



**HAL**  
open science

# Asymptotic Gaussian Fluctuations of Eigenvectors in Spectral Clustering

Hugo Lebeau, Florent Chatelain, Romain Couillet

► **To cite this version:**

Hugo Lebeau, Florent Chatelain, Romain Couillet. Asymptotic Gaussian Fluctuations of Eigenvectors in Spectral Clustering. IEEE Signal Processing Letters, 2024, 31, pp.1920-1924. 10.1109/LSP.2024.3422886 . hal-04673319

**HAL Id: hal-04673319**

**<https://hal.science/hal-04673319>**

Submitted on 20 Aug 2024

**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

# Asymptotic Gaussian Fluctuations of Eigenvectors in Spectral Clustering

Hugo Lebeau, Florent Chatelain, Romain Couillet

**Abstract**—The performance of spectral clustering relies on the fluctuations of the entries of the eigenvectors of a similarity matrix, which has been left uncharacterized until now. In this letter, it is shown that the *signal + noise* structure of a general spike random matrix model is transferred to the eigenvectors of the corresponding Gram kernel matrix and the fluctuations of their entries are Gaussian in the large-dimensional regime. This CLT-like result was the last missing piece to precisely predict the classification performance of spectral clustering. The proposed proof is very general and relies solely on the rotational invariance of the noise. Numerical experiments on synthetic and real data illustrate the universality of this phenomenon.

**Index Terms**—Spectral clustering, central limit theorem, kernel matrix, spike eigenvector, Gaussian fluctuations.

## I. INTRODUCTION

SPECTRAL clustering is a popular unsupervised classification technique which finds applications in many domains, such as image segmentation [1], text mining [2], and as a general purpose method for data analysis [3], [4], [5]. It relies on the spectrum of a suitably chosen similarity matrix to perform dimensionality reduction before applying a standard clustering algorithm such as  $K$ -means. Consider, e.g., the following toy example where  $n$  vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$  are separated in two clusters  $\mathcal{C}^+, \mathcal{C}^-$  centered around  $+\boldsymbol{\mu}, -\boldsymbol{\mu}$  respectively, i.e.,  $\mathbf{x}_i = \pm\boldsymbol{\mu} + \mathbf{w}_i$  where  $\mathbf{w}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Then, the dominant eigenvector  $\hat{\mathbf{v}}$  of the Gram kernel matrix  $\mathbf{K} = \frac{1}{p}[\mathbf{x}_i^\top \mathbf{x}_j]_{1 \leq i, j \leq n}$  is an information-theoretically optimal estimator [6] of the vector  $\frac{1}{\sqrt{n}}\mathbf{j}$  such that  $j_i = \pm 1$  if  $\mathbf{x}_i \in \mathcal{C}^\pm$ . In this case, clustering is achieved with the trivial decision rule  $\mathbf{x}_i \rightarrow \mathcal{C}^\pm$  if  $\hat{v}_i \geq 0$ .

The achievable performances of spectral clustering can be theoretically predicted thanks to the study of random matrix models corresponding to similarity matrices. For this purpose, random matrix theory offers powerful tools [7], [8], [9]. In particular, it allows to derive the limiting spectral distribution of the kernel matrix and to predict the position of isolated eigenvalues in *spiked* random matrix models [10], [11], [12]. The latter are of particular importance as, in a wide range of problems, the information of interest can be modeled as a low-rank signal corrupted with noise. In our previous toy example, the data matrix  $\mathbf{X} = [\mathbf{x}_1 \ \dots \ \mathbf{x}_n]$  is a rank-one perturbation  $\boldsymbol{\mu}\mathbf{j}^\top$  of a noise matrix  $\mathbf{W}$  with i.i.d.  $\mathcal{N}(0, 1)$  entries. In order to theoretically predict the error rate of

spectral clustering for a given signal-to-noise ratio, one must therefore study the behavior of the dominant eigenvectors of the similarity matrix. Tools such as the ones used in [12], [13] allow to express the quality of their alignment with the true underlying signal, i.e.,  $\frac{1}{\sqrt{n}}|\mathbf{j}^\top \hat{\mathbf{v}}|$ . Although this tells us when an estimation of the signal is possible (depending on the signal-to-noise ratio) and its efficiency, a precise characterization of the fluctuations of the entries of spiked eigenvectors still lacks to rigorously predict the error rate of spectral clustering. Indeed, in our toy example, the expected error rate  $\mathbb{P}(\hat{v}_i j_i < 0)$  cannot be expressed unless the law of  $\hat{\mathbf{v}}$  is known.

Yet, it is often stated that the entries of  $\hat{\mathbf{v}}$  have Gaussian fluctuations in the large-dimensional regime, so that  $\mathbb{P}(\hat{v}_i j_i < 0)$  is a Gaussian integral. In [14], this result is formally stated but no proof is given. Hence, we fill this missing gap with a rigorous proof of this phenomenon for a general spiked random matrix model. Although we stick to a simple *signal + noise* model here, the proposed proof is not restricted to Gaussian noise (in fact, the noise only needs to be rotationally invariant) and can easily be adapted to most standard spiked models (such as, notably, the general model considered in [11]). Our result and its proof thus support a wide range of previous works studying the performance of spectral algorithms. The demonstration can be summarized in two simple facts 1) an eigenvector of the kernel matrix can be decomposed into the sum of a deterministic signal part and a random noise part 2) the random part is uniformly distributed on a certain sphere, hence any finite subset of its entries tends to a centered Gaussian vector in the large-dimensional limit.

In this letter, we consider a general *signal + noise* random matrix model and briefly recall known results regarding its limiting spectral distribution and the behavior of its dominant eigenvalues and eigenvectors. Then, we show that the entries of the kernel eigenvectors indeed have Gaussian fluctuations in the large-dimensional regime. We present a short, self-contained and general proof which is our main contribution. Finally, we illustrate this result with numerical experiments on synthetic and real data.

