



Markov Random Geometric Graph (MRGG): A Growth Model for Temporal Dynamic Networks

Quentin Duchemin, Yohann de Castro

► To cite this version:

Quentin Duchemin, Yohann de Castro. Markov Random Geometric Graph (MRGG): A Growth Model for Temporal Dynamic Networks. *Electronic Journal of Statistics* , 2022, 16 (1), pp.671 – 699. 10.1214/21-EJS1969 . hal-02865542v3

HAL Id: hal-02865542

<https://hal.science/hal-02865542v3>

Submitted on 8 Mar 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.

Markov random geometric graph, MRGG: A growth model for temporal dynamic networks

Quentin Duchemin

LAMA, Univ Gustave Eiffel, CNRS, Marne-la-Vallée, France.

quentin.duchemin@univ-eiffel.fr

&

Yohann De Castro

Institut Camille Jordan, École Centrale de Lyon, Lyon, France

yohann.de-castro@ec-lyon.fr

Abstract

We introduce Markov Random Geometric Graphs (MRGGs), a growth model for temporal dynamic networks. It is based on a Markovian latent space dynamic: consecutive latent points are sampled on the Euclidean Sphere using an unknown Markov kernel; and two nodes are connected with a probability depending on a unknown function of their latent geodesic distance.

More precisely, at each stamp-time k we add a latent point X_k sampled by jumping from the previous one X_{k-1} in a direction chosen uniformly Y_k and with a length r_k drawn from an unknown distribution called the *latitude function*. The connection probabilities between each pair of nodes are equal to the *envelope function* of the distance between these two latent points. We provide theoretical guarantees for the non-parametric estimation of the latitude and the envelope functions.

We propose an efficient algorithm that achieves those non-parametric estimation tasks based on an ad-hoc Hierarchical Agglomerative Clustering approach. As a by product, we show how MRGGs can be used to detect dependence structure in growing graphs and to solve link prediction problems.

1 Introduction

In Random Geometric Graphs (RGG), nodes are sampled independently in latent space \mathbb{R}^d . Two nodes are connected if their distance is smaller than a threshold. A thorough probabilistic study of RGGs can be found in [26]. RGGs have been widely studied recently due to their ability to provide a powerful modeling tool for networks with spatial structure. We can mention applications in bioinformatics [16] or analysis of social media [17]. One main feature is to uncover hidden representation of nodes using latent space and to model interactions by relative positions between latent points.

Furthermore, nodes interactions may evolve with time. In some applications, this evolution is given by the arrival of new nodes as in online collection growth [22], online social network growth [3, 19], or outbreak modeling [31] for instance. The network is growing as more nodes are entering. Other time evolution modelings have been studied, we refer to [28] for a review.

A natural extension of RGG consists in accounting this time evolution. In [12], the expected length of connectivity and dis-connectivity periods of the Dynamic Random Geometric Graph is studied: each node choose at random an angle in $[0, 2\pi)$ and make a constant step size move in that direction. In [29], a random walk model for RGG on the hypercube is studied where at each time step a vertex is either appended or deleted from the graph. Their model falls into the class of Geometric Markovian Random Graphs that are generally defined in [7].

As far as we know, there is no extension of RGG to growth model for temporal dynamic networks. For the first time, we consider a Markovian dynamic on the latent space where the new latent point is drawn with respect to the latest latent point and some Markov kernel to be estimated.

This work was supported by grants from Région Ile-de-France.

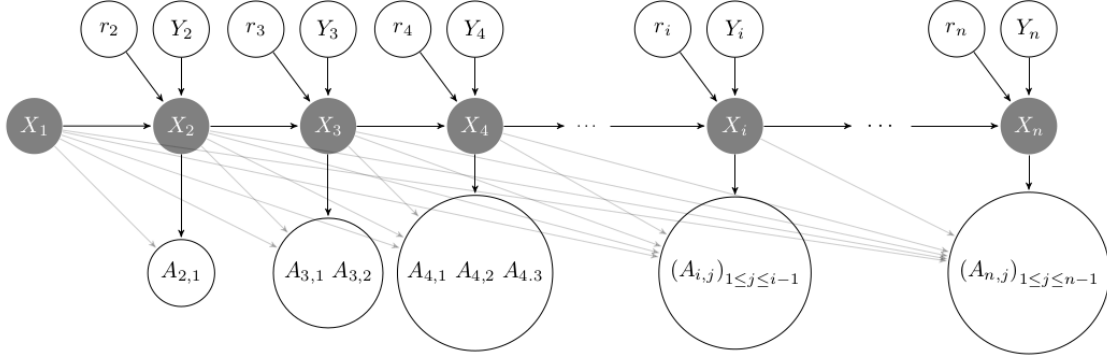


Figure 1: Graphical model of the MRGG model: Markovian dynamics on Euclidean sphere where we jump from X_k onto X_{k+1} . The Y_k encodes direction of jump while r_k encodes its distance, see (1).

Estimation of graphon in RGGs: the Euclidean sphere case Random graphs with latent space can be defined using a *graphon*, cf. [23]. A graphon is a kernel function that defines edge distribution. In [30], Tang and al. prove that spectral method can recover the matrix formed by graphon evaluated at latent points up to an orthogonal transformation, assuming that graphon is a positive definite kernel (PSD). Going further, algorithms have been designed to estimate graphons, as in [20] which provide sharp rates for the Stochastic Block Model (SBM). Recently, the paper [9] provides a non-parametric algorithm to estimate RGGs on Euclidean spheres, without PSD assumption.

We present here RGG on Euclidean sphere. Given n points X_1, X_2, \dots, X_n on the Euclidean sphere \mathbb{S}^{d-1} , we set an edge between nodes i and j (where $i, j \in [n]$, $i \neq j$) with independent probability $\mathbf{p}(\langle X_i, X_j \rangle)$. The unknown function $\mathbf{p} : [-1, 1] \rightarrow [0, 1]$ is called the *envelope function*. This RGG is a graphon model with a symmetric kernel W given by $W(x, y) = \mathbf{p}(\langle x, y \rangle)$. Once the latent points are given, independently draw the random undirected adjacency matrix A by

$$A_{i,j} \sim \text{Ber}(\mathbf{p}(\langle X_i, X_j \rangle)), \quad i < j$$

with Bernoulli r.v. drawn independently (set zero on the diagonal and complete by symmetry), and set

$$T_n := \frac{1}{n} (\mathbf{p}(\langle X_i, X_j \rangle))_{i,j \in [n]} \quad \text{and} \quad \hat{T}_n := \frac{1}{n} A, \quad (1)$$

We do not observe the latent points and we have to estimate the envelope \mathbf{p} from A only. A standard strategy is to remark that \hat{T}_n is a random perturbation of T_n and to dig into T_n to uncover \mathbf{p} .

One important feature of this model is that the interactions between nodes is depicted by a simple object: the envelope function \mathbf{p} . The envelope summarises how individuals connect each others given their latent positions. Standard examples [6] are given by $\mathbf{p}_\tau(t) = \mathbb{1}_{\{t \geq \tau\}}$ where one connects two points as soon as their geodesic distance is below some threshold. The non-parametric estimation of \mathbf{p} is given by [9] where the authors assume that latent points X_i are independently and uniformly distributed on the sphere, which will not be the case in this work.

A new growth model: the latent Markovian dynamic Consider RGGs where latent points are sampled with Markovian jumps, the Graphical Model under consideration can be found in Figure 1. Namely, we sample n points X_1, X_2, \dots, X_n on the Euclidean sphere \mathbb{S}^{d-1} using a Markovian dynamic. We start by sampling randomly X_1 on \mathbb{S}^{d-1} . Then, for any $i \in \{2, \dots, n\}$, we sample

- a unit vector $Y_i \in \mathbb{S}^{d-1}$ uniformly, orthogonal to X_{i-1} .
- a real $r_i \in [-1, 1]$ encoding the distance between X_{i-1} and X_i , see (2). r_i is sampled from a distribution $f_{\mathcal{L}} : [-1, 1] \rightarrow [0, 1]$, called the *latitude function*.

then X_i is defined by

$$X_i = r_i \times X_{i-1} + \sqrt{1 - r_i^2} \times Y_i.$$

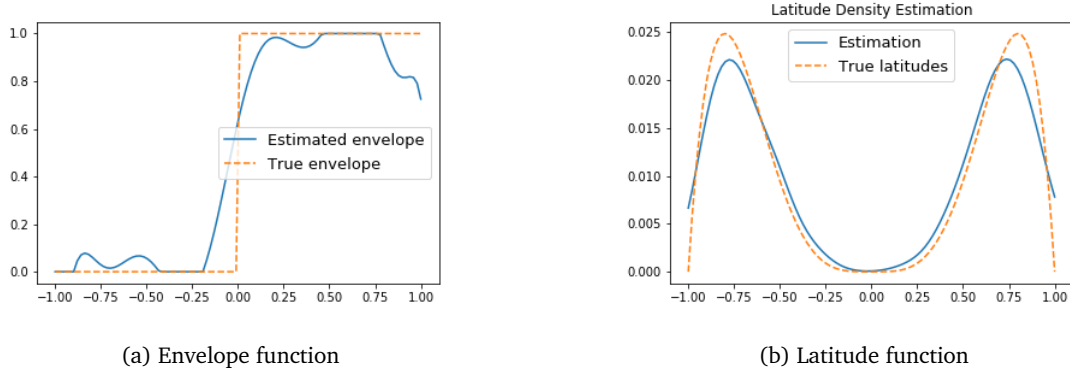


Figure 2: Non-parametric estimation of envelope and latitude functions using algorithms of Sections 2 and 3. We built a graph of 1500 nodes sampled on the sphere \mathbb{S}^2 and using envelope $\mathbf{p}^{(1)}$ and latitude $f_{\mathcal{L}}^{(1)}$ (dot orange curves) defined in Section 5 by Eq.(11). The estimated envelope is thresholded to get a function in $[0, 1]$ and the estimated latitude function is normalized with integral 1 (plain blue lines).

This dynamic can be pictured as follows. Consider that X_{i-1} is the north pole, then chose uniformly a direction (i.e., a longitude) and, in a independent manner, randomly move along the latitudes (the longitude being fixed by the previous step). The geodesic distance γ_i drawn on the latitudes satisfies

$$\gamma_i = \arccos(r_i), \quad (2)$$

where random variable $r_i = \langle X_i, X_{i-1} \rangle$ has density $f_{\mathcal{L}}(r_i)$. The resulting model will be referred to as the Markov Random Geometric Graph (MRGG) and is described with Figure 1.

Temporal Dynamic Networks: MRGG estimation strategy. Seldom growth models exist for temporal dynamic network modeling, see [28] for a review. In our model, we add one node at a time making a Markovian jump from the previous latent position. It results in

the observation of $(A_{i,j})_{1 \leq j \leq i-1}$ at time $T = i$,

as pictured in Figure 1. Namely, we observe how a new node connects to the previous ones. For such dynamic, we aim at estimating the model, namely envelope \mathbf{p} and respectively latitude $f_{\mathcal{L}}$. These functions capture in a simple function on $\Omega = [-1, 1]$ the range of interaction of nodes (represented by \mathbf{p}) and respectively the dynamic of the jumps in latent space (represented by $f_{\mathcal{L}}$), where, in abscissa Ω , values $r = \langle X_i, X_j \rangle$ near 1 corresponds to close point $X_i \simeq X_j$ while values close to -1 corresponds to antipodal points $X_i \simeq -X_j$. These functions may be non-parametric.

From snapshots of the graph at different time steps, can we recover envelope and latitude functions? We prove that it is possible under mild conditions on the Markovian dynamic of the latent points and our approach is summed up with Figure 3.

Define $\lambda(T_n) := (\lambda_1, \dots, \lambda_n)$ and resp. $\lambda(\hat{T}_n) := (\hat{\lambda}_1, \dots, \hat{\lambda}_n)$ the spectrum of T_n and resp. \hat{T}_n , see (1). Building clusters from $\lambda(\hat{T}_n)$, Algorithm 1 (SCCHEi) estimates the spectrum of envelope \mathbf{p} while Algorithm 3 [1] (HEiC, cf. Section F) extracts d eigenvectors of \hat{T}_n to uncover the Gram matrix of the latent positions. Both can then be used to estimate the unknown functions of our model (cf. Figure 2).

Previous works. The latent space approach to model dynamics of network has already been studied in a large span of recent works. Most of them focus on block models with dynamic generalizations covering discrete dynamic evolution via hidden Markov models (cf. [24]) or continuous time analysis via extended Kalman filter (cf. [32]). [33] and [11] use a Gamma Markov process allowing to model evolving mixed membership in graphs using respectively the Bernoulli Poisson link function and the logistic function to generate edges from the latent space representation. While the above mentioned papers consider community based random graphs with fixed size where edges and communities change through time, we

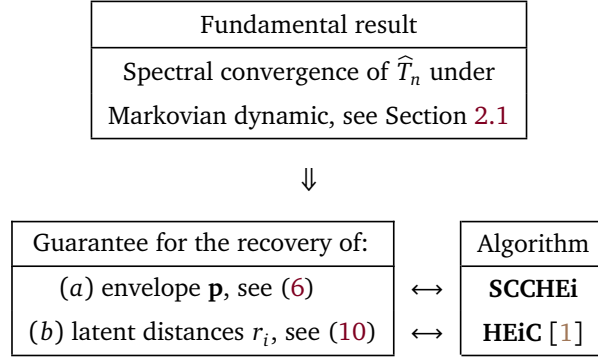


Figure 3: Presentation of our method to recover the envelope and the latitude functions.

focus on growing RGGs on Euclidean sphere where new nodes are added along time.

Non-parametric estimation of RGGs on \mathbb{S}^{d-1} has been investigated in [9] with i.i.d. latent points. Estimation of latent point relative distances with HEiC Algorithm has been introduced in [1] under i.i.d. latent points assumption. Phase transitions on the detection of geometry in RGGs (against Erdős Rényi alternatives) has been investigated in [6].

For the first time, we introduce latitude function and non-parametric estimations of envelope and latitude using new results on kernel matrices concentration with dependent variables.

Outline Sections 2 and 3 present the estimation method with new theoretical results under Markovian dynamic. These new results are random matrices operator norm control and resp. U-statistics control under Markovian dynamic, presented in the Appendix at Section H and resp. Section G. The envelope *adaptive* estimate is built from a size constrained clustering (Algorithm 1) tuned by slope heuristic Eq.(7), and the latitude function estimate (cf. Section 3.1) is derived from estimates of latent distances r_i . Our method can handle random graphs with logarithmic growth node degree (i.e., new comer at time $T = n$ connects to $\mathcal{O}(\log n)$ previous nodes), referred to as *relatively sparse* models, see Section 4. Sections 5 and 6 investigate synthetic data experiments. We propose heuristics to solve link prediction problems and to test for a Markovian dynamic. In a last section (Section 7), we dig deeper into the analysis of our methods by studying their behaviour under model misspecification or under slow mixing conditions. We conclude by presenting final remarks and future research directions. At the end of Section 7, we provide with Figure 14 a synthetic presentation of the estimation methods of this paper.

Notations. Consider a dimension $d \geq 3$. Denote by $\|\cdot\|_2$ (resp. $\langle \cdot, \cdot \rangle$) the Euclidean norm (resp. inner product) on \mathbb{R}^d . Consider the d -dimensional sphere $\mathbb{S}^{d-1} := \{x \in \mathbb{R}^d : \|x\|_2 = 1\}$ and denote by π the uniform probability measure on \mathbb{S}^{d-1} . For any matrix $M = (m_{i,j})_{i,j} \in \mathbb{R}^{D_1 \times D_2}$, we define $\|M\|_F^2 := \sum_{i=1}^{D_1} \sum_{j=1}^{D_2} |m_{i,j}|^2$ and the operator norm of M as $\|M\| := \sup_{x \in \mathbb{S}^{D_2-1}} \|Mx\|_2$. For two real valued sequences $(u_n)_{n \in \mathbb{N}}$ and $(v_n)_{n \in \mathbb{N}}$, denote $u_n \underset{n \rightarrow \infty}{=} \mathcal{O}(v_n)$ if there exist $k_1 > 0$ and $n_0 \in \mathbb{N}$ such that $\forall n > n_0$, $|u_n| \leq k_1 |v_n|$. For any $x, y \in \mathbb{R}$, $x \wedge y := \min(x, y)$ and $x \vee y := \max(x, y)$. Given two sequences x, y of reals—completing finite sequences by zeros—such that $\sum_i x_i^2 + y_i^2 < \infty$, we define the ℓ_2 rearrangement distance $\delta_2(x, y)$ as

$$\delta_2^2(x, y) := \inf_{\sigma \in \mathfrak{S}} \sum_i (x_i - y_{\sigma(i)})^2,$$

where \mathfrak{S} is the set of permutations with finite support. This pseudo-distance is useful to compare two spectra.

2 Nonparametric estimation of the envelope function

One can associate with $W(x, y) = \mathbf{p}(\langle x, y \rangle)$ the integral operator

$$\mathbb{T}_W : L^2(\mathbb{S}^{d-1}) \rightarrow L^2(\mathbb{S}^{d-1}),$$

such that for any $g \in L^2(\mathbb{S}^{d-1})$,

$$\forall x \in \mathbb{S}^{d-1}, \quad (\mathbb{T}_W g)(x) = \int_{\mathbb{S}^{d-1}} g(y) \mathbf{p}(\langle x, y \rangle) \pi(dy),$$

where π is the uniform probability measure on \mathbb{S}^{d-1} . The operator \mathbb{T}_W is Hilbert-Schmidt and it has a countable number of bounded eigenvalues λ_k^* with zero as only accumulation point. The eigenfunctions of \mathbb{T}_W have the remarkable property that they do not depend on \mathbf{p} (cf. [8] Lemma 1.2.3): they are given by the real Spherical Harmonics. We denote \mathcal{H}_l the space of real Spherical Harmonics of degree l with dimension d_l and with orthonormal basis $(Y_{l,j})_{j \in [d_l]}$ where

$$d_l := \dim(\mathcal{H}_l) = \begin{cases} 1 & \text{if } l = 0 \\ d & \text{if } l = 1 \\ \binom{l+d-1}{l} - \binom{l+d-3}{l-2} & \text{otherwise.} \end{cases}$$

We end up with the following spectral decomposition

$$\mathbf{p}(\langle x, y \rangle) = \sum_{l \geq 0} p_l^* \sum_{1 \leq j \leq d_l} Y_{l,j}(x) Y_{l,j}(y) = \sum_{k \geq 0} p_k^* c_k G_k^\beta(\langle x, y \rangle), \quad (3)$$

where $\lambda(\mathbb{T}_W) = \{p_0^*, p_1^*, \dots, p_1^*, \dots, p_l^*, \dots, p_l^*, \dots\}$ meaning that each eigenvalue p_l^* has multiplicity d_l ; and G_k^β is the Gegenbauer polynomial of degree k with parameter $\beta := \frac{d-2}{2}$ and $c_k := \frac{2k+d-2}{d-2}$ (cf. Appendix C). Since \mathbf{p} is bounded, one has $\mathbf{p} \in L^2((-1, 1), w_\beta)$ where the weight function w_β is defined by $w_\beta(t) := (1 - t^2)^{\beta - \frac{1}{2}}$ and

$$L^2((-1, 1), w_\beta) := \{g : [-1, 1] \rightarrow \mathbb{R} \mid \|g\|_2^2 := \int_{-1}^1 |g(t)|^2 w_\beta(t) dt < +\infty\}.$$

\mathbf{p} can be decomposed as $\mathbf{p} \equiv \sum_{k \geq 0} p_k^* c_k G_k^\beta$ and the Gegenbauer polynomials G_k^β are an orthogonal basis of $L^2((-1, 1), w_\beta)$.

We finally introduce for any resolution level $R \in \mathbb{N}$ the truncated graphon W_R which is obtained from W by keeping only the \tilde{R} first eigenvalues, that is

$$\forall x, y \in \mathbb{S}^{d-1}, \quad W_R(x, y) := \sum_{k=0}^R p_k^* \sum_{l=1}^{d_k} Y_{k,l}(x) Y_{k,l}(y).$$

Similarly, we denote for all $t \in [0, 1]$, $\mathbf{p}_R(t) = \sum_{k=0}^R p_k^* c_k G_k^\beta(t)$.

Weighted Sobolev space The space $Z_{w_\beta}^s((-1, 1))$ with regularity $s > 0$ is defined as the set of functions $g = \sum_{k \geq 0} g_k^* c_k G_k^\beta \in L^2((-1, 1), w_\beta)$ such that

$$\|g\|_{Z_{w_\beta}^s((-1, 1))}^* := \left[\sum_{l=0}^{\infty} d_l |g_l^*|^2 (1 + (l(l+2\beta))^s) \right]^{1/2} < \infty.$$

2.1 Integral operator spectrum estimation with dependent variables

One key result is a new control of U -statistics with latent Markov variables (cf. Section G) and it makes use of a Talagrand's concentration inequality for Markov chains. This article follows the hypotheses made on the Markov chain $(X_i)_{i \geq 1}$ by [10]. Namely, we work under the following assumption.

Assumption A The latitude function $f_{\mathcal{L}}$ is such that $\|f_{\mathcal{L}}\|_{\infty} < \infty$ and makes the chain $(X_i)_{i \geq 1}$ uniformly ergodic.

Under **Assumption A**, we prove in Section B that the unique stationary distribution of the Markov chain $(X_i)_{i \geq 1}$ is the uniform probability measure on \mathbb{S}^{d-1} denoted π . Theorem 1 is a theoretical guarantee for a random matrix approximation of the spectrum of integral operator with **dependent** latent variables. Theorem 5 in Section H gives explicitly the constants hidden in the big O below which depend on the absolute spectral gap of the Markov chain $(X_i)_{i \geq 1}$ (cf. Definition 11).

Theorem 1. *We consider that **Assumption A** holds and we assume the envelope \mathbf{p} has regularity $s > 0$. Then, it holds*

$$\mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda(T_n))] = \mathcal{O}\left(\left[\frac{n}{\log^2(n)}\right]^{-\frac{2s}{2s+d-1}}\right).$$

Using this preliminary result and the near optimal error bound for the operator norm of random matrices from [4] we obtain

$$\mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda^{R_{opt}}(\widehat{T}_n))] = \mathcal{O}\left(\left[\frac{n}{\log^2(n)}\right]^{-\frac{2s}{2s+d-1}}\right),$$

with $\lambda^{R_{opt}}(\widehat{T}_n) = (\hat{\lambda}_1, \dots, \hat{\lambda}_{\widehat{R}_{opt}}, 0, 0, \dots)$ and $R_{opt} = \lfloor (n/\log^2(n))^{\frac{1}{2s+d-1}} \rfloor$. $\hat{\lambda}_1, \dots, \hat{\lambda}_n$ are the eigenvalues of \widehat{T}_n sorted in decreasing order of magnitude.

Remark. In Theorem 1 and Theorem 4, note that we recover, up to a log factor, the *minimax rate of non-parametric estimation* of s -regular functions on a space of (Riemannian) dimension $d - 1$. Even with i.i.d. latent variables, it is still an open question to know if this rate is the minimax rate of non-parametric estimation of RGGs.

Eq.(3) shows that one could use an approximation of $(p_k^*)_{k \geq 1}$ to estimate the envelope \mathbf{p} and Theorem 1 states we can recover $(p_k^*)_{k \geq 1}$ up to a permutation. In most cases, the problem of finding such a permutation is NP-hard and we introduce in the next section an efficient algorithm to fix this issue.

2.2 Size Constrained Clustering Algorithm

Note the spectrum of \mathbb{T}_W is given by $(p_l^*)_{l \geq 0}$ where p_l^* has multiplicity d_l . In order to recover envelope \mathbf{p} , we build clusters from eigenvalues of \widehat{T}_n while respecting the dimension d_l of each eigen-space of \mathbb{T}_W . In [9], an algorithm is proposed testing all permutations of $\{0, \dots, R\}$ for a given maximal resolution R . To bypass the high computational cost of such approach, we propose an efficient method based on the tree built from *Hierarchical Agglomerative Clustering* (HAC). In the following, for any $v_1, \dots, v_n \in \mathbb{R}$, we denote by $\text{HAC}(\{v_1, \dots, v_n\}, d_c)$ the tree built by a HAC on the real values v_1, \dots, v_n using the complete linkage function d_c defined by $\forall A, B \subset \mathbb{R}, d_c(A, B) = \max_{a \in A} \max_{b \in B} \|a - b\|_2$. Algorithm 1 describes our approach.

Given some resolution level $R \in \mathbb{N}$, our estimator $\widehat{\mathbf{p}}_R$ of the envelope function \mathbf{p} is obtained from the clustering of the eigenvalues obtained by the SCCHEi algorithm as follows

$$\widehat{\mathbf{p}}_R : t \mapsto \sum_{k=0}^R \widehat{p}_k c_k G_k^\beta(t) \quad \text{where} \quad \forall k \in \{0, \dots, R\}, \quad \widehat{p}_k := \frac{1}{d_k} \sum_{\lambda \in \mathcal{C}_{d_k}} \lambda. \quad (4)$$

Algorithm 1 Size Constrained Clustering for Harmonic Eigenvalues (SCCHEi).

Data: Resolution R , matrix $\widehat{T}_n = \frac{1}{n}A$, dimensions $(d_k)_{k=0}^R$.

