

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

---

Yohann De Castro<sup>1</sup>




*joint with:* Quentin Duchemin<sup>2</sup> & Claire Lacour<sup>3</sup>

November 9th, 2022

<sup>1</sup>Institut Camille Jordan (ICJ) & Institut Universitaire de France (IUF)

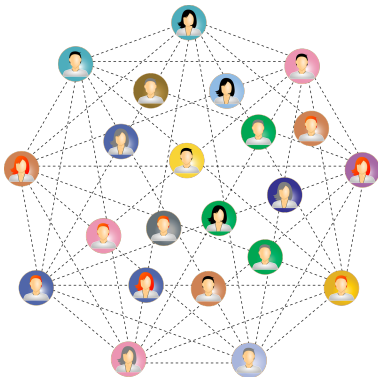
<sup>2</sup>SDSC, EPFL

<sup>3</sup>LAMA, Univ. Gustave Eiffel

- What is common to ...
-  Social networks
  -  Protein-protein interactions
  -  Freight transport between cities
  -  Spread of the virus in a population

?

All these datasets have a **graph structure**.



What can be tackled from graph data?

What can be tackled from graph data?

- \* **Security issues**

*Is the user privacy preserved?*

What can be tackled from graph data?

- \* **Security issues**

*Is the user privacy preserved?*

- \* **Prediction tasks**

*What is the probability for two new users to be connected?*

What can be tackled from graph data?

- \* **Security issues**

*Is the user privacy preserved?*

- \* **Prediction tasks**

*What is the probability for two new users to be connected?*

- \* **Estimation tasks**

*Can we learn the preference of each user?*

What can be tackled from graph data?

- \* **Security issues** → **Graph alignment**

*Is the user privacy preserved?*

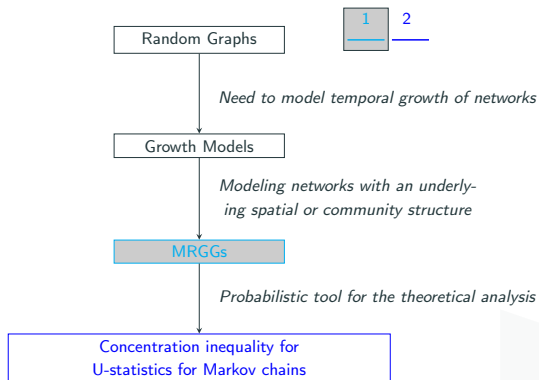
- \* **Prediction tasks** → **Link prediction**

*What is the probability for two new users to be connected?*

- \* **Estimation tasks** → **Community detection**

*Can we learn the preference of each user?*

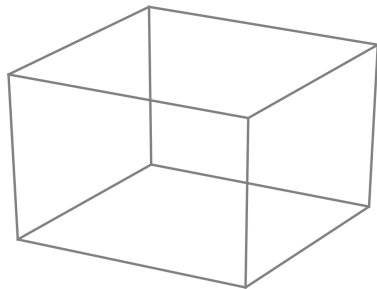
# Outline of the presentation



## I.1 Random Geometric Graphs: the problem of *geometry detection*

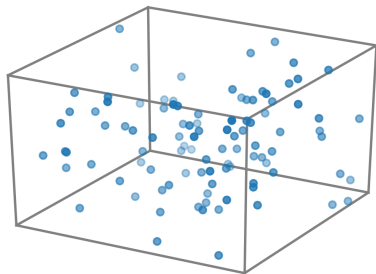
## Presentation of the Random Geometric Graph (RGG)

- 1 Choose a latent space.



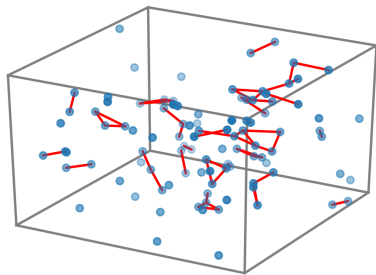
## Presentation of the Random Geometric Graph (RGG)

- ② Sample points.



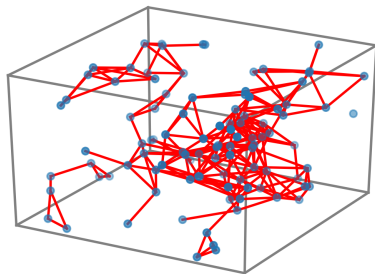
## Presentation of the Random Geometric Graph (RGG)

- 3 Draw edges.



## Presentation of the Random Geometric Graph (RGG)

- 3 Draw edges.



- \* Hard RGGs on  $\mathbb{S}^{d-1}$ :

$$X_i \stackrel{i.i.d.}{\sim} \mathcal{U}(\mathbb{S}^{d-1})$$

$$i \sim j \Leftrightarrow \langle X_i, X_j \rangle \geq t^*$$

---

Y.D.C. & Quentin Duchemin. Random Geometric Graph: Some recent developments and perspectives. *High Dimensional Probability IX, Progress in Probability series*, 2022.

- \* Hard RGGs on  $\mathbb{S}^{d-1}$ :

$$X_i \stackrel{i.i.d.}{\sim} \mathcal{U}(\mathbb{S}^{d-1})$$

$$i \sim j \Leftrightarrow \langle X_i, X_j \rangle \geq t^*$$

- \* Studying the behaviour of Hard RGGs in **the high-dimensional setting** (i.e. when  $d = d(n)$  goes to  $+\infty$  as  $n \rightarrow +\infty$ ):

Given some graph  $G$  with  $n$  points, the *geometry detection problem* reads as

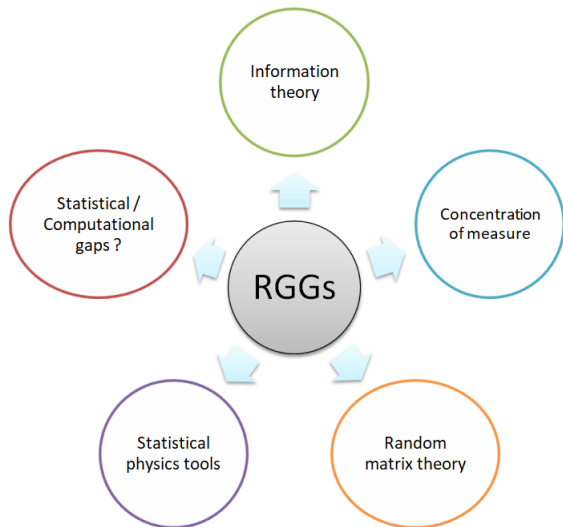
$$\mathbb{H}_0 : "G \sim G(n, p)" \quad \text{VS} \quad \mathbb{H}_1 : "G \sim \text{RGG}(n, p, d)",$$

where  $G(n, p)$  is the Erdős-Renyi distribution and  $\text{RGG}(n, p, d)$  is the distribution of Hard-RGGs on  $\mathbb{S}^{d-1}$  with matching expected degree.

---

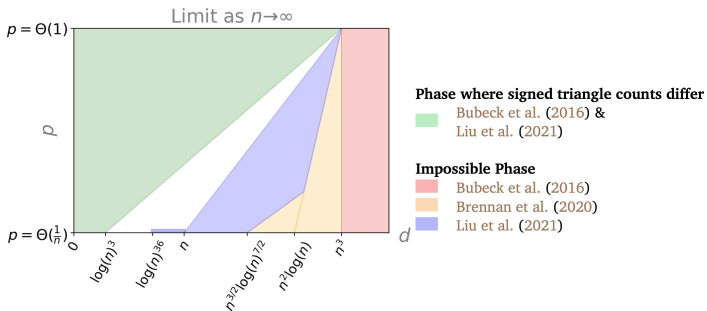
Y.D.C. & Quentin Duchemin. Random Geometric Graph: Some recent developments and perspectives. *High Dimensional Probability IX, Progress in Probability series*, 2022.