**Simulations.** Python codes to reproduce simulations are available in the following GitHub repository [https://github.com/HugoLebeau/asymptotic\\_fluctuations\\_spectral\\_clustering](https://github.com/HugoLebeau/asymptotic_fluctuations_spectral_clustering).

## II. MODEL AND MAIN RESULT

### A. Notations

The symbols  $a$ ,  $\mathbf{a}$  and  $\mathbf{A}$  respectively denote a scalar, a vector and a matrix. Their entries are  $a_i$  and  $A_{i,j}$ . The set of positive integers below  $n$  is  $[n] = \{1, \dots, n\}$ . The

H. Lebeau and R. Couillet are with Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, 38000 Grenoble, France, (e-mail: {hugo.lebeau, romain.couillet}@univ-grenoble-alpes.fr).

F. Chatelain is with Univ. Grenoble Alpes, CNRS, Grenoble INP, GIPSA-lab, 38000 Grenoble, France (e-mail: florent.chatelain@grenoble-inp.fr).

cardinality of a set  $\mathcal{E}$  is  $|\mathcal{E}|$ . Given an ordered set of indices  $\mathcal{I} = (i_1, \dots, i_{|\mathcal{I}|})$ ,  $[\mathbf{a}]_{\mathcal{I}}$  is the vector  $[a_{i_1} \dots a_{i_{|\mathcal{I}|}}]^\top$ . The norm of a vector  $\mathbf{x} \in \mathbb{R}^n$  is  $\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}}$ . The unit sphere in  $\mathbb{R}^n$  is  $\mathbb{S}^{n-1} = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| = 1\}$ . The  $n \times n$  identity matrix is  $\mathbf{I}_n$ . For a real number  $x \in \mathbb{R}$ ,  $[x]^+ = \max(0, x)$  and  $\delta_x$  is the Dirac measure at  $x$ . The imaginary unit is denoted  $i$ . If the random variable  $X$  follows the law  $\mathcal{L}$ , we write  $X \sim \mathcal{L}$ . The convergence in distribution of a sequence of random variables  $(X_n)_{n \geq 0}$  to  $\mathcal{L}$  is denoted  $X_n \xrightarrow[n \rightarrow +\infty]{\mathcal{D}} \mathcal{L}$ . Its almost sure convergence to  $L$  is denoted  $X_n \xrightarrow[n \rightarrow +\infty]{\text{a.s.}} L$ . The multivariate normal distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$  is denoted  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . The set of eigenvalues of a square matrix  $\mathbf{A}$  is its spectrum,  $\text{Sp } \mathbf{A}$ . Given two real-valued sequences  $(u_n)_{n \geq 0}$  and  $(v_n)_{n \geq 0}$ , we write  $u_n = \mathcal{O}(v_n)$  if  $|u_n/v_n|$  is bounded as  $n \rightarrow +\infty$  and  $u_n \asymp v_n$  if  $u_n/v_n \rightarrow 1$ .

## B. Spiked Matrix Model

Consider the following statistical model

$$\mathbf{X} = \mathbf{P} + \mathbf{W} \in \mathbb{R}^{p \times n}, \quad \mathbf{P} = \mathbf{L}\mathbf{V}^\top \quad (1)$$

with  $\mathbf{L} \in \mathbb{R}^{p \times K}$  and  $\mathbf{V} = [\mathbf{v}_1 \dots \mathbf{v}_K] \in \mathbb{R}^{n \times K}$  such that  $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_K$ . It models a *low-rank* signal  $\mathbf{P}$  corrupted by additive Gaussian noise  $W_{i,j}$   $\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ . In a spectral clustering perspective,  $K$  represents the number of classes and  $\mathbf{P} = \mathbf{M}\mathbf{J}^\top$  where  $\mathbf{M} = [\boldsymbol{\mu}_1 \dots \boldsymbol{\mu}_K]$  is a matrix gathering the  $K$  cluster means and  $J_{i,k} = 1$  if  $\mathbf{x}_i$  is in the  $k$ -th cluster (i.e.,  $\mathbf{x}_i = \boldsymbol{\mu}_k + \mathbf{w}_i$ ) and 0 otherwise. This is congruent with model (1): define the  $K \times K$  diagonal matrix  $\mathbf{D}$  such that  $D_{k,k} = n_k$  where  $n_k$  is the number of samples belonging to the  $k$ -th cluster, then  $\mathbf{L} = \mathbf{M}\mathbf{D}^{1/2}$  and  $\mathbf{V} = \mathbf{J}\mathbf{D}^{-1/2}$ .

Given model (1), we are interested in the reconstruction of  $\mathbf{V}$  from the dominant eigenvectors of the Gram kernel matrix  $\mathbf{K} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$ . We study this problem in the regime where  $K$  is fixed and  $p, n \rightarrow +\infty$  at the same rate, i.e.,  $0 < c \stackrel{\text{def}}{=} \lim p/n < +\infty$ . This models the fact that, in practice, the number of samples  $n$  is comparable to the number of features  $p$  and they are both large. Moreover, we make the following assumptions.

**Assumption 1.** All classes are of comparable size, i.e.,  $\liminf n_k/n > 0$  as  $p, n \rightarrow +\infty$  for all  $k \in [K]$ .

**Assumption 2.**  $\lim_{p, n \rightarrow +\infty} \max_{\substack{1 \leq i \leq n \\ 1 \leq k \leq K}} \sqrt{n} V_{i,k}^2 = 0$ .

**Assumption 3.** As  $n \rightarrow +\infty$ , the eigenvalues  $\ell_1 \geq \dots \geq \ell_K > 0$  of  $\frac{1}{n} \mathbf{L}^\top \mathbf{L}$  are not degenerate (i.e., have multiplicity one) and the columns of  $\mathbf{V}$  are ordered accordingly.