- 1: Let $\hat{\lambda}_1, \dots, \hat{\lambda}_n$ be the eigenvalues of \widehat{T}_n sorted in decreasing order of magnitude.
- 2: Set $\mathcal{P} := \{\hat{\lambda}_1, \dots, \hat{\lambda}_{\tilde{R}}\}$ and $\text{dims} = [d_0, d_1, \dots, d_R]$.
- 3: **while** All eigenvalues in \mathcal{P} are not clustered **do**
- 4: $\text{tree} \leftarrow \text{HAC}(\text{nonclustered eigenvalues in } \mathcal{P}, d_c)$
- 5: **for** $d \in \text{dims}$ **do**
- 6: Search for a cluster of size d in tree as close as possible to the root.
- 7: **if** such a cluster \mathcal{C}_d exists **then** $\text{Update}(\text{dims}, \text{tree}, \mathcal{C}_d, d)$.
- 8: **end for**
- 9: **for** $d \in \text{dims}$ **do**
- 10: Search for the group \mathcal{C} in tree with a size larger than d and as close as possible to d .
- 11: **if** such a group exists **then** $\text{Update}(\text{dims}, \text{tree}, \mathcal{C}, d)$ **else** Go to line 3.
- 12: **end for**
- 13: **end while**

Return: $\mathcal{C}_{d_0}, \dots, \mathcal{C}_{d_R}, \{\hat{\lambda}_{\tilde{R}+1}, \dots, \hat{\lambda}_n\}$

Algorithm 2 $\text{Update}(\text{dims}, \text{tree}, \mathcal{C}, d)$.

- 1: Save the subset \mathcal{C}_d consisting of the d eigenvalues in \mathcal{C} with the largest absolute values.
 - 2: Delete from tree all occurrences to eigenvalues in \mathcal{C}_d and delete d from dims .
-

2.3 Theoretical guarantees

Let us recall that for any resolution level $R \geq 0$,

$$\lambda(\mathbb{T}_{W_R}) = (\lambda_1^*, \dots, \lambda_R^*, 0, 0, \dots) \text{ and } \lambda^R(\widehat{T}_n) = (\hat{\lambda}_1, \dots, \hat{\lambda}_{\tilde{R}}, 0, 0, \dots)$$

where $\hat{\lambda}_1, \dots, \hat{\lambda}_n$ are the eigenvalues of \widehat{T}_n sorted in decreasing order of magnitude. We order the eigenvalues $\hat{\lambda}_1, \dots, \hat{\lambda}_{\tilde{R}}$ and in the following we consider that $\lambda^R(\widehat{T}_n)_1 \geq \dots \geq \lambda^R(\widehat{T}_n)_{\tilde{R}}$.

Theorem 2. *Let us consider some resolution level $R \in \mathbb{N}$. We keep the assumptions of Theorem 1. We recall that we consider $\lambda^R(\widehat{T}_n)_1 \geq \dots \geq \lambda^R(\widehat{T}_n)_{\tilde{R}}$.*

Then for n large enough, the clusters $\mathcal{C}_{d_0}, \dots, \mathcal{C}_{d_R}$ obtained from the SCCHEi algorithm satisfy

$$\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\widehat{T}_n)) = \sum_{k=0}^R \sum_{\hat{\lambda} \in \mathcal{C}_{d_k}} (\hat{\lambda} - p_k^*)^2.$$

Proof of Theorem 2. Let us denote

$$\Delta^G = \min_{0 \leq k \neq l \leq R, p_k^* \neq p_l^*} |p_k^* - p_l^*| \wedge \min_{0 \leq k \leq R, p_k^* \neq 0} |p_k^*| > 0.$$

For any $g \in (0, \frac{\Delta^G}{4})$, the proof of Theorem 1 (cf. Section H) ensures that for n large enough it holds

$$\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\widehat{T}_n)) \leq g^2. \quad (5)$$

Let us recall that

$$\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\widehat{T}_n)) = \inf_{\sigma \in \mathfrak{S}} \sum_{i \geq 1} (\lambda(\mathbb{T}_{W_R})_{\sigma(i)} - \lambda^R(\widehat{T}_n)_i)^2.$$

The proof of Theorem 2 relies on the following two Lemmas. The proofs of these Lemmas are postponed to Section D.

Lemma 1. We keep the assumptions of Theorem 2. Then, for n large enough for Eq.(5) to hold, one can choose a permutation σ^* such that

- $\sigma^*({1, \dots, \tilde{R}}) = {1, \dots, \tilde{R}}$.
- $\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\hat{T}_n)) = \sum_{i=1}^{\tilde{R}} (\lambda(\mathbb{T}_{W_R})_{\sigma^*(i)} - \lambda^R(\hat{T}_n)_i)^2$.

Moreover, the function f^* given by

$$\begin{aligned} f^* : {1, \dots, \tilde{R}} &\rightarrow \{p_k^*, 0 \leq k \leq R\} \\ i &\mapsto \lambda(\mathbb{T}_{W_R})_{\sigma^*(i)}, \end{aligned}$$

is non-increasing.

Lemma 2. We keep the assumptions and notations of Lemma 1. A clustering $(\hat{\mathcal{C}}_{d_k})_{0 \leq k \leq R}$ at depth R in the tree of the HAC algorithm applied to $\mathcal{P} := \{\lambda^R(\hat{T}_n)_1, \dots, \lambda^R(\hat{T}_n)_{\tilde{R}}\}$ is said to be of type (\mathcal{S}) if it satisfies:

$$\begin{aligned} \hat{\mathcal{C}}_{d_0} &\subset \{\lambda^R(\hat{T}_n)_i \mid 1 \leq i \leq \tilde{R}, f^*(i) = p_0^*\}, \quad |\hat{\mathcal{C}}_{d_0}| = d_0, \\ \hat{\mathcal{C}}_{d_1} &\subset \{\lambda^R(\hat{T}_n)_i \mid 1 \leq i \leq \tilde{R}, f^*(i) = p_1^*\}, \quad |\hat{\mathcal{C}}_{d_1}| = d_1, \\ &\dots \\ \hat{\mathcal{C}}_{d_R} &\subset \{\lambda^R(\hat{T}_n)_i \mid 1 \leq i \leq \tilde{R}, f^*(i) = p_R^*\}, \quad |\hat{\mathcal{C}}_{d_R}| = d_R. \end{aligned}$$

Then the HAC algorithm with complete linkage applied to \mathcal{P} reaches (after $\tilde{R} - R - 1$ iterations) a state $(\hat{\mathcal{C}}_{d_k})_{0 \leq k \leq R}$ of type (\mathcal{S}) . As a consequence, the SCCHEi algorithm returns the clusters $\mathcal{C}_{d_0} = \hat{\mathcal{C}}_{d_0}, \dots, \mathcal{C}_{d_R} = \hat{\mathcal{C}}_{d_R}$.

Theorem 2 directly follows from the conclusion of Lemma 2 since we get that

$$\begin{aligned} \sum_{k=0}^R \sum_{\hat{\lambda} \in \mathcal{C}_{d_k}} (\hat{\lambda} - p_k^*)^2 &= \sum_{i=1}^{\tilde{R}} (\lambda^R(\hat{T}_n)_i - f^*(i))^2 = \sum_{i=1}^{\tilde{R}} (\lambda^R(\hat{T}_n)_i - \lambda(\mathbb{T}_{W_R})_{\sigma^*(i)})^2 \\ &= \delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\hat{T}_n)), \end{aligned}$$

where the first equality comes from the conclusion of Lemma 2, the second one comes from the definition of f^* from Lemma 1 and the last one comes from the choice of σ^* from Lemma 1. \square

Theorem 2 ensures that under appropriate conditions, the SCCHEi leads to a clustering of the eigenvalues of the adjacency matrix that achieves the δ_2 distance between $\lambda(\mathbb{T}_{W_R})$ and $\lambda^R(\hat{T}_n)$. Nevertheless, this is not a sufficient condition to ensure that the L^2 error between the true envelope function and our plug-in estimator (cf. Eq.(4)) goes to 0 has $n \rightarrow +\infty$. This is due to identifiability issues coming from the δ_2 metric. This was already mentioned in [9, Section 3.6], where the authors present the following example. Consider the case $d = 3$, which implies $\beta = 1/2$, $d_k = 2k + 1$, $c_k = 2k + 1$. For $\mu > 0$, let

$$\begin{aligned} \mathbf{p}_a &= \frac{1}{2} c_0 G_0^\beta + \mu c_1 G_1^\beta + 0 \times c_2 G_2^\beta + 0 \times c_3 G_3^\beta + \mu c_4 G_4^\beta \\ \mathbf{p}_b &= \frac{1}{2} c_0 G_0^\beta + 0 \times c_1 G_1^\beta + \mu c_2 G_2^\beta + \mu c_3 G_3^\beta + 0 \times c_4 G_4^\beta \end{aligned}$$

Then the associated spectrum are

$$\begin{aligned} \lambda_a^* &= (1/2, \underbrace{\mu, \mu, \mu}_3, \underbrace{0, 0, 0}_5, \underbrace{0, 0, 0}_7, \underbrace{0, 0, 0}_9, \mu, \mu, \mu, \mu, \mu, \mu, \mu) \\ \lambda_b^* &= (1/2, \underbrace{0, 0, 0}_3, \underbrace{\mu, \mu, \mu}_5, \underbrace{\mu, \mu, \mu}_7, \underbrace{\mu, \mu, \mu}_9, 0, 0, 0, 0, 0, 0, 0) \end{aligned}$$

which are indistinguishable in δ_2 metric, although $\|\mathbf{p}_a - \mathbf{p}_b\|_2 = \mu\sqrt{24}$.

Nevertheless, we can obtain a theoretical guarantee on the L^2 error between the true envelope function and our plug-in estimate using Theorem 2 if we consider additional conditions on the eigenvalues $(p_k^*)_{k \geq 0}$.

Theorem 3. Assume that the envelope function \mathbf{p} is polynomial of degree $D \in \mathbb{N}$, i.e., $p_k^* = 0$ for any $k > D$ and $p_D^* \neq 0$. Assume also that all nonzeros p_k^* for $k \in \{0, \dots, D\}$ are distinct and that $R \geq D$. Then for n large enough it holds with probability at least $1 - n^{-8}$,

$$\|\widehat{\mathbf{p}}_R - \mathbf{p}\|_2^2 \leq c \frac{\widetilde{R}}{n} \ln(n),$$

where $c > 0$ is a universal numerical constant.

Remarks.

- The question of whether the problem of estimating \mathbf{p} is NP-hard was still completely open. Theorem 3 brings a first partial answer to this question by showing that \mathbf{p} can be estimated in polynomial time in the case where \mathbf{p} is a polynomial with all non-zero eigenvalues distinct.
- The proof of Theorem 3 is strictly analogous to the one of [9, Proposition 9]. In a nutshell, considering that the envelope function \mathbf{p} is a polynomial with all non-zeros eigenvalues p_k^* distinct ensures that (since $R \geq D$)

$$\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\widehat{T}_n)) = \delta_2^2(\lambda(\mathbb{T}_W), \lambda^R(\widehat{T}_n)),$$

which coincides with the L^2 norm of the difference between \mathbf{p} and its estimate

$$\widehat{\mathbf{p}}_{opt,R} := \sum_{k=0}^R \widehat{p}_{opt,k} c_k G_k^\beta \quad \text{with} \quad \widehat{p}_{opt,k} := \frac{1}{d_k} \sum_{i \in (\sigma^*)^{-1}([\widetilde{k}+1, \widetilde{k}+1])} \lambda^R(\widehat{T}_n)_i,$$

where σ^* is a permutation as defined in Lemma 1. Since we proved that for n large enough, the clusters returned by the SCCHEi algorithm correspond to an allocation given by f^* , we deduce that the L^2 norm between \mathbf{p} and our plug-in estimate $\widehat{\mathbf{p}}_R$ is equal to the δ_2 distance between spectra. The result then comes directly using Theorem 1.

2.4 Adaptation: Slope heuristic as model selection of Resolution

A data-driven choice of model size R can be done by *slope heuristic*, see [2] for a nice review. One main idea of slope heuristic is to penalize the empirical risk by $\kappa \text{pen}(\widetilde{R})$ and to calibrate $\kappa > 0$. If the sequence $(\text{pen}(\widetilde{R}))_{\widetilde{R}}$ is equivalent to the sequence of variances of the population risk of empirical risk minimizer (ERM) as model size \widetilde{R} grows, then, penalizing the empirical risk (as done in Eq.(7)), one may ultimately uncover an empirical version of the U -shaped curve of the population risk. Hence, minimizing it, one builds a model size \widehat{R} that balances between bias (*under-fitting* regime) and variance (*over-fitting* regime). First, note that empirical risk is given by the intra-class variance below.

Definition 1. For any output $(\mathcal{C}_{d_0}, \dots, \mathcal{C}_{d_R}, \Lambda)$ of the Algorithm SCCHEi, the thresholded intra-class variance is defined by

$$\mathcal{J}_R := \frac{1}{n} \left[\sum_{k=0}^R \sum_{\lambda \in \mathcal{C}_{d_k}} \left(\lambda - \frac{1}{d_k} \sum_{\lambda' \in \mathcal{C}_{d_k}} \lambda' \right)^2 + \sum_{\lambda \in \Lambda} \lambda^2 \right],$$

and the estimations $(\widehat{p}_k)_{k \geq 0}$ of the eigenvalues $(p_k^*)_{k \geq 0}$ is given by

$$\forall k \in \mathbb{N}, \quad \widehat{p}_k = \begin{cases} \frac{1}{d_k} \sum_{\lambda \in \mathcal{C}_{d_k}} \lambda & \text{if } k \in \{0, \dots, \widehat{R}\} \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

Second, as underlined in the proof of Theorem 1 (see Theorem 5 in Section H), the estimator's variance of our estimator scales linearly in \widetilde{R} .

Hence, we apply Algorithm SCCHEi for R varying from 0 to R_{\max} (with $R_{\max} := \max\{R \geq 0 : \widetilde{R} \leq n\}$) to compute the *thresholded intra-class variance* \mathcal{J}_R (see Definition 1) and given some $\kappa > 0$, we select

$$R(\kappa) \in \arg \min_{R \in \{0, \dots, R_{\max}\}} \left\{ \mathcal{J}_R + \kappa \frac{\widetilde{R}}{n} \right\}. \quad (7)$$

The hyper-parameter κ controlling the bias-variance trade-off is set to $2\kappa_0$ where κ_0 is the value of $\kappa > 0$ leading to the “largest jump” of the function $\kappa \mapsto R(\kappa)$. Once $\hat{R} := R(2\kappa_0)$ has been computed, we approximate the envelope function \mathbf{p} using Eq.(6) (see Eq.(20) for the closed form). We denote this estimator $\hat{\mathbf{p}}$ and with the notations of Eq.(4) it holds $\hat{\mathbf{p}} = \hat{\mathbf{p}}_{\hat{R}}$. In Appendix E, we describe this slope heuristic on a concrete example and our results can be reproduced using the notebook *Experiments*¹ in the Appendix.

3 Nonparametric estimation of the latitude function

3.1 Our approach to estimate the latitude function in a nutshell

In Theorem 4 (see below), we show that we are able to estimate consistently the pairwise distances encoded by the Gram matrix G^* where

$$G^* := \frac{1}{n} (\langle X_i, X_j \rangle)_{i,j \in [n]}.$$

Taking the diagonal just above the main diagonal (referred to as *superdiagonal*) of \hat{G} - an estimate of the matrix G to be specified - we get estimates of the i.i.d. random variables $(\langle X_i, X_{i-1} \rangle)_{2 \leq i \leq n} = (r_i)_{2 \leq i \leq n}$ sampled from $f_{\mathcal{L}}$. Using $(\hat{r}_i)_{2 \leq i \leq n}$ the superdiagonal of $n\hat{G}$, we can build a kernel density estimator of the latitude function $f_{\mathcal{L}}$. In the following, we describe the algorithm used to build our estimator \hat{G} with theoretical guarantees.

3.2 Spectral gap condition and Gram matrix estimation

The Gegenbauer polynomial of degree one is defined by $G_1^\beta(t) = 2\beta t$, $\forall t \in [-1, 1]$. As a consequence, using the *addition theorem* (cf. [8, Lem.1.2.3 and Thm.1.2.6]), the Gram matrix G^* is related to the Gegenbauer polynomial of degree one. More precisely, for any $i, j \in [n]$ it holds

$$G_{i,j}^* = \frac{1}{2\beta n} G_1^\beta(\langle X_i, X_j \rangle) = \frac{1}{nd} \sum_{k=1}^d Y_{1,k}(X_i) Y_{1,k}(X_j). \quad (8)$$

Denoting for all $k \in [d]$ $v_k^* := \frac{1}{\sqrt{n}} (Y_{1,k}(X_1), \dots, Y_{1,k}(X_n)) \in \mathbb{R}^n$, and $V^* = (v_1^*, \dots, v_d^*) \in \mathbb{R}^{n \times d}$, Eq.(8) becomes

$$G^* := \frac{1}{d} V^* (V^*)^\top.$$

We will prove that for n large enough there exists a matrix $\hat{V} \in \mathbb{R}^{n \times d}$ where each column is an eigenvector of \hat{T}_n , such that $\hat{G} := \frac{1}{d} \hat{V} \hat{V}^\top$ approximates G^* well, in the sense that the Frobenius norm $\|G^* - \hat{G}\|_F$ converges to 0. To choose the d eigenvectors of the matrix \hat{T}_n that we will use to build the matrix \hat{V} , we need the following spectral gap condition

$$\Delta^* := \min_{k \in \mathbb{N}, k \neq 1} |p_1^* - p_k^*| > 0. \quad (9)$$

This condition will allow us to apply Davis-Kahan type inequalities.

Now, thanks to Theorem 1, we know that the spectrum of the matrix \hat{T}_n converges towards the spectrum of the integral operator \mathbb{T}_W . Then, based on Eq.(8), one can naturally think that extracting the d eigenvectors of the matrix \hat{T}_n related with the eigenvalues that converge towards p_1^* , we can approximate the Gram matrix G^* of the latent positions. Theorem 4 proves that the latter intuition is true with high probability under the spectral gap condition (9). The algorithm HEiC [1] (cf. Section F for a presentation) aims at identifying the above mentioned d eigenvectors of the matrix \hat{T}_n to build our estimate of the Gram matrix G^* .

¹<https://github.com/quentin-duchemin/Markovian-random-geometric-graph>

Theorem 4. We consider that *Assumption A* holds, we assume $\Delta^* > 0$, and we assume that graphon W has regularity $s > 0$. We denote $\hat{V} \in \mathbb{R}^{n \times d}$ the d eigenvectors of the matrix \hat{T}_n associated with the eigenvalues returned by the algorithm HEiC and we define $\hat{G} := \frac{1}{d} \hat{V} \hat{V}^\top$. Then for n large enough and for some constant $D > 0$, it holds with probability at least $1 - 5/n^2$,

$$\|G^* - \hat{G}\|_F \leq D \left(\frac{n}{\log^2(n)} \right)^{\frac{-s}{2s+d-1}}. \quad (10)$$

Based on Theorem 4, we propose a kernel density approach to estimate the latitude function $f_{\mathcal{L}}$ based on the super-diagonal of the matrix \hat{G} , namely $(\hat{r}_i := n\hat{G}_{i-1,i})_{i \in \{2, \dots, n\}}$. In the following, we denote $\hat{f}_{\mathcal{L}}$ this estimator.

4 Relatively Sparse Regime

Although we deal so far with the so-called *dense* regime (i.e. when the expected number of neighbors of each node scales linearly with n), our results may be generalized to the *relatively sparse* model connecting nodes i and j with probability $W(X_i, X_j) = \zeta_n \mathbf{p}(\langle X_i, X_j \rangle)$ where $\zeta_n \in (0, 1]$ satisfies $\liminf_n \zeta_n n / \log n \geq Z$ for some universal constant $Z > 0$.

In the relatively sparse model, one can show following the proof of Theorem 1 that the resolution should be chosen as $\hat{R} = \left(\frac{n\zeta_n}{1 + \zeta_n \log^2 n} \right)^{\frac{1}{2s+d-1}}$. Specifying that $\lambda^* = (p_0^*, p_1^*, \dots, p_1^*, p_2^*, \dots)$ and $\hat{T}_n = A/n$, Theorem 1 becomes for a graphon with regularity $s > 0$

$$\mathbb{E} \left[\delta_2^2 \left(\lambda^*, \frac{\lambda(\hat{T}_n)}{\zeta_n} \right) \right] = \mathcal{O} \left(\left(\frac{n\zeta_n}{1 + \zeta_n \log^2 n} \right)^{\frac{-2s}{2s+d-1}} \right).$$

Figure 4 illustrates the estimation of the latitude and the envelope functions in some relatively sparse regimes.

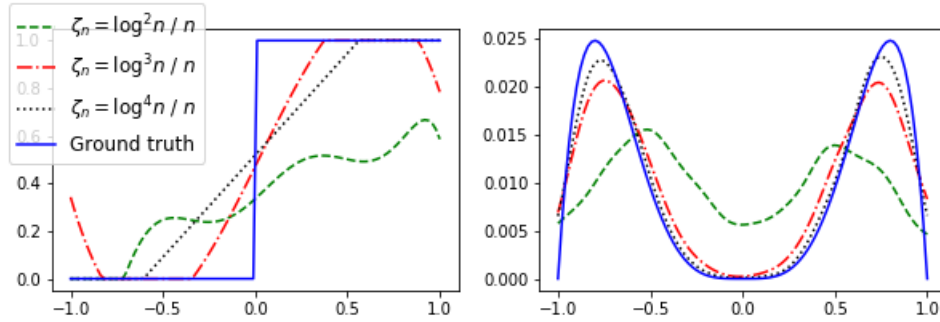


Figure 4: Results of our algorithms for graph of size 2000 with functions $\mathbf{p}^{(1)}$ and $f_{\mathcal{L}}^{(1)}$ of Eq.(11) and sparsity parameter $\zeta_n = \log^k n / n$, $k \in \{2, 3, 4\}$.

5 Experiments

In the following, we test our methods using different envelope and latitude functions. Note that a common choice of connection functions in RGGs are the *Rayleigh fading* activation functions which take the form

$$\mathcal{R}_{\zeta,\eta,r}(\rho) = \exp[-\zeta\rho^\eta], \quad \zeta > 0, \eta > 0.$$

Any Rayleigh function $\mathcal{R}_{\zeta,\eta}$ corresponds to the following envelope function

$$\mathbf{p}_{\zeta,\eta} : t \mapsto \mathcal{R}_{\zeta,\eta}(2(1-t)),$$

so that it holds

$$\forall x, y \in \mathbb{S}^{d-1}, \quad \mathbf{p}_{\zeta,\eta}(\langle x, y \rangle) = \mathcal{R}_{\zeta,\eta}(\|x - y\|_2).$$

Let us also denote for any $\alpha, \beta > 0$ $g(\cdot; \alpha, \beta)$ the density of the beta distribution $\mathcal{B}(\alpha, \beta)$ with parameters (α, β) . In this paper, we will study the numerical results of our methods considering the following envelope and latitude functions

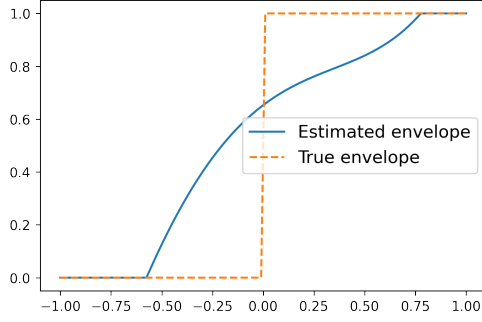
$$\begin{aligned} \mathbf{p}^{(1)} : x &\mapsto \mathbb{1}_{x \geq 0}, & \mathbf{p}^{(2)} &\equiv \mathbf{p}_{0.5,1} \\ f_{\mathcal{L}}^{(1)} : r &\mapsto \begin{cases} \frac{1}{2}g(1-r; 2, 2) & \text{if } r \geq 0 \\ \frac{1}{2}g(1+r; 2, 2) & \text{otherwise} \end{cases}, & f_{\mathcal{L}}^{(2)} : r &\mapsto \frac{1}{2}g\left(\frac{1-r}{2}; 1, 3\right) \\ \text{and } \mathbf{p}^{(3)} &\equiv \mathbf{p}_{0.25,3} \\ f_{\mathcal{L}}^{(3)} : r &\mapsto \frac{1}{2}g\left(\frac{1-r}{2}; 2, 2\right). \end{aligned} \tag{11}$$

Note that considering the latitude function $f_{\mathcal{L}}^{(2)}$ (resp. $f_{\mathcal{L}}^{(3)}$) is equivalent to consider that one fourth of the Euclidean distance between consecutive latent positions is distributed as $Z \sim \mathcal{B}(1, 3)$ (resp. $Z \sim \mathcal{B}(2, 2)$). With Figures 5, 6 and 7, we present the results of our experiments for the three different settings described in Eq.(11). In each case, we work with a latent dimension $d = 4$ and we show:

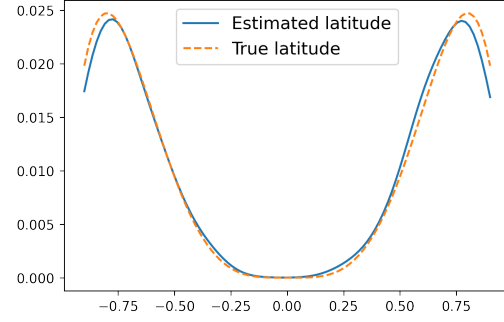
1. the estimates of the envelope and latitude functions obtained with our adaptive procedure working the graph of 1500 nodes (see Figures (a) and (b)).
2. the corresponding clustering obtained by the SCCHEi algorithm for the resolution level R determined by the slope heuristic (see Figures (c)).

