



HAL
open science

A Random Matrix Analysis of Data Stream Clustering: Coping With Limited Memory Resources

Hugo Lebeau, Romain Couillet, Florent Chatelain

► **To cite this version:**

Hugo Lebeau, Romain Couillet, Florent Chatelain. A Random Matrix Analysis of Data Stream Clustering: Coping With Limited Memory Resources. Proceedings of Machine Learning Research, 2022, 162, pp.12253-12281. hal-03755939

HAL Id: hal-03755939

<https://hal.science/hal-03755939v1>

Submitted on 22 Aug 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

A Random Matrix Analysis of Data Stream Clustering: Coping With Limited Memory Resources

Hugo Lebeau¹ Romain Couillet¹ Florent Chatelain²

Abstract

This article introduces a random matrix framework for the analysis of clustering on high-dimensional data streams, 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 banded 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 clustering 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,

¹Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, 38000 Grenoble, France ²Univ. Grenoble Alpes, CNRS, Grenoble INP, GIPSA-lab, 38000 Grenoble, France. Correspondence to: Hugo Lebeau <hugo.lebeau@univ-grenoble-alpes.fr>.

thus inducing a limited memory footprint.

Numerous works have proposed various algorithms to cluster data streams in an unsupervised manner (see, e.g., Ghesmoune et al. (2016); Zubaroğlu & Atalay (2021) 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; Cohen-Addad et al., 2021), 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 grows 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).

From a random matrix theory perspective, spectral clustering is also of particular interest. Following the works of El Karoui (2010) and 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 adaptation of spectral clustering is then performed on \mathbf{K}_L to retrieve the class information.

In technical terms, the present analysis derives the limiting 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.

The remainder of the paper is organized as follows. Section 2 introduces the model and a circulant approximation of the Toeplitz mask \mathbf{T} , which will be used to derive our main results, presented in section 3. The limiting spectral distribution of the kernel matrix is studied first (Theorem 3.1) and a closer look is then given to the behavior of its isolated eigenvalues and associated eigenvectors (Theorem 3.3). Based on the previous results, section 4 presents some theoretical considerations on the classification performance achievable on a data stream and proposes an online kernel spectral clustering algorithm, which is tested on image clustering tasks. Section 5 gives some concluding remarks.

Proofs and simulations All proofs are deferred to the appendix. Python codes to reproduce simulations are available in the following GitHub repository https://github.com/HugoLebeau/online_learning/.

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 and does not vanish asymptotically:

$$0 < \liminf_{p \rightarrow +\infty} \|\mathbf{M}\| \leq \limsup_{p \rightarrow +\infty} \|\mathbf{M}\| < +\infty.$$

Assumption 2.3 (Additive noise model). $\mathbf{X} = \mathbf{P} + \mathbf{Z}$ where $\mathbf{P} = \mathbf{M}\mathbf{J}^\top$ is a deterministic *signal* matrix and \mathbf{Z} is a random standard Gaussian *noise* matrix with independent entries¹.

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 ($\|\mathbf{M}\| \rightarrow +\infty$) nor too hard ($\|\mathbf{M}\| \rightarrow 0$). The classification 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 $\mathbf{K}_{i,j} = \frac{1}{p} \mathbf{x}_i^\top \mathbf{x}_j$ of the Gram kernel matrix can be computed only if $|i - j| < L$. 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

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

As standard (offline) spectral clustering is “optimal”², 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 $(2L - 1)/n \rightarrow \varepsilon \in]0, +\infty[$.³

2.2. The circulant approximation

An important trick to derive our main results lies in the fact that the Toeplitz matrix \mathbf{T} can be approximated to some extent by its circulant “version” $\mathbf{C} = [\mathbf{1}_{|i-j| < L} + \mathbf{1}_{|i-j| > n-L}]_{1 \leq i, j \leq n}$ (Gray, 2006). Indeed, 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

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

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

³The provided results are asymptotic for theoretical convenience, modeling the fact that n, p and L are large. The convergence rates being at least $\mathcal{O}(1/\sqrt{n})$ as $n, p, L \rightarrow +\infty$, they remain valid for a large but finite horizon.

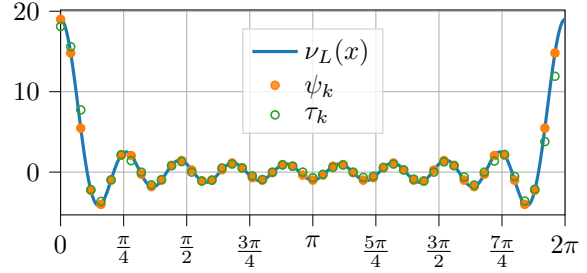


Figure 1. Graph of ν_L on $[0, 2\pi[$ (one period) with a plot of $\psi_k = \nu_L(\frac{2k\pi}{n})$ and τ_k for $0 \leq k < n$ (the eigenvalues of \mathbf{C} and \mathbf{T} respectively). **Experimental setting:** $n = 50, L = 10$.

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 ε is sufficiently small.

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

$$\mathbf{C} = \mathbf{F}\mathbf{\Psi}\mathbf{F}^* \quad \text{with} \quad \mathbf{F}_{i,j} = \frac{1}{\sqrt{n}} e^{-\frac{2i\pi}{n}(i-1)(j-1)},$$

i.e., \mathbf{F} is the $n \times n$ Fourier matrix and $\mathbf{\Psi} = \text{diag}(\psi_k)_{0 \leq k < n}$ 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 ν_L the eigenvalues of \mathbf{C} and \mathbf{T} .⁴ The τ_k ’s roughly follow the graph of ν_L , as if they were noisy versions of the ψ_k ’s.

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

⁴Although 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 corresponding order for the eigenvalues of \mathbf{T} : after numerically computing them in ascending order, we apply the same permutation that maps the eigenvalues of \mathbf{C} in ascending order to $(\psi_0, \dots, \psi_{n-1})$. This yields the corresponding $(\tau_0, \dots, \tau_{n-1})$.

of the empirical spectral measure $\mu_n = \frac{1}{n} \sum_{\xi \in \text{sp}(\mathbf{K}_L)} \delta_\xi$ 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 Γ circling around an eigenvalue ξ 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 ξ .⁵

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 $\tilde{\mathbf{K}}_L = \frac{\mathbf{X}^T \mathbf{X}}{p} \odot \mathbf{C}$. Namely, we find a deterministic matrix $\tilde{\mathbf{Q}}(z)$ 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), $\frac{1}{n} \text{tr} \mathbf{A}_n (\tilde{\mathbf{Q}}(z) - \tilde{\mathbf{Q}}(z)) \rightarrow 0$ and $\mathbf{a}_n^T (\tilde{\mathbf{Q}}(z) - \tilde{\mathbf{Q}}(z)) \mathbf{b}_n \rightarrow 0$ almost surely as $n, p, L \rightarrow +\infty$. This will be simply denoted $\tilde{\mathbf{Q}}(z) \leftrightarrow \tilde{\mathbf{Q}}(z)$.