Assumption 3 is only to simplify the presentation of the results so it is not necessary, but often verified in practice. However, Assumption 2 states that  $\mathbf{V}$  must be *delocalized*, i.e., not sparse. It is naturally verified for spectral clustering as a result of Assumption 1 since  $\mathbf{V} = \mathbf{J}\mathbf{D}^{-1/2}$ , but the results presented below concern the statistical model (1), which is more general and also encompasses PCA for example [13].

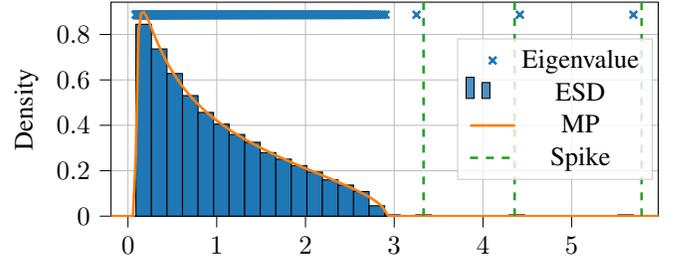


Fig. 1. Empirical Spectral Distribution (ESD) of  $\mathbf{K} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$  and Marčenko-Pastur Distribution (MP). The green dashed lines are the positions  $\xi_k$  of isolated eigenvalues predicted by Theorem 1. **Experimental setting:**  $n = 1000$ ,  $p = 2000$ ,  $K = 3$ ,  $(n_1, n_2, n_3) = (333, 334, 333)$ ,  $(\|\boldsymbol{\mu}_1\|, \|\boldsymbol{\mu}_2\|, \|\boldsymbol{\mu}_3\|) = (3, 4, 5)$ .

## C. Eigenvalue Distribution and Spiked Eigenvalues

We briefly recall known results on model (1) in order to set the ground for our main result in Theorem 2.

Firstly, the empirical spectral distribution of  $\mathbf{K}$ , that is  $\frac{1}{n} \sum_{\lambda \in \text{Sp } \mathbf{K}} \delta_\lambda$ , converges weakly almost surely to the Marčenko-Pastur distribution  $\mu_{\text{MP}} = [1-c]^+ \delta_0 + \nu$  where  $\nu$  has density supported on  $[E_-, E_+]$  with  $E_\pm = (1 \pm \sqrt{c^{-1}})^2$  [15], [7], [8], [9]. In other words, as  $p, n \rightarrow +\infty$ , the histogram of eigenvalues of  $\mathbf{K}$  approaches  $\mu_{\text{MP}}$ , as depicted in Figure 1.

Due to the low-rank perturbation  $\mathbf{P}$ , the  $K$  dominant eigenvalues of  $\mathbf{K}$  may isolate themselves from the bulk characterized by the Marčenko-Pastur distribution if their corresponding signal-to-noise ratios (the eigenvalues of  $\frac{1}{n} \mathbf{L}^\top \mathbf{L}$ ) are large enough — they are then called *spikes*. Their behavior is specified in the following theorem, which is a particular case of Theorem 2 in [13].

**Theorem 1 (Spikes).** Let  $(\lambda_k, \hat{\mathbf{v}}_k)_{k \in [K]}$  denote the dominant eigenvalue-eigenvector pairs of  $\mathbf{K}$  such that  $\lambda_1 \geq \dots \geq \lambda_K$ . Then, for all  $k \in [K]$ ,

$$\lambda_k \xrightarrow[p, n \rightarrow +\infty]{\text{a.s.}} \xi_k \stackrel{\text{def}}{=} \begin{cases} \frac{(\ell_k + c)(\ell_k + 1)}{\ell_k c} & \text{if } \ell_k > \sqrt{c} \\ E_+ & \text{otherwise} \end{cases},$$

$$|\mathbf{v}_k^\top \hat{\mathbf{v}}_k|^2 \xrightarrow[p, n \rightarrow +\infty]{\text{a.s.}} \zeta_k \stackrel{\text{def}}{=} \begin{cases} 1 - \frac{\ell_k + c}{\ell_k (\ell_k + 1)} & \text{if } \ell_k > \sqrt{c} \\ 0 & \text{otherwise} \end{cases}.$$

This result states that  $\lambda_k$  leaves the bulk if, and only if,  $\ell_k > \sqrt{c}$  (this is a well-known phase transition phenomenon [16]) and further gives its almost sure asymptotic position  $\xi_k$ . This is illustrated in Figure 1. Moreover, Theorem 1 indicates the almost sure asymptotic alignment  $\zeta_k$  of the corresponding eigenvector  $\hat{\mathbf{v}}_k$  with the underlying signal  $\mathbf{v}_k$ . This is depicted in Figure 2 (top row).

It should be noted that the previous results are not restricted to Gaussian noise: up to a control on the moments of the distribution, they can be generalized thanks to an “interpolation trick” [17, Corollary 3.1]. In addition, a similar spectral behavior is observed with non-i.i.d. noise following the realistic assumption that it is *concentrated* [18], [19].