Blue crosses represent the \tilde{R} eigenvalues of \hat{T}_n with the largest magnitude, which are used to form clusters corresponding to the $R + 1$ -first spherical harmonic spaces. The red plus are the estimated eigenvalues $(\hat{p}_k)_{0 \leq k \leq R}$ (plotted with multiplicity) defined from the clustering given by our algorithm SCCHEi (see Eq. (6)). Those results show that SCCHEi achieves a relevant clustering of the eigenvalues of \hat{T}_n which allows us to recover the envelope function.

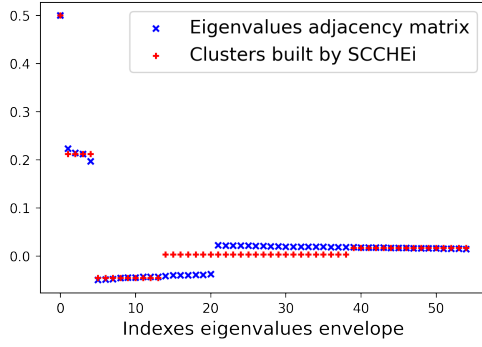
3. the errors between the estimated functions and the true ones in δ_2 metric and in L^2 norm for different size of graphs (see Figures (d) and (e)). We notice that a significant decrease of the δ_2 distance between spectra does not necessarily means that the L^2 norm between the estimated and the true envelope functions shrinks seriously. We refer in particular to Figures 5 and 7. The identifiability issue highlighted in Section 2.3 is one of the possible explanations of this phenomenon. Nevertheless, these experiments show that both the δ_2 and L^2 errors on our estimate of the envelope or the latitude functions are decreasing as the size of the graph is getting larger. Let us also recall that Theorem 3 ensures that the L^2 error on our estimate of the envelope function goes to zero as n grows when \mathbf{p} has a finite number of non zeros eigenvalues that are all distinct.



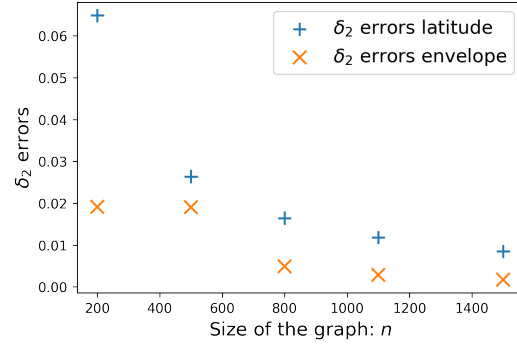
(a) Envelope function



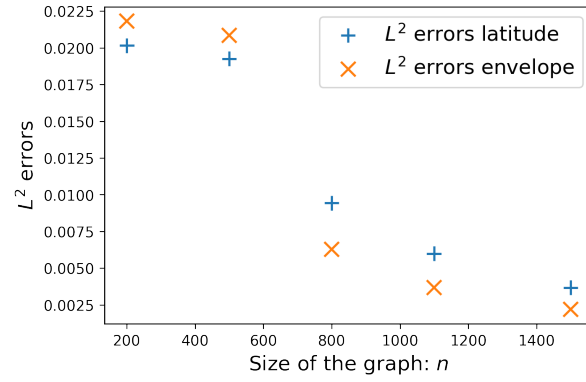
(b) Latitude function



(c) Eigenvalues envelope

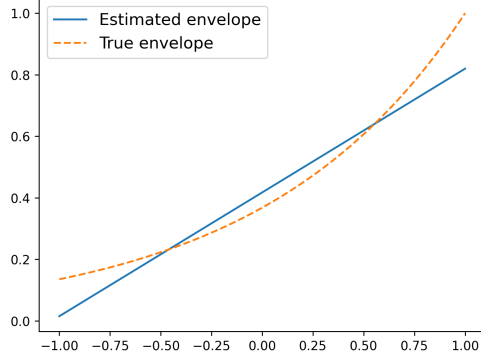


(d) δ_2 errors

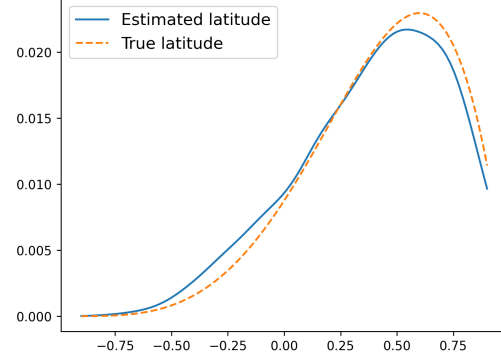


(e) L^2 errors

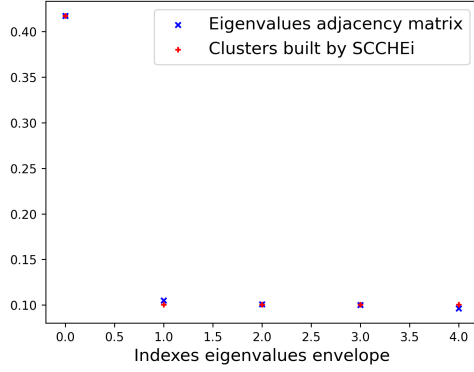
Figure 5: Results for $d = 4$, the envelope $\mathbf{p}^{(1)}$ and the latitude $f_{\mathcal{L}}^{(1)}$ of Eq.(11).



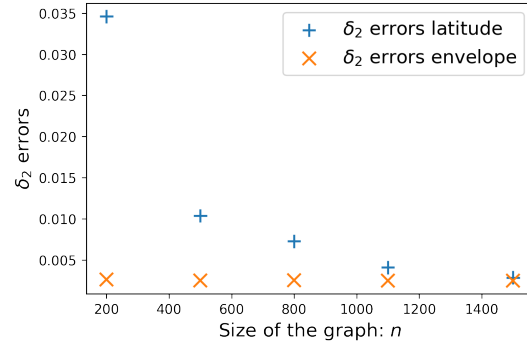
(a) Envelope function



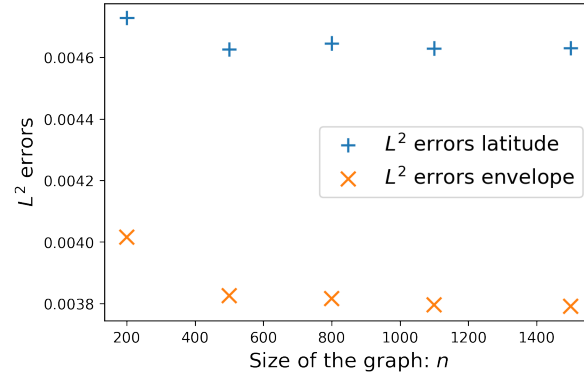
(b) Latitude function



(c) Eigenvalues envelope

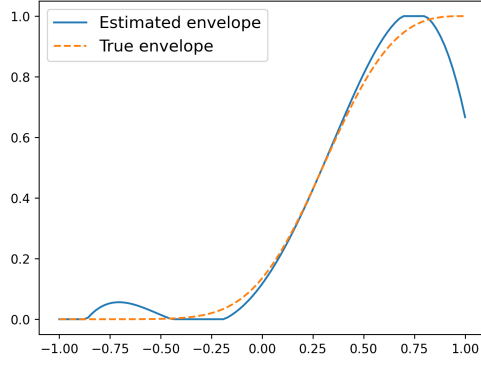


(d) δ_2 errors

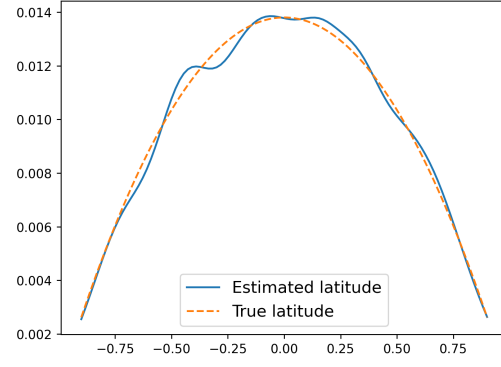


(e) L^2 errors

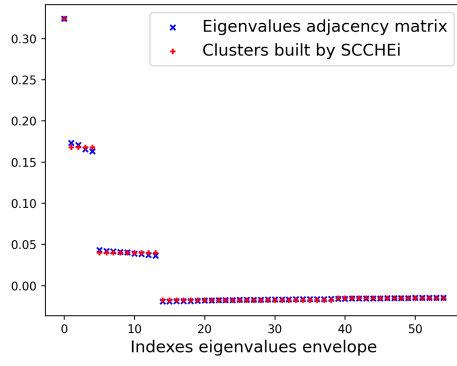
Figure 6: Results for $d = 4$, the envelope $\mathbf{p}^{(2)}$ and the latitude $f_{\mathcal{L}}^{(2)}$ of Eq.(11).



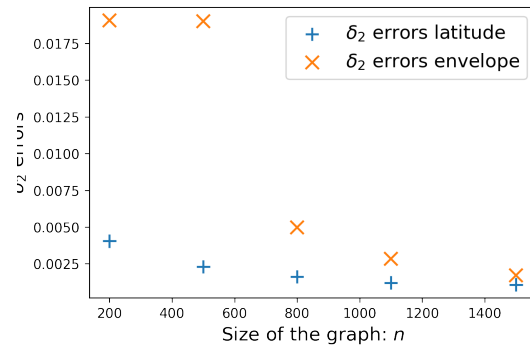
(a) Envelope function



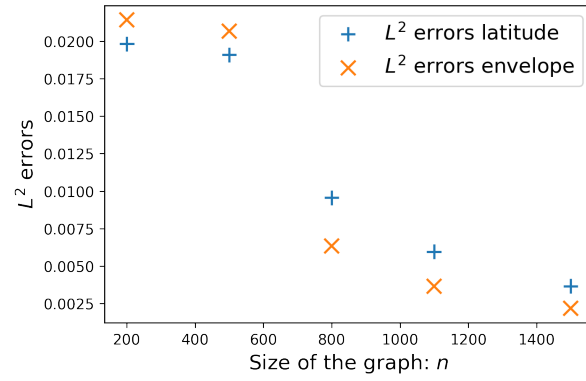
(b) Latitude function



(c) Eigenvalues envelope



(d) δ_2 errors



(e) L^2 errors

Figure 7: Results for $d = 4$, the envelope $\mathbf{p}^{(3)}$ and the latitude $f_{\mathcal{L}}^{(3)}$ of Eq.(11).

6 Applications

In this section, we apply the MRGG model to link prediction and hypothesis testing in order to demonstrate the usefulness of our approach as well as the estimation procedure.

6.1 Markovian Dynamic Testing

As a first application of our model, we propose a hypothesis test to statistically distinguish between an independent sampling the latent positions and a Markovian dynamic. The null is then set to \mathbb{H}_0 : *nodes are independent and uniformly distributed on the sphere* (i.e., *no Markovian dynamic*). Our test is based on estimate $\hat{f}_{\mathcal{L}}$ of latitude and thus the null can be rephrased as \mathbb{H}_0 : $f_{\mathcal{L}} = f_{\mathcal{L}}^0$ where $f_{\mathcal{L}}^0$ is the latitude of uniform law, dynamic is then i.i.d. dynamic.

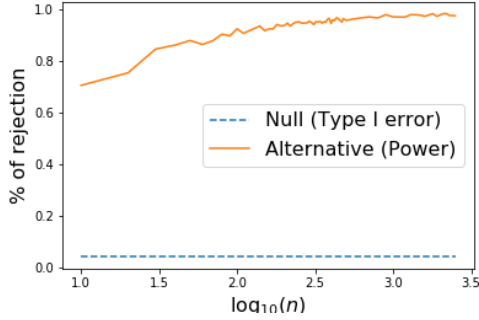


Figure 8: Hypothesis testing.

Figure 8 shows the power of a hypothesis test with level 5% (Type I error). One can use any *black-box goodness-of-fit test* comparing $\hat{f}_{\mathcal{L}}$ to $f_{\mathcal{L}}^0$, and we choose χ^2 -test discretizing $(-1, 1)$ in 70 regular intervals. Rejection region is calibrated (i.e., threshold of the χ^2 -test here) by *Monte Carlo simulations under the null*. It allows us to control Type I error as depicted by dotted blue line. We choose alternative given by Heaviside envelope $\mathbf{p}^{(1)}$ and latitude $f_{\mathcal{L}}^{(1)}$ of Eq.(11). We run our algorithm to estimate latitude from which we sample a batch to compute the χ^2 -test statistic. We see that for graphs of size larger than 1,000, the rejection rate is almost 1 under the alternative (Type II error is almost zero), the test is very powerful.

6.2 Link Prediction

Suppose that we observe a graph with n nodes. Link prediction is the task that consists in estimating the probability of connection between a given node of the graph and the upcoming node.

6.2.1 Bayes Link Prediction

We propose to show the usefulness of our model solving a link prediction problem. Let us recall that we do not estimate the latent positions but only the *pairwise distances* (embedding task is not necessary for our purpose). Denoting by $\text{proj}_{X_n^\perp}(\cdot)$ the orthogonal projection onto the orthogonal complement of $\text{Span}(X_n)$, the decomposition of $\langle X_i, X_{n+1} \rangle$ defined by

$$\begin{aligned} & \langle X_i, X_n \rangle \langle X_n, X_{n+1} \rangle \\ & + \sqrt{1 - \langle X_n, X_{n+1} \rangle^2} \sqrt{1 - \langle X_i, X_n \rangle^2} \left\langle \frac{\text{proj}_{X_n^\perp}(X_i)}{\|\text{proj}_{X_n^\perp}(X_i)\|_2}, Y_{n+1} \right\rangle, \end{aligned} \quad (12)$$

shows that latent distances are enough for link prediction. Indeed, it can be achieved using a *forward step* on our Markovian dynamic, giving the posterior probability (cf. Definition 2) $\eta_i(\mathbf{D}_{1:n})$ defined by

$$\int_{[-1,1]^2} \mathbf{p}(\langle X_i, X_n \rangle r + \sqrt{1-r^2} \sqrt{1-\langle X_i, X_n \rangle^2} u) f_{\mathcal{L}}(r) w_{\frac{d-3}{2}}(u) \frac{\Gamma(\frac{d-1}{2})}{\Gamma(\frac{d-2}{2}) \sqrt{\pi}} dr du, \quad (13)$$

where $w_{\frac{d-3}{2}}(u) := (1-u^2)^{\frac{d-3}{2}-\frac{1}{2}}$ and where $\Gamma : a \in]0, +\infty[\mapsto \int_0^{+\infty} t^{a-1} e^{-t} dt$.

Definition 2. (Posterior probability function)

The posterior probability function η is defined for any latent pairwise distances $\mathbf{D}_{1:n} = (\langle X_i, X_j \rangle)_{1 \leq i, j \leq n} \in [-1, 1]^{n \times n}$ by

$$\forall i \in [n], \quad \eta_i(\mathbf{D}_{1:n}) = \mathbb{P}(A_{i,n+1} = 1 \mid \mathbf{D}_{1:n}),$$

where $A_{i,n+1} \sim \text{Ber}(\mathbf{p}(\langle X_i, X_{n+1} \rangle))$ is a random variable that equals 1 if there is an edge between nodes i and $n+1$, and is zero otherwise.

We consider a classifier g (cf. Definition 3) and an algorithm that, given some latent pairwise distances $\mathbf{D}_{1:n}$, estimates $A_{i,n+1}$ by putting an edge between nodes X_i and X_{n+1} if $g_i(\mathbf{D}_{1:n})$ is 1.

Definition 3. A classifier is a function which associates to any pairwise distances $\mathbf{D}_{1:n} = (\langle X_i, X_j \rangle)_{1 \leq i, j \leq n}$, a label $(g_i(\mathbf{D}_{1:n}))_{i \in [n]} \in \{0, 1\}^n$.

The risk of this algorithm is as in binary classification,

$$\begin{aligned} \mathcal{R}(g, \mathbf{D}_{1:n}) &:= \frac{1}{n} \sum_{i=1}^n \mathbb{P}(g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1} \mid \mathbf{D}_{1:n}) \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ (1 - \eta_i(\mathbf{D}_{1:n})) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=1} + \eta_i(\mathbf{D}_{1:n}) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=0} \right\}, \end{aligned} \quad (14)$$

where we used the independence between $A_{i,n+1}$ and $g_i(\mathbf{D}_{1:n})$ conditionally on $\mathbf{x}(\mathbf{D}_{1:n})$. Pushing further this analogy, we can define the classification error of some classifier g by $L(g) = \mathbb{E}[\mathcal{R}(g, \mathbf{D}_{1:n})]$. Proposition 1 shows that the Bayes estimator - introduced in Definition 4 - is optimal for the risk defined in Eq.(14).

Definition 4. (Bayes estimator)

We keep the notations of Definition 2. The Bayes estimator g^* of $(A_{i,n+1})_{1 \leq i \leq n}$ is defined by

$$\forall i \in [n], \quad g_i^*(\mathbf{D}_{1:n}) = \begin{cases} 1 & \text{if } \eta_i(\mathbf{D}_{1:n}) \geq \frac{1}{2} \\ 0 & \text{otherwise.} \end{cases}$$

Proposition 1. (Optimality of the Bayes classifier for the risk \mathcal{R})

We keep the notations of Definitions 2 and 4. For any classifier g , it holds for all $i \in [n]$,

$$\begin{aligned} &\mathbb{P}(g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1} \mid \mathbf{D}_{1:n}) - \mathbb{P}(g_i^*(\mathbf{D}_{1:n}) \neq A_{i,n+1} \mid \mathbf{D}_{1:n}) \\ &= 2 \left| \eta_i(\mathbf{D}_{1:n}) - \frac{1}{2} \right| \times \mathbb{E} \left\{ \mathbb{1}_{g_i(\mathbf{D}_{1:n}) \neq g_i^*(\mathbf{D}_{1:n})} \mid \mathbf{D}_{1:n} \right\}, \end{aligned}$$

which immediately implies that

$$\mathcal{R}(g, \mathbf{D}_{1:n}) \geq \mathcal{R}(g^*, \mathbf{D}_{1:n}) \text{ and therefore } L(g) \geq L(g^*).$$

6.2.2 Heuristic for Link Prediction

One natural method to approximate the Bayes classifier from the previous section is to use the *plug-in approach*. This leads to the MRGG classifier introduced in Definition 5.

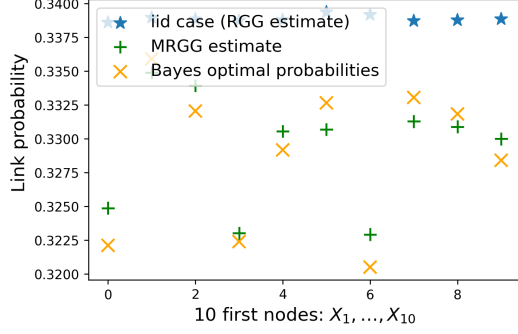
Definition 5. (The MRGG classifier)

For any n and any $i \in [n]$, we define $\hat{\eta}_i(\mathbf{D}_{1:n})$ as

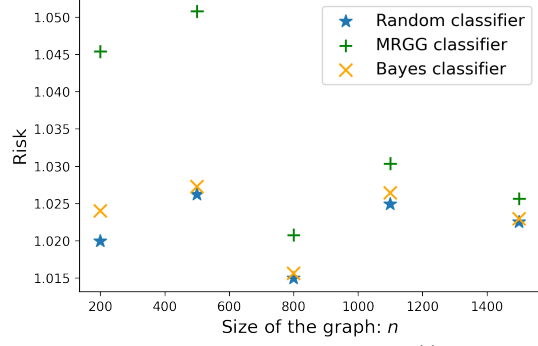
$$\int \hat{\mathbf{p}}(\hat{r}_{i,n}r + \sqrt{1-r^2}\sqrt{1-\hat{r}_{i,n}^2}u) \hat{f}_{\mathcal{L}}(r) w_{\frac{d-3}{2}}(u) \frac{\Gamma(\frac{d-1}{2})}{\Gamma(\frac{d-2}{2})\sqrt{\pi}} dr du, \quad (15)$$

where $\hat{\mathbf{p}}$ and $\hat{f}_{\mathcal{L}}$ denote respectively the estimate of the envelope function and the latitude function with our method and where $\hat{r} := n\hat{G}$. The MRGG classifier is defined by

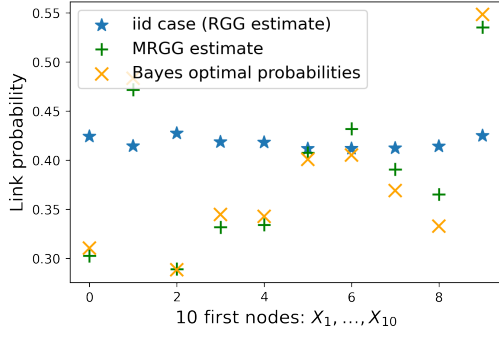
$$\forall i \in [n], \quad g_i^{\text{MRGG}}(\mathbf{D}_{1:n}) = \begin{cases} 1 & \text{if } \hat{\eta}_i(\mathbf{D}_{1:n}) \geq \frac{1}{2} \\ 0 & \text{otherwise.} \end{cases}$$



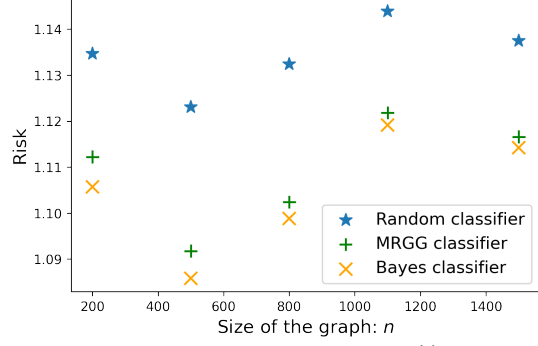
(a) Envelope $\mathbf{p}^{(1)}$, Latitude $f_{\mathcal{L}}^{(1)}$



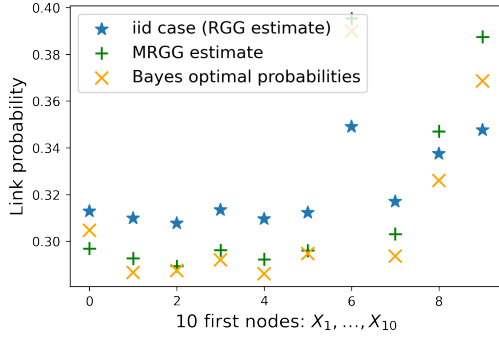
(b) Envelope $\mathbf{p}^{(1)}$, Latitude $f_{\mathcal{L}}^{(1)}$



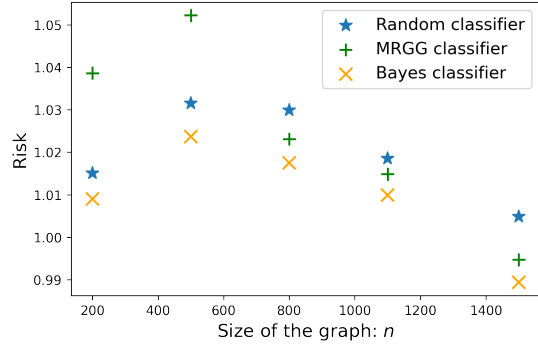
(c) Envelope $\mathbf{p}^{(2)}$, Latitude $f_{\mathcal{L}}^{(2)}$



(d) Envelope $\mathbf{p}^{(2)}$, Latitude $f_{\mathcal{L}}^{(2)}$



(e) Envelope $\mathbf{p}^{(3)}$, Latitude $f_{\mathcal{L}}^{(3)}$



(f) Envelope $\mathbf{p}^{(3)}$, Latitude $f_{\mathcal{L}}^{(3)}$

Figure 9: ← **On the left:** Link predictions between the future node X_{n+1} and the 10 first nodes X_1, \dots, X_{10} . → **On the right:** Comparison between the risk (defined in Eq.(14)) of the MRGG classifier, the random classifier and the risk of the optimal Bayes classifier.

To illustrate our approach we work with a graph of 1500 nodes with $d = 4$, and we consider the envelope and latitude functions defined in Eq.(11). The plots on the left column of Figure 9 show that we are able to recover the probabilities of connection of the nodes already present in the graph with the coming node X_{n+1} . Using the decomposition of $\langle X_i, X_{n+1} \rangle$ given by Eq.(12), orange crosses are computed using Eq.(13). Green plus are computed similarly replacing \mathbf{p} and $f_{\mathcal{L}}$ by their estimations $\hat{\mathbf{p}}$ and $\hat{f}_{\mathcal{L}}$ following Eq.(15). Blue stars are computed using Eq.(13) by replacing $f_{\mathcal{L}}$ by $\frac{w_{\beta}}{\|w_{\beta}\|_1}$ (with $\beta = \frac{d-2}{2}$) which implicitly supposes that the points are sampled uniformly on the sphere.

With the plots on the left column of Figure 9, we compare the risk of the *random* classifier - whose guess $g_i(\mathbf{D}_{1:n})$ is a Bernoulli random variable with parameter given by the ratio of edges compared to complete graph - with the risk of the MRGG classifier (cf. Definition 5). These figures show that for a small number of nodes, the risk estimate provided by the MRGG classifier can be significantly far from the one of the Bayes classifier. However, when the number of nodes is getting larger, the MRGG classifier gives similar results compared to the optimal Bayes classifier. This risk estimate can be significantly smaller than the one of the random classifier (see for example the plots corresponding to the envelope $\mathbf{p}^{(2)}$ and the latitude $f_{\mathcal{L}}^{(2)}$).

7 Discussion

In this section, we want to push the investigation of the performance of our estimation methods as far as possible. In Section 7.1 we study the robustness of our methods under model misspecification before inspecting the influence of the mixing time of the Markov chain $(X_i)_{i \geq 1}$ on the estimation error in Section 7.2.

