xi}_k)}{p} \psi_l \right)^2}{\frac{\|\boldsymbol{\mu}\|^2}{p} m'(\bar{\xi}_k) \psi_l}. \end{aligned}$$

Recalling that  $m(\bar{\xi}_k) = \frac{-1}{\frac{\|\boldsymbol{\mu}\|^2 + 1}{p} \psi_k}$ , we have  $1 + \frac{m(\bar{\xi}_k)}{p} \psi_k = \frac{\|\boldsymbol{\mu}\|^2}{\|\boldsymbol{\mu}\|^2 + 1}$ . Hence,

$$\begin{aligned} \lim_{z \rightarrow \bar{\xi}_k} \frac{(z - \bar{\xi}_k) m(z)}{1 + \frac{\|\boldsymbol{\mu}\|^2}{p} \frac{m(z) \psi_l}{1 + \frac{m(z)}{p} \psi_l}} &= \frac{\|\boldsymbol{\mu}\|^2}{\|\boldsymbol{\mu}\|^2 + 1} \frac{1}{\frac{\|\boldsymbol{\mu}\|^2 + 1}{p} \psi_k} \frac{m(\bar{\xi}_k)}{m'(\bar{\xi}_k)} \\ &= -\frac{\|\boldsymbol{\mu}\|^2}{\|\boldsymbol{\mu}\|^2 + 1} \frac{m^2(\bar{\xi}_k)}{m'(\bar{\xi}_k)}. \end{aligned}$$

Let us calculate an expression of  $\frac{m^2(\bar{\xi}_k)}{m'(\bar{\xi}_k)}$ . Differentiating in  $z$  the fixed-point equation of proposition B.3 yields

$$m(z) + z m'(z) = \frac{1}{n} \sum_{r=0}^{n-1} \frac{m'(z) \psi_r}{(1 + m(z) \psi_r / p)^2}$$

1155 thus,

$$\begin{aligned}
1156 \quad \frac{m^2(\bar{\xi}_k)}{m'(\bar{\xi}_k)} &= -\bar{\xi}_k m(\bar{\xi}_k) + \frac{1}{n} \sum_{r=0}^{n-1} \frac{m(\bar{\xi}_k) \psi_r}{(1 + m(\bar{\xi}_k) \psi_r/p)^2} \\
1157 \quad &= 1 - \frac{1}{n} \sum_{r=0}^{n-1} \frac{m(\bar{\xi}_k) \psi_r}{1 + m(\bar{\xi}_k) \psi_r/p} + \frac{1}{n} \sum_{r=0}^{n-1} \frac{m(\bar{\xi}_k) \psi_r}{(1 + m(\bar{\xi}_k) \psi_r/p)^2} \quad \text{from proposition B.3} \\
1158 \quad & \\
1159 \quad \frac{m^2(\bar{\xi}_k)}{m'(\bar{\xi}_k)} &= 1 - \frac{p}{n} \sum_{r=0}^{n-1} \left[ \frac{m(\bar{\xi}_k) \psi_r/p}{1 + m(\bar{\xi}_k) \psi_r/p} \right]^2
\end{aligned}$$

1165 and we just need to remember that  $m(\bar{\xi}_k) = \frac{-1}{\frac{\|\mu\|^2+1}{p} \psi_k}$  to get the result.  $\square$

1168 We can now state the following proposition which defines the phase transition position as the value of  $\|\mu\|$  at which  $\bar{\zeta}_k$  changes sign.

1170 **Proposition C.4** (Phase transition). For  $0 \leq k < n$ ,

$$1172 \quad \bar{\xi}_k \text{ is a singular point of } \bar{\mathbf{Q}} \iff \bar{\zeta}_k > 0.$$

1174 *Proof.* Let us consider a singular point  $\bar{\xi}_k$  of  $\bar{\mathbf{Q}}$ .

1175 As a Stieljes transform,  $m$  is increasing on all connected components of  $\mathbb{R} \setminus \text{Sp}_\infty(\tilde{\mathbf{K}}_L)$  and the restriction of its functional inverse  $z(\cdot)$  to the real line, here denoted  $x(\cdot)$ , is also growing on every connected component of  $m(\mathbb{R} \setminus \text{Sp}_\infty(\tilde{\mathbf{K}}_L))$ . Then, as  $\bar{\xi}_k$  is outside  $\text{Sp}_\infty(\tilde{\mathbf{K}}_L)$ , it implies  $x' \left( \frac{-1}{\frac{\|\mu\|^2+1}{p} \psi_k} \right) > 0$ .