## D. Fluctuations of Spiked Eigenvectors Entries

The convergence of  $|\mathbf{v}_k^\top \hat{\mathbf{v}}_k|^2$  to  $\zeta_k$  stated in Theorem 1 is an important result which justifies the use of the dominant

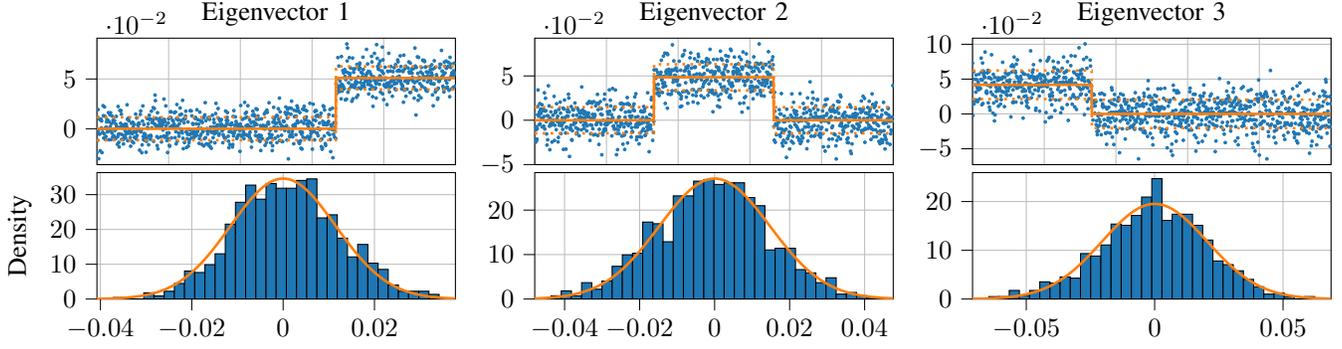


Fig. 2. Dominant eigenvectors of  $\mathbf{K} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$ . **Top:** Coordinates of  $\hat{\mathbf{v}}_k$  (blue) and the underlying signal  $\sqrt{\zeta_k} \mathbf{v}_k$  (orange) with  $\zeta_k$  given in Theorem 1. The dotted orange lines are the  $\pm 1\sigma$ -error curves deduced from Theorem 2. **Bottom:** Histogram of the entries of  $\hat{\mathbf{v}}_k - \sqrt{\zeta_k} \mathbf{v}_k$  (blue) and probability density function of  $\mathcal{N}(0, \frac{1-\zeta_k}{n})$  (orange). **Experimental setting:** like in Figure 1.

eigenvectors of  $\mathbf{K}$  as estimators of the underlying signal  $\mathbf{V}$ . Yet, it is not enough to characterize its reconstruction performance. Indeed, the fluctuations of the entries of  $\hat{\mathbf{v}}_k$  must be known to fully characterize *how* it is aligned with  $\mathbf{v}_k$ .

Consider, e.g., the multi-class spectral clustering problem with  $\mathbf{P} = \mathbf{M}\mathbf{J}^\top$ . Here,  $[\mathbf{v}_k]_i = J_{i,k}/\sqrt{n_k}$ . Hence,  $\mathbf{x}_i$  is classified in the  $k$ -th class if  $[\hat{\mathbf{v}}_k]_i > [\hat{\mathbf{v}}_{k'}]_i$  for all  $k' \neq k$ . The reconstruction performance thus depends on the probability of correct classification  $\mathbb{P}([\hat{\mathbf{v}}_k]_i > [\hat{\mathbf{v}}_{k'}]_i \mid J_{i,k} = 1)$ . In the theorem below, we show that the entries of  $\hat{\mathbf{v}}_k$  asymptotically have Gaussian fluctuations around those of  $\mathbf{v}_k$  with variance  $(1 - \zeta_k)/n$ , as illustrated in the bottom row of Figure 2.

**Theorem 2.** *For all finite ordered set of indices  $\mathcal{I} = (i_1, \dots, i_{|\mathcal{I}|}) \subset [n]$  and  $k \in [K]$ ,*

$$\frac{\sqrt{n} [\hat{\mathbf{v}}_k - \sqrt{\zeta_k} \mathbf{v}_k]_{\mathcal{I}}}{\sqrt{1 - \zeta_k}} \xrightarrow[p, n \rightarrow +\infty]{\mathcal{D}} \mathcal{N}(\mathbf{0}, \mathbf{I}_{|\mathcal{I}|})$$

with  $\hat{\mathbf{v}}_k$  such that  $\mathbf{v}_k^\top \hat{\mathbf{v}}_k \geq 0$  (otherwise, consider  $-\hat{\mathbf{v}}_k$ ).

This result invokes the quantity  $\zeta_k$  introduced in Theorem 1, which quantifies the alignment of  $\hat{\mathbf{v}}_k$  with  $\mathbf{v}_k$ . Theorem 2 specifies that  $[\hat{\mathbf{v}}_k]_{\mathcal{I}}$  behaves like  $\mathcal{N}(\sqrt{\zeta_k} [\mathbf{v}_k]_{\mathcal{I}}, \frac{1-\zeta_k}{n} \mathbf{I}_{|\mathcal{I}|})$  in the large-dimensional regime. That is, the more  $\hat{\mathbf{v}}_k$  is aligned with  $\mathbf{v}_k$  (i.e., the closer  $\zeta_k$  is to 1), the more it concentrates around  $\sqrt{\zeta_k} \mathbf{v}_k$ , since the variance is  $(1 - \zeta_k)/n$ . Furthermore, the entries of  $[\hat{\mathbf{v}}_k]_{\mathcal{I}}$  are *asymptotically independent* for any finite ordered set of indices  $\mathcal{I}$ . In the multi-class spectral clustering problem considered above, since  $\hat{\mathbf{v}}_k$  and  $\hat{\mathbf{v}}_{k'}$  are asymptotically independent if  $k' \neq k$  [20, Theorem 4], Theorem 2 yields

$$\mathbb{P}([\hat{\mathbf{v}}_k]_i > [\hat{\mathbf{v}}_{k'}]_i \mid J_{i,k} = 1) \asymp \Phi \left( \sqrt{\frac{n}{n_k} \frac{\zeta_k}{2 - (\zeta_k + \zeta_{k'})}} \right)$$

where  $\Phi : x \mapsto \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$  is the Gaussian cumulative distribution function.