On a more theoretical side, we show that replacing the use of the complete linkage by the Ward distance in the SCCHEi algorithm, Theorem 2 might not be true anymore. We conclude with some remarks and by highlighting future research directions.

7.1 Robustness to model misspecification

We consider a mixture model for the sampling scheme of the latent position. We fix some $\varepsilon \in (0, 1)$ and we draw X_1 randomly on the sphere. Then at time step $i \geq 2$, the point X_i is sampled as follows:

- with probability $1 - \varepsilon$, X_i is drawn following the Markovian dynamic described in Section 1 (based on X_{i-1}).
- with probability ε , X_i is drawn uniformly on the sphere.

Figure 10 and Figure 11 show the numerical results obtained under this misspecified model. We consider the hypothesis testing question presented in Section 6.1 with the same settings namely $d = 3$ and the envelope and latitude functions $\mathbf{p}^{(1)}$ and $f_{\mathcal{L}}^{(1)}$ of Eq.(11). We can see that when $\varepsilon = 0$, the power of our test is 1 and we always reject the null hypothesis (uniform sampling of the latent positions) under the alternative. On the contrary, when $\varepsilon = 1$, the points are sampled uniformly on the sphere and we obtain a power of the order of the level of our test (i.e. 5%) as expected. The larger the sample size n is, the greater ε can be chosen while keeping a large power. In the case where $n = 1500$, one can afford to sample 75% of latent positions uniformly (and the rest using our Markovian sampling scheme) while keeping a power equal to 1. Figure 11 shows that the larger ε is, the closer the estimated latitude function is to $\frac{w_{\beta}}{\|w_{\beta}\|_1} \equiv \frac{1}{2}$ (since $d = 3$) which corresponds to the density of a one-dimensional marginal of a uniform random point on \mathbb{S}^{d-1} .

7.2 Influence of mixing time on estimation error

In order to assert that the dependence of the latent variables has an influence on the estimation of the unknown functions of our model, we would require a minimax bound. The derivation of such minimax result is still an open problem, even in the independent setting (cf. [9]). Nevertheless, by making explicit the constants involved in concentration inequalities, we can show that the mixing time of the latent

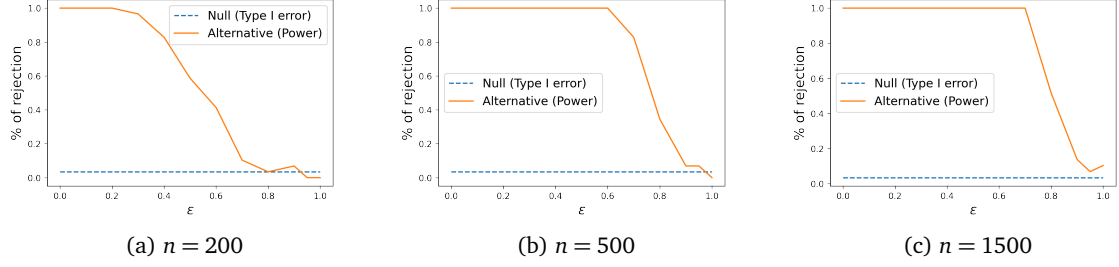


Figure 10: Studying the robustness of our method under model misspecification. We study the evolution of power for Markovian Dynamic Testing when the mixture parameter ε ranges $(0, 1)$. We conduct this analysis for different values of n .

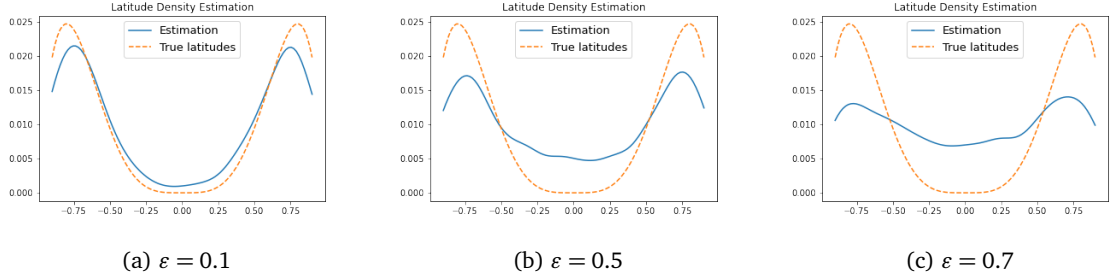


Figure 11: Studying the robustness of the estimation of the latitude function under model misspecification. We plot our kernel density estimator of the latitude function for $n = 1500$, $d = 3$ and for $\varepsilon \in \{0.1, 0.5, 0.7\}$. We use the envelope $\mathbf{p}^{(1)}$ and latitude function $f_{\mathcal{L}}^{(1)}$ defined in Eq.(11).

Markovian dynamic affects our bound on the δ_2 error between spectra. For any $r^* \in (-1, 1)$, let us consider the following latitude function

$$f_{\mathcal{L}}^{r^*}(r) := \frac{1}{I(r^*)}(1-r^2)^{\frac{d-3}{2}} \mathbf{1}_{r \in (r^*, 1)}, \quad I(r^*) := \int_{r^*}^1 (1-r^2)^{\frac{d-3}{2}} dr.$$

Note that the Markov transition kernel P of the chain $(X_i)_{i \geq 1}$ using this latitude function is the one that starting from a point $x \in \mathbb{S}^{d-1}$ samples uniformly a point in the set $\{z \in \mathbb{S}^{d-1} \mid \|x - z\|_2^2 \leq 2(1-r^*)\}$. In particular, when $r^* = -1$, we recover the uniform distribution on the sphere. It is clear that the closer r^* to one, the larger the mixing time of the chain. One can show that for any $r^* \in (-1, 1)$, the chain is uniformly ergodic by proving that there exist an integer $m \geq 1$, a constant $\delta_m > 0$ and a probability measure ν such that

$$\forall x \in \mathbb{S}^{d-1}, \forall A \in \Sigma, \quad P^m(x, A) \geq \delta_m \nu(A) \quad (\text{cf. Definition 9}). \quad (16)$$

Eq.(16) holds by considering for example $\nu = \pi$ the uniform distribution on the sphere. It is straightforward to show that the smallest integer $m(r^*) \geq 1$ satisfying Eq.(16) is larger than $\frac{2}{1-r^*}$.² Taking a closer look at the constants involved in the concentration inequality from [10] (cf. [10, Section 3.1.1]), we get that

$$\mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda(T_n)) \vee \delta_2^2(\lambda(\mathbb{T}_W), \lambda^{R_{opt}}(\hat{T}_n))] < C \left[\frac{n}{\log^2(n)} \right]^{-\frac{2s}{2s+d-1}},$$

where $C > m(r^*)^2 \tau(r^*)^2 \|f_{\mathcal{L}}^{r^*}\|_{\infty}$ and $\tau(r^*) \geq 1$ is the Orlicz norm of some regeneration time. Since for any $0 < r^* < 1$,

$$I(r^*) = \int_{r^*}^1 (1-r^2)^{\frac{d-3}{2}} dr = \int_0^{1-r^*} e^{\frac{d-3}{2} \ln(1-(r+r^*)^2)} dr$$

²Indeed, the latitude function $f_{\mathcal{L}}^{r^*}$ allows to make a jump at each time step of size at most $1-r^*$. Since the length of the shortest arc on \mathbb{S}^{d-1} joining the north pole to the south pole is 2, the result follows.

$$\begin{aligned}
&= (1 - (r^*)^2)^{\frac{d-3}{2}} \int_0^{1-r^*} e^{\frac{d-3}{2} \{\ln(1-(r+r^*)^2) - \ln(1-(r^*)^2)\}} dr \\
&\leq (1 - (r^*)^2)^{\frac{d-3}{2}} \int_0^{1-r^*} e^{-\frac{d-3}{2} \left\{ \frac{2rr^*+r^2}{1-(r^*)^2} \right\}} dr \\
&\leq (1 - (r^*)^2)^{\frac{d-3}{2}} \int_0^{1-r^*} e^{-\frac{d-3}{2} \{2rr^*+r^2\}} dr \\
&\leq (1 - (r^*)^2)^{\frac{d-3}{2}} \int_0^1 e^{-\frac{d-3}{2} \{2rr^*\}} dr \\
&\leq (1 - (r^*)^2)^{\frac{d-3}{2}} \left(1 \wedge \frac{1}{r^*(d-3)} \right),
\end{aligned}$$

we get that $\|f_{\mathcal{L}}^{r^*}\|_{\infty} \geq \frac{1}{l(r^*)} (1 - (r^*)^2)^{\frac{d-3}{2}} \geq r^*(d-3)$. Finally we obtain

$$C > \frac{2r^*}{1-r^*} (d-3),$$

where $r^* \mapsto \frac{2r^*}{1-r^*} (d-3)$ is increasing in r^* and diverges to $+\infty$ when $r^* \rightarrow 1^-$. Hence, the closer r^* is to one, the slower the chain is mixing, and the poorer is our bound.

Figure 12 presents the result of the simulations using the latitude function $f_{\mathcal{L}}^{r^*}$ and the envelope function $\mathbf{p} : t \mapsto \mathbb{1}_{t \geq 0}$. We compute the L^2 error between the true and the estimated envelope functions (respectively the true and the estimated latitude functions). When r^* is getting closer to 1, the chain is mixing slowly and we need to increase the sample size if we want to prevent the L^2 errors from blowing up. Graphs have been generated with a latent dimension $d = 3$ and by sampling the latent positions using our isotropic sampling procedure with latitude function $f_{\mathcal{L}}^{r^*}$.

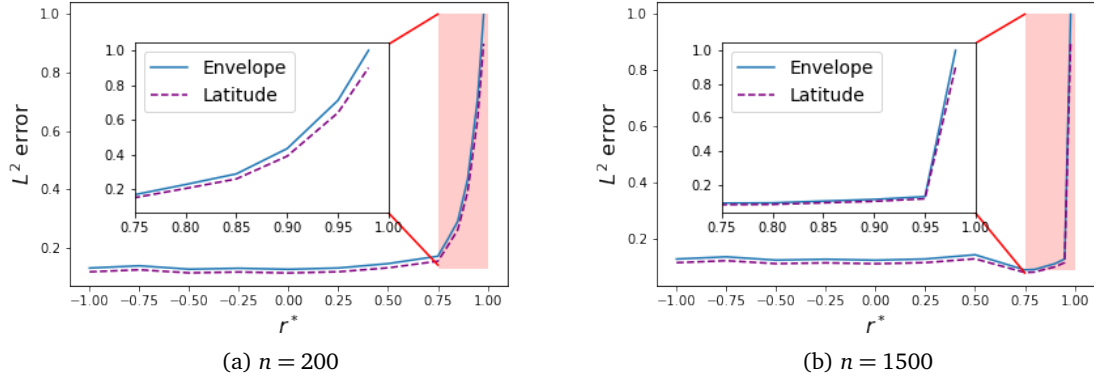


Figure 12: Studying the influence of the mixing time of the chain on the L^2 errors between (i) the envelope function and its estimate (using our adaptive procedure), and (ii) the latitude function and its estimate obtained with a kernel estimator.

7.3 Choice of the clustering algorithm for the SCCHEi

The SCCHEi algorithm relies on the clustering of the eigenvalues of the adjacency matrix provided by the HAC with complete linkage. In this section, we motivate the use of the HAC algorithm with complete linkage by showing that the theoretical results from Section 2.3 could be much more involved to establish by using another clustering procedure. Indeed, if one would consider for example the HAC with the Ward distance, the theoretical result obtained for the correctness of the SCCHEi algorithm (cf. Theorems 2 and 3) is likely to be no longer true (even if the sample size n is chosen arbitrarily large). Let us show this on a simple example.

We fix a resolution level $R = 2$ and we consider some $\Delta^G > 0$. We set $p_0^* = 4\Delta^G$, $p_1^* = 3\Delta^G$, $p_2^* = 2\Delta^G$, and $p_k^* = 0$ for all $k \geq 3$. Let us consider some $g \in (0, \Delta^G/4)$ that can be taken arbitrarily small. Let us denote $\lambda^R(\hat{T}_n) = (\hat{\lambda}_1, \dots, \hat{\lambda}_{\tilde{R}}, 0, 0, \dots)$ and assume that it holds $\hat{\lambda}_1 = p_0^*$, $\hat{\lambda}_2 = \dots = \hat{\lambda}_{d+1} = p_1^*$ (we recall that $d_1 = d$), $\hat{\lambda}_{d+2} = \dots = \hat{\lambda}_{d+1+\lfloor d_2/2 \rfloor} = p_2^* + g$ and $\hat{\lambda}_{d+2+\lfloor d_2/2 \rfloor} = \dots = \hat{\lambda}_{1+d+d_2} = p_2^* - g$. To simplify the presentation, we will assume in the following that $d_2 = \frac{(d+1)d}{2} - 1$ is even (which holds for example if $d = 2k$ for any $k \geq 1$ odd). Figure 13 gives a visualization of this example.

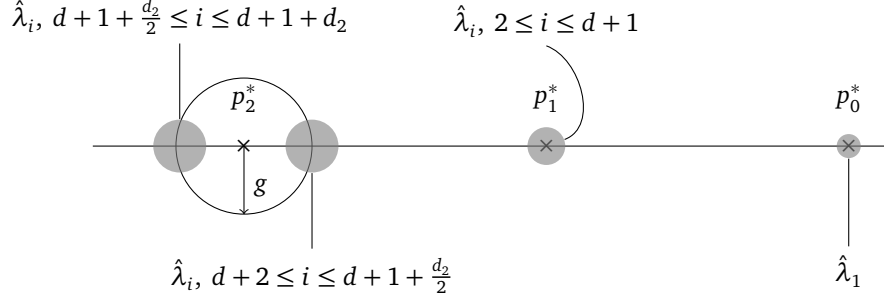


Figure 13: Visualization of the eigenvalues of the envelope function of our example.

Applying the HAC algorithm (with the Ward distance) to the eigenvalues $(\hat{\lambda}_1, \dots, \hat{\lambda}_{\tilde{R}})$, it is obvious that the state reached after $\tilde{R} - 4 = 1 + d + d_2 - 4$ iterations in the HAC procedure will be

$$\begin{aligned}\widehat{\mathcal{G}}_0 &:= \{\hat{\lambda}_1\} \\ \widehat{\mathcal{G}}_1 &:= \{\hat{\lambda}_i \mid 2 \leq i \leq d\} \\ \widehat{\mathcal{G}}_2 &:= \{\hat{\lambda}_i \mid d+2 \leq i \leq d+1+d_2/2\} \\ \widehat{\mathcal{G}}_3 &:= \{\hat{\lambda}_i \mid d+2+d_2/2 \leq i \leq 1+d+d_2\}\end{aligned}$$

Hence, in order to understand which clusters will be merged at the next step of the HAC algorithm, we compute the Ward distance between the different clusters.

Let us recall that for two finite and non-empty sets $S, S' \subset \mathbb{R}$ with respective cardinality $|S|$ and $|S'|$, the Ward distance between S and S' is given by

$$d_W(S, S') := \frac{|S| \times |S'|}{|S| + |S'|} \left(\frac{1}{|S|} \sum_{x_s \in S} x_s - \frac{1}{|S'|} \sum_{x'_s \in S'} x'_s \right)^2.$$

Ward distances between clusters

	$\widehat{\mathcal{G}}_1$	$\widehat{\mathcal{G}}_2$	$\widehat{\mathcal{G}}_3$
$\widehat{\mathcal{G}}_0$	$\frac{d}{d+1}(\Delta^G)^2$	$\frac{d_2}{d_2+2}(2\Delta^G - g)^2$	$\frac{d_2}{d_2+2}(2\Delta^G + g)^2$
$\widehat{\mathcal{G}}_1$		$\frac{d \times d_2}{2d+d_2}(\Delta^G - g)^2$	$\frac{d \times d_2}{2d+d_2}(\Delta^G + g)^2$
$\widehat{\mathcal{G}}_2$			$d_2 \times g^2$

We deduce that all Ward distances between pair of clusters are scaling at least linearly with d except the Ward distances between $\widehat{\mathcal{G}}_0$ and the other three clusters $\widehat{\mathcal{G}}_1$, $\widehat{\mathcal{G}}_2$ and $\widehat{\mathcal{G}}_3$. Indeed, for any $i \in \{1, 2, 3\}$, $d_W(\widehat{\mathcal{G}}_0, \widehat{\mathcal{G}}_i)$ remains bounded independently of the latent dimension d . Hence, for any $g \in (0, \Delta^G/4)$ which can be chosen arbitrarily small, one can take d large enough to ensure that

$$\max \{d_W(\widehat{\mathcal{G}}_0, \widehat{\mathcal{G}}_i), i \in \{1, 2, 3\}\} < d_W(\widehat{\mathcal{G}}_2, \widehat{\mathcal{G}}_3). \quad (17)$$

We deduce that for any $g \in (0, \Delta^G/4)$, we can choose d large enough to ensure that Eq.(17) holds and thus the clusters merged between depths 4 and 3 from the root of the HAC's tree will not be $\widehat{\mathcal{G}}_2$ and $\widehat{\mathcal{G}}_3$.

This means that the state obtained at depth 3 from the root is not of type (\mathcal{S}) (in the sense defined in Lemma 2).

If this is not a sufficient condition to state that the SCCHEi will fail to recover the correct clusters, this example shows that the use of Ward distance can lead to some unexpected clustering of the eigenvalues. Our example proves that using the HAC algorithm with the Ward distance, the result of Lemma 2 does not hold anymore. Namely, regardless of how large the sample size is chosen, there are situations (in particular for a large latent dimension) where the states of type (\mathcal{S}) (cf. Lemma 2) are never reached in the HAC tree with the Ward distance. Hence obtaining a theoretical guarantee for the clustering provided by the SCCHEi in this framework may be impossible or at least much more involved.

7.4 Concluding remarks

7.4.1 Estimation of the latent dimension

The proposed methods implicitly assume that the latent dimension d is known. [1] proved that the latent dimension d can be easily recovered in practice for n large enough provided that the spectral gap condition (9) holds. In the following, we briefly describe their approach.

Given some matrix \hat{T}_n as input and some set of candidates \mathcal{D} for the dimension d (typically $\mathcal{D} = \{2, 3, \dots, d_{\max}\}$), apply the Algorithm HEiC (cf. Algorithm 3 in Section F) for any $d_c \in \mathcal{D}$ and store the returned value $\text{gap} := \text{gap}(d_c)$. Let us recall that $\text{gap}(d_c)$ corresponds to the largest gap between a bulk of d_c eigenvalues of \hat{T}_n and the rest of the spectrum (see the definition of Gap_1 in Section F for details). Once we have computed the different gaps, we pick the candidate d_c that led to the largest one. Given the guarantees provided by Proposition 4, the previously described procedure will find the correct dimension, with high probability (on the event \mathcal{E} with the notations of Proposition 4), if the true dimension of the latent space is in the candidate set \mathcal{D} .

7.4.2 Future research directions

Our work encourages the development of growth model in random graphs and in particular the derivation of similar results in MRGGs with other latent spaces. It would be also desirable to extend our methods to the case where we consider more complex Markovian sampling of the latent positions, typically one that is not isotropic. Our work leaves open the question of getting a theoretical guarantee for the estimation of the latitude function. If we proved (with Theorem 4) that we can consistently estimate the Gram matrix of the latent positions in Frobenius norm, this is not sufficient to ensure that our kernel density estimator is consistent since we cannot ensure that $\frac{1}{n-1} \sqrt{\sum_{i=2}^n (r_i - \hat{r}_i)^2}$ tends to 0 as n goes to $+\infty$. Deriving a theoretical result regarding the estimation of the latitude function seems challenging and we believe that it would require significantly different proof techniques.

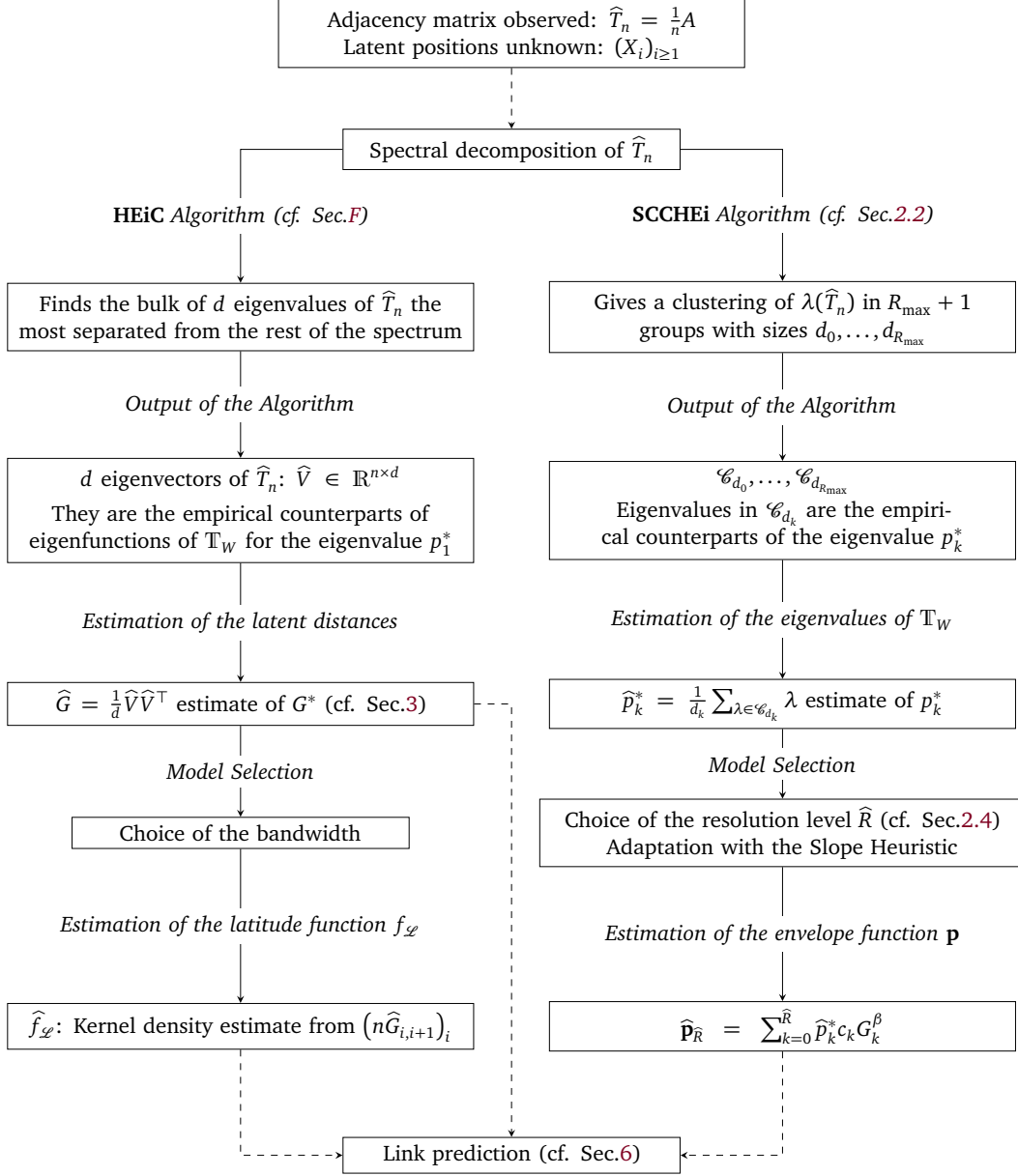


Figure 14: Synthetic presentation of the different estimation procedures.

References

- [1] E. Araya and Y. De Castro. Latent Distance Estimation for Random Geometric Graphs. In *Advances in Neural Information Processing Systems*, pages 8721–8731, 2019.
- [2] S. Arlot. Minimal penalties and the slope heuristics: a survey. *Journal de la Société Française de Statistique*, 160(3):1–106, 2019.
- [3] L. Backstrom, D. Huttenlocher, J. Kleinberg, and X. Lan. Group formation in large social networks: membership, growth, and evolution. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 44–54, 2006.
- [4] A. S. Bandeira and R. van Handel. Sharp nonasymptotic bounds on the norm of random matrices with independent entries. *Ann. Probab.*, 44(4):2479–2506, 07 2016.
- [5] R. Bhatia. *Matrix Analysis*. Graduate Texts in Mathematics. Springer New York, 1996.
- [6] S. Bubeck, J. Ding, R. Eldan, and M. Z. Rácz. Testing for high-dimensional geometry in random graphs. *Random Structures & Algorithms*, 49(3):503–532, Jan 2016.
- [7] A. E. Clementi, F. Pasquale, A. Monti, and R. Silvestri. Information spreading in stationary Markovian evolving graphs. *2009 IEEE International Symposium on Parallel & Distributed Processing*, May 2009.
- [8] Y. Dai, F. and Xu. *Approximation theory and harmonic analysis on spheres and balls*. Springer, 2013.
- [9] Y. De Castro, C. Lacour, and T. M. P. Ngoc. Adaptive Estimation of Nonparametric Geometric Graphs. *Mathematical Statistics and Learning*, 2020.
- [10] Q. Duchemin, Y. de Castro, and C. Lacour. Concentration inequality for U-statistics of order two for uniformly ergodic markov chains, 2020.
- [11] D. Durante and D. Dunson. Bayesian logistic gaussian process models for dynamic networks. In *Artificial Intelligence and Statistics*, pages 194–201. PMLR, 2014.
- [12] J. Díaz, D. Mitsche, and X. Pérez-Giménez. On the connectivity of dynamic random geometric graphs, 01 2008.
- [13] J. Fan, B. Jiang, and Q. Sun. Hoeffding’s lemma for Markov Chains and its applications to statistical learning, 2018.
- [14] D. Ferré, L. Hervé, and J. Ledoux. Limit theorems for stationary Markov processes with L2-spectral gap. *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, 48(2):396–423, May 2012.
- [15] P. Gilles. *The volume of convex bodies and Banach space geometry*. Cambridge University Press, 1989.
- [16] D. J. Higham, M. Rašajski, and N. Pržulj. Fitting a geometric graph to a protein–protein interaction network. *Bioinformatics*, 24(8):1093–1099, 03 2008.
- [17] P. D. Hoff, A. E. Raftery, and M. S. Handcock. Latent Space Approaches to Social Network Analysis. *Journal of the American Statistical Association*, 97(460):1090–1098, 2002.
- [18] B. Jiang, Q. Sun, and J. Fan. Bernstein’s inequality for general Markov Chains, 2018.
- [19] E. M. Jin, M. Girvan, and M. E. Newman. Structure of growing social networks. *Physical review E*, 64(4):046132, 2001.
- [20] O. Klopp, A. B. Tsybakov, and N. Verzelen. Oracle inequalities for network models and sparse graphon estimation. *The Annals of Statistics*, 45(1):316–354, Feb 2017.
- [21] V. Koltchinskii and E. Giné. Random Matrix Approximation of Spectra of Integral Operators. *Bernoulli*, 6, 02 2000.