Theorem 3.1 (Deterministic equivalent of $\tilde{\mathbf{Q}}(z)$). *Under assumptions 2.1 – 2.3, $\tilde{\mathbf{K}}_L$ admits a limiting spectral distribution μ as $n, p, L \rightarrow +\infty$. Its Stieltjes transform m is solution to*

$$1 + zm(z) = \frac{p}{n} \sum_{k=0}^{n-1} \frac{m(z) \frac{\psi_k}{p}}{1 + m(z) \frac{\psi_k}{p}} \quad z \in \mathbb{C} \setminus \text{supp } \mu. \quad (1)$$

Moreover, if $\text{dist}(z, \text{supp } \mu) > \frac{2L-1}{p}$, then

$$\tilde{\mathbf{Q}}(z) \leftrightarrow \tilde{\mathbf{Q}}(z) \equiv m(z) (\mathbf{I}_n + \mathbf{P}^T \mathbf{P} \odot \mathbf{F} \Lambda(z) \mathbf{F}^*)^{-1}$$

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

Proof. See appendix B. \square

A first observation from Theorem 3.1 is that $\tilde{\mathbf{Q}}(z)$ 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 $\mathbf{P}^T \mathbf{P} \odot \mathbf{F} \Lambda(z) \mathbf{F}^*$ are

⁵In fact, this is only true if ξ has multiplicity 1. In the general case, the integral equals the projection matrix on the eigenspace associated to ξ .

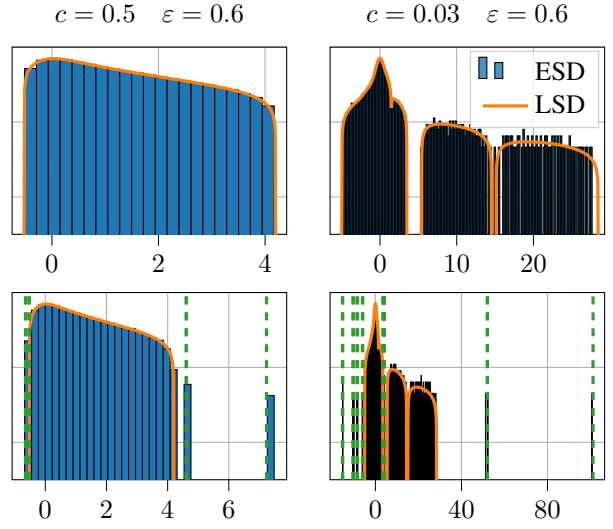


Figure 2. Empirical spectral distribution (ESD) and limiting spectral distribution (LSD) of $\tilde{\mathbf{K}}_L$. **The y-axis is in log scale.** **Top:** noise only, $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$. **Bottom:** two-class mixture, $\mathbf{x}_i \sim \mathcal{N}(\pm\boldsymbol{\mu}, \mathbf{I}_p)$ 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).

small enough for only a few number of corresponding isolated eigenvalues in the spectrum of $\tilde{\mathbf{K}}_L$ to appear.

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

$$zc^{-1}m^2(z) - (1 - c^{-1} - z)m(z) + 1 = 0$$

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

In practice, rather than computing $m(z)$ directly from equation 1, it is easier to solve numerically the following fixed-point equation in η_0

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

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 settings⁶ 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

⁶Recall that $c = \lim p/n$ and $\varepsilon = \lim (2L - 1)/n$.

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 ε . 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 pops 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 focus back on our original clustering objective. We consider two classes \mathcal{C}^\pm whose centroids are $\pm\boldsymbol{\mu}$, i.e.⁷, $\mathbf{P} = \boldsymbol{\mu}\mathbf{j}^\top$ 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 empirically 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 (Theorem 3.1) has a much simpler expression:

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

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 $\mathbf{I}_n + \|\boldsymbol{\mu}\|^2 \boldsymbol{\Lambda}$ 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 signs 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*

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

and

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

The following propositions are equivalent.

1. $\psi_k \neq 0$ and $\bar{\xi}_k \notin \text{supp } \mu$.
2. $\bar{\zeta}_k > 0$.

⁷Consistently with the previous setting, $\mathbf{M} = [+ \boldsymbol{\mu} \quad - \boldsymbol{\mu}]$ and $\mathbf{J}_{i \cdot} = [\mathbf{1}_{\mathbf{x}_i \in \mathcal{C}^+} \quad \mathbf{1}_{\mathbf{x}_i \in \mathcal{C}^-}]$.

3. Almost surely, $\bar{\xi}_k$ is the asymptotic position of an isolated eigenvalue of $\tilde{\mathbf{K}}_L$.

Then, in this case, the matrix $\mathbf{U}_k = [\mathbf{u}_l]_{\substack{\psi_k = \psi_l \\ 0 \leq l < n}}$ gathers all the eigenvectors of $\tilde{\mathbf{K}}_L$ whose associated eigenvalues converge a.s. 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 $\mathcal{D}_k = \text{diag} (\mathbf{1}_{\psi_k = \psi_l})_{0 \leq l < n}$.

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 for index k 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 ψ_k , is rarely simple⁸. However, for finite values of n, p and L , the corresponding eigenvalues of $\tilde{\mathbf{K}}_L$ are not necessarily degenerate (with probability one, they are not), but they have the same limit⁹.

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

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

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 ψ_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 ψ_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

⁸In fact, the only simple eigenvalues of \mathbf{C} are ψ_0 , and $\psi_{n/2}$ when n is even.

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

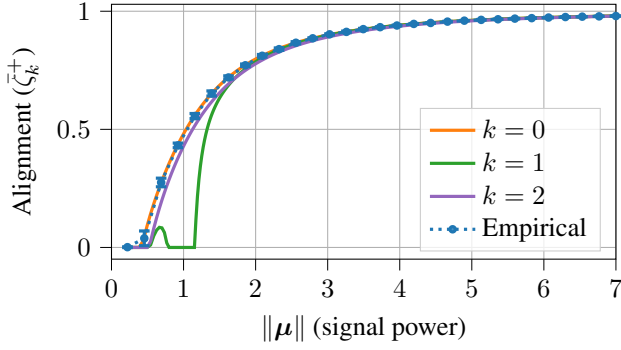


Figure 3. Asymptotic alignment $\bar{\zeta}_k^+$ versus $\|\mu\|$ for three values of k . The empirical alignment is computed as the mean of $|\mathbf{u}_0^* \mathbf{v}_0|^2$ on 10 realizations (error bars indicate the standard deviation). **Experimental setting:** $n = 2500$, $p = 75$, $L = 750$.