We prove Theorem 2 in Section III below. The proof hinges on the rotational invariance of the noise (Lemma 1). In fact, it does not need the entries of  $\mathbf{W}$  to be distributed according to the Gaussian law, but only that its distribution be invariant under isometries. This makes it a very general proof, which can easily be adapted to most standard spiked models as those discussed, e.g., in [9, §2.5.4].

### III. PROOF OF MAIN RESULT

Consider the tangent-normal decomposition

$$\hat{\mathbf{v}}_k = \sum_{\kappa=1}^K \tau_{\kappa} \mathbf{v}_{\kappa} + \sqrt{1 - \|\boldsymbol{\tau}\|^2} \hat{\mathbf{v}}_k^{\sharp} \quad (2)$$

where  $\hat{\mathbf{v}}_k^{\sharp} = (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\top) \frac{\hat{\mathbf{v}}_k}{\sqrt{1 - \|\boldsymbol{\tau}\|^2}}$  is a unit-norm vector orthogonal to the span of  $\mathbf{V}$  and  $\boldsymbol{\tau} = [\tau_1 \dots \tau_K]^\top$  with  $\tau_{\kappa} = \mathbf{v}_{\kappa}^\top \hat{\mathbf{v}}_k$  measuring the cosine between  $\mathbf{v}_{\kappa}$  and  $\hat{\mathbf{v}}_k$ . Let  $\mathbf{O}$  be an  $n \times n$  orthogonal matrix such that  $\mathbf{O}\mathbf{V} = \mathbf{V}$  — i.e., a rotational symmetry about the span of  $\mathbf{V}$  — and  $\tilde{\mathbf{K}} \stackrel{\text{def}}{=} \mathbf{O}\mathbf{K}\mathbf{O}^\top$ . Then, since  $\mathbf{K} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$  and  $\mathbf{X} = \mathbf{L}\mathbf{V}^\top + \mathbf{W}$ ,

$$\tilde{\mathbf{K}} = \frac{1}{p} \left( [\mathbf{O}\mathbf{V}] \mathbf{L}^\top \mathbf{L} [\mathbf{O}\mathbf{V}]^\top + [\mathbf{O}\mathbf{V}] \mathbf{L}^\top [\mathbf{W}\mathbf{O}^\top] + [\mathbf{W}\mathbf{O}^\top]^\top \mathbf{L} [\mathbf{O}\mathbf{V}]^\top + [\mathbf{W}\mathbf{O}^\top]^\top [\mathbf{W}\mathbf{O}^\top] \right).$$

**Lemma 1.**  *$\mathbf{W}$  and  $\mathbf{W}\mathbf{O}^\top$  are identically distributed.*

*Proof.* The distribution of  $[\mathbf{W}\mathbf{O}^\top]_{i,j} = \sum_{k=1}^n W_{i,k} O_{j,k}$  is  $\mathcal{N}(0, 1)$  and  $\text{Cov}([\mathbf{W}\mathbf{O}^\top]_{i,j}, [\mathbf{W}\mathbf{O}^\top]_{i',j'})$  is 1 if  $(i, j) = (i', j')$  and 0 otherwise. Hence  $[\mathbf{W}\mathbf{O}^\top]_{i,j} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ .  $\square$

According to the previous lemma,  $\tilde{\mathbf{K}}$  follows the same model as  $\mathbf{K}$  since  $\mathbf{O}\mathbf{V} = \mathbf{V}$ . Therefore, its  $k$ -th dominant eigenvector can likewise be decomposed as

$$\tilde{\mathbf{v}}_k = \sum_{\kappa=1}^K \tilde{\tau}_{\kappa} \mathbf{v}_{\kappa} + \sqrt{1 - \|\tilde{\boldsymbol{\tau}}\|^2} \tilde{\mathbf{v}}_k^{\sharp}$$

with  $\tilde{\tau}_{\kappa} = \mathbf{v}_{\kappa}^\top \tilde{\mathbf{v}}_k$  and  $\tilde{\mathbf{v}}_k^{\sharp} = (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\top) \frac{\tilde{\mathbf{v}}_k}{\sqrt{1 - \|\tilde{\boldsymbol{\tau}}\|^2}}$  identically distributed to  $\hat{\mathbf{v}}_k^{\sharp}$ . Yet,  $\tilde{\mathbf{v}}_k = \mathbf{O}\hat{\mathbf{v}}_k$ . Thus,  $\hat{\mathbf{v}}_k^{\sharp}$  and  $\mathbf{O}\hat{\mathbf{v}}_k^{\sharp}$  are identically distributed for all  $n \times n$  orthogonal matrix  $\mathbf{O}$  such that  $\mathbf{O}\mathbf{V} = \mathbf{V}$ . Consequently, denoting  $\eta$  the probability distribution of  $\hat{\mathbf{v}}_k^{\sharp}$  and  $\mathbf{V}^\perp = \{\mathbf{w} \in \mathbb{R}^n \mid \mathbf{V}^\top \mathbf{w} = \mathbf{0}\}$ , then, for all  $\mathbf{x}, \mathbf{y} \in \mathbb{S}^{n-1} \cap \mathbf{V}^\perp$ , we have  $d\eta(\mathbf{x}) = d\eta(\mathbf{O}\mathbf{x}) = d\eta(\mathbf{y})$  with  $\mathbf{O} = \mathbf{I}_n - \frac{(\mathbf{x}-\mathbf{y})(\mathbf{x}-\mathbf{y})^\top}{1-\mathbf{x}^\top \mathbf{y}}$  satisfying  $\mathbf{O}\mathbf{V} = \mathbf{V}$ . This shows that  $\hat{\mathbf{v}}_k^{\sharp}$  is uniformly distributed on  $\mathbb{S}^{n-1} \cap \mathbf{V}^\perp$ .