- [22] C. Lo, J. Cheng, and J. Leskovec. Understanding Online Collection Growth Over Time: A Case Study of Pinterest. In *Proceedings of the 26th International Conference on World Wide Web Companion*, WWW '17 Companion, page 545–554, Republic and Canton of Geneva, CHE, 2017. International World Wide Web Conferences Steering Committee.
- [23] L. Lovász. *Large Networks and Graph Limits*. American Mathematical Society, 01 2012.
- [24] C. Matias and V. Miele. Statistical clustering of temporal networks through a dynamic stochastic block model. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(4):1119–1141, 2017.
- [25] S. P. Meyn and R. L. Tweedie. *Markov chains and stochastic stability*. Communications and Control Engineering Series. Springer-Verlag London, Ltd., London, 1993.
- [26] M. Penrose. *Random Geometric Graphs*, 01 2003.
- [27] G. O. Roberts and J. S. Rosenthal. General state space Markov chains and MCMC algorithms. *Probability Surveys*, 1(0):20–71, 2004.
- [28] G. Rossetti and R. Cazabet. Community discovery in dynamic networks: a survey. *ACM Computing Surveys (CSUR)*, 51(2):1–37, 2018.
- [29] R. S. . G. Staples. *Dynamic Geometric Graph Processes : Adjacency Operator Approach*, 2009.
- [30] M. Tang, D. L. Sussman, and C. E. Priebe. Universally consistent vertex classification for latent positions graphs. *The Annals of Statistics*, 41(3):1406–1430, Jun 2013.
- [31] J. Ugander, L. Backstrom, C. Marlow, and J. Kleinberg. Structural diversity in social contagion. *Proceedings of the National Academy of Sciences*, 109(16):5962–5966, 2012.
- [32] K. S. Xu and A. O. Hero. Dynamic stochastic blockmodels for time-evolving social networks. *IEEE Journal of Selected Topics in Signal Processing*, 8(4):552–562, 2014.
- [33] S. Yang and H. Koepl. Dependent relational gamma process models for longitudinal networks. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5551–5560. PMLR, 10–15 Jul 2018.
- [34] Y. Yu, T. Wang, and R. J. Samworth. A useful variant of the Davis–Kahan theorem for statisticians. *Biometrika*, 102(2):315–323, 04 2014.

Guidelines for the Appendix

Sections A to C: Basic definitions and Complements

In Section A we recall basic definitions on Markov chains which are required for Section B where we describe some properties verified by the Markov chain $(X_i)_{i \geq 1}$. Section C provides complementary results on the Harmonic Analysis on \mathbb{S}^{d-1} which will be useful for our proofs.

Sections D to F: Algorithms and Experiments

In Section D, we give the proof of Lemma 1 and Lemma 2. They are the cornerstones of the proof of Theorem 2 that provides a theoretical guarantee for the correctness of the algorithm SCCHEi. Section E describes precisely the slope heuristic used to perform the adaptive selection of the model dimension \hat{R} . Section F provides a complete description of the HEiC algorithm used to extract d -eigenvectors of the adjacency matrix that will be used to estimate the Gram matrix of the latent positions.

Sections G to I: Proofs of theoretical results

Thereafter, we dig into the most theoretical part of the Appendix. In Section G, we discuss the assumptions we made on the Markov chain $(X_i)_{i \geq 1}$. Section G is also dedicated to the presentation of a concentration result for a particular U-statistic of the Markov chain $(X_i)_{i \geq 1}$ that is an essential element of the proof of Theorem 1 which is provided in Section H. Finally, the proof of Theorem 4 can be found in Section I.

A Definitions for general Markov chains

We consider a state space E and a sigma-algebra Σ on E which is a standard Borel space. We denote by $(X_i)_{i \geq 1}$ a time homogeneous Markov chain on the state space (E, Σ) with transition kernel P .

A.1 Ergodic and reversible Markov chains

Definition 6. [27, section 3.2] (φ -irreducible Markov chains)

The Markov chain $(X_i)_{i \geq 1}$ is said φ -irreducible if there exists a non-zero σ -finite measure φ on E such that for all $A \in \Sigma$ with $\varphi(A) > 0$, and for all $x \in E$, there exists a positive integer $n = n(x, A)$ such that $P^n(x, A) > 0$ (where $P^n(x, \cdot)$ denotes the distribution of X_{n+1} conditioned on $X_1 = x$).

Definition 7. [27, section 3.2] (Aperiodic Markov chains)

The Markov chain $(X_i)_{i \geq 1}$ with invariant distribution π is aperiodic if there do not exist $m \geq 2$ and disjoint subsets $A_1, \dots, A_m \subset E$ with $P(x, A_{i+1}) = 1$ for all $x \in A_i$ ($1 \leq i \leq m-1$), and $P(x, A_1) = 1$ for all $x \in A_m$, such that $\pi(A_1) > 0$ (and hence $\pi(A_i) > 0$ for all i).

Definition 8. [27, section 3.4] (Geometric ergodicity)

The Markov chain $(X_i)_{i \geq 1}$ is said geometrically ergodic if there exists an invariant distribution π , $\rho \in (0, 1)$ and $C : E \rightarrow [1, \infty)$ such that

$$\|P^n(x, \cdot) - \pi\|_{TV} \leq C(x)\rho^n, \quad \forall n \geq 0, \pi\text{-a.e } x \in E,$$

where $\|\mu\|_{TV} := \sup_{A \in \Sigma} |\mu(A)|$.

Definition 9. [27, section 3.3] and [25, Chapter 16] (Uniform ergodicity)

The Markov chain $(X_i)_{i \geq 1}$ is said uniformly ergodic if there exists an invariant distribution π and constants $0 < \rho < 1$ and $L > 0$ such that

$$\|P^n(x, \cdot) - \pi\|_{TV} \leq L\rho^n, \quad \forall n \geq 0, \pi\text{-a.e } x \in E,$$

where $\|\mu\|_{TV} := \sup_{A \in \Sigma} |\mu(A)|$.

Equivalently, the Markov chain $(X_i)_{i \geq 1}$ is uniformly ergodic if the whole space \mathbb{S}^{d-1} is a small set, namely if there exist an integer $m \geq 1$, $\delta_m > 0$ and a probability measure ν such that

$$\forall x \in \mathbb{S}^{d-1}, \forall A \in \Sigma, \quad P^m(x, A) \geq \delta_m \nu(A).$$

Remark. A Markov chain geometrically or uniformly ergodic admits a unique invariant distribution and is aperiodic.

Definition 10. A Markov chain is said reversible if there exists a distribution π satisfying

$$\pi(dx)P(x, dy) = \pi(dy)P(y, dx).$$

A.2 Spectral gap

This section is largely inspired from [13]. Let us consider that the Markov chain $(X_i)_{i \geq 1}$ admits a unique invariant distribution π on \mathbb{S}^{d-1} .

For any real-valued, Σ -measurable function $h : E \rightarrow \mathbb{R}$, we define $\pi(h) := \int h(x)\pi(dx)$. The set

$$\mathcal{L}_2(E, \Sigma, \pi) := \{h : \pi(h^2) < \infty\}$$

is a Hilbert space endowed with the inner product

$$\langle h_1, h_2 \rangle_\pi = \int h_1(x)h_2(x)\pi(dx), \quad \forall h_1, h_2 \in \mathcal{L}_2(E, \Sigma, \pi).$$

The map

$$\|\cdot\|_\pi : h \in \mathcal{L}_2(E, \Sigma, \pi) \mapsto \|h\|_\pi = \sqrt{\langle h, h \rangle_\pi},$$

is a norm on $\mathcal{L}_2(E, \Sigma, \pi)$. $\|\cdot\|_\pi$ naturally allows to define the norm of a linear operator T on $\mathcal{L}_2(E, \Sigma, \pi)$ as

$$N_\pi(T) = \sup\{\|Th\|_\pi : \|h\|_\pi = 1\}.$$

To each transition probability kernel $P(x, B)$ with $x \in E$ and $B \in \Sigma$ invariant with respect to π , we can associate a bounded linear operator $h \mapsto \int h(y)P(\cdot, dy)$ on $\mathcal{L}_2(E, \Sigma, \pi)$. Denoting this operator P , we get

$$Ph(x) = \int h(y)P(x, dy), \quad \forall x \in E, \quad \forall h \in \mathcal{L}_2(E, \Sigma, \pi).$$

Let $\mathcal{L}_2^0(\pi) := \{h \in \mathcal{L}_2(E, \Sigma, \pi) : \pi(h) = 0\}$. We define the absolute spectral gap of a Markov operator.

Definition 11. (Spectral gap) A Markov operator P reversible admits an absolute spectral gap $1 - \lambda$ if

$$\lambda := \sup \left\{ \frac{\|Ph\|_\pi}{\|h\|_\pi} : h \in \mathcal{L}_2^0(\pi), h \neq 0 \right\} < 1.$$

The next result provides a connection between spectral gap and geometric ergodicity for reversible Markov chains.

Proposition 2. [14, section 2.3]

A uniformly ergodic Markov chain admits a spectral gap.

B Properties of the Markov chain

In the following, we denote $\lambda_{Leb} \equiv \lambda_{Leb,d}$ the Lebesgue measure on \mathbb{S}^{d-1} and $\lambda_{Leb,d-1}$ the Lebesgue measure on \mathbb{S}^{d-2} . Using [8, Section 1.1], it holds $b_d := \int_{x \in \mathbb{S}^{d-1}} \lambda_{Leb,d}(dx) = \frac{2\pi^{d/2}}{\Gamma(d/2)}$. Let P be the Markov operator of the Markov chain $(X_i)_{i \geq 1}$. By abuse of notation, we will also denote $P(x, \cdot)$ the density of the measure $P(x, dz)$ with respect to $\lambda_{Leb}(dz)$. For any $x, z \in \mathbb{S}^{d-1}$, we denote $R_x^z \in \mathbb{R}^{d \times d}$ a rotation matrix sending x to z (i.e. $R_x^z x = z$) and keeping $\text{Span}(x, z)^\perp$ fixed. In the following, we denote $e_d := (0, 0, \dots, 0, 1) \in \mathbb{R}^d$.

B.1 Invariant distribution and reversibility for the Markov chain

Reversibility of the Markov chain $(X_i)_{i \geq 1}$.

Lemma 3. For all $x, z \in \mathbb{S}^{d-1}$, $P(x, z) = P(z, x) = P(e_d, R_x^{e_d} x)$.

Proof of Lemma 3. Using our model described in Section 2, we get $X_2 = rX_1 + \sqrt{1-r^2}Y$ where conditionally on X_1 , Y is uniformly sampled on $\mathcal{S}(X_1) := \{q \in \mathbb{S}^{d-1} : \langle q, X_1 \rangle = 0\}$, and where r has density $f_{\mathcal{S}}$ on $[-1, 1]$. Let us consider a Gaussian vector $W \sim \mathcal{N}(0, I_d)$. Using the Cochran's theorem and Lemma 4, we know that conditionally on X_1 , the random variable $\frac{W - \langle W, X_1 \rangle X_1}{\|W - \langle W, X_1 \rangle X_1\|_2}$ is distributed uniformly on $\mathcal{S}(X_1)$.

Lemma 4. Let $W \sim \mathcal{N}(0, I_d)$. Then, $\frac{W}{\|W\|_2}$ is distributed uniformly on the sphere \mathbb{S}^{d-1} .

In the following, we denote $\stackrel{(d)}{=}$ the equality in distribution sense. We have conditionally on X_1

$$R_{X_1}^{e_d} \frac{W - \langle W, X_1 \rangle X_1}{\|W - \langle W, X_1 \rangle X_1\|_2} = \frac{\hat{W} - \langle \hat{W}, e_d \rangle e_d}{\|\hat{W} - \langle \hat{W}, e_d \rangle e_d\|_2},$$

where $\hat{W} = R_{X_1}^{e_d} W \sim \mathcal{N}(0, I_d)$. Using Cochran's theorem, we know that $\hat{W} - \langle \hat{W}, e_d \rangle e_d$ is a centered normal vector with covariance matrix the orthographic projection matrix onto the space $\text{Span}(e_d)^\perp$, leading to

$$\hat{W} - \langle \hat{W}, e_d \rangle e_d \stackrel{(d)}{=} \begin{bmatrix} Y \\ 0 \end{bmatrix},$$

where $Y \sim \mathcal{N}(0, I_{d-1})$. Using Lemma 4, we conclude that conditionally on X_1 , the random variable $\frac{W - \langle W, X_1 \rangle X_1}{\|W - \langle W, X_1 \rangle X_1\|_2}$ is distributed uniformly on $\mathcal{S}(X_1)$ (because the distribution of Y is invariant by rotation).

We deduce that

$$\begin{aligned} X_2 &\stackrel{(d)}{=} rX_1 + \sqrt{1-r^2} \frac{W - \langle W, X_1 \rangle X_1}{\|W - \langle W, X_1 \rangle X_1\|_2} \\ &\stackrel{(d)}{=} rX_1 + \sqrt{1-r^2} \frac{R_{X_2}^{X_1} W' - \langle R_{X_2}^{X_1} W', X_1 \rangle X_1}{\|R_{X_2}^{X_1} W' - \langle R_{X_2}^{X_1} W', X_1 \rangle X_1\|_2}, \end{aligned}$$

where $W' := R_{X_1}^{X_2} W$. Note that $W' \in \mathbb{R}^d$ is also a standard centered Gaussian vector because this distribution is invariant by rotation. Since $\langle R_{X_2}^{X_1} W', X_1 \rangle = \langle W', X_2 \rangle$ and $\|R_{X_2}^{X_1} q\|_2 = \|q\|_2$, $\forall q \in \mathbb{S}^{d-1}$, we deduce that

$$X_2 - rX_1 \stackrel{(d)}{=} R_{X_2}^{X_1} \left[\sqrt{1-r^2} \frac{W' - \langle W', X_2 \rangle X_2}{\|W' - \langle W', X_2 \rangle X_2\|_2} \right]. \quad (18)$$

$R_{X_1}^{X_2}$ is the rotation that sends X_1 to X_2 keeping the other dimensions fixed. Let us denote $a_1 := X_1$, $a_2 := \frac{X_2 - rX_1}{\|X_2 - rX_1\|_2}$ and complete the linearly independent family (a_1, a_2) in an orthonormal basis of \mathbb{R}^d given by $a := (a_1, a_2, \dots, a_d)$. Then, the matrix of $R_{X_1}^{X_2}$ in the basis a is

$$\begin{bmatrix} r & -\sqrt{1-r^2} & 0_{d-2}^\top \\ \sqrt{1-r^2} & r & 0_{d-2}^\top \\ 0_{d-2} & 0_{d-2} & I_{d-2} \end{bmatrix}.$$

We deduce that

$$\begin{aligned} (R_{X_2}^{X_1})^{-1} (X_2 - rX_1) &= R_{X_1}^{X_2} (X_2 - rX_1) \\ &= \|X_2 - rX_1\|_2 R_{X_1}^{X_2} \left(\frac{X_2 - rX_1}{\|X_2 - rX_1\|_2} \right) \\ &= \|X_2 - rX_1\|_2 R_{X_1}^{X_2} a_2 \\ &= \|X_2 - rX_1\|_2 \left[-\sqrt{1-r^2} a_1 + r a_2 \right] \\ &= -\sqrt{1-r^2} \|X_2 - rX_1\|_2 X_1 + r X_2 - r^2 X_1 \\ &= -(1-r^2) X_1 + r X_2 - r^2 X_1 \\ &= -X_1 + r X_2. \end{aligned}$$

Going back to Eq.(18), we deduce that

$$X_1 \stackrel{(d)}{=} rX_2 + \sqrt{1-r^2} \frac{\tilde{W} - \langle \tilde{W}, X_2 \rangle X_2}{\|\tilde{W} - \langle \tilde{W}, X_2 \rangle X_2\|_2}, \quad (19)$$

where $\tilde{W} = -W'$ is also a standard centered Gaussian vector in \mathbb{R}^d . Thus, we proved the first equality of Lemma 3. Based on Eq.(19) we have,

$$\begin{aligned} R_{X_2}^{e_d} X_1 &\stackrel{(d)}{=} r R_{X_2}^{e_d} X_2 + \sqrt{1-r^2} \frac{R_{X_2}^{e_d} \tilde{W} - \langle \tilde{W}, X_2 \rangle R_{X_2}^{e_d} X_2}{\|\tilde{W} - \langle \tilde{W}, X_2 \rangle X_2\|_2} \\ &= r e_d + \sqrt{1-r^2} \frac{R_{X_2}^{e_d} \tilde{W} - \langle R_{X_2}^{e_d} \tilde{W}, e_d \rangle e_d}{\|R_{X_2}^{e_d} \tilde{W} - \langle R_{X_2}^{e_d} \tilde{W}, e_d \rangle e_d\|_2}, \end{aligned}$$

which proves that $P(e_d, R_{X_2}^{e_d} x_1) = P(x_2, x_1)$ for any $x_1, x_2 \in \mathbb{S}^{d-1}$ because $R_{X_2}^{e_d} \tilde{W}$ is again a standard centered Gaussian vector in \mathbb{R}^d . \square

Stationary distribution of the Markov chain.

Proposition 3. *The uniform distribution on the sphere \mathbb{S}^{d-1} is a stationary distribution of the Markov chain $(X_i)_{i \geq 1}$.*

Proof of Proposition 3. Let us consider $z \in \mathbb{S}^{d-1}$. We have using Lemma 3,

$$\int_{x \in \mathbb{S}^{d-1}} P(x, z) \lambda_{Leb}(dx) = \int_{x \in \mathbb{S}^{d-1}} P(z, x) \lambda_{Leb}(dx) = 1,$$

which proves that the uniform distribution on the sphere is a stationary distribution of the Markov chain. \square

B.2 Ergodicity of the Markov chain

Our results hold under the condition that the Markov chain $(X_i)_{i \geq 1}$ is uniformly ergodic (cf. [Assumption A](#)). In this section, we provide a sufficient condition on the latitude function $f_{\mathcal{L}}$ for uniform ergodicity to hold.

Lemma 5. *We consider that $f_{\mathcal{L}}$ is bounded away from zero. Then, the Markov chain $(X_i)_{i \geq 1}$ is π -irreducible and aperiodic.*

Lemma 6. *We consider that $f_{\mathcal{L}}$ is bounded away from zero. Then the Markov chain $(X_i)_{i \geq 1}$ is uniformly ergodic.*

Proof of Lemmas 5 and 6. Considering for π the uniform distribution on \mathbb{S}^{d-1} , we get that for any $x \in \mathbb{S}^{d-1}$ and any $A \subset \mathbb{S}^{d-1}$ with $\pi(A) > 0$,

$$\begin{aligned} P(x, A) &= \int_{z \in A} P(x, z) \frac{\lambda_{Leb, d}(dz)}{b_d} \\ &= \int_{z \in A} P(e_d, R_x^{e_d} z) \frac{\lambda_{Leb, d}(dz)}{b_d} \quad (\text{Using Lemma 3}) \\ &= \int_{z \in R_x^{e_d} A} P(e_d, z) \frac{\lambda_{Leb, d}(dz)}{b_d} \\ & \quad (\text{Using the change of variable } z \mapsto R_x^{e_d} z \text{ with } R_x^{e_d} A = \{R_x^{e_d} a : a \in A\}) \\ &= \int_{r \in [-1, 1]} \int_{\xi \in \mathbb{S}^{d-2}} f_{\mathcal{L}}(r) \mathbf{1}_{(\xi^\top \sqrt{1-r^2}, r)^\top \in R_x^{e_d} A} dr \frac{\lambda_{Leb, d-1}(d\xi)}{b_{d-1} b_d} \\ &\geq \inf_{s \in [-1, 1]} f_{\mathcal{L}}(s) \int_{r \in [-1, 1]} \int_{\xi \in \mathbb{S}^{d-2}} \mathbf{1}_{(\xi^\top \sqrt{1-r^2}, r)^\top \in R_x^{e_d} A} dr \frac{\lambda_{Leb, d-1}(d\xi)}{b_{d-1} b_d} \\ &\geq \inf_{s \in [-1, 1]} f_{\mathcal{L}}(s) \int_{r \in [-1, 1]} \int_{\xi \in \mathbb{S}^{d-2}} \mathbf{1}_{(\xi^\top \sqrt{1-r^2}, r)^\top \in R_x^{e_d} A} (1-r^2)^{\frac{d-3}{2}} \frac{dr \lambda_{Leb, d-1}(d\xi)}{b_{d-1} b_d} \end{aligned}$$

$$= \frac{1}{b_{d-1}} \inf_{s \in [-1,1]} f_{\mathcal{L}}(s) \pi(R_x^{e_d} A) = \frac{1}{b_{d-1}} \inf_{s \in [-1,1]} f_{\mathcal{L}}(s) \pi(A),$$

since π is invariant by rotation and $f_{\mathcal{L}}$ is bounded away from zero. We also used that $\int_{-1}^1 (1-r^2)^{\frac{d-3}{2}} dr = \frac{b_d}{b_{d-1}}$. This result means that the whole space \mathbb{S}^{d-1} is a small set. Hence, the Markov chain is uniformly ergodic (cf. [25, Theorem 16.0.2]) and thus aperiodic and π -irreducible. \square

B.3 Computation of the absolute spectral gap of the Markov chain

Thanks to Proposition 2 (in Appendix A), we know that if $f_{\mathcal{L}}$ is such that $(X_i)_{i \geq 1}$ is uniformly ergodic, the Markov chain has an absolute spectral gap (cf. Definition 11). In the following, we show that this absolute spectral gap is equal to 1.

Keeping notations of Appendix A, let us consider $h \in L_0^2(\pi)$ such that $\|h\|_{\pi} = 1$. Then

$$\begin{aligned} \|Ph\|_{\pi}^2 &= \int_{x \in \mathbb{S}^{d-1}} \left(\int_{y \in \mathbb{S}^{d-1}} P(x, dy) h(y) \right)^2 \pi(dx) \\ &= \int_{x \in \mathbb{S}^{d-1}} \left(\int_{y \in \mathbb{S}^{d-1}} P(x, y) h(y) \pi(dy) \right)^2 \pi(dx) \\ &= \int_{x \in \mathbb{S}^{d-1}} \left(\int_{y \in \mathbb{S}^{d-1}} P(e_d, R_y^{e_d} x) h(y) \pi(dy) \right)^2 \pi(dx) \quad (\text{Using Lemma 3}) \\ &= \int_{x \in \mathbb{S}^{d-1}} \left(\int_{y \in \mathbb{S}^{d-1}} P(e_d, x) h(y) \pi(dy) \right)^2 \pi(dx) \\ &\quad (\text{Using the rotational invariance of } \pi) \\ &= \int_{x \in \mathbb{S}^{d-1}} P(e_d, x)^2 \left(\int_{y \in \mathbb{S}^{d-1}} h(y) \pi(dy) \right)^2 \pi(dx) \\ &= 0, \end{aligned}$$

where the last equality comes from $h \in L_0^2(\pi)$. Hence, the Markov chain $(X_i)_{i \geq 1}$ has 1 for absolute spectral gap.

C Complement on Harmonic Analysis on the sphere

This section completes the brief introduction to Harmonic Analysis on the sphere \mathbb{S}^{d-1} provided in Section 2. We will need in our proof the following result which states that fixing one variable and integrating with respect to the other one with the uniform measure on \mathbb{S}^{d-1} gives $\|W - W_R\|_2^2$.

Lemma 7. For any $x \in \mathbb{S}^{d-1}$,

$$\mathbb{E}_{X \sim \pi}[(W - W_R)^2(x, X)] = \|W - W_R\|_2^2,$$

where π is the uniform measure on the \mathbb{S}^{d-1} .