# Mathematical Tools to study RGGs in high dimensions



# Phase diagram for the problem of geometry detection

Task	Current state of knowledge	Ref.
Recognizing if a graph can be realized as a RGG	NP-hard	1998
Testing between $G(n, p, d)$ and $G(n, p)$ in high-dimension for $p \in (0, 1)$ fixed	$0$ — Polynomial time test — $n^3$ — Undistinguishable — $d$	2016
Testing between $G(n, \frac{c}{n}, d)$ and $G(n, \frac{c}{n})$ in high-dimension for $c > 0$	$0$ — Polynomial time test — $\log^3 n$ — ? — $\log^{36} n$ — Undistinguishable — $d$	2016 & 2021



- \* First motivated to model **wireless communication systems**.
- \* An active line of research is studying the behaviour of RGGs in **the high-dimensional setting** (i.e. when  $d = d(n)$  goes to  $+\infty$  as  $n \rightarrow +\infty$ ).

- \* First motivated to model **wireless communication systems**.
- \* An active line of research is studying the behaviour of RGGs in **the high-dimensional setting** (i.e. when  $d = d(n)$  goes to  $+\infty$  as  $n \rightarrow +\infty$ ).

In this work, we adopt another point of view and we focus on **non-parametric estimation in a new *growth model* for RGGs**.

## I.2 Non-parametric estimation in RGGs

## Non-parametric estimation in RGGs (1/2)

We work with the Euclidean d-dimensional sphere

$$\mathbb{S}^{d-1} := \{x \in \mathbb{R}^d : \|x\| = 1\}.$$

## Non-parametric estimation in RGGs (1/2)

We work with the Euclidean  $d$ -dimensional sphere

$$\mathbb{S}^{d-1} := \{x \in \mathbb{R}^d : \|x\| = 1\}.$$

The graph is built by

- ✦ sampling  $n$  points  $X_1, \dots, X_n$  on  $\mathbb{S}^{d-1}$ .

## Non-parametric estimation in RGGs (1/2)

We work with the Euclidean  $d$ -dimensional sphere

$$\mathbb{S}^{d-1} := \{x \in \mathbb{R}^d : \|x\| = 1\}.$$

The graph is built by

- \* sampling  $n$  points  $X_1, \dots, X_n$  on  $\mathbb{S}^{d-1}$ .
- \* setting an edge between nodes  $i$  and  $j$  with  $i \neq j$  with probability  $p(\langle X_i, X_j \rangle)$  where  $p : [-1, 1] \rightarrow [0, 1]$  is called the *envelope function*.

## Non-parametric estimation in RGGs (1/2)

We work with the Euclidean  $d$ -dimensional sphere

$$\mathbb{S}^{d-1} := \{x \in \mathbb{R}^d : \|x\| = 1\}.$$

The graph is built by

- \* sampling  $n$  points  $X_1, \dots, X_n$  on  $\mathbb{S}^{d-1}$ .
- \* setting an edge between nodes  $i$  and  $j$  with  $i \neq j$  with probability  $p(\langle X_i, X_j \rangle)$  where  $p : [-1, 1] \rightarrow [0, 1]$  is called the *envelope function*.

The adjacency matrix  $A$  is such that

$$\forall i, j \in [n], \quad \begin{cases} A_{i,j} = 0 & \text{if } i = j \\ A_{i,j} \sim \text{Ber}(p(\langle X_i, X_j \rangle)) & \text{if } i \neq j \end{cases}.$$

This RGG corresponds to a **graphon-type model** where the graphon  $W$  is

$$\forall x, y \in \mathbb{S}^{d-1}, \quad W(x, y) := p(\langle x, y \rangle).$$

Typical example:  $p : t \mapsto \mathbb{1}_{t \geq t^*}$ ,  $t^* \in (-1, 1)$ .

Can we recover  $p$  from the observation of the graph (i.e.  $A$ )?

**Previous work** In some work<sup>1</sup>, the authors address the non-parametric estimation of the envelope function  $p$  in the **independent setting**.

---

<sup>1</sup>Y.D.C., C. Lacour and T.M Pham Ngoc, Adaptive Estimation of Nonparametric Geometric Graphs. *Mathematical Statistics and Learning*, 2020.

Can we recover  $p$  from the observation of the graph (i.e.  $A$ )?

**Previous work** In some work<sup>1</sup>, the authors address the non-parametric estimation of the envelope function  $p$  in the **independent setting**.

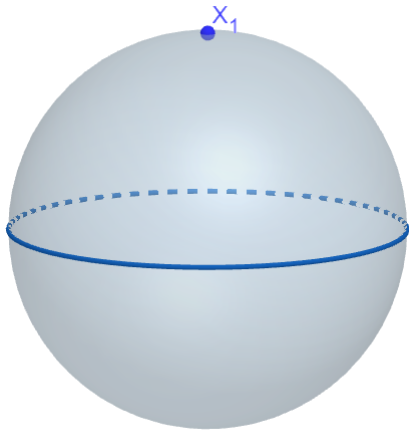
### Goal

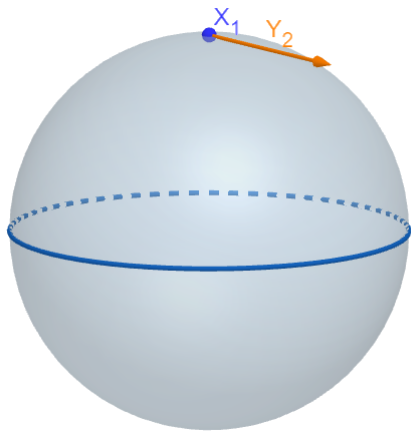
- \* Extend the previous work to a dependent framework.  
↪ We sample the latent space using a Markovian dynamic.
- \* Estimate the Markov transition kernel.  
↪ We propose a heuristic to solve link prediction problems.

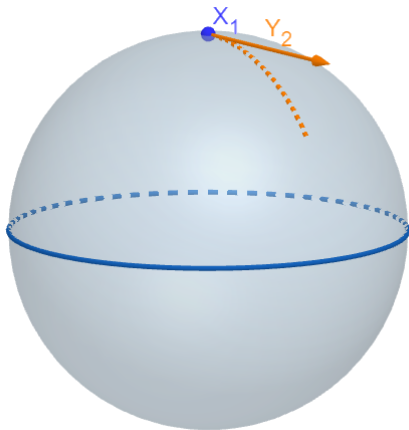
---

<sup>1</sup>Y.D.C., C. Lacour and T.M Pham Ngoc, Adaptive Estimation of Nonparametric Geometric Graphs. *Mathematical Statistics and Learning*, 2020.