1180 We have

$$\begin{aligned}
1181 \quad x(m) &= -\frac{1}{m} + \frac{1}{n} \sum_{l=0}^{n-1} \frac{\psi_l}{1 + m \psi_l/p} \\
1182 \quad x'(m) &= \frac{1}{m^2} - \frac{p}{n} \sum_{l=0}^{n-1} \left[ \frac{\psi_l/p}{1 + m \psi_l/p} \right]^2
\end{aligned}$$

1187 thus

$$1188 \quad x' \left( \frac{-1}{\frac{\|\mu\|^2+1}{p} \psi_k} \right) > 0 \iff 1 - \frac{p}{n} \sum_{l=0}^{n-1} \left[ \frac{\psi_l}{(\|\mu\|^2 + 1) \psi_k - \psi_l} \right]^2 > 0.$$

1192 Therefore, if  $\bar{\xi}_k$  is a singular point of  $\bar{\mathbf{Q}}$ , then  $\bar{\zeta}_k > 0$ .

1193 Conversely, if  $\bar{\xi}_k$  is not a singular point of  $\bar{\mathbf{Q}}$ , then either  $\psi_k = 0$  or  $\bar{\xi}_k \in \text{Sp}_\infty(\tilde{\mathbf{K}}_L)$ . If  $\psi_k = 0$ , we immediately see that  $\bar{\zeta}_k = \frac{\|\mu\|^2}{\|\mu\|^2+1} (1-p) < 0$ .

1197 On the other hand, if  $\bar{\xi}_k \in \text{Sp}_\infty(\tilde{\mathbf{K}}_L)$  and  $\psi_k \neq 0$  then  $x' \left( \frac{-1}{\frac{\|\mu\|^2+1}{p} \psi_k} \right) \leq 0$  (otherwise  $\bar{\xi}_k$  would be a spike) and  $\bar{\zeta}_k \leq 0$ .  $\square$

## 1200 D. Predictions with a Toeplitz mask

1201 Figures 8a and 8b compare simulations with a Toeplitz mask and the predictions of theorems 3.1 and 3.3 with the  $\psi_k$ 's replaced by the  $\tau_k$ 's and  $\mathbf{F}$  replaced by  $\mathbf{G}$ .

1204 Apart from extra mass around 0 in the second setting ( $c = 0.03$  and  $\varepsilon = 0.6$ ), the shape of the limiting spectral distribution is very well predicted, as well as the position of the isolated eigenvalues. Empirical alignments  $|\mathbf{u}_0^* \mathbf{v}_0|^2$  are also fit well the predicted curve.

1207 Note that, contrary to the circulant mask, the eigenvalues of  $\mathbf{T}$  are mostly simple (see theorem 5 of (Trench, 1994)). Thus, we also represent  $\bar{\zeta}_{n-1}^+$  in Figure 8b, which was confounded with  $\bar{\zeta}_1^+$  in Figure 3 ( $\psi_1 = \psi_{n-1}$  but  $\tau_1 \neq \tau_{n-1}$ ).