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 the vector $\mathbf{v}_k = \mathbf{F}_{\cdot, k} \odot \mathbf{j}$ and the noise level is indicated by $0 \leq 1 - \bar{\zeta}_k \leq 1$. In fact, \mathbf{v}_k is the vector \mathbf{j} — the information sought — modulated by the $(k + 1)$ -th Fourier mode.¹⁰

Figure 3 displays the value of $\bar{\zeta}_k^+ = \max(\bar{\zeta}_k, 0)$ as a function of $\|\mu\|$ 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 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.

¹⁰Recall that Fourier modes are the eigenvectors of \mathbf{C} .

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 .¹¹ 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, as observed empirically.

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

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

4. Online spectral clustering of large data

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

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 classification error is given by $\mathcal{Q}(\sqrt{\bar{\zeta}_0^+ / (1 - \bar{\zeta}_0^+)})$, 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 $\varepsilon = \frac{2L-1}{n}$ (green curve), with the asymptotic classification error of online kernel spectral clustering (orange density map), 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 (red curve).
- *Punctured kernel spectral clustering* (Zarrouk et al., 2020; Couillet et al., 2021), i.e., *offline* clustering performed with a sparsified kernel matrix $\mathbf{K}_\varepsilon = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{B}$,

¹¹One can notice that removing the first and last $L - 1$ rows and columns of \mathbf{C} and \mathbf{T} yields the same two Toeplitz matrices.

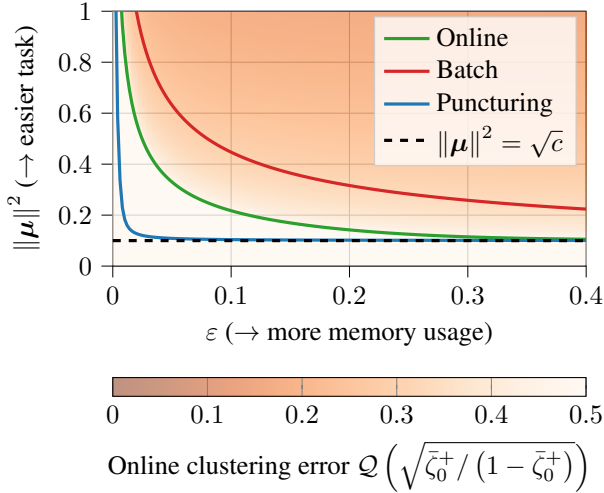


Figure 4. Phase transition position ($\|\boldsymbol{\mu}\|^2$) of the dominant eigenvector of the kernel matrix against the sparsity parameter (ε) with $n/p = 100$. Classification is only possible above the curve corresponding to the method used. The black dashed line is the optimal phase transition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

where $\mathbf{B}_{i,j} = \mathbf{B}_{j,i} \sim \text{Bern}(\varepsilon)$ and $\mathbf{B}_{i,i} = 1$ (blue curve).¹²

As ε grows, the phase transition position of online spectral clustering 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). This is expected, since increasing the memory size allows to encapsulate more information. Still, the green curve specifies *how* memory limitations impair performance. Although we lack some information-theoretic result, the distance between the green curve and the black dashed line yields an upper bound on the difference between the performances of our method and an optimistic optimum (which, as we see, can get very close to 0). Moreover, with ε fixed (fixed memory size), the classification error vanishes as $\|\boldsymbol{\mu}\|$ increases (the signal becomes more powerful). In order to keep L — the memory usage — small without impairing too much the performance under this setting, a good compromise appears to be $0.1 \lesssim \varepsilon \lesssim 0.2$, i.e., $\frac{n}{20} \lesssim L \lesssim \frac{n}{10}$.

Our method performs better (i.e., the phase transition occurs earlier) than any method based on batches of the L points available in memory. It is also able to classify the n previ-

¹²This of course is not doable with a memory bank of size L since the computation of \mathbf{K}_ε requires the (almost) full knowledge of \mathbf{X} . Still, the comparison is interesting as the number of entries in \mathbf{K}_ε and \mathbf{K}_L is the same.

ous 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 (i.e., the same value of ε), the puncturing method performs better. Yet, this requires the access to n data points to compute \mathbf{K}_ε , which is not possible in an online fashion.

4.2. Online clustering algorithm

Before diving into the simulations, we detail a clustering algorithm based on our previous results. We now use the banded version of the kernel matrix: $\mathbf{K}_L = \frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{T}$ (the circulant mask is 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 while \mathbf{x}_{t-L} is discarded. The kernel matrix is then 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 (Memory management policy). A different memory management policy — not restricted to only choosing the previous L points to keep in memory — could be considered. However, we found that having points spread over a greater period of time (i.e., discarding newer ones to keep older ones) does not bring more information. To get a grasp, remark that the mean leaving time of the pipeline cannot be different than L , whatever the policy.

Remark 4.2 (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 , 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 ($n \gg L$ or $\varepsilon \rightarrow 0$), its eigenvectors tend to localize (Hata & Nakao, 2017), making classification more challenging.

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, standard clustering algorithms such as k -means — which are usually employed on spectral embeddings — perform poorly here, because of the particular shape of the eigenvectors caused by the

Toeplitz mask¹³ (see Figure 6).

Remark 4.3. 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**

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

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

end if

end for

The careful reader may wonder here whether the performance of the algorithm could be improved by using eigenvectors other than just the top one. In fact, the top eigenvector already contains all the information that can be retrieved. Since the classification is performed very easily with the signs of the coordinates in the binary setting $\mathbf{x}_i \sim \mathcal{N}(\pm\boldsymbol{\mu}, \mathbf{I}_p)$, the use of other spike eigenvectors does not bring more information. However, in a general setting $\mathbf{x}_i \sim \sum_{k=1}^K \pi_k \mathcal{N}(\boldsymbol{\mu}_k, \mathbf{I}_p)$, we no longer have an alignment result such as Theorem 3.3 and it can become much harder to distinguish the classes from just the top eigenvector. In this case, the combination of several spike eigenvectors can make the classification easier. The interested reader is referred to appendix E, where we propose a — more complex and heuristic — online spectral clustering algorithm capable of handling K -class mixtures and test it on Fashion-MNIST images.