Then,  $\hat{\mathbf{v}}_k^{\sharp}$  can be written as  $\mathbf{U}\mathbf{u}$  where  $\mathbf{u}$  is uniformly distributed on  $\mathbb{S}^{n-1-K} \subset \mathbb{R}^{n-K}$  and  $\mathbf{U} \in \mathbb{R}^{n \times (n-K)}$  is such that  $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_{n-K}$  and  $\mathbf{U}^\top \mathbf{V} = \mathbf{0}$  (the columns of  $\mathbf{U}$  form

an orthonormal basis of  $\mathbf{V}^\perp$  in  $\mathbb{R}^n$ ). We use the following theorem to identify the asymptotic distribution of  $\sqrt{n}[\hat{\mathbf{v}}_k^\#]_{\mathcal{I}}$ .

**Theorem 3** ([21], [22]). *The characteristic function of a vector  $\mathbf{w}$  uniformly distributed on  $\mathbb{S}^{n-1}$  is given by  $\varphi_{\mathbf{w}}(\mathbf{t}) \stackrel{\text{def}}{=} \mathbb{E}[e^{i\mathbf{t}^\top \mathbf{w}}] = \Omega_n(\|\mathbf{t}\|)$  where  $\Omega_n$  is such that  $r \mapsto \Omega_n(r\sqrt{n})$  converges uniformly in  $r \geq 0$  to  $r \mapsto e^{-\frac{r^2}{2}}$  as  $n \rightarrow +\infty$ .*

Let  $\mathbf{t} \in \mathbb{R}^n$  be such that  $t_i = 0$  if  $i \notin \mathcal{I}$ . The characteristic function of  $\sqrt{n}[\hat{\mathbf{v}}_k^\#]_{\mathcal{I}}$  is

$$\begin{aligned} \varphi_{\sqrt{n}[\hat{\mathbf{v}}_k^\#]_{\mathcal{I}}}(t_{i_1}, \dots, t_{i_{|\mathcal{I}|}}) &= \mathbb{E} \left[ e^{i\sqrt{n}\mathbf{t}^\top \mathbf{U}\mathbf{u}} \right] \\ &= \Omega_{n-K}(\sqrt{n} \|\mathbf{U}^\top \mathbf{t}\|) \end{aligned}$$

and  $\|\mathbf{U}^\top \mathbf{t}\| = \sqrt{\|\mathbf{t}\|^2 - \|\mathbf{V}^\top \mathbf{t}\|^2} = \|\mathbf{t}\| + \mathcal{O}(\|\mathbf{V}^\top \mathbf{t}\|^2)$ . According to Assumption 2,  $\sqrt{n}\|\mathbf{V}^\top \mathbf{t}\|^2 \rightarrow 0$  as  $p, n \rightarrow +\infty$ , thus  $\Omega_{n-K}(\sqrt{n}\|\mathbf{U}^\top \mathbf{t}\|) = \Omega_{n-K}(\sqrt{n-K}\|\mathbf{t}\| + \epsilon_n)$  with  $\epsilon_n \rightarrow 0$  as  $p, n \rightarrow +\infty$  and

$$\begin{aligned} \left| \Omega_{n-K}(\sqrt{n-K}\|\mathbf{t}\| + \epsilon_n) - e^{-\frac{1}{2}\|\mathbf{t}\|^2} \right| &\leq \\ \left| \Omega_{n-K}(\sqrt{n-K}\|\mathbf{t}\| + \epsilon_n) - e^{-\frac{1}{2}[\|\mathbf{t}\| + \epsilon_n/\sqrt{n-K}]^2} \right| &+ \\ \left| e^{-\frac{1}{2}[\|\mathbf{t}\| + \epsilon_n/\sqrt{n-K}]^2} - e^{-\frac{1}{2}\|\mathbf{t}\|^2} \right|. \end{aligned}$$

As  $p, n \rightarrow +\infty$ , the first term vanishes from the uniform convergence given in Theorem 3 and the second term vanishes by continuity. Therefore,  $\varphi_{\sqrt{n}[\hat{\mathbf{v}}_k^\#]_{\mathcal{I}}}(t_{i_1}, \dots, t_{i_{|\mathcal{I}|}}) \rightarrow e^{-\frac{\|\mathbf{t}\|^2}{2}}$  and, by Lévy's continuity theorem [23], we can conclude that  $\sqrt{n}[\hat{\mathbf{v}}_k^\#]_{\mathcal{I}} \xrightarrow{\mathcal{D}} \mathcal{N}(\mathbf{0}, \mathbf{I}_{|\mathcal{I}|})$  as  $p, n \rightarrow +\infty$ . And, finally, given decomposition (2), Theorem 1 and the independence of  $\hat{\mathbf{v}}_k$  and  $\hat{\mathbf{v}}_{k'}$  if  $k' \neq k$  [20, Theorem 4],

$$\hat{\mathbf{v}}_k = \sqrt{\zeta_k} \mathbf{v}_k + \sqrt{1 - \zeta_k} \hat{\mathbf{v}}_k^\# + \varepsilon \quad \text{with} \quad \|\varepsilon\| \xrightarrow[p, n \rightarrow +\infty]{\text{a.s.}} 0.$$

This concludes the proof of Theorem 2.

#### IV. NUMERICAL EXPERIMENTS

To illustrate this result, we conduct a first experiment on synthetic data following model (1) with  $K = 3$  classes of equal size and  $(\|\boldsymbol{\mu}_1\|, \|\boldsymbol{\mu}_2\|, \|\boldsymbol{\mu}_3\|) = (3, 4, 5)$ . The  $\mathbf{x}_i$ 's are ordered by class. The spectral distribution of  $\mathbf{K}$  is plotted in Figure 1 and Figure 2 shows the dominant eigenvectors with the histograms of residuals  $\hat{\mathbf{v}}_k - \sqrt{\zeta_k} \mathbf{v}_k$ . We observe a very good fit of the latter to the probability density function of  $\mathcal{N}(0, \frac{1-\zeta_k}{n})$  — the  $\hat{\mathbf{v}}_k$ 's exactly correspond to a deterministic signal  $\sqrt{\zeta_k} \mathbf{v}_k$  corrupted by additive centered Gaussian noise. The *signal + noise* structure of model (1) has been transferred to the spectral estimator of  $\mathbf{V}$ .