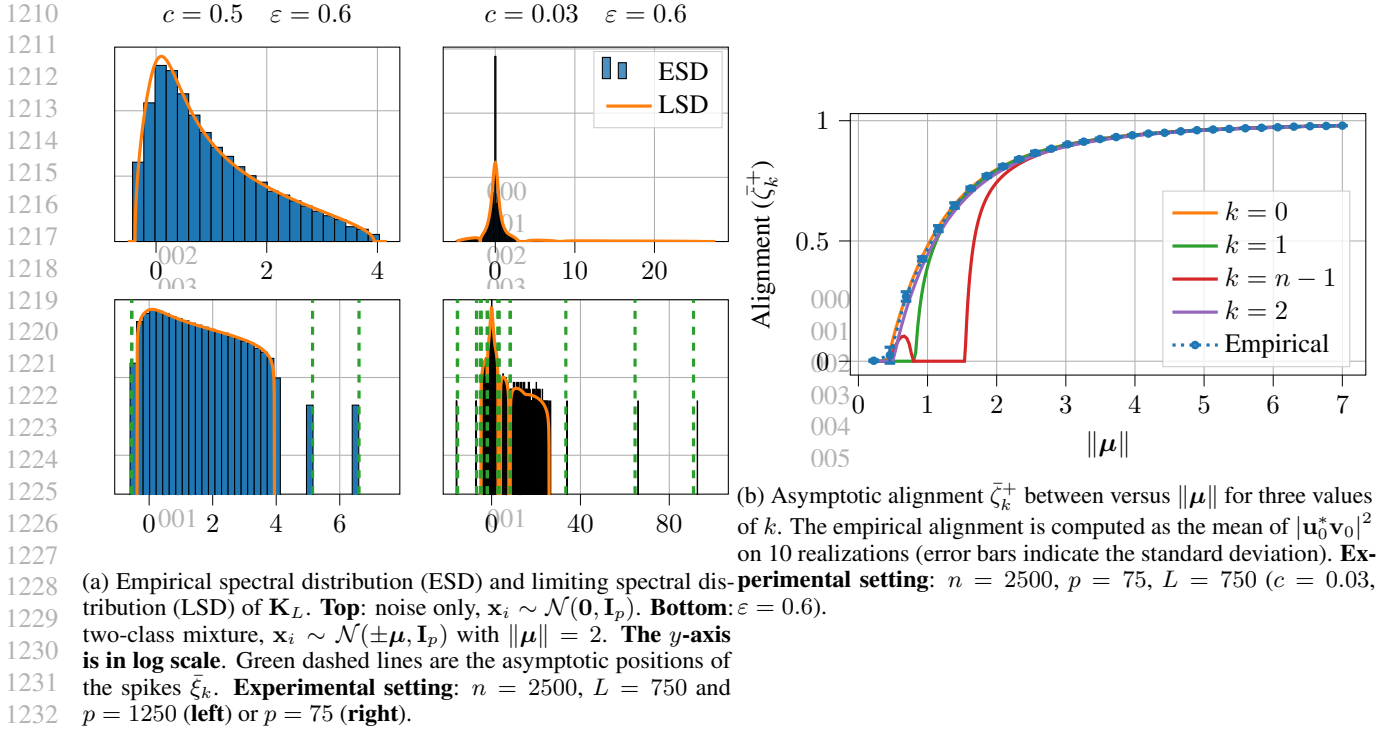


Figure 8. Predictions of theorems 3.1 and 3.3 adapted for a Toeplitz mask.

## E. $K$ -classes online kernel spectral clustering algorithm

### E.1. General presentation and simulations

We use a set of spike eigenvectors  $\{\mathbf{u}_k^{(t)}\}_{k \in \mathcal{K}}$  (with a set of indices  $\mathcal{K}$ ) to estimate the  $|\mathcal{K}|$ -dimensional “trend” of each class. That is, denoting  $\mathcal{C}[t]$  the class of  $\mathbf{x}_t$ , we consider the following model

$$\left[\mathbf{u}_k^{(t)}\right]_i = \left[\mathbf{h}_{k, \mathcal{C}[t-n+i]}^{(t)} + \boldsymbol{\epsilon}_k^{(t)}\right]_i$$

where, for  $k \in \mathcal{K}$ ,  $\mathbf{h}_{k, \mathcal{C}}^{(t)} \in \mathbb{R}^n$  is the “trend” of class  $\mathcal{C}$  and  $\boldsymbol{\epsilon}_k^{(t)}$  is a centered noise vector. A deeper analysis of the deterministic equivalent of theorem 3.1 is needed to properly understand the behavior of the vectors  $\mathbf{h}_{k, \mathcal{C}}^{(t)}$ . From our general understanding so far, it is expected that they are linear combinations of a few dominant eigenvectors of  $\mathbf{T}$ . Using this approach, we are able to estimate the trends from  $\{\mathbf{u}_k^{(t)}\}_{k \in \mathcal{K}}$  (see the left part of Figure 9). Each point is then associated to the class whose curve is the nearest. The details of this algorithm are given in the following subsection.

This algorithm is tested on a stream made of  $T = 21\,000$  centered raw-images from the Fashion-MNIST dataset (Xiao et al., 2017). Their dimension is  $p = 784$  and we want to discriminate between `trouser`, `coat` and `ankle boot` images in an online fashion. We choose  $n = 1\,000$  and  $L = 100$  and we use the 5 dominant eigenvectors of  $\mathbf{K}_L^{(t)}$  for the estimation.