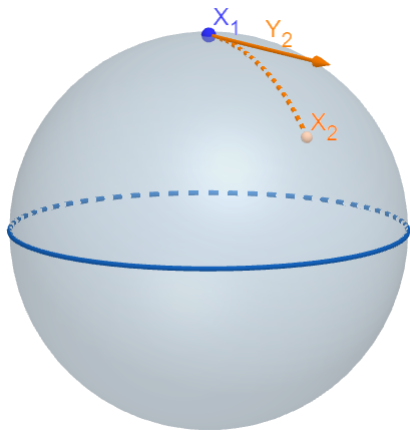
## I.3 The Markovian dynamic

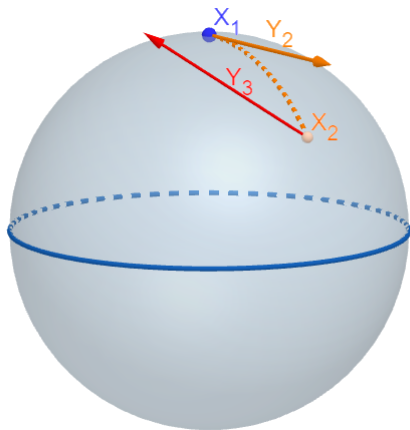




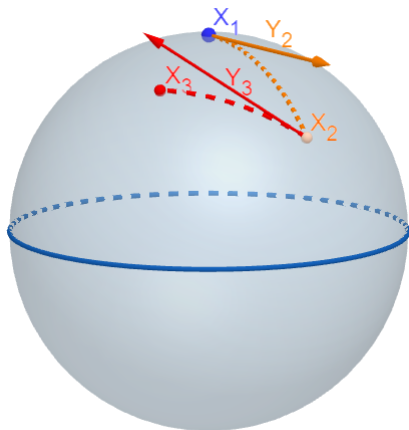


## Latent Markovian dynamic: Isotropic sampling

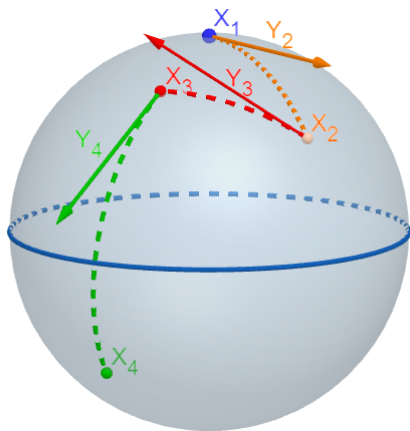




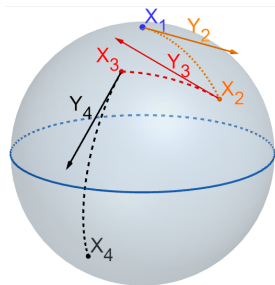
## Latent Markovian dynamic: Isotropic sampling



# Latent Markovian dynamic: Isotropic sampling



Sampling procedure: Isotropic Markovian dynamic

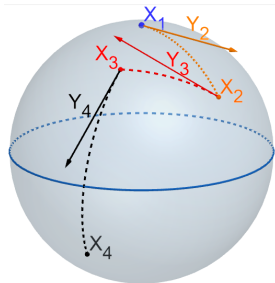


## Sampling procedure: Isotropic Markovian dynamic

- \*  $X_1 \sim \mathcal{U}(\mathbb{S}^{d-1})$
- \*  $\forall i \in \{2, \dots, n\}$ ,
  - \*  $Y_i \sim \mathcal{U}(\mathbb{S}^{d-1}) \mid Y_i \perp X_{i-1}$ .
  - \*  $r_i \sim f_{\mathcal{L}} : [-1, 1] \rightarrow [0, 1]$ .  
 $f_{\mathcal{L}}$  is called the *latitude function*.

Then  $X_i$  is defined by

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

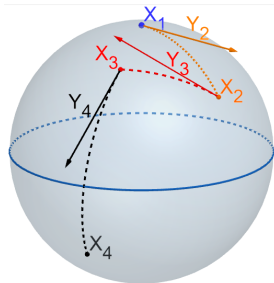


## Sampling procedure: Isotropic Markovian dynamic

- \*  $X_1 \sim \mathcal{U}(\mathbb{S}^{d-1})$
- \*  $\forall i \in \{2, \dots, n\}$ ,
  - \*  $Y_i \sim \mathcal{U}(\mathbb{S}^{d-1}) \mid Y_i \perp X_{i-1}$ .
  - \*  $r_i \sim f_{\mathcal{L}} : [-1, 1] \rightarrow [0, 1]$ .  
 $f_{\mathcal{L}}$  is called the *latitude function*.

Then  $X_i$  is defined by

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



## Construction of the graph

$i \leftrightarrow j$  (i.e.  $A_{i,j} = 1$ ) with proba.  $W(X_i, X_j) := p(\langle X_i, X_j \rangle)$ .

## I.3 Graphon estimation via Harmonic decomposition and Matrix concentration

## Estimation of $\rho$ : A plug-in approach

- \*  $W$  can be viewed as a translation-invariant integral operator.

$$\mathbb{T}_W : f \in L^2(\mathbb{S}^{d-1}) \mapsto \int_{\mathbb{S}^{d-1}} \rho(\langle x, \cdot \rangle) f(x) \sigma(dx) \in L^2(\mathbb{S}^{d-1}).$$

## Estimation of $\rho$ : A plug-in approach

- \*  $W$  can be viewed as a translation-invariant integral operator.

$$\mathbb{T}_W : f \in L^2(\mathbb{S}^{d-1}) \mapsto \int_{\mathbb{S}^{d-1}} p(\langle x, \cdot \rangle) f(x) \sigma(dx) \in L^2(\mathbb{S}^{d-1}).$$

- \*  $\mathbb{T}_W$  is Hilbert-Schmidt.

- Spectrum of  $\mathbb{T}_W$ :  $\lambda^* = \{\rho_0^*, \rho_1^*, \dots, \rho_1^*, \dots, \rho_l^*, \dots, \rho_l^*, \dots\}$ , with known multiplicities that **depend only on the dimension of the sphere**.

- Decomposition of  $p$ :  $p(t) = \sum_{k \geq 0} \rho_k^* \phi_k(t)$  where  $(\phi_k)_k$  **do not depend on** the graphon  $W$ , they are the **Gegenbauer polynomials**.

- $k = 1$  corresponds to linear kernel:  $p(\langle x, y \rangle) = c_0 \rho_0^* + c_1 \rho_1^* \langle x, y \rangle + \dots$  (with  $c_k$  explicit constants) and  $\rho_1^*$  **has multiplicity  $d$**  in  $\lambda^*$ .

## Estimation of $\rho$ : A plug-in approach

- \*  $W$  can be viewed as a translation-invariant integral operator.

$$\mathbb{T}_W : f \in L^2(\mathbb{S}^{d-1}) \mapsto \int_{\mathbb{S}^{d-1}} p(\langle x, \cdot \rangle) f(x) \sigma(dx) \in L^2(\mathbb{S}^{d-1}).$$

- \*  $\mathbb{T}_W$  is Hilbert-Schmidt.

- Spectrum of  $\mathbb{T}_W$ :  $\lambda^* = \{\rho_0^*, \rho_1^*, \dots, \rho_1^*, \dots, \rho_l^*, \dots, \rho_l^*, \dots\}$ , with known multiplicities that **depend only on the dimension of the sphere**.

- Decomposition of  $p$ :  $p(t) = \sum_{k \geq 0} \rho_k^* \phi_k(t)$  where  $(\phi_k)_k$  **do not depend on** the graphon  $W$ , they are the **Gegenbauer polynomials**.

- $k = 1$  corresponds to linear kernel:  $p(\langle x, y \rangle) = c_0 \rho_0^* + c_1 \rho_1^* \langle x, y \rangle + \dots$  (with  $c_k$  explicit constants) and  $\rho_1^*$  **has multiplicity  $d$**  in  $\lambda^*$ .

- \* **Goal:**

Show that  $\lambda(T_n) \rightarrow \lambda^*$  where  $T_n = \frac{1}{n} ((1 - \delta_{i,j}) p(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n}$ .

Show that  $\lambda(\widehat{T}_n) \rightarrow \lambda(T_n)$  where  $\widehat{T}_n = \frac{1}{n} A$ .