Note that these algorithms 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). This will nonetheless modify the structure of the kernel matrix \mathbf{K}_L and additional work may be necessary to recover theoretical grounds.

¹³The 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).

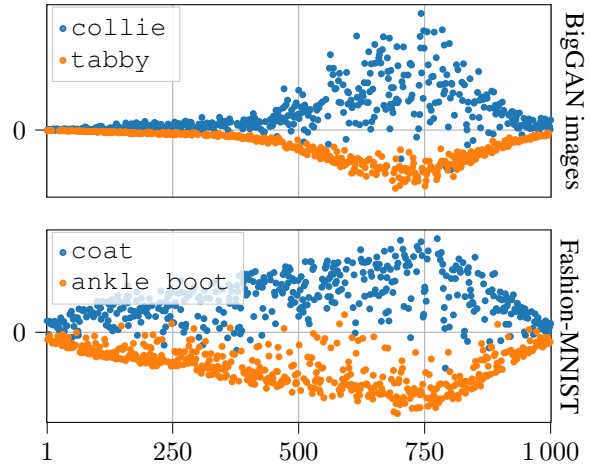


Figure 6. Dominant eigenvector of $\mathbf{K}_L^{(t)}$ with BigGAN-generated images (top) and Fashion-MNIST images (bottom). **Experimental setting:** $T = 20\,000$, $n = 1\,000$, $p = 4\,096$, $L = 100$ (BigGAN images) and $T = 14\,000$, $n = 1\,000$, $p = 784$, $L = 100$ (Fashion-MNIST).

4.3. Simulations on real-world images

We illustrate our findings with two applications on image clustering 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 $p = 4\,096$ 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 $L = 100$. This means that, at each time step, 100 images are kept in memory and, from the $n \times n$ kernel matrix, we are able to classify the previous 1 000 images. This is a realistic choice of parameters (it can easily be run on most standard laptops) from which good performances are expected ($\varepsilon \simeq 0.2$).

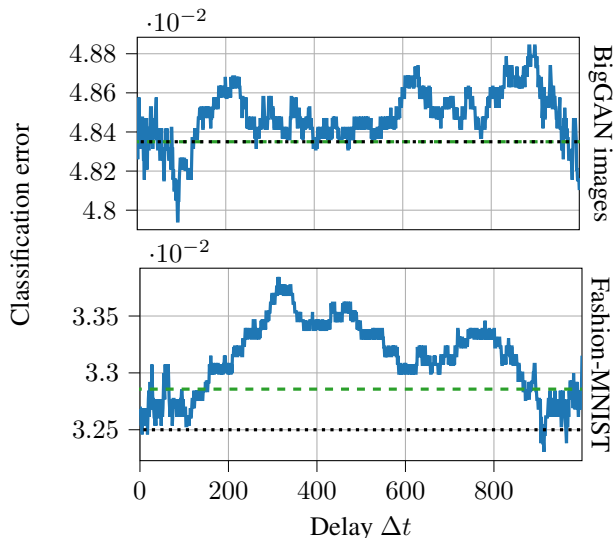


Figure 7. Classification error against delay Δt on BigGAN-generated images (**top**) and Fashion-MNIST images (**bottom**). This is the mean classification error at time $t_0 + \Delta t$ of a point arrived at t_0 . The green dashed line indicates the overall classification error when the class is chosen 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.

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 classification error at $t_0 + \Delta t$ of a data point seen at t_0 ¹⁴, 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 standard $T \times T$ offline kernel spectral clustering¹⁵ (black dotted line). The mean classification error remains constant with Δt , thus showing that our algorithm does not lose any discriminative power between t_0 and $t_0 + n - 1$. The classification performances of our algorithm are very close to those of the standard (offline and costly) 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

¹⁴Recall that a data point arriving at t_0 is classified at each time step between t_0 and $t_0 + n - 1$.

¹⁵For which optimality results are known.

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 does not only introduce a new algorithm for on-line clustering, but also paves the path towards the question of large-dimensional learning on data streams with theoretical guarantees. Still, here we miss an information-theoretic result of optimality for the proposed approach (which exists in the standard unbanded case), a key direction we currently investigate.

Acknowledgment

We thank our reviewers for the time and effort they have devoted to our work. Their precious remarks have allowed us to greatly improve the content of the present article.

References

- 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., Kossmann, 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.
- Cohen-Addad, V., Guedj, B., Kanade, V., and Rom, G. Online k-means Clustering. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, pp. 1126–1134. PMLR, March 2021. URL <https://proceedings.mlr.press/v130/cohen-addad21a.html>. ISSN: 2640-3498.
- 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 Cl  men  on, 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., Kriegel, 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: 10.1561/01000000006. URL <https://www.nowpublishers.com/article/Details/CIT-006>. Publisher: Now Publishers, Inc.
- Gribonval, R., Chatalic, A., Keriven, N., Schellekens, V., Jacques, L., and Schniter, P. Sketching Data Sets for Large-Scale Learning: Keeping only what you need. *IEEE Signal Processing Magazine*, 38(5):12–36, September 2021. doi: 10.1109/MSP.2021.3092574. URL <https://hal.inria.fr/hal-03350599>. Publisher: Institute of Electrical and Electronics Engineers.
- Hata, S. and Nakao, H. Localization of Laplacian eigenvectors on random networks. *Scientific Reports*, 7(1):1121, April 2017. ISSN 2045-2322. doi: 10.1038/s41598-017-01010-0. URL <https://www.nature.com/articles/s41598-017-01010-0>. Bandiera_abtest: a Cc.license.type: cc_by Cg.type: Nature Research Journals Number: 1 Primary.atype: Research Publisher: Nature Publishing Group Subject.term: Complex networks;Nonlinear phenomena Subject.term.id: complex-networks;nonlinear-phenomena.
- Kadavankandy, A. and Couillet, R. Asymptotic Gaussian Fluctuations of Spectral Clustering Eigenvectors. In *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, pp. 694–698, December 2019. doi: 10.1109/CAMSAP45676.2019.9022474.
- Keriven, N., Bourrier, A., Gribonval, R., and P  rez, P. Sketching for Large-Scale Learning of Mixture Models. *arXiv:1606.02838 [cs, stat]*, May 2017. URL