In figure 9 are displayed the shape of the dominant eigenvector  $\mathbf{u}_0^{(t)}$  at a given time during the execution of the algorithm with the estimated trends of each class<sup>14</sup> (left) and the mean clustering error at  $t_0 + \Delta t$  of a data point seen at  $t_0$  with the overall classification error obtained after a majority vote (right). The classification error curve is U-shaped: classes are better estimated around  $t_0 + \frac{n}{2}$  than  $t_0$  or  $t_0 + n - 1$ . This can be understood by the slightly-localized shape of  $\mathbf{u}_0^{(t)}$  (Figure 6, bottom) — it is easier to discriminate between the trends in the middle of the eigenvector than on its edges. Nevertheless, the majority vote counteract this weakness and the overall classification error touches the bottom of the U-shape.

*Remark E.1.* In a binary setting, Algorithm 1 does not suffer this limitation as class estimates are directly given by the sign

<sup>14</sup>This is only the first dimension of a 5-dimensional trend.

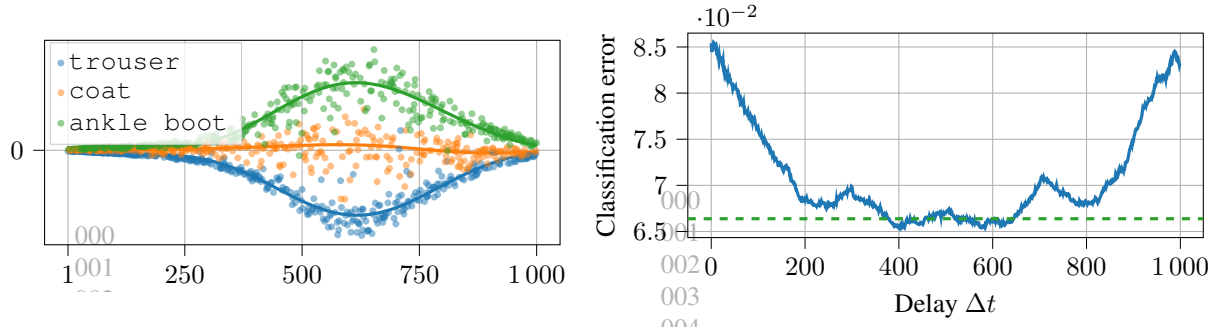


Figure 9. Clustering on Fashion-MNIST images (trouser vs. coat vs. ankle boot). **Left:** dominant eigenvector of  $\mathbf{K}_L^{(t)}$ . Solid curves are the estimated trend of each class  $\mathbf{h}_{k,C}^{(t)}$ . **Right:** Classification error against delay  $\Delta t$ . This is the mean classification error at time  $t_0 + \Delta t$  of a point arrived at  $t_0$ . The green dashed line indicate the overall classification error when the class is chosen by a majority vote. **Experimental setting:**  $T = 21\,000$ ,  $n = 1\,000$ ,  $p = 784$ ,  $L = 100$ .

of the coordinates of  $\mathbf{u}_0$  (no trend needs to be estimated).

Here, the overall classification error is 6.638% while a standard  $T \times T$  offline kernel spectral clustering has only a 3.662% error rate.

## E.2. Details of the algorithm

We consider a set  $\mathcal{K}$  of indices of spikes and the following model for  $\mathbf{u}_k^{(t)}$ ,  $k \in \mathcal{K}$ ,

$$\left[ \mathbf{u}_k^{(t)} \right]_i = \left[ \mathbf{h}_{k,C}^{(t)} + \boldsymbol{\epsilon}_k^{(t)} \right]_i \quad 1 \leq i \leq n$$

where  $\mathbf{h}_{k,C}^{(t)} \in \mathbb{R}^n$  is the trend of class  $C$  and  $\boldsymbol{\epsilon}_k^{(t)}$  is a centered noise vector.

Our goal is to estimate the trend  $\mathbf{h}_{k,C}^{(t)}$  from the eigenvectors  $\left\{ \mathbf{u}_k^{(t)} \right\}_{k \in \mathcal{K}}$ . Since we assume they are linear combinations of a few dominant eigenvectors of  $\mathbf{T}$ , we define a set of indices  $\mathcal{K}_*$  specifying the eigenvectors  $\left\{ \mathbf{g}_k \right\}_{k \in \mathcal{K}_*}$  which we expect the  $\mathbf{h}_{k,C}^{(t)}$ 's being linear combinations of.