## Estimation of $\rho$ : Main result

### Definition (Distance between spectra)

Given two sequences  $x, y$  of reals such that  $\sum_i x_i^2 + y_i^2 < \infty$ , we define

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

where  $S$  is the set of permutations of  $\mathbb{N}$ .



## Estimation of $\rho$ : Main result

### Definition (Distance between spectra)

Given two sequences  $x, y$  of reals such that  $\sum_i x_i^2 + y_i^2 < \infty$ , we define

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

where  $S$  is the set of permutations of  $\mathbb{N}$ .



## Estimation of $\rho$ : Main result

### Definition (Distance between spectra)

Given two sequences  $x, y$  of reals such that  $\sum_i x_i^2 + y_i^2 < \infty$ , we define

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

where  $S$  is the set of permutations of  $\mathbb{N}$ .

- \*  $\delta_2^2(\lambda(\widehat{T}_n), \lambda(T_n)) \rightarrow 0$ ,  $\rightsquigarrow$  follows from one nice work<sup>2</sup>.
- \* It remains to prove that  $\delta_2^2(\lambda(T_n), \lambda^*) \rightarrow 0$ .

---

<sup>2</sup>Bandeira & Van Handel, Sharp nonasymptotic bounds on the norm of random matrices with independent entries. *AoP*, 2016.

## Estimation of $\rho$ : Main result

### Definition (Distance between spectra)

Given two sequences  $x, y$  of reals such that  $\sum_i x_i^2 + y_i^2 < \infty$ , we define

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

where  $S$  is the set of permutations of  $\mathbb{N}$ .

- \*  $\delta_2^2(\lambda(\hat{T}_n), \lambda(T_n)) \rightarrow 0$ ,  $\rightsquigarrow$  follows from one nice work<sup>2</sup>.
- \* It remains to prove that  $\delta_2^2(\lambda(T_n), \lambda^*) \rightarrow 0$ .

### Theorem (EJS: 2022)

If \*  $\rho$  has regularity  $s > 0$

\*  $\|f_{\mathcal{L}}\|_{\infty} < \infty$  and  $\inf_{r \in [-1, 1]} f_{\mathcal{L}}(r) > 0$

then for  $n$  large enough it holds

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

where  $\lambda^{R_{opt}}(\hat{T}_n) = (\hat{\lambda}_1^{sort}, \dots, \hat{\lambda}_{\hat{R}_{opt}}^{sort}, 0, \dots)$  and  $R_{opt} = \lfloor \left( \frac{n}{\log^2(n)} \right)^{\frac{1}{2s+d-1}} \rfloor$ .

<sup>2</sup>Bandeira & Van Handel, Sharp nonasymptotic bounds on the norm of random matrices with independent entries. *AoP*, 2016.

## I.4 Finding eigenvalues clusters

The previous theorem proves that we can recover  $(\rho_k^*)_{k \geq 1}$  up to a permutation! But ...

The previous theorem proves that we can recover  $(\rho_k^*)_{k \geq 1}$  up to a permutation! But ...

- \* ... finding such a permutation is NP-hard.

We develop an algorithm (the SCCHEi) based on a **Hierarchical Agglomerative Clustering** of the eigenvalues  $\lambda(\hat{T}_n)$  to estimate the true partition.

The previous theorem proves that we can recover  $(\rho_k^*)_{k \geq 1}$  up to a permutation! But ...

- \* ... finding such a permutation is NP-hard.

We develop an algorithm (the SCCHEi) based on a **Hierarchical Agglomerative Clustering** of the eigenvalues  $\lambda(\hat{T}_n)$  to estimate the true partition.

- \* ... the optimal resolution level  $R^{opt}$  is unknown.

We use a data-driven choice of model size  $R$  based on the *slope heuristic*.

## I.5 Link prediction

## The characters of the Gold Linear Harmonic Rush

- \* Adjacency matrix observed:  $\widehat{T}_n = \frac{1}{n}A$ ;
- \* Latent positions (unknown):  $(X_i)_{i \geq 1}$ ;

## The characters of the Gold Linear Harmonic Rush

- \* Adjacency matrix observed:  $\widehat{T}_n = \frac{1}{n}A$ ;
- \* Latent positions (unknown):  $(X_i)_{i \geq 1}$ ;
- \* Latent distances (unknown):  $G^* := \frac{1}{n}(\langle X_i, X_j \rangle)_{1 \leq i, j \leq n}$  (has rank  $d$ );
- \* Linear part of the kernel (unknown):  $p(\langle x, y \rangle) = c_0 p_0^* + c_1 p_1^* \langle x, y \rangle + \dots$   
(with  $c_k$  explicit constants) and  $p_1^*$  has multiplicity  $d$  in  $\lambda^*$ ;
- \* The empirical eigenvectors corresponding to the linear Harmonic space (to be found):  $\widehat{V} \in \mathbb{R}^{n \times d}$ ;

## The characters of the Gold Linear Harmonic Rush

- \* Adjacency matrix observed:  $\widehat{T}_n = \frac{1}{n}A$ ;
- \* Latent positions (unknown):  $(X_i)_{i \geq 1}$ ;
- \* Latent distances (unknown):  $G^* := \frac{1}{n}(\langle X_i, X_j \rangle)_{1 \leq i, j \leq n}$  (has rank  $d$ );
- \* Linear part of the kernel (unknown):  $p(\langle x, y \rangle) = c_0 p_0^* + c_1 p_1^* \langle x, y \rangle + \dots$   
(with  $c_k$  explicit constants) and  $p_1^*$  has multiplicity  $d$  in  $\lambda^*$ ;
- \* The empirical eigenvectors corresponding to the linear Harmonic space (to be found):  $\widehat{V} \in \mathbb{R}^{n \times d}$ ;
- \* The key remark (our hope):  $p_1^*$  is the ONLY eigenvalue with multiplicity  $d$ . Note that the multiplicities are pairwise disjoint as soon as  $d \geq 3$  (otherwise we resort to classical Fourier analysis on the torus);

Adjacency matrix observed:  $\widehat{T}_n = \frac{1}{n}A$   
Latent positions unknown:  $(X_i)_{i \geq 1}$

Spectral decomposition of  $\widehat{T}_n$

**HEiC Algorithm (NeurIPS 2019)**

Finds the bulk of  $d$  eigenvalues of  $\widehat{T}_n$  the most separated from the rest of the spectrum

Output of the Algorithm

$d$  eigenvectors of  $\widehat{T}_n$ :  $\widehat{V} \in \mathbb{R}^{n \times d}$

Estimation of the latent distances

$\widehat{G} = \frac{1}{d} \widehat{V} \widehat{V}^T$  estimate of  $G^*$

Estimation of the latitude function  $f_{\mathcal{L}}$

$\widehat{f}_{\mathcal{L}}$ : Kernel density estimate from  $(n\widehat{G}_{i,i+1})_i$

**SCCHei Algorithm (EJS 2022)**

Gives a clustering of  $\lambda(\widehat{T}_n)$  in  $R_{\max} + 1$  groups with sizes  $d_0, \dots, d_{R_{\max}}$

Output of the Algorithm

$\mathcal{C}_{d_0}, \dots, \mathcal{C}_{d_{R_{\max}}}$

Estimation of the eigenvalues of  $\mathbb{T}_W$

$\widehat{p}_k^* = \frac{1}{d_k} \sum_{\lambda \in \mathcal{C}_{d_k}} \lambda$  estimate of  $p_k^*$

Estimation of the envelope function  $p$

$\widehat{p}_{\widehat{R}} = \sum_{k=0}^{\widehat{R}} \widehat{p}_k^* \phi_k$

Link prediction

## From latent distances to link prediction

Denoting  $\text{proj}_{X_n^\perp}(\cdot)$  the orthogonal projection onto  $\text{Span}(X_n)^\perp$ , it holds

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

$\rightsquigarrow$  Latent distances  $D_{1:n} = (\langle X_i, X_j \rangle)_{1 \leq i, j \leq n} \in [-1, 1]^{n \times n}$  are enough for link prediction.