- <http://arxiv.org/abs/1606.02838>. arXiv: 1606.02838.
- Liao, Z., Couillet, R., and Mahoney, M. W. Sparse Quantized Spectral Clustering. *arXiv:2010.01376 [cs, math, stat]*, October 2020. URL <http://arxiv.org/abs/2010.01376>. arXiv: 2010.01376.
- Liberty, E. Simple and Deterministic Matrix Sketching. *arXiv:1206.0594 [cs]*, July 2012. URL <http://arxiv.org/abs/1206.0594>. arXiv: 1206.0594.
- Liberty, E., Sriharsha, R., and Sviridenko, M. An Algorithm for Online K-Means Clustering. In *2016 Proceedings of the Meeting on Algorithm Engineering and Experiments (ALENEX)*, Proceedings, pp. 81–89. Society for Industrial and Applied Mathematics, December 2015. doi: 10.1137/1.9781611974317.7. URL <https://epubs.siam.org/doi/abs/10.1137/1.9781611974317.7>.
- Lytova, A. and Pastur, L. Central limit theorem for linear eigenvalue statistics of random matrices with independent entries. *The Annals of Probability*, 37(5):1778–1840, 2009. Publisher: Institute of Mathematical Statistics.
- 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 dimensional data. *arXiv:1711.03404 [cs, stat]*, November 2017. URL <http://arxiv.org/abs/1711.03404>. arXiv: 1711.03404.
- Marčenko, V. A. and Pastur, L. A. Distribution of eigenvalues 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 efficiently updating the eigen-system. *Pattern Recognition*, 43(1):113–127, January 2010. ISSN 0031-3203. doi: 10.1016/j.patcog.2009.06.001. URL <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/journals/annals-of-statistics/volume-41/issue-3/Asymptotic-power-of-sphericity-tests-for-high-dimensional-data/10.1214/13-AOS1100.full>. Publisher: Institute of Mathematical Statistics.
- Simonyan, K. and Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv:1409.1556 [cs]*, April 2015. URL <http://arxiv.org/abs/1409.1556>. arXiv: 1409.1556.
- Stein, C. M. Estimation of the Mean of a Multivariate 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/journals/annals-of-statistics/volume-9/issue-6/Estimation-of-the-Mean-of-a-Multivariate-Normal-Distribution/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 Computer Science, pp. 81–92, Berlin, Heidelberg, 2007. Springer. ISBN 978-3-540-74825-0. doi: 10.1007/978-3-540-74825-0.8.
- Trench, W. F. Some Spectral Properties of Hermitian Toeplitz Matrices. *SIAM Journal on Matrix Analysis and Applications*, 15(3):938–942, July 1994. ISSN 0895-4798. doi: 10.1137/S0895479892239007. URL <https://epubs.siam.org/doi/10.1137/S0895479892239007>. Publisher: Society for Industrial and Applied Mathematics.
- von Luxburg, U. A tutorial on spectral clustering. *Statistics and Computing*, 17(4):395–416, December 2007. ISSN 1573-1375. doi: 10.1007/s11222-007-9033-z. URL <https://doi.org/10.1007/s11222-007-9033-z>.
- Xiao, H., Rasul, K., and Vollgraf, R. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. *arXiv:1708.07747 [cs, stat]*, September 2017. URL <http://arxiv.org/abs/1708.07747>. arXiv: 1708.07747.
- Yoo, S., Huang, H., and Kasiviswanathan, S. P. Streaming spectral clustering. In *2016 IEEE 32nd International Conference on Data Engineering (ICDE)*, pp. 637–648, May 2016. doi: 10.1109/ICDE.2016.7498277.
- Zarrouk, T., Couillet, R., Chatelain, F., and Le Bihan, N. Performance-Complexity Trade-Off in Large Dimensional Statistics. In *2020 IEEE 30th International Workshop on Machine Learning for Signal Processing (MLSP)*, pp. 1–6, September 2020. doi: 10.1109/MLSP49062.2020.9231568. ISSN: 1551-2541.

Zhang, T., Ramakrishnan, R., and Livny, M. BIRCH: an efficient data clustering method for very large databases. In *Proceedings of the 1996 ACM SIGMOD international conference on Management of data*, SIGMOD '96, pp. 103–114, New York, NY, USA, June 1996. Association for Computing Machinery. ISBN 978-0-89791-794-0. doi: 10.1145/233269.233324. URL <https://doi.org/10.1145/233269.233324>.

Zubaroglu, A. and Atalay, V. Data stream clustering: a review. *Artificial Intelligence Review*, 54(2):1201–1236, February 2021. ISSN 1573-7462. doi: 10.1007/s10462-020-09874-x. URL <https://doi.org/10.1007/s10462-020-09874-x>.

A. Toolbox

A.1. Useful results

For the reader's convenience, we recall here some standard results.

Proposition A.1 (Jensen's inequality). *Let*

- φ be a real convex function,
- $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ be n points in the domain of φ ,
- $\lambda_1, \dots, \lambda_n \in [0, 1]$ be such that $\sum_{i=1}^n \lambda_i = 1$.

Then,

$$\varphi \left(\sum_{i=1}^n \lambda_i \mathbf{x}_i \right) \leq \sum_{i=1}^n \lambda_i \varphi(\mathbf{x}_i).$$

Proposition A.2 (Resolvent identity). *Given two invertible matrices \mathbf{A} and \mathbf{B} ,*

$$\mathbf{A}^{-1} - \mathbf{B}^{-1} = \mathbf{A}^{-1} (\mathbf{B} - \mathbf{A}) \mathbf{B}^{-1}.$$

Proposition A.3 (Cauchy's integral formula). *Let*

- U be a simply connected open subset of \mathbb{C} ,
- $f : U \rightarrow \mathbb{C}$ be a holomorphic function,
- Γ be a positively-oriented closed contour inside U ,
- $\mathring{\Gamma}$ be the open domain enclosed by Γ .

For $z_0 \in \mathbb{C}$,

$$-\frac{1}{2i\pi} \oint_{\Gamma} \frac{f(z)}{z_0 - z} dz = \begin{cases} f(z_0) & \text{if } z_0 \in \mathring{\Gamma} \\ 0 & \text{otherwise} \end{cases}.$$

A.2. Singular value inequalities

We recall here two useful inequalities which can be found in Corollary A.12 and Theorem A.14 of Bai & Silverstein (2010).

For a matrix \mathbf{A} , we denote by $s_i(\mathbf{A})$ its i -th singular value in descending order.

Proposition A.4. *For any $n \times n$ complex matrix \mathbf{A} ,*

$$|\text{tr } \mathbf{A}| \leq \sum_{i=1}^n s_i(\mathbf{A}).$$

Proposition A.5. *Let \mathbf{A} and \mathbf{B} be two complex matrices of size $p \times n$ and $n \times m$ respectively. For any integer $1 \leq k \leq \min(m, n, p)$,*

$$\sum_{i=1}^k s_i(\mathbf{AB}) \leq \sum_{i=1}^k s_i(\mathbf{A}) s_i(\mathbf{B}).$$

A.3. Puncturing 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.6.