Proof of Lemma 7.

$$\begin{aligned} &\mathbb{E}_{X \sim \pi}[(W - W_R)^2(x, X)] \\ &= \int_y (W - W_R)^2(x, y) \pi(dy) \\ &= \int_y \left(\sum_{r > R} p_r^* \sum_{l=1}^{d_r} Y_{r,l}(x) Y_{r,l}(y) \right)^2 \pi(dy) \end{aligned}$$

$$\begin{aligned}
&= \int_{\mathcal{Y}} \sum_{r_1, r_2 > R} p_{r_1}^* p_{r_2}^* \sum_{l_1=1}^{d_{r_1}} \sum_{l_2=1}^{d_{r_2}} Y_{r_1, l_1}(x) Y_{r_1, l_1}(y) Y_{r_2, l_2}(x) Y_{r_2, l_2}(y) \pi(dy) \\
&= \sum_{r_1, r_2 > R} p_{r_1}^* p_{r_2}^* \sum_{l_1=1}^{d_{r_1}} \sum_{l_2=1}^{d_{r_2}} Y_{r_1, l_1}(x) Y_{r_2, l_2}(x) \int_{\mathcal{Y}} Y_{r_1, l_1}(y) Y_{r_2, l_2}(y) \pi(dy).
\end{aligned}$$

Since $\int_{\mathcal{Y}} Y_{r, l}(y) Y_{r', l'}(y) \pi(dy)$ is 1 if $r = r'$ and $l = l'$ and 0 otherwise, we have that

$$\begin{aligned}
\mathbb{E}_{X \sim \pi}[(W - W_R)^2(x, X)] &= \sum_{r > R} (p_r^*)^2 \sum_{l=1}^{d_r} Y_{r, l}(x)^2 \\
&= \sum_{r > R} (p_r^*)^2 d_r \quad (\text{Using [8, Eq.(1.2.9)]}) \\
&= \|W - W_R\|_2^2.
\end{aligned}$$

□

Let us consider $\beta := \frac{d-2}{2}$ and the weight function $w_\beta(t) := (1-t^2)^{\beta-\frac{1}{2}}$. As highlighted in section 2, any envelope function $\mathbf{p} \in L^2([-1, 1], w_\beta)$ can be decomposed as $\mathbf{p} \equiv \sum_{k=0}^R p_k^* c_k G_k^\beta$ where G_l^β is the Gegenbauer polynomial of degree l with parameter β and where $c_k := \frac{2k+d-2}{d-2}$. The Gegenbauer polynomials are orthonormal polynomials on $[-1, 1]$ associated with the weight function w_β . The eigenvalues $(p_k^*)_{k \geq 0}$ of the envelope function can be computed numerically through the formula

$$\forall l \geq 0, \quad p_l^* = \left(\frac{c_l b_d}{d_l} \right) \int_{-1}^1 p(t) G_l^\beta(t) w_\beta(t) dt,$$

where $b_d := \frac{\Gamma(\frac{d}{2})}{\Gamma(\frac{1}{2})\Gamma(\frac{d}{2}-\frac{1}{2})}$ with Γ the Gamma function. Hence, it is possible to recover the envelope function \mathbf{p} thanks to the identity

$$\mathbf{p} = \sum_{l \geq 0} \sqrt{d_l} p_l^* \frac{G_l^\beta}{\|G_l^\beta\|_{L^2([-1, 1], w_\beta)}} = \sum_{l \geq 0} p_l^* c_l G_l^\beta. \quad (20)$$

D Proofs of the two key lemmas for Theorem 2

In the proofs of Lemma 1 and Lemma 2 provided in this section, we keep the notations and the assumptions used in the proof of Theorem 2. To ease the reading of this section, we recall here important notations. We denoted

$$\Delta^G = \min_{0 \leq k \neq l \leq R, p_k^* \neq p_l^*} |p_k^* - p_l^*| \wedge \min_{0 \leq k \leq R, p_k^* \neq 0} |p_k^*| > 0.$$

For any $g \in (0, \frac{\Delta^G}{4})$, the proof of Theorem 1 (cf. Section H) ensures that for n large enough it holds

$$\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\widehat{T}_n)) \leq g^2. \quad (21)$$

Let us finally recall (cf. Section 1) that

$$\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\widehat{T}_n)) = \inf_{\sigma \in \mathfrak{S}} \sum_{i \geq 1} ((\lambda(\mathbb{T}_{W_R})_{\sigma(i)} - \lambda^R(\widehat{T}_n)_i)^2. \quad (22)$$

D.1 Proof of Lemma 1

We denote σ^* a permutation achieving the minimum in Eq.(22).

- First we show that we can choose σ^* such that $\sigma^*({1, \dots, \tilde{R}}) = {1, \dots, \tilde{R}}$. We recall that

$$\lambda(\mathbb{T}_{W_R}) = \left(\underbrace{p_0^*, p_1^*, \dots, p_1^*}_{d_0=1}, \underbrace{p_1^*, \dots, p_1^*}_{d_1=d}, \dots, \underbrace{p_R^*, \dots, p_R^*}_{d_R}, 0, 0, \dots \right),$$

$$\text{and } \lambda^R(\hat{T}_n) = \left(\underbrace{\lambda^R(\hat{T}_n)_1, \dots, \lambda^R(\hat{T}_n)_{\tilde{R}}}_{\tilde{R}}, 0, 0, \dots \right),$$

with $\lambda^R(\hat{T}_n)_1 \geq \dots \geq \lambda^R(\hat{T}_n)_{\tilde{R}}$.

\rightsquigarrow If $p_k^* \neq 0$ for all $0 \leq k \leq R$, then it is clear that $\sigma^*({1, \dots, \tilde{R}}) = {1, \dots, \tilde{R}}$. Otherwise, there would exist some $i \in {1, \dots, \tilde{R}}$ such that $\sigma^*(j) \neq i$ for all $j \in {1, \dots, \tilde{R}}$. Hence, we would obtain that $\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\hat{T}_n)) \geq |\lambda(\mathbb{T}_{W_R})_i|^2 \geq (\Delta^G)^2$, which would contradict Eq.(21).

\rightsquigarrow If $p_k^* = 0$ for all $0 \leq k \leq R$, it is clear that we can take $\sigma^* = \text{Id}$.

\rightsquigarrow Otherwise, let us denote *Null* the list of all indexes $i \in {1, \dots, \tilde{R}}$ such that $\lambda(\mathbb{T}_{W_R})_i = 0$. It holds that $N_0 = |\text{Null}| = \sum_{0 \leq k \leq R, s.t. p_k^* = 0} d_k$. We also denote *NoNull* the complement of *Null* in ${1, \dots, \tilde{R}}$ (i.e. the list of indexes in ${1, \dots, \tilde{R}}$ that are not in *Null*).

For any $1 \leq i \leq \tilde{R}$ such that $\lambda(\mathbb{T}_{W_R})_i \neq 0$, it must exist some $j \in {1, \dots, \tilde{R}}$ such that $\sigma^*(j) = i$. Otherwise, we would have

$$\delta_2^2(\lambda(\mathbb{T}_{W_R}), \lambda^R(\hat{T}_n)) \geq |\lambda(\mathbb{T}_{W_R})_i|^2 \geq (\Delta^G)^2,$$

which would contradict Eq.(21). Hence, we get that

$$(\sigma^*)^{-1}(\text{NoNull}) \subset {1, \dots, \tilde{R}}.$$

We deduce that for any $i \in {1, \dots, \tilde{R}} \setminus (\sigma^*)^{-1}(\text{NoNull})$, $\lambda(\mathbb{T}_{W_R})_{\sigma^*(i)} = 0$. Hence, we can define σ^* such that this permutation sends the N_0 indexes in ${1, \dots, \tilde{R}} \setminus (\sigma^*)^{-1}(\text{NoNull})$ to the N_0 indexes in *Null*. Such σ^* still achieves the minimum in Eq.(22). In the following, we thus consider that $\sigma^*({1, \dots, \tilde{R}}) = {1, \dots, \tilde{R}}$.

- Let us recall that the function f^* is defined by

$$f^* : {1, \dots, \tilde{R}} \rightarrow \{p_k^*, 0 \leq k \leq R\}$$

$$i \mapsto \lambda(\mathbb{T}_{W_R})_{\sigma^*(i)}.$$

Note that for any $1 \leq i \leq \tilde{R}$, $\sigma^*(i) \leq \tilde{R}$ thanks to the previous paragraph. We denote $p_{(0)}^* \geq \dots \geq p_{(R)}^*$ the ordered sequence of p_0^*, \dots, p_R^* and $d_{(k)}$ is the multiplicity of the eigenvalue $p_{(k)}^*$ of the operator \mathbb{T}_W . We show that f^* is such that $f^*(1) = \dots = f^*(d_{(0)}) = p_{(0)}^*$, $f^*(d_{(0)} + 1) = \dots = f^*(d_{(0)} + d_{(1)}) = p_{(1)}^*$, $f^*(d_{(0)} + d_{(1)} + 1) = \dots = f^*(d_{(0)} + d_{(1)} + d_{(2)}) = p_{(2)}^*$, This is equivalent to say that the function f^* is non-increasing. If this was not true, it would mean that there exist $1 \leq j < i \leq \tilde{R}$ such that $f^*(j) < f^*(i)$. Since $\lambda^R(\hat{T}_n)_j \geq \lambda^R(\hat{T}_n)_i$ (because $j < i$), we would get that

$$\begin{aligned} \Delta^G &< f^*(i) - f^*(j) \\ &= \underbrace{f^*(i) - \lambda^R(\hat{T}_n)_i}_{\leq g} + \underbrace{\lambda^R(\hat{T}_n)_i - \lambda^R(\hat{T}_n)_j}_{\leq 0} + \underbrace{\lambda^R(\hat{T}_n)_j - f^*(j)}_{\leq g} \\ \text{i.e. } \Delta^G &\leq 2g. \end{aligned}$$

Since we chose g such that $\Delta^G > 4g$, this previous inequality is absurd. This concludes the proof.

D.2 Proof of Lemma 2

We prove our result by induction. In the following, we say that an intermediate state of the HAC algorithm is *valid* if it is still possible to reach state (\mathcal{S}) in the next iterations of the algorithm. Stated otherwise, a state is *valid* if it does not exist $1 \leq i \neq j \leq \tilde{R}$ such that $f^*(i) \neq f^*(j)$ with $\lambda^R(\hat{T}_n)_i$ and $\lambda^R(\hat{T}_n)_j$ in the

same cluster. It is obvious that the initial state of the HAC algorithm is *valid* since all eigenvalues are alone in their respective clusters.

Suppose now that we are at iteration $2 \leq t \leq \tilde{R} - R - 2$ of the HAC algorithm and that our procedure is *valid* until step t . We are sure that we did not reach a state of type (\mathcal{S}) before step t because only the state at depth R from the root of the HAC's tree contains exactly $R + 1$ clusters. For any cluster S formed at step t by the HAC algorithm, we denote by abuse of notation $f^*(S) := f^*(i)$ for any i such that $\lambda^R(\hat{T}_n)_i \in S$ (which is licit since step t is *valid*). By contradiction, assume that the algorithm does not make a valid merging at step $t + 1$. This means that the two merged clusters S_a and S_b at step $t + 1$ are such that $f^*(S_a) \neq f^*(S_b)$. Since at step t we did not reach a state of type (\mathcal{S}) , this means that there are two clusters S_i and S_j with $i \neq j$ such that $f^*(S_i) = f^*(S_j)$.

For any $\lambda^R(\hat{T}_n)_i \in S_i$ and $\lambda^R(\hat{T}_n)_j \in S_j$,

$$|\lambda^R(\hat{T}_n)_i - \lambda^R(\hat{T}_n)_j| \leq |\lambda^R(\hat{T}_n)_i - f^*(S_i)| + \underbrace{|f^*(S_i) - \lambda^R(\hat{T}_n)_j|}_{=|f^*(S_j) - \lambda^R(\hat{T}_n)_j|} \leq 2g,$$

and for any $\lambda^R(\hat{T}_n)_a \in S_a$ and $\lambda^R(\hat{T}_n)_b \in S_b$,

$$\begin{aligned} & |\lambda^R(\hat{T}_n)_a - \lambda^R(\hat{T}_n)_b| \\ & \geq -|\lambda^R(\hat{T}_n)_a - f^*(S_a)| + |f^*(S_a) - \lambda^R(\hat{T}_n)_b| \\ & \geq |f^*(S_a) - f^*(S_b)| - |\lambda^R(\hat{T}_n)_a - f^*(S_a)| - |\lambda^R(\hat{T}_n)_b - f^*(S_b)| \\ & \geq \Delta^G - 2g. \end{aligned}$$

Since we chose $\Delta^G > 4g$, we get

$$d_c(S_a, S_b) > d_c(S_i, S_j).$$

This is a contradiction since at step t , the HAC algorithm merges the two clusters with the smallest complete linkage distance. Hence, the algorithm performs a valid merging at step $t + 1$.

We proved that a state of type (\mathcal{S}) is reached by the HAC algorithm with complete linkage at iteration $\tilde{R} - R - 1$. Since $d \geq 3$, it holds $d_0 < d_1 < d_2 < \dots$ and since the SCCHEi starts by selecting the cluster of size d_0 in the tree as close as possible to the root, we get $\mathcal{C}_{d_0} = \hat{\mathcal{C}}_{d_0}$. Continuing the process of the "for loop" in the SCCHEi algorithm, the SCCHEi algorithm then selects the cluster of size d_1 in the remaining tree (where we removed all eigenvalues in $\hat{\mathcal{C}}_{d_0}$ in the tree of the HAC). Hence, the SCCHEi algorithm sets $\mathcal{C}_{d_1} = \hat{\mathcal{C}}_{d_1}$. Following this procedure, it is straightforward to see that the SCCHEi returns the partition $\mathcal{C}_{d_0} = \hat{\mathcal{C}}_{d_0}, \dots, \mathcal{C}_{d_R} = \hat{\mathcal{C}}_{d_R}$.

E Slope heuristic

We propose a detailed analysis of the slope heuristic described in Section 2.4 on simulated data using $d = 3$, the envelope function $\mathbf{p}^{(1)}$ and the latitude function $f_{\mathcal{S}}^{(1)}$ presented in Eq.(11). We recall that $R(\kappa)$ represents the optimal value of R to minimize the bias-variance decomposition defined by Eq.(7) for a given hyperparameter κ . Figure 15 shows the evolution of $\tilde{R}(\kappa)$ with respect to κ which is sampled on a logscale. $\tilde{R}(\kappa)$ is the dimension of the space of Spherical Harmonics with degree at most $R(\kappa)$. Our slope heuristic consists in choosing the value κ_0 leading to the larger jump of the function $\kappa \mapsto \tilde{R}(\kappa)$. In our case, Figure 15 shows that $\kappa_0 = 10^{-3.9}$. As described in Section 2.2, the resolution level \hat{R} selected to cluster the eigenvalues of the matrix \hat{T}_n is given by $R(2\kappa_0)$.

F Reminder on Harmonic EigenCluster(HEiC)

Before presenting the algorithm HEiC, let us define for a given set of indices $i_1, \dots, i_d \in [n]$

$$\text{Gap}_1(\hat{T}_n; i_1, \dots, i_d) := \min_{i \notin \{i_1, \dots, i_d\}} \max_{j \in \{i_1, \dots, i_d\}} |\hat{\lambda}_i - \hat{\lambda}_j|.$$

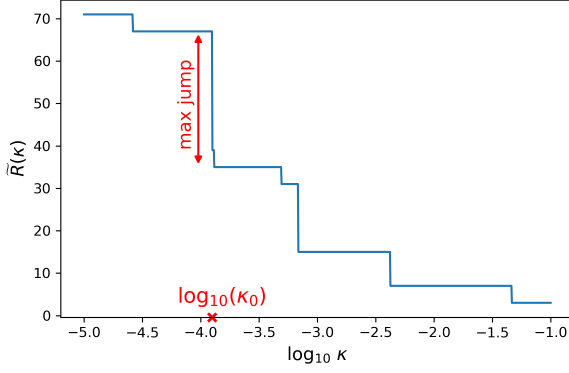


Figure 15: We sample the parameter κ on a logscale between 10^{-5} and 10^{-1} and we compute the corresponding $R(\kappa)$ defined in Eq. (7). We plot the values of $\tilde{R}(\kappa)$ with respect to κ . The larger jump allows us to define κ_0 .

Algorithm 3 Harmonic EigenCluster(HEiC) algorithm.

Data: Adjacency matrix A . Dimension d .

- 1: $(\hat{\lambda}_1^{sort}, \dots, \hat{\lambda}_n^{sort}) \leftarrow$ eigenvalues of \hat{T}_n sorted in decreasing order.
- 2: $\Lambda_1 \leftarrow \{\hat{\lambda}_1^{sort}, \dots, \hat{\lambda}_d^{sort}\}$.
- 3: Initialize $i = 2$ and $\text{gap} = \text{Gap}_1(\hat{T}_n; 1, 2, \dots, d)$.
- 4: **while** $i \leq n - d + 1$ **do**
- 5: **if** $\text{Gap}_1(\hat{T}_n; i, i + 1, \dots, i + d - 1) > \text{gap}$ **then**
- 6: $\Lambda_1 \leftarrow \{\hat{\lambda}_i^{sort}, \dots, \hat{\lambda}_{i+d-1}^{sort}\}$
- 7: **end if**
- 8: $i = i + 1$
- 9: **end while**

Return: Λ_1, gap .

G Concentration inequality for U-statistics with Markov chains

In this section, we present a recent concentration inequality for a U-statistic of the Markov chain $(X_i)_{i \geq 1}$ from [10] which is a key result to prove Theorem 1. In the first subsection, we remind the assumptions made on the Markovian dynamic, namely **Assumption A**.

G.1 Assumptions and notations for the Markov chain

Let us recall that **Assumption A** states that the latitude function $f_{\mathcal{L}}$ is such that $\|f_{\mathcal{L}}\|_{\infty} < \infty$ and makes the chain $(X_i)_{i \geq 1}$ uniformly ergodic. **Assumption A** guarantees in particular that there exists $\delta_M > 0$ such that

$$\forall x \in \mathbb{S}^{d-1}, \forall A \in \mathcal{B}(\mathbb{S}^{d-1}), \quad P(x, A) \leq \delta_M \nu(A),$$

for some probability measure ν (e.g. the uniform measure on the sphere π).

In Section B.2, we provide a sufficient condition on the latitude function $f_{\mathcal{L}}$ ensuring the uniform ergodicity of the chain with associated constants $L > 0$ and $0 < \rho < 1$ (cf. Definition 9). In Section B.3, we explain why **Assumption A** ensures that the Markov chain $(X_i)_{i \geq 1}$ has a spectral gap and we show that this spectral gap is equal to 1.

G.2 Concentration inequality of U-statistic for Markov chain

One key result to prove Theorem 1 is the concentration of the following U-statistic

$$U_{stat}(n) = \frac{1}{n^2} \sum_{1 \leq i < j \leq n} [(W - W_R)^2(X_i, X_j) - \|W - W_R\|_2^2].$$

Note that $\|W - W_R\|_2^2$ corresponds to the expectation of the kernel $(W - W_R)^2(\cdot, \cdot)$ under the uniform distribution on \mathbb{S}^{d-1} which is known to be the unique invariant distribution π of the Markov chain $(X_i)_{i \geq 1}$

(cf. Appendix B). More precisely, for any $x \in \mathbb{S}^{d-1}$, it holds

$$\|W - W_R\|_2^2 = \mathbb{E}_{X \sim \pi}[(W - W_R)^2(x, X)] = \mathbb{E}_{(X, X') \sim \pi \otimes \pi}[(W - W_R)^2(X, X')],$$

see Lemma 7 for a proof. Applying [10, Theorem 2] in our framework leads to the following result.

Lemma 8. *Let us consider $\gamma \in (0, 1)$ satisfying $\log(e \log(n)/\gamma) \leq n$. Then it holds with probability at least $1 - \gamma$,*

$$U_{\text{stat}}(n) \leq M \frac{\|\mathbf{p} - \mathbf{p}_R\|_\infty^2 \log n}{n} \log(e \log(n)/\gamma),$$

where $M > 0$ only depends on constants related to the Markov chain $(X_i)_{i \geq 1}$.

H Proof of Theorem 1

The proof of Theorem 1 mainly lies in the following result which is proved in Section H.1. Coupling the convergence of the spectrum of the matrix of probability T_n with a concentration result on the spectral norm of random matrices with independent entries (cf. [4]), we show the convergence in metric δ_2 of the spectrum of \hat{T}_n towards the spectrum of the integral operator \mathbb{T}_W .

Theorem 5. *Let us consider $\gamma \in (0, 1)$ satisfying $\log(e \log(n)/\gamma) \leq n/(13\tilde{R})$. Then it holds with probability at least $1 - \gamma$,*

$$\begin{aligned} & \delta_2(\lambda(\mathbb{T}_W), \lambda(T_n)) \\ & \leq 2\|\mathbf{p} - \mathbf{p}_R\|_2 + 8\sqrt{\frac{\tilde{R}}{n} \ln(e/\gamma)} + M\|\mathbf{p} - \mathbf{p}_R\|_\infty \sqrt{\frac{\log n}{n}} (\log(e \log(n)/\gamma))^{1/2}, \end{aligned}$$

where $M > 0$ only depends on constants related to the Markov chain $(X_i)_{i \geq 1}$ (cf. Lemma 8).

First part of the proof for Theorem 1 We start by establishing the convergence rate for $\delta_2(\lambda(\mathbb{T}_W), \lambda(T_n))$. We keep notations of Theorem 5. Let us consider $\gamma \in (0, 1)$ satisfying $\log(e \log(n)/\gamma) \leq n/(13\tilde{R})$, and assume that $p \in Z_{w_\beta}^s((-1, 1))$ with $s > 0$.

Let us define the event

$$\begin{aligned} \Omega(\gamma) := & \left\{ \delta_2(\lambda(\mathbb{T}_W), \lambda(T_n)) \leq 2\|\mathbf{p} - \mathbf{p}_R\|_2 + 8\sqrt{\frac{\tilde{R}}{n} \ln(e/\gamma)} \right. \\ & \left. + M\|\mathbf{p} - \mathbf{p}_R\|_\infty \sqrt{\frac{\log n}{n}} (\log(e \log(n)/\gamma))^{1/2} \right\}. \end{aligned}$$

Using Theorem 5, it holds $\mathbb{P}(\Omega(\gamma)) \geq 1 - \gamma$. Remarking further that

$$\delta_2(\lambda(\mathbb{T}_W), \lambda(T_n)) \leq \delta_2(\lambda(\mathbb{T}_W), 0) + \delta_2(0, \lambda(T_n)) \leq \|\mathbf{p}\|_2 + \sqrt{n} \leq \sqrt{2} + \sqrt{n},$$

we have

$$\begin{aligned} & \mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda(T_n))] \\ & = \mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda(T_n)) \mathbb{1}_{\Omega(\gamma)}] + (1 + \sqrt{2})^2 n \mathbb{P}(\Omega(\gamma)^c) \\ & \leq c\|\mathbf{p} - \mathbf{p}_R\|_2^2 + c\frac{\tilde{R}}{n} \log(e/\gamma) + c\|\mathbf{p} - \mathbf{p}_R\|_\infty^2 \frac{\log n}{n} \log(e \log(n)/\gamma) \\ & \quad + (1 + \sqrt{2})^2 n \gamma, \end{aligned}$$

where $c > 0$ is a constant that does not depend on R , d nor n . Since for some constant $C(\mathbf{p}, s, d) > 0$ (depending only on \mathbf{p} , s and d)

$$\|\mathbf{p} - \mathbf{p}_R\|_2^2 = \sum_{k > R} (p_k^*)^2 d_k \frac{(1 + k(k + 2\beta))^s}{(1 + k(k + 2\beta))^s} \leq C(\mathbf{p}, s, d) R^{-2s}, \quad (23)$$

and since

$$\tilde{R} = O(R^{d-1}), \quad (24)$$

we have choosing $\gamma = 1/n^2$

$$\mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda(T_n))] \leq D' \left[R^{-2s} + R^{d-1} \frac{\log(n)}{n} + \|\mathbf{p} - \mathbf{p}_R\|_\infty^2 \frac{\log^2(n)}{n} \right], \quad (25)$$

where $D' > 0$ is a constant independent of n and R . Let us show that choosing $R = \lfloor (n/\log^2(n))^{\frac{1}{2s+d-1}} \rfloor$ concludes the proof. Since $\|G_k^\beta\|_\infty = G_k^\beta(1) = d_k/c_k$, we get that

$$\|\mathbf{p}_R\|_\infty \leq \sum_{k=0}^R |p_k^*| c_k G_k^\beta(1) = \sum_{k=0}^R |p_k^*| d_k \leq \sqrt{\tilde{R}} \|\mathbf{p}_R\|_2,$$

and using Eq.(30), we deduce that

$$\|\mathbf{p} - \mathbf{p}_R\|_\infty \leq \|\mathbf{p}\|_\infty + \|\mathbf{p}_R\|_\infty \leq 1 + \sqrt{2\tilde{R}}. \quad (26)$$

Hence, Eq.(25) becomes

$$\mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda(T_n))] \leq D'' \left[R^{-2s} + R^{d-1} \frac{\log(n)}{n} + \tilde{R} \frac{\log^2(n)}{n} \right],$$

where D'' is a constant that does not depend on n nor R . Choosing $R = \lfloor (n/\log^2(n))^{\frac{1}{2s+d-1}} \rfloor$ and using Eq.(24) we get

$$\begin{aligned} & \mathbb{E}[\delta_2^2(\lambda(\mathbb{T}_W), \lambda(T_n))] \\ & \leq D'' \left[\left(\frac{n}{\log^2(n)} \right)^{\frac{-2s}{2s+d-1}} + 2 \left(\frac{n}{\log^2(n)} \right)^{\frac{d-1}{2s+d-1}} \frac{\log^2(n)}{n} \right] \\ & \leq 3D'' \left(\frac{n}{\log^2(n)} \right)^{\frac{-2s}{2s+d-1}}. \end{aligned}$$

Second part of the proof for Theorem 1 Let us recall that in the statement of Theorem 1, $\lambda^{R_{opt}}(\hat{T}_n)$ is the sequence of the \tilde{R}_{opt} first eigenvalues (sorted in decreasing absolute values) of the matrix \hat{T}_n where R_{opt} is the value of the parameter R leading to the optimal bias-variance trade off, namely

$$\lambda^{R_{opt}}(\hat{T}_n) = (\hat{\lambda}_1, \dots, \hat{\lambda}_{\tilde{R}_{opt}}, 0, 0, \dots).$$