Indeed,  $\forall i \in [n]$ ,

$$\begin{aligned}\eta_i(D_{1:n}) &= \mathbb{P}(A_{i,n+1} = 1 \mid D_{1:n}) \\ &= \int_{r, u \in (-1, 1)} \mathbf{p} \left( \langle X_i, X_n \rangle r + \sqrt{1 - r^2} \sqrt{1 - \langle X_i, X_n \rangle^2} u \right) f_{\mathcal{L}}(r) dr (1 - u^2)^{\frac{d-4}{2}} \frac{du}{b_d},\end{aligned}$$

where  $A_{i,n+1} \in \{0, 1\}$  is one if and only if node  $n + 1$  is connected to node  $i$ .

## Theorem

Assume that

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

Then, we can compute  $\widehat{G} \in \mathbb{R}^{n \times n}$  s.t. with high probability

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

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

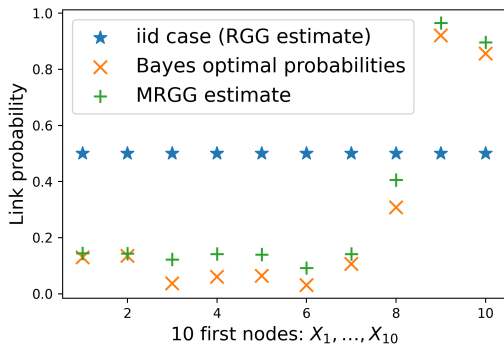
**Heuristic from link prediction** Denoting  $\widehat{r} = n\widehat{G}$ ,

- \* estimate  $p$  with  $\widehat{p} \equiv \sum_{i=1}^{\widehat{R}_{opt}} \widehat{p}_i \phi_i$ .
- \* let  $\widehat{f}_{\mathcal{L}}$  be a kernel density estimator of  $f_{\mathcal{L}}$  based on  $(\widehat{r}_{i,i+1})_{1 \leq i \leq n-1}$ .

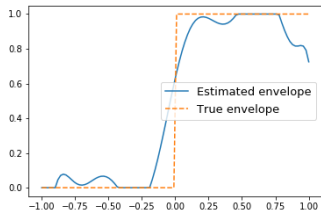
Use the plug-in estimate

$$\widehat{\eta}_i(D_{1:n}) = \int_{r, u \in (-1, 1)} \widehat{p} \left( \widehat{r}_{i,n} r + \sqrt{1-r^2} \sqrt{1-\widehat{r}_{i,n}^2} u \right) \widehat{f}_{\mathcal{L}}(r) dr \frac{(1-u^2)^{\frac{d-4}{2}} du}{b_d},$$

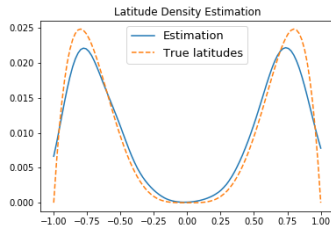
## Link prediction: experiments



**Figure 1:** Link predictions between the future node  $X_{n+1}$  and the 10 first nodes  $X_1, \dots, X_{10}$ . We work with  $n = 2\,000$ ,  $d = 3$ ,  $p(t) = \mathbb{1}_{t \geq \frac{1}{2}}$  and  $f_{\mathcal{L}}(r) = \frac{1}{2}f_{(5,1)}(\frac{r+1}{2})$  where  $f_{(5,1)}$  is the pdf of the Beta distribution with parameter  $(5, 1)$ .



(a) Envelope function



(b) Latitude function

**Figure 2:** Non-parametric estimation of envelope and latitude functions.

Transition towards the concentration toolbox

### Theorem

If \*  $\rho$  has regularity  $s > 0$ ,

\*  $\|f_{\mathcal{L}}\|_{\infty} < \infty$  and  $\inf_{r \in [-1,1]} f_{\mathcal{L}}(r) > 0$ ,

then for  $n$  large enough it holds

$$\mathbb{E}[\delta_2^2(\lambda(T_n), \lambda^*)] \lesssim \left[ \frac{n}{\log^2 n} \right]^{-\frac{2s}{2s+(d-1)}}.$$

## Transition: Proof and use of $U$ -statistic concentration

We recall that  $T_n = \frac{1}{n} ((1 - \delta_{i,j})\rho(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n}$  and we define

$$T_{R,n} = \frac{1}{n} ((1 - \delta_{i,j})\rho_R(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n},$$

where

$$\rho_R(t) := \sum_{k=0}^R \rho_k^* \phi_k(t), \quad R \geq 0.$$

## Transition: Proof and use of $U$ -statistic concentration

We recall that  $T_n = \frac{1}{n} ((1 - \delta_{i,j})p(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n}$  and we define

$$T_{R,n} = \frac{1}{n} ((1 - \delta_{i,j})p_R(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n},$$

where

$$p_R(t) := \sum_{k=0}^R p_k^* \phi_k(t), \quad R \geq 0.$$

Then

$$\delta_2(\lambda(T_n), \lambda^*)^2 \lesssim \delta_2(\lambda(T_n), \lambda(T_{R,n}))^2 + \delta_2(\lambda(T_{R,n}), \lambda^*)^2.$$

## Transition: Proof and use of $U$ -statistic concentration

We recall that  $T_n = \frac{1}{n} ((1 - \delta_{i,j})p(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n}$  and we define

$$T_{R,n} = \frac{1}{n} ((1 - \delta_{i,j})p_R(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n},$$

where

$$p_R(t) := \sum_{k=0}^R p_k^* \phi_k(t), \quad R \geq 0.$$

Then

$$\delta_2(\lambda(T_n), \lambda^*)^2 \lesssim \delta_2(\lambda(T_n), \lambda(T_{R,n}))^2 + \delta_2(\lambda(T_{R,n}), \lambda^*)^2.$$

The first term involves a  $U$ -statistic.

$$\delta_2(\lambda(T_{R,n}), \lambda(T_n))^2 \stackrel{\substack{\leq \\ \text{Hoffman-Wielandt} \\ \text{inequality}}}{\leq} \|T_{R,n} - T_n\|_F^2 = \frac{1}{n^2} \sum_{i \neq j} (p - p_R)^2(\langle X_i, X_j \rangle).$$

## Transition: Proof and use of $U$ -statistic concentration

We recall that  $T_n = \frac{1}{n} ((1 - \delta_{i,j})p(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n}$  and we define

$$T_{R,n} = \frac{1}{n} ((1 - \delta_{i,j})p_R(\langle X_i, X_j \rangle))_{1 \leq i, j \leq n},$$

where

$$p_R(t) := \sum_{k=0}^R p_k^* \phi_k(t), \quad R \geq 0.$$

Then

$$\delta_2(\lambda(T_n), \lambda^*)^2 \lesssim \delta_2(\lambda(T_n), \lambda(T_{R,n}))^2 + \delta_2(\lambda(T_{R,n}), \lambda^*)^2.$$

The first term involves a  $U$ -statistic.

$$\delta_2(\lambda(T_{R,n}), \lambda(T_n))^2 \underbrace{\leq}_{\substack{\text{Hoffman-Wielandt} \\ \text{inequality}}} \|T_{R,n} - T_n\|_F^2 = \frac{1}{n^2} \sum_{i \neq j} (p - p_R)^2(\langle X_i, X_j \rangle).$$

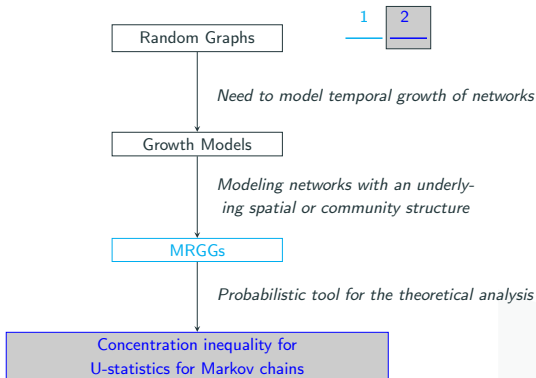
### Definition (U-statistic of order 2)

A  $U$ -statistic of order 2 is a sum of the form

$$\sum_{1 \leq i < j \leq n} h_{i,j}(X_i, X_j),$$

where  $X_1, \dots, X_n$  are r.v. taking values in a measurable space  $(E, \Sigma)$  and where  $h_{i,j} : E^2 \rightarrow \mathbb{R}$ .