We denote  $\hat{C}^{(t)}[s]$  the class of  $\mathbf{x}_{t-n+s}$  estimated at time  $t$ .

In order to compute an estimation  $\left\{ \hat{C}^{(t)}[i] \right\}_{1 \leq i \leq n}$  of the classes at a given time  $t$ , we propose a two-step algorithm. Firstly, we compute a rough estimation  $\left\{ \hat{C}_0^{(t)}[i] \right\}_{1 \leq i \leq n}$  of the classes by following the  $K$  paths with an exponential smoothing in the coordinates of the eigenvectors  $\left\{ \mathbf{u}_k^{(t)} \right\}_{k \in \mathcal{K}}$ , this is called the *pre-classification* step. Then, we refine this estimation with projections on  $\text{span} \left\{ \mathbf{g}_k \right\}_{k \in \mathcal{K}_*}$ , this is the *classification* step.

In the following, we drop the time dependency when it is not needed to ease notations.

### E.2.1. PRE-CLASSIFICATION STEP

Given the number of classes  $K$  and the eigenvectors  $\left\{ \mathbf{u}_k \right\}_{k \in \mathcal{K}}$ , we consider the set of  $n$  points in  $\mathbb{R}^{|\mathcal{K}|}$  defined by the coordinates of each eigenvector:  $\left[ \mathbf{u}_{\mathcal{K}} \right]_i = \left( \left[ \mathbf{u}_k \right]_i \right)_{k \in \mathcal{K}}$  for  $1 \leq i \leq n$ . As  $i$  goes from 1 to  $n$ , these points draw  $K$  paths. The goal is to guess which path (and therefore which class) each point belong to.

Let us suppose we have already estimated  $\hat{C}_0[1], \dots, \hat{C}_0[i-1]$  and the first  $i-1$  coordinates of the vectors  $\left\{ \tilde{\mathbf{h}}_k \right\}_{k \in \mathcal{K}}$  such that  $\left[ \tilde{\mathbf{h}}_k \right]_j$  is an estimation of  $\left[ \mathbf{h}_{k, \hat{C}_0[j]} \right]_j$  (initialization is discussed later). As for  $\left\{ \mathbf{u}_k \right\}_{k \in \mathcal{K}}$ , we see  $\left\{ \tilde{\mathbf{h}}_k \right\}_{k \in \mathcal{K}}$  as a set of  $n$  points in  $\mathbb{R}^{|\mathcal{K}|}$ , which have to be estimated. The estimation of the  $i$ -th point  $\left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_i$  is induced by the class estimate  $\hat{C}_0[i]$  —



the corresponding path is updated with an exponential smoothing:

$$\begin{aligned} \left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_i &= \mathcal{E}_\alpha(i, \mathbf{u}_{\mathcal{K}}, \tilde{\mathbf{h}}_{\mathcal{K}}, \hat{\mathcal{C}}_0[i]) \equiv \frac{\alpha [\mathbf{u}_{\mathcal{K}}]_i + M \left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_{I[\hat{\mathcal{C}}_0[i], i]}}{\alpha + M} \\ \text{where } M &= \frac{1 - \alpha}{i - I[\hat{\mathcal{C}}_0[i], i]} \left[ 1 + \frac{1 - \alpha}{\alpha} \left( 1 - (1 - \alpha)^{i - I[\hat{\mathcal{C}}_0[i], i] - 1} \right) \right], \end{aligned}$$

$I[\hat{\mathcal{C}}_0[i], i] = \max \{ 1 \leq j \leq i - 1 \mid \hat{\mathcal{C}}_0[j] = \hat{\mathcal{C}}_0[i] \}$  is the index of the last seen point in  $\hat{\mathcal{C}}_0[i]$  and  $\alpha \in [0, 1]$  is the smoothing parameter. The reasons for such a formula are detailed in appendix F.