From the computations of the first part of the proof, we know that $R_{opt} = \lfloor (n/\log^2(n))^{\frac{1}{2s+d-1}} \rfloor$. That corresponds to the situation where we choose optimally R and it is in practice possible to approximate this best model dimension using e.g. the slope heuristic. Therefore, $\delta_2(\lambda(\mathbb{T}_W), \lambda^{R_{opt}}(\hat{T}_n))$ is the quantity of interest since it represents the distance between the eigenvalues used to build our estimates $(\hat{p}_k)_k$ and the true spectrum of the envelope function \mathbf{p} . Since $\tilde{R} = \mathcal{O}(R^{d-1})$ for all integer $R \geq 0$, we have $\tilde{R}_{opt} = \mathcal{O}\left((n/\log^2(n))^{\frac{d-1}{2s+d-1}}\right)$. We deduce that for n large enough $2\tilde{R}_{opt} \leq n$ and using [9, Proposition 15] we obtain

$$\begin{aligned} & \delta_2(\lambda^{R_{opt}}(\hat{T}_n), \lambda(\mathbb{T}_{W_{R_{opt}}})) \\ & \leq \delta_2(\lambda(T_n), \lambda(\mathbb{T}_{W_{R_{opt}}})) + \sqrt{2\tilde{R}_{opt}} \|\hat{T}_n - T_n\| \\ & \leq \delta_2(\lambda(T_n), \lambda(\mathbb{T}_W)) + \delta_2(\lambda(\mathbb{T}_W), \lambda(\mathbb{T}_{W_{R_{opt}}})) + \sqrt{2\tilde{R}_{opt}} \|\hat{T}_n - T_n\|, \end{aligned} \quad (27)$$

where $\lambda(\mathbb{T}_{W_{R_{opt}}}) = (\lambda_1^*, \dots, \lambda_{\tilde{R}_{opt}}^*, 0, 0, \dots)$. Let us consider $\gamma \in (0, 1)$. Using Theorem 5, we know that with probability at least $1 - \gamma$ it holds for n large enough

$$\begin{aligned} \delta_2(\lambda(T_n), \lambda(\mathbb{T}_W)) &\leq 2\|\mathbf{p} - \mathbf{p}_{R_{opt}}\|_2 + 8\sqrt{\frac{\tilde{R}_{opt}}{n} \ln(e/\gamma)} \\ &\quad + M\|\mathbf{p} - \mathbf{p}_{R_{opt}}\|_\infty \sqrt{\frac{\log n}{n}} (\log(e \log(n)/\gamma))^{1/2}. \end{aligned}$$

Using Eq.(23), Eq.(26) and the fact that $\tilde{R} = \mathcal{O}(R^{d-1})$, it holds with probability at least $1 - 1/n^2$,

$$\begin{aligned} \delta_2^2(\lambda(T_n), \lambda(\mathbb{T}_W)) &\leq c \left[R_{opt}^{-2s} + R_{opt}^{d-1} \frac{\log n}{n} + M R_{opt}^{d-1} \frac{\log^2 n}{n} \right] \\ &\leq (M')^2 (n/\log^2 n)^{\frac{-2s}{2s+d-1}}, \end{aligned}$$

where $c > 0$ is a numerical constant and $M' > 0$ depends on constants related to the Markov chain $(X_i)_{i \geq 1}$ (see Theorem 5 for details). Moreover,

$$\begin{aligned} \delta_2^2(\lambda(\mathbb{T}_W), \lambda(\mathbb{T}_{W_{R_{opt}}})) &= \|\mathbf{p} - \mathbf{p}_{R_{opt}}\|_2^2 \\ &\leq C(\mathbf{p}, s, d) R_{opt}^{-2s} = \mathcal{O}\left((n/\log^2 n)^{\frac{-2s}{2s+d-1}}\right), \end{aligned} \quad (28)$$

where we used Eq.(23). Finally, using the concentration of spectral norm for random matrices with independent entries from [4], there exists a universal constant $C_0 > 0$ such that conditionally on $(X_i)_{i \geq 1}$, it holds with probability at least $1 - 1/n^2$,

$$\|T_n - \hat{T}_n\| \leq \frac{3}{\sqrt{2n}} + C_0 \frac{\sqrt{\log(n^3)}}{n}.$$

Using again $\tilde{R} = \mathcal{O}(R^{d-1})$, this implies that for n large enough, it holds conditionally on $(X_i)_{i \geq 1}$ with probability at least $1 - 1/n^2$,

$$\sqrt{2\tilde{R}_{opt}} \|T_n - \hat{T}_n\| \leq D(n/\log^2 n)^{\frac{-s}{2s+d-1}},$$

where $D > 0$ is a numerical constant. From Eq.(27), we deduce that $\mathbb{P}(\Omega) \geq 1 - 2/n^2$ where the event Ω is defined by

$$\Omega = \left\{ \delta_2^2(\lambda^{R_{opt}}(\hat{T}_n), \lambda(\mathbb{T}_{W_{R_{opt}}})) \leq (C(\mathbf{p}, s, d)^{1/2} + D + M')^2 (n/\log^2 n)^{\frac{-2s}{2s+d-1}} \right\}.$$

Remarking finally that

$$\begin{aligned} \delta_2(\lambda^{R_{opt}}(\hat{T}_n), \lambda(\mathbb{T}_{W_{R_{opt}}})) &\leq \delta_2(\lambda(\mathbb{T}_{W_{R_{opt}}}), 0) + \delta_2(0, \lambda(\hat{T}_n)) \\ &\leq \|\mathbf{p}\|_2 + \sqrt{n} \leq \sqrt{2} + \sqrt{n}, \end{aligned}$$

we obtain

$$\begin{aligned} &\mathbb{E} \left[\delta_2^2(\lambda^{R_{opt}}(\hat{T}_n), \lambda(\mathbb{T}_{W_{R_{opt}}})) \right] \\ &\leq \mathbb{E} \left[\delta_2^2(\lambda^{R_{opt}}(\hat{T}_n), \lambda(\mathbb{T}_{W_{R_{opt}}})) \mid \Omega \right] + \mathbb{P}(\Omega^c) (\sqrt{2} + \sqrt{n})^2 \\ &\leq (C(\mathbf{p}, s, d)^{1/2} + D + M')^2 (n/\log^2 n)^{\frac{-2s}{2s+d-1}} + 2 \frac{(\sqrt{2} + \sqrt{n})^2}{n^2} \\ &= \mathcal{O}\left((n/\log^2 n)^{\frac{-2s}{2s+d-1}}\right). \end{aligned} \quad (29)$$

Using the triangle inequality, Eq.(28) and Eq.(29) lead to

$$\begin{aligned} &\mathbb{E} \left[\delta_2^2(\lambda^{R_{opt}}(\hat{T}_n), \lambda(\mathbb{T}_W)) \right] \\ &\leq 3\mathbb{E} \left[\delta_2^2(\lambda^{R_{opt}}(\hat{T}_n), \lambda(\mathbb{T}_{W_{R_{opt}}})) \right] + 3\delta_2^2(\lambda(\mathbb{T}_{W_{R_{opt}}}), \lambda(\mathbb{T}_W)) \\ &= \mathcal{O}\left((n/\log^2 n)^{\frac{-2s}{2s+d-1}}\right), \end{aligned}$$

which concludes the proof of Theorem 1.

H.1 Proof of Theorem 5

We follow the same sketch of proof as in [9]. Let $R \geq 1$ and define,

$$\begin{aligned}
\Phi_{k,l} &= \frac{1}{\sqrt{n}} [Y_{k,l}(X_1), \dots, Y_{k,l}(X_n)] \in \mathbb{R}^n, \\
E_{R,n} &= (\langle \Phi_{k,l}, \Phi_{k',l'} \rangle - \delta_{(k,l),(k',l')})_{(k,k') \in [R], l \in \{1, \dots, d_k\}, l' \in \{1, \dots, d_{k'}\}} \in \mathbb{R}^{\tilde{R} \times \tilde{R}}, \\
X_{R,n} &= [\Phi_{0,1}, \Phi_{1,1}, \Phi_{1,2}, \dots, \Phi_{R,d_R}] \in \mathbb{R}^{n \times \tilde{R}}, \\
A_{R,n} &= (X_{R,n}^\top X_{R,n})^{1/2} \text{ with } A_{R,n}^2 = \text{Id}_{\tilde{R}} + E_{R,n}, \\
K_R &= \text{Diag}(\lambda_1(\mathbb{T}_W), \dots, \lambda_{\tilde{R}}(\mathbb{T}_W)), \\
T_{R,n} &= \sum_{k=0}^R p_k^* \sum_{l=1}^{d_k} \Phi_{k,l}(\Phi_{k,l})^\top = X_{R,n} K_R X_{R,n}^\top \in \mathbb{R}^{n \times n} \\
\tilde{T}_{R,n} &= ((1 - \delta_{i,j}) T_{R,n})_{i,j \in [n]} \in \mathbb{R}^{n \times n}, \\
T_{R,n}^* &= A_{R,n} K_R A_{R,n}^\top \in \mathbb{R}^{\tilde{R} \times \tilde{R}}, \\
W_R(x, y) &= \sum_{k=0}^R p_k^* \sum_{l=1}^{d_k} Y_{k,l}(x) Y_{k,l}(y).
\end{aligned}$$

It holds

$$\delta_2(\lambda(\mathbb{T}_W), \lambda(\mathbb{T}_{W_R})) = \left(\sum_{k \geq R} d_k (p_k^*)^2 \right)^{1/2}.$$

We point out the equality between spectra of the operator \mathbb{T}_{W_R} and the matrix K_R . Using the SVD decomposition of $X_{R,n}$, one can also easily prove that $\lambda(T_{R,n}) = \lambda(T_{R,n}^*)$. We deduce that

$$\begin{aligned}
\delta_2(\lambda(\mathbb{T}_{W_R}), \lambda(T_{R,n})) &= \delta_2(\lambda(K_R), \lambda(T_{R,n}^*)) \\
&\leq \|T_{R,n}^* - K_R\|_F \\
&= \|A_{R,n} K_R A_{R,n} - K_R\|_F,
\end{aligned}$$

with the Hoffman-Wielandt inequality. Using equation (4.8) at ([21] p.127) gives

$$\delta_2(\lambda(\mathbb{T}_{W_R}), \lambda(T_{R,n})) \leq \sqrt{2} \|K_R\|_F \|E_{R,n}\| = \sqrt{2} \|W_R\|_2 \|E_{R,n}\|.$$

Using again the Hoffman-Wielandt inequality we get

$$\delta_2(\lambda(T_{R,n}), \lambda(\tilde{T}_{R,n})) \leq \|\tilde{T}_{R,n} - T_{R,n}\|_F = \left[\frac{1}{n^2} \sum_{i=1}^n W_R(X_i, X_i)^2 \right]^{1/2},$$

and

$$\delta_2(\lambda(\tilde{T}_{R,n}), \lambda(T_n)) \leq \|\tilde{T}_{R,n} - T_n\|_F = \left[\frac{1}{n^2} \sum_{i \neq j} (W - W_R)^2(X_i, X_j) \right]^{1/2}.$$

Now, we invoke Lemmas 8, 9 and 10 to conclude the proof. The proofs of these last two lemmas are provided in Section H.2 and Section H.3 respectively.

Lemma 9. *Let us consider $\gamma > 0$ and assume that $13\tilde{R} \ln(e/\gamma) \leq n$. Then it holds with probability at least $1 - \gamma$*

$$\|E_{R,n}\| \leq 4 \sqrt{\frac{\tilde{R}}{n} \ln(2/\gamma)}.$$

Lemma 10. *Let $R \geq 1$. We have*

$$\frac{1}{n^2} \sum_{i=1}^n W_R(X_i, X_i)^2 = \frac{1}{n} \left(\sum_{k=0}^R p_k^* d_k \right)^2.$$

For any $\gamma \in (0, 1)$ with $\log(e \log(n)/\gamma) \leq (n/(13\tilde{R}))$, it holds with probability at least $1 - \gamma$,

$$\begin{aligned} & \delta_2(\lambda(\mathbb{T}_W), \lambda(T_n)) \\ & \leq \delta_2(\lambda(\mathbb{T}_W), \lambda(\mathbb{T}_{W_R})) + \delta_2(\lambda(\mathbb{T}_{W_R}), \lambda(T_{R,n})) + \delta_2(\lambda(T_{R,n}), \lambda(\tilde{T}_{R,n})) \\ & \quad + \delta_2(\lambda(\tilde{T}_{R,n}), \lambda(T_n)) \\ & \leq 4\sqrt{\frac{\tilde{R}}{n} \ln(2/\gamma)} + \sqrt{2} \left(\sum_{k=0}^R d_k (p_k^*)^2 \right)^{1/2} + \frac{1}{\sqrt{n}} \left| \sum_{k=0}^R p_k^* d_k \right| + 2\|\mathbf{p} - \mathbf{p}_R\|_2 \\ & \quad + M\|\mathbf{p} - \mathbf{p}_R\|_\infty \sqrt{\frac{\log n}{n} (\log(e \log(n)/\gamma))^{1/2}}, \end{aligned}$$

where $M > 0$ depends only on constants related to the Markov chain $(X_i)_{i \geq 1}$. Now remark that

$$\left| \sum_{k=0}^R p_k^* d_k \right| \leq \left(\sum_{k=0}^R d_k \right)^{1/2} \left(\sum_{k=0}^R d_k (p_k^*)^2 \right)^{1/2} = \sqrt{\tilde{R}} \|\mathbf{p}_R\|_2,$$

and that

$$\|\mathbf{p}_R\|_2^2 \leq \|\mathbf{p}\|_2^2 \leq 2, \quad (30)$$

because \mathbf{p}_R is the orthogonal projection of \mathbf{p} , and $|\mathbf{p}| \leq 1$. We deduce that

$$\begin{aligned} & \delta_2(\lambda(\mathbb{T}_W), \lambda(T_n)) \\ & \leq 2\|\mathbf{p} - \mathbf{p}_R\|_2 + 4\sqrt{\frac{\tilde{R}}{n} \ln(2/\gamma)} + \sqrt{\frac{2\tilde{R}}{n}} \\ & \quad + M\|\mathbf{p} - \mathbf{p}_R\|_\infty \sqrt{\frac{\log n}{n} (\log(e \log(n)/\gamma))^{1/2}} \\ & \leq 2\|\mathbf{p} - \mathbf{p}_R\|_2 + 8\sqrt{\frac{\tilde{R}}{n} \ln(e/\gamma)} \\ & \quad + M\|\mathbf{p} - \mathbf{p}_R\|_\infty \sqrt{\frac{\log n}{n} (\log(e \log(n)/\gamma))^{1/2}}. \end{aligned}$$

H.2 Proof of Lemma 9

Observe that $nE_{R,n} = \sum_{i=1}^n (Z_i Z_i^\top - \text{Id}_{\tilde{R}})$ where for all $i \in [n]$, $Z_i \in \mathbb{R}^{\tilde{R}}$ is defined by

$$Z_i := Z(X_i) := (Y_{0,1}(X_i), Y_{1,1}(X_i), Y_{1,2}(X_i), \dots, Y_{1,d_1}(X_i), \dots, Y_{R,1}(X_i), \dots, Y_{R,d_R}(X_i)).$$

By definition of the spectral norm for a Hermitian matrix,

$$\left\| \frac{1}{n} \sum_{i=1}^n Z_i Z_i^\top - \text{Id}_{\tilde{R}} \right\| = \max_{x, \|x\|_2=1} \left| x^\top \left(\frac{1}{n} \sum_{i=1}^n Z_i Z_i^\top \right) x - 1 \right|.$$

We use a covering set argument based on the following Lemma.

Lemma 11. *(cf. [15, Lemma 4.10])*

Let us consider an integer $D \geq 2$. For any $\varepsilon_0 > 0$, there exists a set $Q \subset \mathbb{S}^{D-1}$ of cardinality at most $(1 + 2/\varepsilon_0)^D$ such that

$$\forall \alpha \in \mathbb{S}^{D-1}, \quad \exists q \in Q, \quad \|\alpha - q\|_2 \leq \varepsilon_0.$$

We consider Q the set given by Lemma 11 with $D = d$ and $\varepsilon_0 \in (0, 1/2)$. Let us define $x_0 \in \mathbb{S}^{d-1}$ such that $|x_0^\top E_{R,n} x_0| = \|E_{R,n}\|$ and $q_0 \in Q$ such that $\|x_0 - q_0\|_2 \leq \varepsilon_0$. Then,

$$\begin{aligned} |x_0^\top E_{R,n} x_0| - |q_0^\top E_{R,n} q_0| &\leq |x_0^\top E_{R,n} x_0 - q_0^\top E_{R,n} q_0| \text{ (by triangle inequality)} \\ &= |x_0^\top E_{R,n} (x_0 - q_0) - (q_0 - x_0)^\top E_{R,n} q_0| \\ &\leq \|x_0\|_2 \|E_{R,n}\| \|x_0 - q_0\|_2 + \|q_0 - x_0\|_2 \|E_{R,n}\| \|q_0\|_2 \\ &\leq 2\varepsilon_0 \|E_{R,n}\|. \end{aligned}$$

which leads to

$$|x_0^\top E_{R,n} x_0| = \|E_{R,n}\| \leq |q_0^\top E_{R,n} q_0| + 2\varepsilon_0 \|E_{R,n}\|.$$

Hence,

$$\|E_{R,n}\| \leq \frac{1}{1 - 2\varepsilon_0} \max_{q \in Q} |q^\top E_{R,n} q|.$$

We introduce for any $q \in Q$ the function

$$F_q : x = (x_1, \dots, x_n) \mapsto \frac{1}{n} \sum_{i=1}^n q^\top (Z_i Z_i^\top - 1) q := \frac{1}{n} \sum_{i=1}^n f_q(x_i),$$

where $f_q(x) = q^\top (Z(x)Z(x)^\top - 1)q$.

Let us consider $t > 0$. We want to apply Bernstein's inequality for Markov chains from [18, Theorem 1.1]. In the following, we denote $\mathbb{E}_\pi[\cdot]$ the expectation with respect to the measure π . We remark that $\mathbb{E}_\pi[f_q(X)] = 0$ and that $\|f_q\|_\infty \leq \tilde{R} - 1$. For all $m \in [\tilde{R}]$, we denote $\varphi_m = Y_{r,l}$ with $r \in \{0, \dots, R\}$ and $l \in [d_r]$ such that $m = l + \sum_{i=0}^r d_i - 1$. Then, for any $x \in \mathbb{S}^{d-1}$, and for all $k, l \in [\tilde{R}]$, $((Z(x)^\top Z(x))^2)_{k,l} = \sum_{m=1}^{\tilde{R}} \varphi_l(x) \varphi_m(x)^2 \varphi_k(x) = \tilde{R} \varphi_l(x) \varphi_k(x) = \tilde{R} (Z(x)Z(x)^\top)_{k,l}$ where we used [8, Eq.(1.2.9)]. We deduce that

$$\begin{aligned} \mathbb{E}_\pi[f_q(X)^2] &= \mathbb{E}_\pi[q^\top Z(X)Z(X)^\top q q^\top Z(X)Z(X)^\top q] - 2\mathbb{E}_\pi[q^\top Z(X)Z(X)^\top q] + 1 \\ &= \mathbb{E}_\pi[q^\top \underbrace{(Z(X)Z(X)^\top)^2}_{=\tilde{R}Z(X)Z(X)^\top} q] - 2q^\top \underbrace{\mathbb{E}_\pi[Z(X)Z(X)^\top]}_{=\text{Id}} q + 1 \\ &= \tilde{R} \cdot q^\top \mathbb{E}_\pi[Z(X)Z(X)^\top] q - 1 \\ &= \tilde{R} - 1. \end{aligned}$$

Using that the Markov chain $(X_i)_{i \geq 1}$ has an absolute spectral gap equals to 1 (cf. Section B.3), we get from [18, Eq. (1.6)] that

$$\mathbb{P}(|F_q(X)| \geq t) = \mathbb{P}(|q^\top E_{R,n} q| \geq t) \leq 2 \exp\left(\frac{-nt^2}{4(\tilde{R}-1) + 10(\tilde{R}-1)t}\right),$$

which leads to

$$\begin{aligned} \mathbb{P}\left(\max_{q \in Q} |q^\top E_{R,n} q| \geq t\right) &\leq \mathbb{P}\left(\bigcup_{q \in Q} |q^\top E_{R,n} q| \geq t\right) \\ &\leq 2 \exp\left(\frac{-nt^2/(\tilde{R}-1)}{4 + 10t}\right) (1 + 2/\varepsilon_0)^{\tilde{R}}. \end{aligned}$$

Choosing $\varepsilon_0 = 2\left(\exp\left(\frac{nt^2/2}{(\tilde{R}-1)\tilde{R}(4+10t)}\right) - 1\right)^{-1}$ in order to satisfy $(1+2/\varepsilon_0)^{\tilde{R}} = \exp(nt^2(\tilde{R}-1)^{-1}(4+10t)^{-1}/2)$, we get

$$\mathbb{P}\left(\max_{q \in Q} |q^\top E_{R,n} q| \geq t\right) \leq 2 \exp\left(\frac{-nt^2}{(\tilde{R}-1)(8+20t)}\right).$$

We deduce that if $\frac{25}{2} \ln(2/\alpha) \tilde{R} \leq n$, it holds with probability at least $1 - \alpha$,

$$\max_{q \in Q} |q^\top E_{R,n} q| \leq 16 \sqrt{\frac{\tilde{R}}{n} \ln(2/\alpha)}.$$

Assuming that $200 \ln(7) \tilde{R}^3 \ln(2/\alpha) \leq n^3$ in order to have $1/(1 - 2\varepsilon_0) \leq 4$, it holds with probability at least $1 - \alpha$

$$\|E_{R,n}\| \leq \frac{1}{1 - 2\varepsilon_0} \max_{q \in Q} |q^\top E_{R,n} q| \leq 4 \sqrt{\frac{\tilde{R}}{n} \ln(2/\alpha)}.$$

H.3 Proof of Lemma 10

Reminding that for all $x \in \mathbb{S}^{d-1}$ and for all $k \geq 0$, $\sum_{l=1}^{d_k} Y_{k,l}(x)^2 = d_k$ (cf. Corollary 1.2.7 from [8]), we get

$$\begin{aligned} \frac{1}{n^2} \sum_{i=1}^n W_R(X_i, X_i)^2 &= \frac{1}{n^2} \sum_{i=1}^n \left(\sum_{k=0}^R p_k^* \sum_{l=1}^{d_k} Y_{k,l}(X_i)^2 \right)^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \left(\sum_{k=0}^R p_k^* d_k \right)^2 \\ &= \frac{1}{n} \left(\sum_{k=0}^R p_k^* d_k \right)^2. \end{aligned}$$

I Proof of Theorem 4

Proposition 4 is the counterpart of Proposition 1 in [1] in our dependent framework. This result is the cornerstone of Theorem 4 and is proved in Section I.1.

Proposition 4. *We assume that $\Delta^* > 0$. Let us consider $\gamma > 0$ and define the event*

$$\mathcal{E} := \left\{ \delta_2(\lambda(T_n), \lambda(\mathbb{T}_W)) \vee \frac{2^{\frac{9}{2}} \sqrt{d}}{\Delta^*} \|T_n - \hat{T}_n\| \leq \frac{\Delta^*}{4} \right\}.$$

Then for n large enough,

$$\mathbb{P}(\mathcal{E}) \geq 1 - \gamma/2.$$

Moreover, on the event \mathcal{E} , there exists one and only one set Λ_1 , consisting of d eigenvalues of \hat{T}_n , whose diameter is smaller than $\Delta^*/2$ and whose distance to the rest of the spectrum of \hat{T}_n is at least $\Delta^*/2$. Furthermore, on the event \mathcal{E} , the algorithm HEiC returns the matrix $\hat{G} = \frac{1}{d} \hat{V} \hat{V}^\top$, where \hat{V} has by columns the eigenvectors corresponding to the eigenvalues in Λ_1 .

In the following, we work on the event \mathcal{E} . Let us consider $\gamma \in (0, 1)$.