## II. Concentration Inequality for U-Statistics of order 2 for Uniformly Ergodic Markov Chains



- ✦ Arcones & Giné 1993

$$\|h_{i,j}\|_{\infty} < \infty \quad \text{and}$$
$$\forall i, j \in [n], \forall x, \mathbb{E}_X [h_{i,j}(x, X)] = \mathbb{E}_X [h_{i,j}(X, x)] = 0.$$

- \* Arcones & Giné 1993

$$\|h_{i,j}\|_{\infty} < \infty \quad \text{and}$$

$$\forall i, j \in [n], \forall x, \mathbb{E}_X [h_{i,j}(x, X)] = \mathbb{E}_X [h_{i,j}(X, x)] = 0.$$

- \* Houdré & Reynaud-Bouret 2002

Improved the constants in the exponential inequality of Arcones & Giné for U-statistics of order two.

- \* Arcones & Giné 1993

$$\|h_{i,j}\|_{\infty} < \infty \quad \text{and}$$
$$\forall i, j \in [n], \forall x, \mathbb{E}_X [h_{i,j}(x, X)] = \mathbb{E}_X [h_{i,j}(X, x)] = 0.$$

- \* Houdré & Reynaud-Bouret 2002

Improved the constants in the exponential inequality of Arcones & Giné for U-statistics of order two.

- \* Giné, Latala & Zinn 2000

Proved that exponential inequality can still be obtained for U-statistics with sufficiently light tails of arbitrary order.

- \* Arcones & Giné 1993

$$\|h_{i,j}\|_{\infty} < \infty \quad \text{and} \\ \forall i, j \in [n], \forall x, \mathbb{E}_X [h_{i,j}(x, X)] = \mathbb{E}_X [h_{i,j}(X, x)] = 0.$$

- \* Houdré & Reynaud-Bouret 2002

Improved the constants in the exponential inequality of Arcones & Giné for U-statistics of order two.

- \* Giné, Latala & Zinn 2000

Proved that exponential inequality can still be obtained for U-statistics with sufficiently light tails of arbitrary order.

- \* Joly & Lugosi in 2016

Working with kernels that have finite  $p$ -th moment for some  $p \in (1, 2]$ , they construct an estimator of the mean of the U-process using the median-of-means technique.

## Assumption 1 (Uniform ergodicity)

*The Markov chain  $(X_i)_{i \geq 1}$  valued in the measurable space  $(E, \Sigma)$  is uniformly ergodic with transition kernel  $P : E^2 \rightarrow \mathbb{R}_+$  and with invariant distribution  $\pi$ .*

## Assumption 1 (Uniform ergodicity)

*The Markov chain  $(X_i)_{i \geq 1}$  valued in the measurable space  $(E, \Sigma)$  is uniformly ergodic with transition kernel  $P : E^2 \rightarrow \mathbb{R}_+$  and with invariant distribution  $\pi$ .*

$\Leftrightarrow \exists \rho \in (0, 1)$  and  $\exists L > 0$  such that

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

## Assumption 1 (Uniform ergodicity)

*The Markov chain  $(X_i)_{i \geq 1}$  valued in the measurable space  $(E, \Sigma)$  is uniformly ergodic with transition kernel  $P : E^2 \rightarrow \mathbb{R}_+$  and with invariant distribution  $\pi$ .*

## Assumption 2 ("Upper-bounded" transition kernel)

$\exists \delta_M > 0, \exists \nu$  proba. meas. on  $E$  s.t.  $\forall x \in E, \forall A \in \Sigma, P(x, A) \leq \delta_M \nu(A)$ .

## Assumption 1 (Uniform ergodicity)

The Markov chain  $(X_i)_{i \geq 1}$  valued in the measurable space  $(E, \Sigma)$  is uniformly ergodic with transition kernel  $P : E^2 \rightarrow \mathbb{R}_+$  and with invariant distribution  $\pi$ .

## Assumption 2 ("Upper-bounded" transition kernel)

$\exists \delta_M > 0, \exists \nu$  proba. meas. on  $E$  s.t.  $\forall x \in E, \forall A \in \Sigma, P(x, A) \leq \delta_M \nu(A)$ .

This holds for

- \* Aperiodic and irreducible Markov chains on finite state space.
- \* AR(1) process with mild conditions.
- \* ARCH process with mild conditions.

## Assumption 1 (Uniform ergodicity)

The Markov chain  $(X_i)_{i \geq 1}$  valued in the measurable space  $(E, \Sigma)$  is uniformly ergodic with transition kernel  $P : E^2 \rightarrow \mathbb{R}_+$  and with invariant distribution  $\pi$ .

## Assumption 2 ("Upper-bounded" transition kernel)

$\exists \delta_M > 0, \exists \nu$  proba. meas. on  $E$  s.t.  $\forall x \in E, \forall A \in \Sigma, P(x, A) \leq \delta_M \nu(A)$ .

## Assumption 3 (Bounded and $\pi$ -canonical kernels)

$\forall i, j, h_{i,j} : (E^2, \Sigma \otimes \Sigma) \rightarrow (\mathbb{R}, \mathcal{B}(\mathbb{R}))$  is measurable, bounded and  $\pi$ -canonical, i.e.

$$\forall x \in E, \mathbb{E}_\pi h_{i,j}(X, x) = \mathbb{E}_\pi h_{i,j}(x, X) = 0.$$

## Theorem (Bernoulli: 2022)

We consider a **stationary** Markov chain satisfying **Assumptions 1, 2 and 3**. Then there exist  $\beta, \kappa > 0$  s.t.  $\forall u > 0$ , it holds with probability at least  $1 - \beta e^{-u} \log n$ ,

$$U_{\text{stat}}(n) := \sum_{1 \leq i < j \leq n} (h_{i,j}(X_i, X_j) - \mathbb{E}[h_{i,j}(X_i, X_j)])$$
$$\leq \kappa \log(n) \left( [C_n + A \log(n) \sqrt{n}] \sqrt{u} + [A + B_n \sqrt{n}] u + [2A \sqrt{n}] u^{3/2} + A [u^2 + \log n] \right).$$