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

Proof.

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

□

Lemma A.7 ((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.8. *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.

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

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.6, 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,

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

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.

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

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}_{r,\cdot}} \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}(\frac{1}{p})} \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.8, 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 η

Now taking $\mathbf{A} = \mathbf{C}_{\cdot,r}\mathbf{C}_{s,\cdot}$ in the first equation of Proposition A.8, 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 $\eta_0 = \frac{1}{n} \text{tr} \bar{\eta}$,

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

Remark B.1. To perform the previous matrix inversion, we must check that $\frac{1}{1-z-\eta_0} \neq -\frac{p}{\psi_k}$ for any given $0 \leq k < n$ such that $\psi_k \neq 0$. This can be proven by contradiction: assuming that $\frac{1}{1-z-\eta_0} = -\frac{p}{\psi_k}$, we have $\mathbf{C}^2 = (\mathbf{C} - \psi_k \mathbf{I}_n) \bar{\boldsymbol{\eta}}$. Since $\text{rank}(\mathbf{C} - \psi_k \mathbf{I}_n) = \text{rank } \mathbf{C} - 1$, we conclude that $\text{rank } \mathbf{C}^2 < \text{rank } \mathbf{C}$, which is a contradiction.

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

Theorem B.2 (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 η_0 is solution to the fixed-point equation

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

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.3. 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 η_0 from going to $\pm\infty$.

Proposition B.4 (Fixed-point equation for $m(z)$). *Under the setting of Theorem B.2, $m(z)$ is also solution to a fixed-point equation:*

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

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

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

therefore,

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

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

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.5.

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

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 $\mathbf{Q}_0 = \left(\frac{\mathbf{Z}^\top \mathbf{Z}}{p} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1}$.

$$\begin{aligned} \left| \frac{1}{n} \operatorname{tr} \mathbf{Q}_A - \frac{1}{n} \operatorname{tr} \mathbf{Q}_0 \right| &= \frac{1}{n} |\operatorname{tr} \mathbf{Q}_0 (\mathbf{A} \odot \mathbf{C}) \mathbf{Q}_A| \quad \text{from the resolvent identity (Proposition A.2)} \\ &\leq \frac{1}{n} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \sum_{i=1}^n s_i(\mathbf{A} \odot \mathbf{C}) \quad \text{from Proposition A.4 and A.5} \\ &\leq \frac{1}{n} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \sqrt{n \sum_{i=1}^n s_i^2(\mathbf{A} \odot \mathbf{C})} \quad \text{from Jensen's inequality (Proposition A.1)} \\ &\leq \frac{1}{n} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \sqrt{n \sum_{i=1}^n s_i^2(\mathbf{A})} \quad \text{since } \|\mathbf{A} \odot \mathbf{C}\|_F \leq \|\mathbf{A}\|_F \\ \left| \frac{1}{n} \operatorname{tr} \mathbf{Q}_A - \frac{1}{n} \operatorname{tr} \mathbf{Q}_0 \right| &\leq \sqrt{\frac{3K}{n}} \|\mathbf{Q}_0\| \|\mathbf{Q}_A\| \|\mathbf{A}\| = \mathcal{O}_{n,p,L \rightarrow +\infty} \left(\frac{1}{\sqrt{n}} \right) \quad \text{since } \mathbf{A} \text{ has rank at most } 3K. \end{aligned}$$

□

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.2, Proposition B.5 justifies that the limiting spectral distribution is unchanged by the presence of signal.

Let us now seek a deterministic equivalent of $\mathbf{Q} = \left(\frac{\mathbf{X}^\top \mathbf{X}}{p} \odot \mathbf{C} - z\mathbf{I}_n \right)^{-1}$.

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 $\mathbf{Q}_{i,j}$,

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

$\mathbf{P}^\top \mathbf{P}$ is deterministic so there is no work to do on the last term of the sum. Expanding the other terms yields (see Proposition A.8)

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

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.6.

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

Proof. Similarly to the proof of Proposition B.5, 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.

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

Hence,

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

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

So far, we have

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

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.7.

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

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.5.

Now, using the third equation of Proposition A.8, 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 $T \geq 1$,

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

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

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

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

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

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

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

which completes the proof. \square

We can now state the following theorem.

Theorem B.8 (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*

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

Remark B.9. This is coherent with the result of Theorem B.2 when $\mathbf{P} = \mathbf{0}$.

Remark B.10. From Proposition B.4, we see that $1 + zm(z) = \frac{p}{n} \text{tr} \mathbf{\Lambda}$.

Remark B.11. As m is a Stieltjes transform, we know that $|m(z)| \leq \frac{1}{\text{dist}(z, \text{supp} \mu)}$ for all $z \in \mathbb{C} \setminus \text{supp} \mu$. Therefore $\text{dist}(z, \text{supp} \mu) > \frac{2L-1}{p} \implies \left| \frac{2L-1}{p} m(z) \right| < 1$.

Remark B.12. It is strongly believed that the condition $\left| \frac{2L-1}{p} m(z) \right| < 1$ can be removed by means of analytic continuation.

C. Proof of Theorem 3.3

In this section, we use the following notation:

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

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*

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

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

Proof. From Theorem B.8,

$$\begin{aligned}\mathbf{Q} &\leftrightarrow m(z) \left(\mathbf{I}_n + \|\boldsymbol{\mu}\|^2 \mathbf{j}\mathbf{j}^\top \odot \mathbf{F}\boldsymbol{\Lambda}\mathbf{F}^* \right)^{-1} \\ &\leftrightarrow m(z) \left(\mathbf{I}_n + \|\boldsymbol{\mu}\|^2 \mathbf{D}_j \mathbf{F}\boldsymbol{\Lambda}\mathbf{F}^* \mathbf{D}_j \right)^{-1} \\ \mathbf{Q} &\leftrightarrow m(z) [\mathbf{D}_j \mathbf{F}] \left(\mathbf{I}_n + \|\boldsymbol{\mu}\|^2 \boldsymbol{\Lambda} \right)^{-1} [\mathbf{D}_j \mathbf{F}]^* \quad \text{since } \mathbf{D}_j \mathbf{F} \text{ is a unitary matrix.}\end{aligned}$$

□

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

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

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 $\text{sp}_\infty(\tilde{\mathbf{K}}_L)$.