We choose $R = (n/\log^2 n)^{\frac{1}{2s+d-1}}$. Reminding that W_R is the rank R approximation of W , the Gram matrix associated with the kernel W_R is

$$T_{R,n} = \sum_{k=0}^R p_k^* \sum_{l=1}^{d_k} \Phi_{k,l}(\Phi_{k,l})^\top = X_{R,n} K_R X_{R,n}^\top \in \mathbb{R}^{n \times n}$$

where

$$\begin{aligned} \Phi_{k,l} &= \frac{1}{\sqrt{n}} [Y_{k,l}(X_1), \dots, Y_{k,l}(X_n)] \in \mathbb{R}^n, \\ X_{R,n} &= [\Phi_{0,1}, \Phi_{1,1}, \Phi_{1,2}, \dots, \Phi_{R,d_R}] \in \mathbb{R}^{n \times \tilde{R}} \text{ and} \\ K_R &= \text{Diag}(\lambda_1(\mathbb{T}_W), \dots, \lambda_{\tilde{R}}(\mathbb{T}_W)). \end{aligned}$$

Let us denote now \tilde{V} (resp. \tilde{V}_R) the orthonormal matrix formed by the eigenvectors of the matrix T_n (resp. $T_{R,n}$). We have the following eigenvalue decompositions

$$T_n = \tilde{V} \Lambda \tilde{V}^\top \text{ and } T_{R,n} = \tilde{V}_R \Lambda_R \tilde{V}_R^\top,$$

where $\Lambda = \text{Diag}(\lambda_1, \dots, \lambda_n)$ are the eigenvalues of the matrix T_n and where $\Lambda_R = (p_0^*, p_1^*, \dots, p_1^*, \dots, p_R^*, \dots, p_R^*, 0, \dots, 0) \in \mathbb{R}^n$ where each p_k^* has multiplicity d_k . Then, we note by $V \in \mathbb{R}^{n \times d}$ (resp. V_R) the matrix formed by the columns 1, ..., d of the matrix \tilde{V} (resp. \tilde{V}_R). The matrix $V^* \in \mathbb{R}^{n \times d}$ is the orthonormal matrix with i -th column $\frac{1}{\sqrt{n}}(Y_{1,1}(X_i), \dots, Y_{1,d}(X_i))$. The matrices G^*, G, G_R and G_{proj}^* are defined as follows

$$\begin{aligned} G^* &:= \frac{1}{c_1} V^* (V^*)^\top, & G &:= \frac{1}{c_1} V V^\top \\ G_R &:= \frac{1}{c_1} V_R V_R^\top, & G_{proj}^* &:= V^* ((V^*)^\top V^*)^{-1} (V^*)^\top. \end{aligned}$$

G_{proj}^* is the projection matrix for the columns span of the matrix V^* . Using the triangle inequality we have

$$\|G^* - G\|_F \leq \|G^* - G_{proj}^*\|_F + \|G_{proj}^* - G_R\|_F + \|G_R - G\|_F.$$

Step 1: Bounding $\|G - G_R\|_F$. Since the columns of the matrices V and V_R correspond respectively to the eigenvectors of the matrices T_n and $T_{R,n}$, applying the Davis Kahan sinus Theta Theorem (cf. Theorem 6) gives that there exists $O \in \mathbb{R}^{d \times d}$ such that

$$\|VO - V_R\|_F \leq \frac{2^{3/2} \|T_n - T_{R,n}\|_F}{\Delta},$$

where $\Delta := \min_{k \in \{0,2,3,\dots,R\}} |p_1^* - p_k^*| \geq \Delta^* = \min_{k \in \mathbb{N}, k \neq 1} |p_1^* - p_k^*|$. Using Lemma 12 and $c_1 = \frac{d}{d-2}$, we get that

$$\|G - G_R\|_F = \frac{d-2}{d} \|VO(VO)^\top - V_R V_R^\top\|_F \leq 2 \|VO - V_R\|_F.$$

Hence, using the proof of Theorem 1, we get that with probability at least $1 - 1/n^2$,

$$\|G - G_R\|_F \leq 2 \|VO - V_R\|_F \leq \frac{C}{\Delta^*} \left(\frac{n}{\log^2 n} \right)^{-\frac{s}{2s+d-1}},$$

where $C > 0$ is a constant.

Step 2: Bounding $\|G^* - G_{proj}^*\|_F$. To bound $\|G^* - G_{proj}^*\|_F$, we apply first Lemma 13 with $B = V^*$. This leads to

$$\|G^* - G_{proj}^*\|_F \leq \|\text{Id}_d - (V^*)^\top V^*\|_F \leq \sqrt{d} \|\text{Id}_d - (V^*)^\top V^*\|.$$

Using a proof rigorously analogous to the proof of Lemma 9, it holds with probability at least $1 - \gamma$ and for n large enough,

$$\|\text{Id}_d - (V^*)^\top V^*\| \leq 4 \sqrt{\frac{d \log(e/\gamma)}{n}}.$$

We get by choosing $\gamma = 1/n^2$ that it holds with probability at least $1 - 1/n^2$,

$$\|\text{Id}_d - (V^*)^\top V^*\| \leq C' \sqrt{\frac{d \log(n)}{n}},$$

where $C' > 0$ is a universal constant.

Step 3: Bounding $\|G_{proj}^* - G_R\|_F$. We proceed exactly like in [1] but we provide here the proof for completeness. Since G_{proj}^* and G_R are projectors we have, using for example [5, p.202],

$$\|G_{proj}^* - G_R\|_F = 2 \|G_{proj}^* G_R^\perp\|_F. \quad (31)$$

We use Theorem 7 with $E = G_{proj}^*$, $F = G_R^\perp$, $B = T_{R,n}$ and $A = T_{R,n} + H$ where

$$H = \tilde{X}_{R,n} K_R \tilde{X}_{R,n}^\top - X_{R,n} K_R X_{R,n},$$

where the columns of the matrix $\tilde{X}_{R,n}$ are obtained using a Gram-Schmidt orthonormalization process on the columns of $X_{R,n}$. Hence there exists a matrix L such that $\tilde{X}_{R,n} = X_{R,n}(L^{-1})^\top$. This matrix L is such that a Cholesky decomposition of $X_{R,n}^\top X_{R,n}$ reads as LL^\top .

A and B are symmetric matrices thus we can apply Theorem 7. On the event \mathcal{E} , we can take $S_1 = (\lambda_1 - \frac{\Delta^*}{8}, \lambda_1 + \frac{\Delta^*}{8})$ and $S_2 = \mathbb{R} \setminus (\lambda_1 - \frac{7\Delta^*}{8}, \lambda_1 + \frac{7\Delta^*}{8})$. By Theorem 7 we get

$$\|G_{proj}^* G_R^\perp\|_F \leq \frac{\|A - B\|_F}{\Delta^*} = \frac{\|H\|_F}{\Delta^*}. \quad (32)$$

We only need to bound $\|H\|_F$.

$$\begin{aligned} \|H\|_F &\leq \|L^{-\top} K_R L^{-1} - K_R\|_F \|X_{R,n}^\top X_{R,n}\| \\ &\leq \|K_R\|_F \|L^{-1} L^{-\top} - \text{Id}_{\tilde{R}}\| \|X_{R,n}^\top X_{R,n}\|, \end{aligned} \quad (33)$$

where the last inequality comes from Lemma 14. From the previous remarks on the matrix L , we directly get

$$\|L^{-1} L^{-\top} - \text{Id}_{\tilde{R}}\| = \|(X_{R,n}^\top X_{R,n})^{-1} - \text{Id}_{\tilde{R}}\|.$$

Using the notations of the proof of Theorem 5 which is provided in Section H.1, we get

$$\|L^{-1} L^{-\top} - \text{Id}_{\tilde{R}}\| \|X_{R,n}^\top X_{R,n}\| = \|X_{R,n}^\top X_{R,n} - \text{Id}_{\tilde{R}}\| = \|E_{R,n}\|.$$

Noticing further that $\|K_R\|_F^2 \leq \sum_{k \geq 0} (p_k^*)^2 d_k = \|\mathbf{p}\|_2^2 \leq 2$ (because $|\mathbf{p}| \leq 1$), Eq.(33) becomes

$$\|H\|_F \leq \sqrt{2} \|E_{R,n}\|. \quad (34)$$

Using Lemma 9, it holds with probability at least $1 - \gamma$ and for n large enough,

$$\|E_{R,n}\| \leq 4 \sqrt{\frac{\tilde{R}}{n} \ln(2/\gamma)}. \quad (35)$$

Since $\tilde{R} = \mathcal{O}(R^{d-1})$ and $R = \mathcal{O}\left((n/\log^2 n)^{\frac{1}{2s+d-1}}\right)$, we obtain using Eqs.(31), (32), (34) and (35) that with probability at least $1 - 1/n^2$ it holds

$$\|G_{proj}^* - G_R\|_F = 2\|G_{proj}^* G_R^\perp\|_F \leq \frac{C_d}{\Delta^*} \left(\frac{n}{\log^2(n)}\right)^{\frac{-s}{2s+d-1}},$$

where $C_d > 0$ is a constant that may depend on d and on constants related to the Markov chain $(X_i)_{i \geq 1}$.

Conclusion. We proved that on the event \mathcal{E} , it holds with probability at least $1 - 3/n^2$,

$$\|G^* - G\|_F \leq D_1 \left(\frac{n}{\log^2(n)}\right)^{\frac{-s}{2s+d-1}},$$

where $D_1 > 0$ is a constant that depends on Δ^* , d and on constants related to the Markov chain $(X_i)_{i \geq 1}$. Moreover, Eq.(39) from the proof of Proposition 4 gives that on the event \mathcal{E} , we have

$$\|G - \hat{G}\|_F = \frac{d-2}{d} \|V V^\top - \hat{V} \hat{V}^\top\|_F \leq \frac{2^{\frac{9}{2}} \sqrt{d} \|T_n - \hat{T}_n\|}{3\Delta^*}.$$

Using the concentration result from [4] on spectral norm of centered random matrix with independent entries we get that there exists some constant $D_2 > 0$ such that with probability at least $1 - 1/n^2$ it holds

$$\|G - \hat{G}\|_F \leq D_2 \frac{\sqrt{\log n}}{n}.$$

Using again Proposition 4, we know that for n large enough, $\mathbb{P}(\mathcal{E}) \geq 1 - 1/n^2$. We conclude that for n large enough, it holds with probability at least $1 - 5/n^2$,

$$\|G^* - \hat{G}\|_F \leq D_3 \left(\frac{n}{\log^2(n)} \right)^{\frac{-s}{2s+d-1}},$$

for some constant $D_3 > 0$ that depends on Δ^* , d and on constants related to the Markov chain $(X_i)_{i \geq 1}$ (see Theorem 5 for details).

I.1 Proof of Proposition 4

First part of the proof Let us consider $\gamma > 0$.

Using the concentration of spectral norm for random matrices with independent entries from [4], there exists a universal constant C_0 such that

$$\mathbb{P} \left(\|T_n - \hat{T}_n\| \leq \frac{3\sqrt{2D_0}}{n} + C_0 \frac{\sqrt{\log n/\gamma}}{n} \right) \leq \gamma,$$

where denoting $Y = T_n - \hat{T}_n$, we define $D_0 := \max_{1 \leq i \leq n} \sum_{j=1}^n Y_{i,j} (1 - Y_{i,j})$. We deduce that for n large enough, it holds with probability at least $1 - \gamma/4$,

$$\|T_n - \hat{T}_n\| \leq \frac{(\Delta^*)^2}{2^{\frac{13}{2}} \sqrt{d}}. \quad (36)$$

Using now Theorem 1, it holds with probability at least $1 - \gamma/4$ for n large enough

$$\delta_2(\lambda(T_n), \lambda(\mathbb{T}_W)) \leq C \left(\frac{\log^2 n}{n} \right)^{\frac{s}{2s+d-1}} \leq \frac{\Delta^*}{8}. \quad (37)$$

Putting together Eq.(36) and Eq.(37), we deduce that for n large enough,

$$\mathbb{P}(\mathcal{E}) \geq 1 - \gamma/2.$$

Second part of the proof In the following, we work on the event \mathcal{E} . Since $\Delta^* > 0$ by assumption, we get that $p_1^* = \lambda_1^* = \dots = \lambda_d^*$ is the only eigenvalue of \mathbb{T}_W with multiplicity d . Indeed, all eigenvalue p_k^* with $k > d$ has multiplicity $d_k > d$ and p_0^* has multiplicity 1. Moreover, from Eq.(37), we have that there exists a unique set of d eigenvalues of T_n , denoted $\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_d}$, such that they are at a distance least $3\Delta^*/4$ away from the other eigenvalues, i.e.

$$\Delta := \min_{v_1 \in \lambda(T_n) \setminus \{\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_d}\}} \max_{v_2 \in \{\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_d}\}} |v_1 - v_2| \geq \frac{3\Delta^*}{4}. \quad (38)$$

Let us form the matrix $V \in \mathbb{R}^{n \times d}$ where the k -th column is the eigenvector of T_n associated with the eigenvalue λ_{i_k} . We denote further $G := VV^\top/d$. Let $\hat{V} \in \mathbb{R}^{n \times d}$ be the matrix with columns corresponding to the eigenvectors associated to eigenvalues $\hat{\lambda}_{i_1}, \hat{\lambda}_{i_2}, \dots, \hat{\lambda}_{i_d}$ of \hat{T}_n and $\hat{G} := \hat{V}\hat{V}^\top/d$. Using Theorem 6 there exists some orthonormal matrix $O \in \mathbb{R}^{d \times d}$ such that

$$\|VO - \hat{V}\|_F \leq \frac{2^{\frac{3}{2}} \min\{\sqrt{d}\|T_n - \hat{T}_n\|, \|T_n - \hat{T}_n\|_F\}}{\Delta}.$$

Denoting $\lambda_{i_1}^{sort} \geq \lambda_{i_2}^{sort} \geq \dots \geq \lambda_{i_d}^{sort}$ (resp. $\hat{\lambda}_{i_1}^{sort} \geq \hat{\lambda}_{i_2}^{sort} \geq \dots \geq \hat{\lambda}_{i_d}^{sort}$) the sorted version of the eigenvalues $\lambda_{i_1}, \lambda_{i_2}, \dots, \lambda_{i_d}$ (resp. $\hat{\lambda}_{i_1}, \hat{\lambda}_{i_2}, \dots, \hat{\lambda}_{i_d}$), we have

$$\begin{aligned}
& \left[\sum_{k=1}^d (\lambda_{i_k}^{sort} - \hat{\lambda}_{i_k}^{sort})^2 \right]^{1/2} \\
& \leq \|VV^\top - \hat{V}\hat{V}^\top\|_F \quad (\text{Hoffman-Wielandt inequality [5, Thm VI.4.1]}) \\
& \leq 2\|VO - \hat{V}\|_F \quad (\text{using Lemma 12}) \\
& \leq \frac{2^{\frac{5}{2}} \min\{\sqrt{d}\|T_n - \hat{T}_n\|, \|T_n - \hat{T}_n\|_F\}}{\Delta} \\
& \leq \frac{2^{\frac{9}{2}} \min\{\sqrt{d}\|T_n - \hat{T}_n\|, \|T_n - \hat{T}_n\|_F\}}{3\Delta^*} \quad (\text{using Eq.(38)}) \\
& \leq \Delta^*/8. \quad (\text{using Eq.(36)})
\end{aligned} \tag{39}$$

Using the triangle inequality, we get that

$$\hat{\Delta} := \min_{v_1 \in \lambda(\hat{T}_n) \setminus \{\hat{\lambda}_{i_1}, \hat{\lambda}_{i_2}, \dots, \hat{\lambda}_{i_d}\}} \max_{v_2 \in \{\hat{\lambda}_{i_1}, \hat{\lambda}_{i_2}, \dots, \hat{\lambda}_{i_d}\}} |v_1 - v_2| \geq \frac{\Delta^*}{2}. \tag{40}$$

We proved that on the event \mathcal{E} , the eigenvalues in $\Lambda_1 := \{\hat{\lambda}_{i_1}, \dots, \hat{\lambda}_{i_d}\}$ are at distance at least $\Delta^*/2$ from the other eigenvalues of \hat{T}_n (cf. Eq.(40)) and are at distance at most $\Delta^*/8$ of the eigenvalues $\lambda_{i_1}, \dots, \lambda_{i_d}$ of T_n . We could have done this analysis for different eigenvalues. Let us consider some $k \geq 0$. Eq.(37) shows that on the event \mathcal{E} , there exists a set of d_k eigenvalues of T_n which concentrate around p_k^* and such that it has diameter at most $\Delta^*/4$. Weyl's inequality (cf. [5, p.63]) proves that there exist d_k eigenvalues of \hat{T}_n that are at distance at most $\Delta^*/4$ from p_k^* . If we consider now a subset $L \neq \Lambda_1$ of d eigenvalues of \hat{T}_n , then the previous analysis shows that there exists some eigenvalue $\hat{\lambda}$ of \hat{T}_n which is not in L and that is at distance at most $\Delta^*/4$ from one eigenvalue in L . Using Eq.(38), we deduce that Algorithm (HEiC) returns $\hat{G} = \hat{V}\hat{V}^\top/d$ where the columns of \hat{V} correspond to the eigenvectors of \hat{T}_n associated to the eigenvalues in Λ_1 .

I.2 Useful results

Lemma 12. *Let A, B be two matrices in $\mathbb{R}^{n \times d}$ then*

$$\|AA^\top - BB^\top\|_F \leq (\|A\| + \|B\|)\|A - B\|_F.$$

If $A^\top A = B^\top B = \text{Id}$ then

$$\|AA^\top - BB^\top\|_F \leq 2\|A - B\|_F.$$

Proof of Lemma 12.

$$\begin{aligned}
\|AA^\top - BB^\top\|_F &= \|(A - B)A^\top + B(A^\top - B^\top)\|_F \\
&\leq \|A(A - B)^\top\|_F + \|(B - A)B^\top\|_F \\
&\leq \|(A \otimes \text{Id}_n) \text{vec}(A - B)\|_2 + \|(\text{Id}_d \otimes B) \text{vec}(A - B)^\top\|_2 \\
&\leq (\|A \otimes \text{Id}_n\| + \|\text{Id}_d \otimes B\|)\|A - B\|_F \\
&= (\|A\| + \|B\|)\|A - B\|_F,
\end{aligned}$$

where $\text{vec}(\cdot)$ represents the vectorization of a matrix that is its transformation into a column vector and \otimes is the notation for the Kronecker product between two matrices. \square

Theorem 6. (Davis-Kahan Theorem, cf. [34]) *Let Σ and $\hat{\Sigma}$ be two symmetric $\mathbb{R}^{n \times n}$ matrices with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ and $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_n$ respectively. For $1 \leq r \leq s \leq n$ fixed, we assume that $\min\{\lambda_{r-1} - \lambda_r, \lambda_s - \lambda_{s+1}\} > 0$ where $\lambda_0 := \infty$ and $\lambda_{n+1} = -\infty$. Let $d = s - r + 1$ and V and \hat{V} two*

matrices in $\mathbb{R}^{n \times d}$ with columns $(v_r, v_{r+1}, \dots, v_s)$ and $(\hat{v}_r, \hat{v}_{r+1}, \dots, \hat{v}_s)$ respectively, such that $\Sigma v_j = \lambda_j v_j$ and $\hat{\Sigma} \hat{v}_j = \lambda_j \hat{v}_j$. Then there exists an orthogonal matrix \hat{O} in $\mathbb{R}^{d \times d}$ such that

$$\|\hat{V}\hat{O} - V\|_F \leq \frac{2^{3/2} \min\{\sqrt{d}\|\Sigma - \hat{\Sigma}\|, \|\Sigma - \hat{\Sigma}\|_F\}}{\min\{\lambda_{r-1} - \lambda_r, \lambda_s - \lambda_{s+1}\}}.$$

Lemma 13. Let B be a $n \times d$ matrix with full column rank. Then we have

$$\|BB^\top - B(B^\top B)^{-1}B^\top\|_F = \|\text{Id}_d - B^\top B\|_F.$$

Proof of Lemma 13. Using the cyclic property of the trace, we have

$$\begin{aligned} & \|BB^\top - B(B^\top B)^{-1}B^\top\|_F^2 \\ &= \|B(\text{Id}_d - (B^\top B)^{-1})B^\top\|_F^2 \\ &= \text{Tr}(B(\text{Id}_d - (B^\top B)^{-1})B^\top B(\text{Id}_d - (B^\top B)^{-1})B^\top) \\ &= \text{Tr}(B^\top B(\text{Id}_d - (B^\top B)^{-1})B^\top B(\text{Id}_d - (B^\top B)^{-1})) \\ &= \text{Tr}((B^\top B - \text{Id}_d)(B^\top B - \text{Id}_d)) \\ &= \|\text{Id}_d - B^\top B\|_F^2. \end{aligned}$$

□

Theorem 7. (cf. [5, ThmVII.3.4]) Let A and B be two normal operators and S_1 and S_2 two sets separated by a strip of size δ . Let E be the orthogonal projection matrix of the eigenspaces of A with eigenvalues inside S_1 and F be the orthogonal projection matrix of the eigenspaces of B with eigenvalues inside S_2 . Then

$$\|EF\|_F \leq \frac{1}{\delta} \|E(A-B)F\|_F \leq \frac{1}{\delta} \|A-B\|_F.$$

Lemma 14. (Ostrowski's inequality) Let $A \in \mathbb{R}^{n \times n}$ be a Hermitian matrix and $S \in \mathbb{R}^{d \times n}$ be a general matrix then

$$\|SAS^\top - A\|_F \leq \|A\|_F \times \|S^\top S - \text{Id}_n\|.$$

J Proof of Proposition 1

Notice that for any $i \in [n]$,

$$\mathbb{P}(g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1}) = \mathbb{E}[\mathbb{1}_{g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1}}] = \mathbb{E}\mathbb{E}[\mathbb{1}_{g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1}} \mid \mathbf{D}_{1:n}],$$

and that

$$\begin{aligned} & \mathbb{E}[\mathbb{1}_{g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1}} \mid \mathbf{D}_{1:n}] \\ &= \mathbb{E}[\mathbb{1}_{g_i(\mathbf{D}_{1:n})=1} \mathbb{1}_{A_{i,n+1}=0} \mid \mathbf{D}_{1:n}] + \mathbb{E}[\mathbb{1}_{g_i(\mathbf{D}_{1:n})=0} \mathbb{1}_{A_{i,n+1}=1} \mid \mathbf{D}_{1:n}] \\ &= \eta_i(\mathbf{D}_{1:n}) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=0} + (1 - \eta_i(\mathbf{D}_{1:n})) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=1}, \end{aligned}$$

which leads to

$$\mathbb{P}(g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1}) = \mathbb{E}[\eta_i(\mathbf{D}_{1:n}) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=0} + (1 - \eta_i(\mathbf{D}_{1:n})) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=1}].$$

By definition of the Bayes classifier g^* , we have for any $i \in [n]$,

$$\begin{aligned} & \mathbb{P}(g_i^*(\mathbf{D}_{1:n}) \neq A_{i,n+1}) \\ &= \mathbb{E}[\eta_i(\mathbf{D}_{1:n}) \mathbb{1}_{\eta_i(\mathbf{D}_{1:n}) < \frac{1}{2}} + (1 - \eta_i(\mathbf{D}_{1:n})) \mathbb{1}_{\eta_i(\mathbf{D}_{1:n}) \geq \frac{1}{2}}] \\ &= \mathbb{E}[\min\{\eta_i(\mathbf{D}_{1:n}), 1 - \eta_i(\mathbf{D}_{1:n})\} (\mathbb{1}_{\eta_i(\mathbf{D}_{1:n}) \geq \frac{1}{2}} + \mathbb{1}_{\eta_i(\mathbf{D}_{1:n}) < \frac{1}{2}})] \\ &= \mathbb{E}[\min\{\eta_i(\mathbf{D}_{1:n}), 1 - \eta_i(\mathbf{D}_{1:n})\}] \end{aligned}$$

Given another classifier g , we have for any $i \in [n]$,

$$\begin{aligned}
& \mathbb{P}(g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1}) - \mathbb{P}(g_i^*(\mathbf{D}_{1:n}) \neq A_{i,n+1}) \\
&= \mathbb{E} \left[\eta_i(\mathbf{D}_{1:n}) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=0} + (1 - \eta_i(\mathbf{D}_{1:n})) \mathbb{1}_{g_i(\mathbf{D}_{1:n})=1} \right. \\
&\quad \left. - \left(\eta_i(\mathbf{D}_{1:n}) \mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=0} + (1 - \eta_i(\mathbf{D}_{1:n})) \mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=1} \right) \right] \\
&= \mathbb{E} \left[\eta_i(\mathbf{D}_{1:n}) \left(\mathbb{1}_{g_i(\mathbf{D}_{1:n})=0} - \mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=0} \right) \right. \\
&\quad \left. + (1 - \eta_i(\mathbf{D}_{1:n})) \left(\mathbb{1}_{g_i(\mathbf{D}_{1:n})=1} - \mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=1} \right) \right] \\
&= \mathbb{E} \left[(2\eta_i(\mathbf{D}_{1:n}) - 1) \left(\mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=1} - \mathbb{1}_{g_i(\mathbf{D}_{1:n})=1} \right) \right],
\end{aligned}$$

where we used that $g(\mathbf{D}_{1:n})$ takes only the values 0 and 1, so that

$$\mathbb{1}_{g_i(\mathbf{D}_{1:n})=0} - \mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=0} = \left(\mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=1} - \mathbb{1}_{g_i(\mathbf{D}_{1:n})=1} \right).$$

Since

$$\begin{aligned}
\mathbb{1}_{g_i^*(\mathbf{D}_{1:n})=1} - \mathbb{1}_{g_i(\mathbf{D}_{1:n})=1} &= \begin{cases} 1 & \text{if } g_i^*(\mathbf{D}_{1:n}) = 1 \text{ and } g_i(\mathbf{D}_{1:n}) = 0 \\ 0 & \text{if } g_i^*(\mathbf{D}_{1:n}) = g_i(\mathbf{D}_{1:n}) \\ -1 & \text{if } g_i^*(\mathbf{D}_{1:n}) = 0 \text{ and } g_i(\mathbf{D}_{1:n}) = 1 \end{cases} \\
&= \mathbb{1}_{g_i^*(\mathbf{D}_{1:n}) \neq g_i(\mathbf{D}_{1:n})} \operatorname{sgn}(\eta_i(\mathbf{D}_{1:n}) - 1/2),
\end{aligned}$$

we deduce that

$$\begin{aligned}
& \mathbb{P}(g_i(\mathbf{D}_{1:n}) \neq A_{i,n+1}) - \mathbb{P}(g_i^*(\mathbf{D}_{1:n}) \neq A_{i,n+1}) \\
&= 2\mathbb{E} \left[\left| \eta_i(\mathbf{D}_{1:n}) - \frac{1}{2} \right| \times \mathbb{1}_{g_i(\mathbf{D}_{1:n}) \neq g_i^*(\mathbf{D}_{1:n})} \right],
\end{aligned}$$

which concludes the proof.