## Theorem (Bernoulli: 2022)

We consider a **stationary** Markov chain satisfying **Assumptions 1, 2 and 3**. Then there exist  $\beta, \kappa > 0$  s.t.  $\forall u > 0$ , it holds with probability at least  $1 - \beta e^{-u} \log n$ ,

$$U_{\text{stat}}(n) := \sum_{1 \leq i < j \leq n} (h_{i,j}(X_i, X_j) - \mathbb{E}[h_{i,j}(X_i, X_j)])$$

$$\leq \kappa \log(n) \left( [C_n + A \log(n) \sqrt{n}] \sqrt{u} + [A + B_n \sqrt{n}] u + [2A \sqrt{n}] u^{3/2} + A [u^2 + \log n] \right).$$

with for some  $t_n$  scaling with  $\log n$ ,

$$A := 2 \max_{i,j} \|h_{i,j}\|_{\infty}, \quad C_n^2 := \sum_{j=2}^n \sum_{i=1}^{j-1} \mathbb{E} [\mathbb{E}_{X' \sim \nu} [h_{i,j}^2(X_i, X')]]$$

$$B_n^2 := \max \left[ \max_{0 \leq k \leq t_n} \max_i \sup_x \sum_{j=i+1}^n \mathbb{E}_{X' \sim \nu} \left( \mathbb{E}_{X \sim P^k(X', \cdot)} h_{i,j}(X, X') \right)^2, \right.$$

$$\left. \max_{0 \leq k \leq t_n} \max_j \sup_y \sum_{i=1}^{j-1} \mathbb{E}_{\tilde{X} \sim \pi} \left( \mathbb{E}_{X \sim P^k(y, \cdot)} h_{i,j}(\tilde{X}, X) \right)^2 \right]$$

## Theorem (Bernoulli: 2022)

We consider a **stationary** Markov chain satisfying **Assumptions 1, 2 and 3**. Then there exist  $\beta, \kappa > 0$  s.t.  $\forall u > 0$ , it holds with probability at least  $1 - \beta e^{-u} \log n$ ,

$$U_{\text{stat}}(n) := \sum_{1 \leq i < j \leq n} (h_{i,j}(X_i, X_j) - \mathbb{E}[h_{i,j}(X_i, X_j)])$$

$$\leq \kappa \log(n) \left( [C_n + A \log(n)\sqrt{n}] \sqrt{u} + [A + B_n\sqrt{n}] u + [2A\sqrt{n}] u^{3/2} + A [u^2 + \log n] \right).$$

with for some  $t_n$  scaling with  $\log n$ ,

$$A := 2 \max_{i,j} \|h_{i,j}\|_{\infty}, \quad C_n^2 := \sum_{j=2}^n \sum_{i=1}^{j-1} \mathbb{E} [\mathbb{E}_{X' \sim \nu} [h_{i,j}^2(X_i, X')]] \leq A^2 n^2,$$

$$B_n^2 := \max \left[ \max_{0 \leq k \leq t_n} \max_i \sup_x \sum_{j=i+1}^n \mathbb{E}_{X' \sim \nu} \left( \mathbb{E}_{X \sim P^k(X', \cdot)} h_{i,j}(X, X') \right)^2, \right.$$

$$\left. \max_{0 \leq k \leq t_n} \max_j \sup_y \sum_{i=1}^{j-1} \mathbb{E}_{\tilde{X} \sim \pi} \left( \mathbb{E}_{X \sim P^k(y, \cdot)} h_{i,j}(\tilde{X}, X) \right)^2 \right] \leq A^2 n.$$

- \* "Constants look pretty ugly" ✗

\* "Constants look pretty ugly" ✗

In the specific case where  $\nu = \pi$  (which includes the independent setting), we get that

$$C_n^2 = \sum_{i < j} \mathbb{E} \left\{ \text{Var}_{\tilde{X} \sim \pi} \left[ h_{i,j}(X_i, \tilde{X}) \mid X_i \right] \right\},$$

and using Jensen inequality that

$$B_n^2 \leq \max \left[ \sup_{x,i} \sum_{j=i+1}^n \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(x, \tilde{X})], \sup_{y,j} \sum_{i=1}^{j-1} \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(\tilde{X}, y)] \right].$$

- \* "Constants look pretty ugly" ✗

In the specific case where  $\nu = \pi$  (which includes the independent setting), we get that

$$C_n^2 = \sum_{i < j} \mathbb{E} \left\{ \text{Var}_{\tilde{X} \sim \pi} \left[ h_{i,j}(X_i, \tilde{X}) \mid X_i \right] \right\},$$

and using Jensen inequality that

$$B_n^2 \leq \max \left[ \sup_{x,i} \sum_{j=i+1}^n \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(x, \tilde{X})], \sup_{y,j} \sum_{i=1}^{j-1} \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(\tilde{X}, y)] \right].$$

- \* "Not too bad but **stationarity is restrictive**"

- \* "Constants look pretty ugly" ✗

In the specific case where  $\nu = \pi$  (which includes the independent setting), we get that

$$C_n^2 = \sum_{i < j} \mathbb{E} \left\{ \text{Var}_{\tilde{X} \sim \pi} \left[ h_{i,j}(X_i, \tilde{X}) \mid X_i \right] \right\},$$

and using Jensen inequality that

$$B_n^2 \leq \max \left[ \sup_{x,i} \sum_{j=i+1}^n \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(x, \tilde{X})], \sup_{y,j} \sum_{i=1}^{j-1} \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(\tilde{X}, y)] \right].$$

- \* "Not too bad but stationarity is restrictive" ✓

- \* "Constants look pretty ugly" ✗

In the specific case where  $\nu = \pi$  (which includes the independent setting), we get that

$$C_n^2 = \sum_{i < j} \mathbb{E} \left\{ \text{Var}_{\tilde{X} \sim \pi} \left[ h_{i,j}(X_i, \tilde{X}) | X_i \right] \right\},$$

and using Jensen inequality that

$$B_n^2 \leq \max \left[ \sup_{x,j} \sum_{i=1}^n \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(x, \tilde{X})], \sup_{y,j} \sum_{i=1}^{j-1} \text{Var}_{\tilde{X} \sim \pi} [h_{i,j}(\tilde{X}, y)] \right].$$

- \* "Not too bad but stationarity is restrictive" ✓

### Theorem (Bernoulli: 2022)

Suppose Assumptions 1, 2 and 3. Assume further that

- \* either  $h_{i,j} \equiv h_{1,j}, \forall i, j$ .
- \* or a mild condition on the initial distribution of the chain.

Then there exist  $\beta, \kappa > 0$  s.t.  $\forall u > 0$ , it holds w.p.  $\geq 1 - \beta e^{-u} \log n$ ,

$$\frac{2}{n(n-1)} U_{\text{stat}}(n) \leq \kappa \max_{i,j} \|h_{i,j}\|_{\infty} \log n \left\{ \frac{u}{n} + \left[ \frac{u}{n} \right]^2 \right\}.$$

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} (\mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)]) + \sum_{i < j} (\mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)])$$

$U_{\text{stat}}^{(k)}(n)$

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} (\mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)]) + \sum_{i < j} (\mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)])$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

## Lemma

Let  $(U_m)_{m \in \mathbb{N}}$  be a martingale such that  $U_0 = U_1 = 0$ .

For each  $m \geq 1$  and  $\ell \geq 2$ , let

$$A_m^\ell = \sum_{i=1}^m \mathbb{E}_{i-1}[(U_i - U_{i-1})^\ell] \quad \text{and} \quad A_1^\ell = 0.$$

If  $A_m^\ell \leq w_m^\ell$  for  $w_m^\ell \geq 0$  then

$$\mathbb{E}[e^{\alpha U_m}] \leq e^{\sum_{\ell \geq 2} \alpha^\ell w_m^\ell / \ell!}.$$

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} (\mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)]) + \sum_{i < j} (\mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)])$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

## Lemma

Let  $(U_m)_{m \in \mathbb{N}}$  be a martingale such that  $U_0 = U_1 = 0$ .