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

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

and, supposing $z \in \mathbb{C} \setminus \text{sp}_\infty(\tilde{\mathbf{K}}_L)$, we have, from Proposition B.4,

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

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

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

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

C.2. Alignments and phase transition

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

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

Therefore, with Cauchy's integral formula and a positively-oriented closed contour Γ_k circling around ξ_k and leaving the other eigenvalues outside, we can have access to the quantity

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

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 ξ_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 Γ_k be a positively-oriented closed contour circling around $\bar{\xi}_k$ and leaving all the $\bar{\xi}_l \neq \bar{\xi}_k$ 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

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

Proof. By residue calculus,

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

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

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

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

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

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

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

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.4 yields

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

thus,

$$\begin{aligned} \frac{m^2(\bar{\xi}_k)}{m'(\bar{\xi}_k)} &= -\bar{\xi}_k m(\bar{\xi}_k) + \frac{p}{n} \sum_{r=0}^{n-1} \frac{m(\bar{\xi}_k) \frac{\psi_r}{p}}{\left(1 + m(\bar{\xi}_k) \frac{\psi_r}{p} \right)^2} \\ &= 1 - \frac{p}{n} \sum_{r=0}^{n-1} \frac{m(\bar{\xi}_k) \frac{\psi_r}{p}}{1 + m(\bar{\xi}_k) \frac{\psi_r}{p}} + \frac{p}{n} \sum_{r=0}^{n-1} \frac{m(\bar{\xi}_k) \frac{\psi_r}{p}}{\left(1 + m(\bar{\xi}_k) \frac{\psi_r}{p} \right)^2} \quad \text{from Proposition B.4} \\ \frac{m^2(\bar{\xi}_k)}{m'(\bar{\xi}_k)} &= 1 - \frac{p}{n} \sum_{r=0}^{n-1} \left[\frac{m(\bar{\xi}_k) \frac{\psi_r}{p}}{1 + m(\bar{\xi}_k) \frac{\psi_r}{p}} \right]^2 \end{aligned}$$

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

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

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

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

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

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 $x' \left(\frac{-1}{(\|\boldsymbol{\mu}\|^2+1)^{\frac{\psi_k}{p}}} \right) > 0$.

We have

$$x(m) = -\frac{1}{m} + \frac{p}{n} \sum_{l=0}^{n-1} \frac{\psi_l/p}{1 + m\psi_l/p}$$

$$x'(m) = \frac{1}{m^2} - \frac{p}{n} \sum_{l=0}^{n-1} \left[\frac{\psi_l/p}{1 + m\psi_l/p} \right]^2$$

thus

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

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

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 $\bar{\zeta}_k = \frac{\|\boldsymbol{\mu}\|^2}{\|\boldsymbol{\mu}\|^2+1} (1-p) < 0$.

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

D. Predictions with a Toeplitz mask

Figures 8a and 8b compare simulations with a Toeplitz mask and the predictions of Theorems 3.1 and 3.3 with the ψ_k 's replaced by the τ_k 's and \mathbf{F} replaced by \mathbf{G} .

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 $|\mathbf{u}_0^* \mathbf{v}_0|^2$ also fit well the predicted curve.

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 $\tau_1 \neq \tau_{n-1}$).

E. K -classes online kernel spectral clustering algorithm

E.1. General presentation and simulations

We use a set of spike eigenvectors $\left\{ \mathbf{u}_k^{(t)} \right\}_{k \in \mathcal{K}}$ (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

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

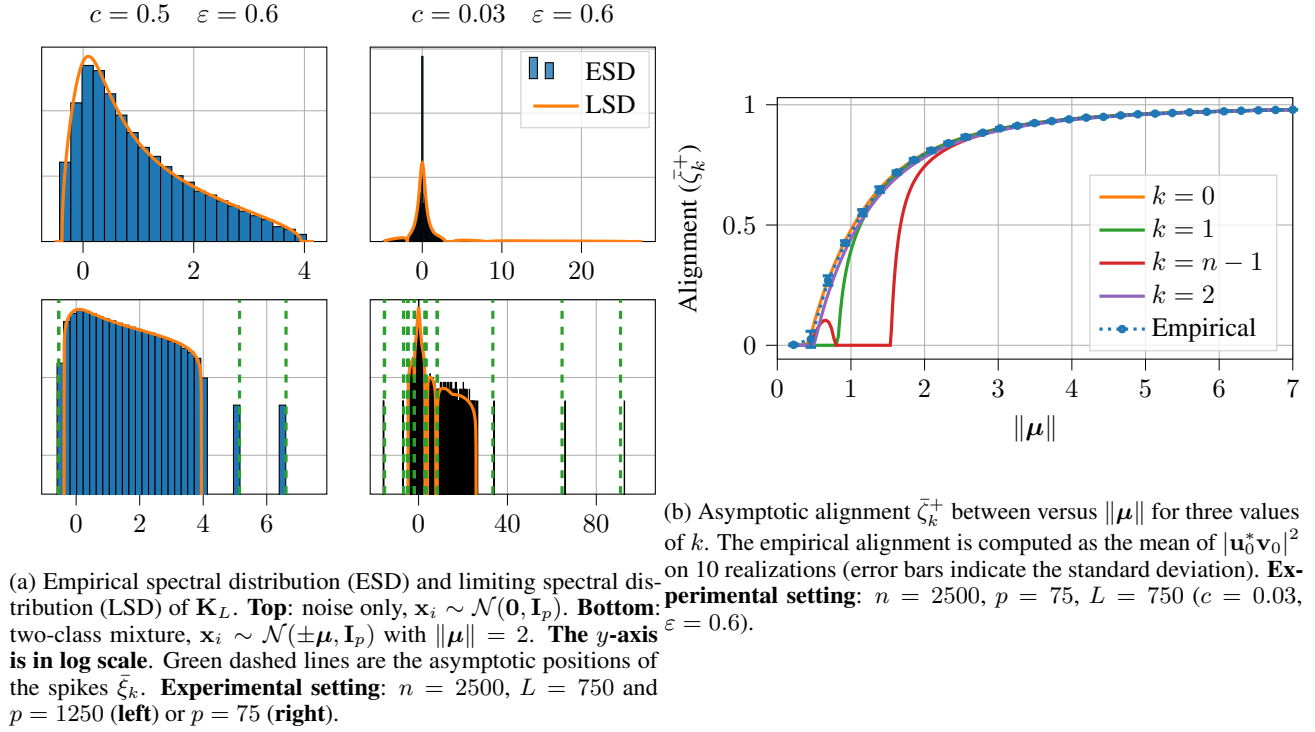


Figure 8. Predictions of Theorems 3.1 and 3.3 adapted for a Toeplitz mask.