Then, we conduct a second experiment on the Fashion-MNIST dataset [24] consisting of  $28 \times 28$  images of clothes separated in 10 classes of size 7000 each. We select two classes  $k_1, k_2$  and perform binary spectral clustering using the dominant eigenvector of  $\mathbf{K} = \frac{1}{p} \mathbf{X}^\top \mathbf{X}$  where the columns of  $\mathbf{X}$  are the 784 pixels of the images from classes  $k_1$  and  $k_2$ . The dimension of  $\mathbf{X}$  is thus  $784 \times 14000$ . Here, we assume a similar model as our toy example in the introduction:  $\mathbf{X} = \boldsymbol{\mu} \mathbf{j}^\top + \mathbf{W}$  where  $j_i = \pm 1$  depending on the class of the  $i$ -th image. Thus, according to Theorem 2,

		Observed										
Predicted	T-shirt/top		0.64	0.54	0.57	0.61	0.89	0.53	0.85	0.48	0.54	T-shirt/top
	Trouser	0.61		0.72	0.54	0.81	0.95	0.58	0.86	0.69	0.71	Trouser
	Pullover	0.55	0.69		0.63	0.55	0.88	0.57	0.85	0.56	0.61	Pullover
	Dress	0.56	0.54	0.63		0.7	0.9	0.53	0.82	0.59	0.58	Dress
	Coat	0.61	0.76	0.54	0.69		0.93	0.63	0.91	0.64	0.7	Coat
	Sandal	0.9	0.97	0.9	0.92	0.95		0.85	0.76	0.88	0.92	Sandal
	Shirt	0.52	0.58	0.57	0.54	0.62	0.85		0.79	0.53	0.52	Shirt
	Sneaker	0.83	0.86	0.84	0.81	0.91	0.76	0.76		0.8	0.85	Sneaker
	Bag	0.5	0.67	0.56	0.59	0.61	0.89	0.53	0.79		0.51	Bag
	Ankle boot	0.52	0.71	0.59	0.59	0.66	0.93	0.52	0.85	0.51		Ankle boot

Fig. 3. Observed (upper right, blue) and predicted (lower left, orange) classification accuracies of binary spectral clustering on the Fashion-MNIST dataset [24].

the  $i$ -th entry of the dominant eigenvector  $\hat{\mathbf{v}}$  asymptotically follows  $\mathcal{N}(\sqrt{\zeta} \frac{j_i}{\sqrt{n}}, \frac{1-\zeta}{n})$  and, given  $j, \zeta$  can be estimated as  $(\sum_{i=1}^n \frac{j_i}{\sqrt{n}} \hat{v}_i)^2$ . We can then compare the observed accuracy  $\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\hat{v}_i j_i > 0}$  to the one expected from Theorem 2,  $\mathbb{P}(\hat{v}_i j_i > 0) \asymp \Phi(\sqrt{\frac{\zeta}{1-\zeta}})$ . The results are presented in Figure 3 for each pair of classes  $k_1, k_2$ .

We find a very good agreement between the observed and predicted accuracies, regardless of whether the problem is easy (e.g., Trouser vs Sandal) or hard (e.g., Bag vs Ankle Boot). This observation confirms the general scope of Theorem 2: starting from real data  $\mathbf{X}$  which is clearly *not* Gaussian, the normal distribution naturally emerges in the fluctuations of the entries of the large-dimensional eigenvector  $\hat{\mathbf{v}}$ .

#### V. CONCLUSION

After recalling known results on spectral clustering under a general *signal + noise* random matrix model, we have shown that the entries of spiked eigenvectors have Gaussian fluctuations in the large-dimensional regime. This formalizes and clearly states a result which is often implicitly assumed in many problems, without ever being actually proven. The proposed proof relies solely on the rotational invariance of the noise. It is thus very general and can easily be extended to most standard spike models. Numerical experiments have demonstrated the universality of this phenomenon: the Gaussian behavior of the entries of spike eigenvectors can even be observed on real unprocessed data. This allows to accurately predict the classification performance of spectral clustering. An interesting problem for future work is to understand how these results extend to more exotic spike models such as [25].