However,  $\hat{\mathcal{C}}_0[i]$  is chosen as the class which minimizes the growth of the corresponding path:

$$\hat{\mathcal{C}}_0[i] = \arg \min_{\hat{\mathcal{C}} \in \{\hat{\mathcal{C}}_1, \dots, \hat{\mathcal{C}}_K\}} \frac{\left\| \mathcal{E}_\alpha(i, \mathbf{u}_{\mathcal{K}}, \tilde{\mathbf{h}}_{\mathcal{K}}, \hat{\mathcal{C}}) - \left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_{I[\hat{\mathcal{C}}, i]} \right\|}{i - I[\hat{\mathcal{C}}, i]}.$$

Indeed, by doing so, we minimize the Lipschitz constant of the estimated trend and ensure some regularity.

From the regularity of the true trend,  $\mathbf{h}_{k, \mathcal{C}}$  is almost flat on its very first coordinates. Therefore, we can initialize the values  $\hat{\mathcal{C}}_0[i]$  for  $1 \leq i \leq H$  with a standard clustering algorithm applied to  $\{[\mathbf{u}_{\mathcal{K}}]_i\}_{1 \leq i \leq H}$ .  $H$  is a parameter which should be taken as small as possible to stay in a domain where the trends are almost flat while still having a few representatives of each class.

The computation of  $\left\{ \left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_i \right\}_{1 \leq i \leq H}$  follows from the class estimates, as presented above.

We found that a hierarchical clustering algorithm and  $H \approx 10K$  worked well for the initialization. As for the smoothing parameter, a good value is  $\alpha \approx 0.15$ .

The pre-classification step is summarized in Algorithm 2.

---

**Algorithm 2** Pre-classification

---

**Input:**  $K, \{\mathbf{u}_k\}_{k \in \mathcal{K}}$ .

**Parameters:**  $H, \alpha$ .

**Output:**  $\{\hat{\mathcal{C}}_0[i]\}_{1 \leq i \leq n}$ .

Set  $\hat{\mathcal{C}}_0[i]$  for  $i = 1$  to  $H$  with agglomerative clustering.

**for**  $i = 1$  to  $H$  **do**

$\left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_i \leftarrow \mathcal{E}_\alpha(i, \mathbf{u}_{\mathcal{K}}, \tilde{\mathbf{h}}_{\mathcal{K}}, \hat{\mathcal{C}}_0[i])$

**end for**

**for**  $i = H + 1$  to  $n$  **do**

$\hat{\mathcal{C}}_0[i] \leftarrow \arg \min_{\hat{\mathcal{C}} \in \{\hat{\mathcal{C}}_1, \dots, \hat{\mathcal{C}}_K\}} \frac{\left\| \mathcal{E}_\alpha(i, \mathbf{u}_{\mathcal{K}}, \tilde{\mathbf{h}}_{\mathcal{K}}, \hat{\mathcal{C}}) - \left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_{I[\hat{\mathcal{C}}, i]} \right\|}{i - I[\hat{\mathcal{C}}, i]}$

$\left[ \tilde{\mathbf{h}}_{\mathcal{K}} \right]_i \leftarrow \mathcal{E}_\alpha(i, \mathbf{u}_{\mathcal{K}}, \tilde{\mathbf{h}}_{\mathcal{K}}, \hat{\mathcal{C}}_0[i])$

**end for**

---

### E.2.2. CLASSIFICATION STEP

The class estimates obtained after the pre-classification step usually are not very satisfying but still are a good basis to estimate  $\mathbf{h}_{k, \mathcal{C}}$  with regressions.

In the second step of the algorithm, we are given a set  $\{\mathbf{g}_k\}_{k \in \mathcal{K}_*}$  of eigenvectors of  $\mathbf{T}$ . It is supposed that the trends  $\{\mathbf{h}_{k, \mathcal{C}}\}_{k \in \mathcal{K}}$  are mixtures of these eigenvectors.

From class estimates  $\{\hat{\mathcal{C}}[i]\}_{1 \leq i \leq n}$ , we can compute an estimation  $\hat{\mathbf{h}}_{\mathcal{K}, \hat{\mathcal{C}}}$  of the trend of each estimated class  $\hat{\mathcal{C}}$  with a linear regression

$$\hat{\mathbf{h}}_{k, \hat{\mathcal{C}}} = \mathbf{v}_{\mathcal{K}_*} \boldsymbol{\beta}_{k, \hat{\mathcal{C}}} \quad \text{where} \quad \boldsymbol{\beta}_{k, \hat{\mathcal{C}}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^{|\mathcal{K}_*|}} \left\| [\mathbf{u}_k]_{\hat{\mathcal{C}}} - [\mathbf{v}_{\mathcal{K}_*}]_{\hat{\mathcal{C}}} \boldsymbol{\beta} \right\|^2$$

1375 where we use the notation  $[\cdot]_{\hat{C}}$  to represent the restriction to  $\hat{C}$ .