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 $\left\{ \mathbf{u}_k^{(t)} \right\}_{k \in \mathcal{K}}$ (see the left part of Figure 9). Each point is then associated to the class whose curve is the nearest in the $|\mathcal{K}|$ -dimensional space. 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 (thus $|\mathcal{K}| = 5$).

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 class¹⁶ (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 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.

¹⁶This is only the first dimension of a 5-dimensional trend.

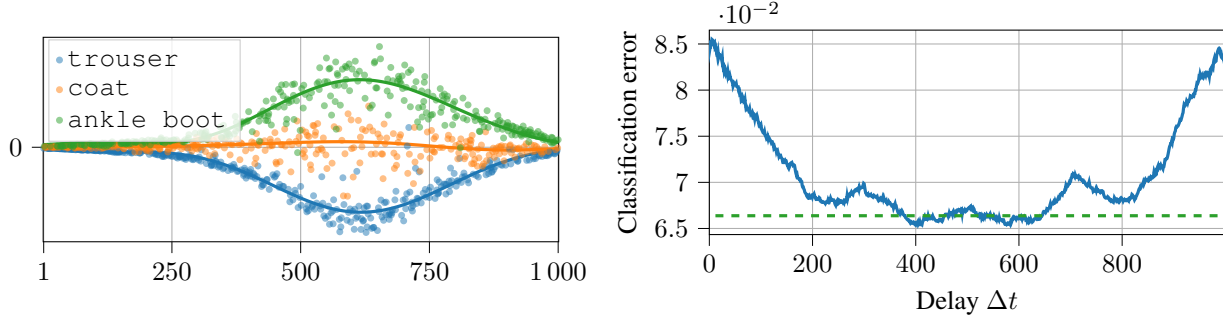


Figure 9. Clustering on Fashion-MNIST images (trouser vs. coat vs. ankle boot). **Left:** dominant eigenvector of $\mathbf{K}_L^{(t)}$. Solid curves are the estimated trend of each class $\mathbf{h}_{k,C}^{(t)}$. **Right:** Classification error against delay Δt . This is the mean classification error 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:** $T = 21\,000$, $n = 1\,000$, $p = 784$, $L = 100$.

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 to be linear combinations of.

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

In order to compute an estimation $\left\{ \hat{\mathcal{C}}^{(t)}[i] \right\}_{1 \leq i \leq n}$ of the classes at a given time t , we propose a two-step algorithm. Firstly, we compute a rough estimation $\left\{ \hat{\mathcal{C}}_0^{(t)}[i] \right\}_{1 \leq i \leq n}$ 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.

Iteration Let us suppose we have already estimated $\hat{\mathcal{C}}_0[1], \dots, \hat{\mathcal{C}}_0[i-1]$ 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 $\left[\mathbf{h}_{k, \hat{\mathcal{C}}_0[j]} \right]_j$ (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 $\hat{\mathcal{C}}_0[i]$ — 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.

Initialization 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, α .

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 are usually not very satisfying but still remain 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

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

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

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

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

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

Algorithm 3 Classification

Input: $K, \{\hat{C}_0[i]\}_{1 \leq i \leq n}, \{\mathbf{u}_k\}_{k \in \mathcal{K}}, \{\mathbf{v}_k\}_{k \in \mathcal{K}_*}$.
Output: $\{\hat{C}[i]\}_{1 \leq i \leq n}$.
for $i = 1$ to n **do**
 $\hat{C}[i] \leftarrow \hat{C}_0[i]$
end for
repeat
 for $\hat{C} \in \{\hat{C}_1, \dots, \hat{C}_K\}$ **do**
 $\hat{\mathbf{h}}_{\mathcal{K}, \hat{C}} \leftarrow \mathbf{g}_{\mathcal{K}_*} \left([\mathbf{g}_{\mathcal{K}_*}]_{\hat{C}}^\top [\mathbf{g}_{\mathcal{K}_*}]_{\hat{C}} \right)^{-1} [\mathbf{g}_{\mathcal{K}_*}]_{\hat{C}}^\top [\mathbf{u}_{\mathcal{K}}]_{\hat{C}}$
 end for
 for $i = 1$ to n **do**
 $\hat{C}[i] \leftarrow \arg \min_{\hat{C} \in \{\hat{C}_1, \dots, \hat{C}_K\}} \left\| [\mathbf{u}_{\mathcal{K}}]_i - [\hat{\mathbf{h}}_{\mathcal{K}}]_i \right\|$
 end for
until convergence

E.2.3. FINAL ALGORITHM

Algorithm 4 Online kernel spectral clustering

Input: $K, \mathcal{K}, \{\mathbf{g}_k\}_{k \in \mathcal{K}_*}$.
Parameters: H, α .
Output: $\{\hat{C}_t[s]\}_{\substack{1 \leq s \leq n \\ n \leq t \leq T}}$.
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 $\mathbf{K}_L^{(t)}$.
 $\mathbf{u}_{\mathcal{K}}^{(t)} \leftarrow \text{PowerIteration}(\mathbf{K}_L^{(t)}, \mathbf{u}_{\mathcal{K}}^{(t-1)})$.
 if $1 \leq t \leq n$ **then**
 Do an iteration as in Algorithm 2.
 end if
 if $t \geq n$ **then**
 Compute $\{\hat{C}^{(t)}[s]\}_{1 \leq s \leq n}$ according to Algorithm 3 with $\{\hat{C}^{(t-1)}[s]\}_{1 \leq s \leq n-1}$.
 end if
end for

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 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 good basis to estimate the classes at time t (both $\hat{C}^{(t-1)}[s + 1]$ and $\hat{C}^{(t)}[s]$ are estimates of the class of \mathbf{x}_{t-n+s}). Moreover, the few 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 not differ much from one time step to another). The final algorithm is presented in Algorithm 4.

F. Exponential smoothing with missing data

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 is to perform an exponential smoothing:

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

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

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} and \mathbf{s}_t only. Expanding the recurrence relation, we have

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

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:

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

Using the following formulae,

$$\begin{aligned} \sum_{k=1}^{h-1} k (1 - \alpha)^k &= \frac{1 - \alpha}{\alpha} \left(1 - h(1 - \alpha)^{h-1} \right) + \left(\frac{1 - \alpha}{\alpha} \right)^2 \left(1 - (1 - \alpha)^{h-1} \right) \\ \text{and} \quad \sum_{k=1}^{h-1} k (1 - \alpha)^{h-k} &= \frac{1 - \alpha}{\alpha} (h - 1) - \left(\frac{1 - \alpha}{\alpha} \right)^2 \left(1 - (1 - \alpha)^{h-1} \right) \end{aligned}$$

we have

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