## REFERENCES

- [1] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 8, pp. 888–905, Aug. 2000, conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence. [Online]. Available: <https://ieeexplore.ieee.org/document/868688>
- [2] C. Brew and S. S. im Walde, "Spectral clustering for german verbs," in *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing (EMNLP 2002)*, 2002, pp. 117–124. [Online]. Available: <https://aclanthology.org/W02-1016.pdf>
- [3] A. Ng, M. Jordan, and Y. Weiss, "On Spectral Clustering: Analysis and an algorithm," in *Advances in Neural Information Processing Systems*, vol. 14. MIT Press, 2002. [Online]. Available: <https://papers.nips.cc/paper/2001/hash/801272ee79cfd7fa5960571fee36b9b-Abstract.html>
- [4] U. von Luxburg, "A tutorial on spectral clustering," *Statistics and Computing*, vol. 17, no. 4, pp. 395–416, Dec. 2007. [Online]. Available: <https://doi.org/10.1007/s11222-007-9033-z>
- [5] C. Ding, X. He, H. Zha, M. Gu, and H. Simon, "A min-max cut algorithm for graph partitioning and data clustering," in *Proceedings 2001 IEEE International Conference on Data Mining*, Nov. 2001, pp. 107–114. [Online]. Available: <https://ieeexplore.ieee.org/document/989507>
- [6] A. Onatski, M. J. Moreira, and M. Hallin, "Asymptotic power of sphericity tests for high-dimensional data," *The Annals of Statistics*, vol. 41, no. 3, pp. 1204–1231, Jun. 2013, publisher: Institute of Mathematical Statistics. [Online]. Available: <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>
- [7] Z. Bai and J. W. Silverstein, *Spectral analysis of large dimensional random matrices*. Springer, 2010, vol. 20.
- [8] L. A. Pastur and M. Shcherbina, *Eigenvalue Distribution of Large Random Matrices*, ser. Mathematical Surveys and Monographs. American Mathematical Society, 2011, no. 171.
- [9] R. Couillet and Z. Liao, *Random Matrix Methods for Machine Learning*. Cambridge: Cambridge University Press, 2022. [Online]. Available: <https://www.cambridge.org/core/books/random-matrix-methods-for-machine-learning/6B681EB69E58B5F888EDB689C160C682>
- [10] J. Baik and J. W. Silverstein, "Eigenvalues of large sample covariance matrices of spiked population models," *Journal of Multivariate Analysis*, vol. 97, no. 6, pp. 1382–1408, 2006. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0047259X0500134X>
- [11] F. Benaych-Georges and R. R. Nadakuditi, "The eigenvalues and eigenvectors of finite, low rank perturbations of large random matrices," *Advances in Mathematics*, vol. 227, no. 1, pp. 494–521, 2011. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0001870811000570>
- [12] R. Couillet and F. Benaych-Georges, "Kernel spectral clustering of large dimensional data," *Electronic Journal of Statistics*, vol. 10, no. 1, pp. 1393–1454, 2016, publisher: Shaker Heights, OH : Institute of Mathematical Statistics. [Online]. Available: <https://hal.archives-ouvertes.fr/hal-01215343>
- [13] R. Couillet, F. Chatelain, and N. L. Bihan, "Two-way kernel matrix puncturing: towards resource-efficient PCA and spectral clustering," in *Proceedings of the 38th International Conference on Machine Learning*. PMLR, Jul. 2021, pp. 2156–2165, iSSN: 2640-3498. [Online]. Available: <https://proceedings.mlr.press/v139/couillet21a.html>
- [14] A. Kadavankandy and R. Couillet, "Asymptotic Gaussian Fluctuations of Spectral Clustering Eigenvectors," in *2019 IEEE 8th International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*, Dec. 2019, pp. 694–698.
- [15] V. A. Marčenko and L. A. Pastur, "Distribution of eigenvalues for some sets of random matrices," *Mathematics of the USSR-Sbornik*, vol. 1, no. 4, p. 457, 1967, publisher: IOP Publishing.
- [16] J. Baik, G. Ben Arous, and S. Péché, "Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices," *The Annals of Probability*, vol. 33, no. 5, pp. 1643–1697, Sep. 2005, publisher: Institute of Mathematical Statistics. [Online]. Available: <https://projecteuclid.org/journals/annals-of-probability/volume-33/issue-5/Phase-transition-of-the-largest-eigenvalue-for-nonnul-complex-sample/10.1214/00911790500000233.full>
- [17] A. Lytova and L. Pastur, "Central limit theorem for linear eigenvalue statistics of random matrices with independent entries," *The Annals of Probability*, vol. 37, no. 5, pp. 1778–1840, 2009, publisher: Institute of Mathematical Statistics.
- [18] N. El Karoui, "Concentration of measure and spectra of random matrices: Applications to correlation matrices, elliptical distributions and beyond," *The Annals of Applied Probability*, vol. 19, no. 6, pp. 2362–2405, Dec. 2009, publisher: Institute of Mathematical Statistics. [Online]. Available: <https://projecteuclid.org/journals/annals-of-applied-probability/volume-19/issue-6/Concentration-of-measure-and-spectra-of-random-matrices--Applications/10.1214/08-AAP548.full>
- [19] C. Louart and R. Couillet, "Concentration of Measure and Large Random Matrices with an application to Sample Covariance Matrices," Jan. 2021, arXiv:1805.08295 [math]. [Online]. Available: <http://arxiv.org/abs/1805.08295>
- [20] R. Couillet and W. Hachem, "Fluctuations of Spiked Random Matrix Models and Failure Diagnosis in Sensor Networks," *IEEE Transactions on Information Theory*, vol. 59, no. 1, pp. 509–525, Jan. 2013, conference Name: IEEE Transactions on Information Theory. [Online]. Available: <https://ieeexplore.ieee.org/document/6320635>
- [21] I. J. Schoenberg, "Metric Spaces and Completely Monotone Functions," *Annals of Mathematics*, vol. 39, no. 4, pp. 811–841, 1938, publisher: Annals of Mathematics. [Online]. Available: <https://www.jstor.org/stable/1968466>
- [22] A. G. M. Steerneman and F. van Perlo-ten Kleij, "Spherical distributions: Schoenberg (1938) revisited," *Expositiones Mathematicae*, vol. 23, no. 3, pp. 281–287, Sep. 2005. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0723086905000034>
- [23] P. Billingsley, *Probability and Measure*, ser. Wiley Series in Probability and Statistics. John Wiley & Sons, Inc., 2012.
- [24] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms," *arXiv:1708.07747 [cs, stat]*, Sep. 2017, arXiv: 1708.07747. [Online]. Available: <http://arxiv.org/abs/1708.07747>
- [25] H. Lebeau, R. Couillet, and F. Chatelain, "A Random Matrix Analysis of Data Stream Clustering: Coping With Limited Memory Resources," in *Proceedings of the 39th International Conference on Machine Learning*. PMLR, 2022, pp. 12 253–12 281, iSSN: 2640-3498. [Online]. Available: <https://proceedings.mlr.press/v162/lebeau22a.html>