1376 Then, new class estimates can be computed by associating each point to the class whose trend is the closest:

$$1378 \quad \hat{C}[i] = \arg \min_{\hat{C} \in \{\hat{C}_1, \dots, \hat{C}_K\}} \left\| [\mathbf{u}_{\mathcal{K}}]_i - [\hat{\mathbf{h}}_{\mathcal{K}, \hat{C}}]_i \right\|.$$

1381 We repeat this process until convergence of the class estimates. The classification step is summarized in Algorithm 3.

---

1383 **Algorithm 3** Classification

---

1385 **Input:**  $K, \{\hat{C}_0[i]\}_{1 \leq i \leq n}, \{\mathbf{u}_k\}_{k \in \mathcal{K}}, \{\mathbf{v}_k\}_{k \in \mathcal{K}_*}$ .

1386 **Output:**  $\{\hat{C}[i]\}_{1 \leq i \leq n}$ .

1387 **for**  $i = 1$  to  $n$  **do**

1388      $\hat{C}[i] \leftarrow \hat{C}_0[i]$

1389 **end for**

1390 **repeat**

1391     **for**  $\hat{C} \in \{\hat{C}_1, \dots, \hat{C}_K\}$  **do**

1392          $\hat{\mathbf{h}}_{\mathcal{K}, \hat{C}} \leftarrow \mathbf{v}_{\mathcal{K}_*} \left( [\mathbf{v}_{\mathcal{K}_*}]_{\hat{C}}^\top [\mathbf{v}_{\mathcal{K}_*}]_{\hat{C}} \right)^{-1} [\mathbf{v}_{\mathcal{K}_*}]_{\hat{C}}^\top [\mathbf{u}_{\mathcal{K}}]_{\hat{C}}$

1393     **end for**

1394     **for**  $i = 1$  to  $n$  **do**

1395          $\hat{C}[i] \leftarrow \arg \min_{\hat{C} \in \{\hat{C}_1, \dots, \hat{C}_K\}} \left\| [\mathbf{u}_{\mathcal{K}}]_i - [\hat{\mathbf{h}}_{\mathcal{K}, \hat{C}}]_i \right\|$

1396     **end for**

1397 **until** convergence

---

1402 E.2.3. FINAL ALGORITHM

---

1404 **Algorithm 4** Online Kernel Spectral Clustering

---

1405 **Input:**  $K, \mathcal{K}, \{\mathbf{g}_k\}_{k \in \mathcal{K}_*}$ .

1406 **Parameters:**  $H, \alpha$ .

1407 **Output:**  $\{\hat{C}_t[s]\}_{\substack{1 \leq s \leq n \\ n \leq t \leq T}}$ .

1408 **for**  $t = 1$  to  $T$  **do**

1409     Get a new point  $\mathbf{x}_t$  into the pipeline.

1410     Compute  $\mathbf{x}_l^* \mathbf{x}_{t-l}$  for  $l = 0$  to  $L - 1$ .

1411     Update  $\mathbf{K}_L^{(t-1)}$  into  $\mathbf{K}_L^{(t)}$ .

1412      $\mathbf{u}_{\mathcal{K}}^{(t)} \leftarrow \text{PowerIteration}(\mathbf{K}_L^{(t)}, \mathbf{u}_{\mathcal{K}}^{(t-1)})$ .

1413     **if**  $1 \leq t \leq n$  **then**

1414         Do an iteration as in Algorithm 2.

1415     **end if**

1416     **if**  $t \geq n$  **then**

1417         Compute  $\{\hat{C}_t[s]\}_{1 \leq s \leq n}$  according to Algorithm 3 with  $\{\hat{C}_{t-1}[s]\}_{1 \leq s \leq n-1}$ .