For each  $m \geq 1$  and  $\ell \geq 2$ , let

$$A_m^\ell = \sum_{i=1}^m \mathbb{E}_{i-1}[(U_i - U_{i-1})^\ell] \quad \text{and} \quad A_1^\ell = 0.$$

If  $A_m^\ell \leq w_m^\ell$  for  $w_m^\ell \geq 0$  then

$$\mathbb{E}[e^{\alpha U_m}] \leq e^{\sum_{\ell \geq 2} \alpha^\ell w_m^\ell / \ell!}.$$

$U_{\text{stat}}^{(k)}(n)$  is a **martingale** !

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} \left( \mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)] \right) + \sum_{i < j} \left( \mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)] \right)$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

$$A_n^\ell = \sum_{j=2}^n \mathbb{E}_{j-1} \left[ (U_{\text{stat}}^{(1)}(j) - U_{\text{stat}}^{(1)}(j-1))^\ell \right] \lesssim \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X_j) \right|^\ell$$
$$p_{i,j}(x, z) := h_{i,j}(x, z) - \mathbb{E}_{X' \sim \nu} [h_{i,j}(X_i, X')]$$

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} \left( \mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)] \right) + \sum_{i < j} \left( \mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)] \right)$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

$$A_n^\ell = \sum_{j=2}^n \mathbb{E}_{j-1} \left[ (U_{\text{stat}}^{(1)}(j) - U_{\text{stat}}^{(1)}(j-1))^\ell \right] \lesssim \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X_j) \right|^\ell$$

$$\leq \delta_M \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X'_j) \right|^\ell, \text{ with } p_{i,j}(x, z) := h_{i,j}(x, z) - \mathbb{E}_{X' \sim \nu} [h_{i,j}(X_i, X')]$$

$(X'_j)_j \stackrel{i.i.d.}{\sim} \nu$

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} \left( \mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)] \right) + \sum_{i < j} \left( \mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)] \right)$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

$$A_n^\ell = \sum_{j=2}^n \mathbb{E}_{j-1} \left[ (U_{\text{stat}}^{(1)}(j) - U_{\text{stat}}^{(1)}(j-1))^\ell \right] \lesssim \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X_j) \right|^\ell$$

$$\leq \delta_M \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X'_j) \right|^\ell, \text{ with } p_{i,j}(x, z) := h_{i,j}(x, z) - \mathbb{E}_{X' \sim \nu} [h_{i,j}(X_i, X')]$$

$(X'_j)_j \text{ i.i.d. } \nu$

$$\stackrel{\text{duality lemma}}{\lesssim} \left[ \sup_{(\xi_j)_{j \in \mathcal{L}_\ell}} \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{E}_{j-1} [p_{i,j}(X_i, X'_j) \xi_j(X'_j)] \right]^\ell$$

# Main proof arguments

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} (\mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)]) + \sum_{i < j} (\mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)])$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

$$A_n^\ell = \sum_{j=2}^n \mathbb{E}_{j-1} [(U_{\text{stat}}^{(1)}(j) - U_{\text{stat}}^{(1)}(j-1))^\ell] \lesssim \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X_j) \right|^\ell$$

$$\leq \delta_M \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X'_j) \right|^\ell, \text{ with } p_{i,j}(x, z) := h_{i,j}(x, z) - \mathbb{E}_{X' \sim \nu} [h_{i,j}(X_i, X')]$$

$(X'_j)_j \text{ i.i.d. } \nu$

duality lemma  $\curvearrowright$

$$\lesssim \left[ \sup_{(\xi_j)_j \in \mathcal{L}^\ell} \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{E}_{j-1} [p_{i,j}(X_i, X'_j) \xi_j(X'_j)] \right]^\ell \rightarrow \text{Use of a Talagrand inequality for Markov chains}$$

# Main proof arguments

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} \left( \mathbb{E}_{j-k+1} [h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k} [h_{i,j}(X_i, X_j)] \right) + \sum_{i < j} \left( \mathbb{E}_{j-t_n} [h_{i,j}(X_i, X_j)] - \mathbb{E} [h_{i,j}(X_i, X_j)] \right)$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

$$A_n^\ell = \sum_{j=2}^n \mathbb{E}_{j-1} \left[ (U_{\text{stat}}^{(1)}(j) - U_{\text{stat}}^{(1)}(j-1))^\ell \right] \lesssim \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X_j) \right|^\ell$$

$$\leq \delta_M \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X'_j) \right|^\ell, \text{ with } p_{i,j}(x, z) := h_{i,j}(x, z) - \mathbb{E}_{X' \sim \nu} [h_{i,j}(X_i, X')]$$

$(X'_j)_j \stackrel{i.i.d.}{\sim} \nu$

$$\stackrel{\text{duality lemma}}{\lesssim} \left[ \sup_{(\xi_j)_{j \in \mathcal{L}_\ell} \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{E}_{j-1} [p_{i,j}(X_i, X'_j) \xi_j(X'_j)] \right]^\ell \rightarrow \text{Use of a Talagrand inequality for Markov chains}$$

converse duality lemma

$$\stackrel{\text{w.p. } \geq 1 - e^{-u}}{\lesssim} \sum_{j=2}^n \mathbb{E} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X'_j) \right|^\ell + n(A\ell u)^\ell + B_n^2 (A^2 n \ell u)^{\ell/2}.$$

# Main proof arguments

Telescopic decomposition of  $U_{\text{stat}}(n)$ .

$$U_{\text{stat}}(n) = \sum_{k=1}^{t_n} \sum_{i < j} \left( \mathbb{E}_{j-k+1}[h_{i,j}(X_i, X_j)] - \mathbb{E}_{j-k}[h_{i,j}(X_i, X_j)] \right) + \sum_{i < j} \left( \mathbb{E}_{j-t_n}[h_{i,j}(X_i, X_j)] - \mathbb{E}[h_{i,j}(X_i, X_j)] \right)$$

$U_{\text{stat}}^{(k)}(n)$

Bounding  $U_{\text{stat}}^{(k)}(n)$ .

$$A_n^\ell = \sum_{j=2}^n \mathbb{E}_{j-1} \left[ (U_{\text{stat}}^{(1)}(j) - U_{\text{stat}}^{(1)}(j-1))^\ell \right] \lesssim \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X_j) \right|^\ell$$

$$\leq \delta_M \sum_{j=2}^n \mathbb{E}_{j-1} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X'_j) \right|^\ell, \text{ with } p_{i,j}(x, z) := h_{i,j}(x, z) - \mathbb{E}_{X' \sim \nu}[h_{i,j}(X_i, X')]$$

$(X'_j)_j \text{ i.i.d. } \nu$

$$\stackrel{\text{duality lemma}}{\lesssim} \left[ \sup_{(\xi_j)_j \in \mathcal{L}_\ell} \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbb{E}_{j-1} [p_{i,j}(X_i, X'_j) \xi_j(X'_j)] \right]^\ell \rightarrow \text{Use of a Talagrand inequality for Markov chains}$$

converse duality lemma

$$\stackrel{\text{w.p. } \geq 1 - e^{-u}}{\lesssim} \sum_{j=2}^n \mathbb{E} \left| \sum_{i=1}^{j-1} p_{i,j}(X_i, X'_j) \right|^\ell + n(A\ell u)^\ell + B_n^2(A^2 n \ell u)^{\ell/2}.$$

## RGGs

- \* Extension of the MRGG to **more general Markovian sampling schemes**.
- \* Study of the **robustness** of the estimation methods in MRGGs.
- \* A large set of open questions related to **geometry detection** in high dimensional RGGs.

## Concentration for U-statistics in a dependent framework and applications

- \* Extension to U-statistics of **higher order** ( $m > 2$ ).
- \* Replacing the Markovian assumption with a **mixing condition**.

**Thanks for your attention!**