1418     **end if**

1419 **end for**

---

1423 In an online fashion, pre-classification can be performed as a warm-up during the first  $n$  time steps. Then, as  $t \geq n$ , only the

1424 classification step is needed: the classes  $\{\hat{C}_{t-1}[s]\}_{1 \leq s \leq n}$  estimated at  $t - 1$  (or during pre-classification if  $t = n$ ) serve as a

1425 good basis to estimate the classes at time  $t$  (both  $\hat{C}_{t-1}[s]$  and  $\hat{C}_t[s + 1]$  are estimates of the class of  $\mathbf{x}_{t-s}$ ). Moreover, the few

1426 interesting eigenvectors  $\mathbf{u}_{\mathcal{K}}^{(t)}$  of  $\mathbf{K}_L^{(t)}$  can be quickly computed with a power iteration algorithm starting at  $\mathbf{u}_{\mathcal{K}}^{(t-1)}$  (they do

1427 not differ much from one time step to another). The final algorithm is presented in Algorithm 4.

1428

1429

1430 **F. Exponential smoothing with missing data**

1431 Let  $(\mathbf{y}_t)_{t \in \mathbb{N}}$  be a time series. Assume we want to compute its trend  $(\mathbf{s}_t)_{t \in \mathbb{N}}$ . A common technique to do so is to perform an  
 1432 exponential smoothing:  
 1433

$$1434 \quad \mathbf{s}_0 = \mathbf{y}_0 \quad \text{and} \quad \mathbf{s}_{t+1} = \alpha \mathbf{y}_{t+1} + (1 - \alpha) \mathbf{s}_t \quad \forall t \in \mathbb{N}$$

1435 where  $\alpha \in [0, 1]$  is the smoothing parameter. It acts as a low-pass filter which removes high-frequency noise.

1436 Let us now assume that we do not have access to  $(\mathbf{y}_t)_{t \in \mathbb{N}}$  at each time step and we want to compute  $\mathbf{s}_{t+h}$  ( $h \geq 1$ ) with  $\mathbf{y}_{t+h}$   
 1437 and  $\mathbf{s}_t$  only. Expanding the recurrence relation, we have  
 1438

$$1439 \quad \mathbf{s}_{t+h} = \alpha \mathbf{y}_{t+h} + \alpha \sum_{k=1}^{h-1} (1 - \alpha)^k \mathbf{y}_{t+h-k} + (1 - \alpha)^h \mathbf{s}_t.$$

1440 We propose to replace the unknown values  $\mathbf{y}_{t+h-k}$  for  $1 \leq k \leq h - 1$  by the linear interpolation of the trend:

$$1441 \quad \begin{aligned} 1442 \quad \mathbf{s}_{t+h} &= \alpha \mathbf{y}_{t+h} + \alpha \sum_{k=1}^{h-1} (1 - \alpha)^k \left[ \frac{k}{h} \mathbf{s}_t + \frac{h-k}{h} \mathbf{s}_{t+h} \right] + (1 - \alpha)^h \mathbf{s}_t \\ 1443 &= \alpha \mathbf{y}_{t+h} + \frac{\alpha}{h} \left( \mathbf{s}_t \sum_{k=1}^{h-1} k (1 - \alpha)^k + \mathbf{s}_{t+h} \sum_{k=1}^{h-1} k (1 - \alpha)^{h-k} \right) + (1 - \alpha)^h \mathbf{s}_t. \end{aligned}$$

1444 Using the following formulae,

$$1445 \quad \sum_{k=1}^{h-1} k (1 - \alpha)^k = \frac{1 - \alpha}{\alpha} \left( 1 - h(1 - \alpha)^{h-1} \right) + \left( \frac{1 - \alpha}{\alpha} \right)^2 \left( 1 - (1 - \alpha)^{h-1} \right)$$

$$1446 \quad \text{and} \quad \sum_{k=1}^{h-1} k (1 - \alpha)^{h-k} = \frac{1 - \alpha}{\alpha} (h - 1) - \left( \frac{1 - \alpha}{\alpha} \right)^2 \left( 1 - (1 - \alpha)^{h-1} \right)$$

1447 we have

$$1448 \quad \left( \alpha + \frac{1 - \alpha}{h} \left[ 1 + \frac{1 - \alpha}{\alpha} \left( 1 - (1 - \alpha)^{h-1} \right) \right] \right) \mathbf{s}_{t+h} = \alpha \mathbf{y}_{t+h} + \frac{1 - \alpha}{h} \left[ 1 + \frac{1 - \alpha}{\alpha} \left( 1 - (1 - \alpha)^{h-1} \right) \right] \mathbf{s}_t.$$

1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484