

## EXPLORATION OF RANDOM GRAPHS BY RESPONDENT DRIVEN SAMPLING (RDS) METHOD

Presented by THUY VO

Co-work with Anthony Cousien (INSERM), Jean-Stéphane Dhersin (Univ Sorbonne Paris Nord) and Viet Chi Tran (Univ Gustave Eiffel)

Biennale SMAI 21-24 Juin 2021

## Motivations

The motivation of this work is to discover the structure and the topology of a hidden network: drug users, MSM,...





Figure 1: Sexual contacts in a population in  $Cuba^1$ 

Figure 2: RDS on the HCV population<sup>2</sup>

- ➡ detect the identities of hidden individuals by exploring the graphs.
- ▶ Proposed methods: Respondent Driven Sampling (RDS)<sup>3</sup>,...

<sup>1.</sup> Clémençon et al. (2015)

<sup>2.</sup> Jauffret-Roustide et al. en cours (2020)

<sup>3.</sup> Respondent Driven Sampling: a new approach to the study of hidden populations; Heckathorn (1997)

## **Respondent Driven Sampling**

There are *c* coupons distributed at each turn of the interview.

- interviewed
- having coupon but have not been interviewed yet
- have been named but without coupon





Step 0



Step 1

Step 3



2

## Mathematical framework

At step  $n \in \mathbb{N}^*$ :

- *n* = # individuals interviewed;
- A<sub>n</sub> = # individuals having coupons but have not been interviewed yet;
- $B_n = \#$  mentioned but have not get any coupon;
- N-  $(n + A_n + B_n) = \#$  vertices not-explored.



 $A_{n+1} = A_n - \mathbf{1}_{\{A_n \ge 1\}} + Z_{n+1} \wedge c,$  $B_{n+1} = B_n + H_{n+1} - Z_{n+1} \wedge c$ 

Let  $X_n = (A_n, B_n)$  be the RDS process to be studied.

When N tends to  $+\infty$ , what is the behavior of  $(X_n)_{n\geq 1}$ ?

## Exploration of sparse random graphs: Erdös-Rényi and SBM (1/2)

The renormalized process:  $X_t^N = \frac{1}{N}(A_{\lfloor Nt \rfloor}, B_{\lfloor Nt \rfloor}).$ 

On the sparse supercritical Erdös-Rényi graphs  $ER(N, \lambda/N)$ , and more general, on the sparse SBM:  $(A^N, B^N) \xrightarrow{(d)} (a, b)$ .



**Figure 3:** Comparing of  $X^N$  with the solution of ODEs  $x_t$  (dashed line) for c = 3,  $\lambda = 2$ .

★ TCL associated to the convergence is established.

## Exploration of a sparse random graphs: Erdös-Rényi and SBM (2/2)

 $\star$  Study the size of graph explored in function of the number of coupons c.



**Figure 4:** Computation of the probability  $\mathbb{P}(\tau > n_0 | A_0 = 1)$  of obtaining a sample of size at least  $n_0$  starting from 1, with c varying from 1 to 10 (colors) and  $n_0$  varying between 1 and 100 (abscissa).

## Random walk: RDS with c=1

★ Denote  $X^{(n)} = (X_1, ..., X_n)$  the explored nodes after *n* steps.



★  $H_n = (V_n, E_n)$  the path of nodes visited by the random walk:  $V_n = \{X_1, ..., X_n\}$  and  $E_n = \bigcup_{i=1}^{n-1} \{X_i, X_{i+1}\}$ 

★  $G_n = G(X^{(n)}, H_n, \kappa)$ : the subgraph discovered.

## Stochastic Block Model<sup>1</sup>



Q blocks (classes)  $\alpha = (\alpha_1, ..., \alpha_Q)$  proportions of blocks  $\pi = (\pi_{qr})_{q,r \in [1,Q]}$  probabilities of connection

★ The observations:

- the random walk  $X^{(n)}$ ;
- the types *Z* = (*Z*<sub>1</sub>, ..., *Z*<sub>n</sub>);
- the adjacency matrix:  $Y = (Y_{ij})_{i,j \in \{1,...,n\}}$ .

**★** The parameter to estimate:  $\theta = (\alpha, \pi)$ .

<sup>1.</sup> The graph of SBM is draw by Julien Chiquet and Catherine Matias.

## The form of explored subgraph - Graphon

Graphon is a symmetric function:  $\kappa : [0,1]^2 \mapsto [0,1]$ .



**\star** Associate to  $G_n$  a graphon  $\kappa$ :

$$\kappa : (x, y) \mapsto \mathbf{1}_{Y_{\lceil nx \rceil, \lceil ny \rceil} = 1}.$$

★ "Infinite" graphs:

- Erdös-Rényi  $\kappa \equiv p$ ;
- SBM( $Q, \alpha, \pi$ ):  $I = (I_1, ..., I_Q)$  a partition of [0,1] such that  $|I_q| = \alpha_q$ . Then If  $x \in I_q$ ,  $y \in I_r$ ,  $\kappa(x, y) = \pi_{qr}$ .

★ For a graph G of n vertices and F a graph of  $k \leq n$  vertices, we define:

$$t(F,G) = \frac{|\operatorname{inj}(F,G)|}{(n)_k}$$

and for the graphon  $\kappa$ :

$$t(F,\kappa) = \int_{[0,1]^k} \prod_{\{i,j\}\in E(F)} \kappa(x_i,x_j) dx_1 \dots dx_k.$$

★ Let  $(F_i)_{i \in \mathbb{N}^*}$  be an enumeration of all the finite graphs. We define:

$$d_{sub}(G,\kappa) = \sum_{i\geq 1} rac{1}{2^i} |t(F_i,G) - t(F_i,\kappa)|$$

**Prop:** When the size of graph  $G_n$  tends to infinity, the SBM graph converges to an SBM graphon for the distance  $d_{sub}$ .

★  $X^{(n)} = (X_1, ..., X_n)$  a RW on  $\kappa$ ,  $X_i \in [0, 1]$  with the transition kernel:

$$P(x, dy) = \frac{\kappa(x, y)dy}{\int_0^1 \kappa(x, v)dv}$$



 $\star$   $X^{(n)}$  admits a stationary measure:

$$m(dx) = \frac{\int_0^1 \kappa(x, v) dv}{\int_0^1 \int_0^1 \kappa(u, v) du \, dv} \, dx.$$

★  $G_n = G(X^{(n)}, H_n, \kappa)$  constructed by  $X^{(n)}$  and the graphon  $\kappa$ .

For the SBM( $Q, \alpha, \pi$ ), the associated graphon is:

$$\kappa(x,y) = \sum_{q=1}^{Q} \sum_{r=1}^{Q} \pi_{qr} \ \mathbf{1}_{l_q}(x) \mathbf{1}_{l_r}(y).$$



#### Prop:

The random walk  $X^{(n)}$  on the graphon  $\kappa$  admits unique invariant measure:

$$m(dx) = \frac{\int_0^1 \kappa(x, v) dv}{\int_0^1 \int_0^1 \kappa(u, v) du \, dv} \, dx = \frac{\sum_{q=1}^Q \left(\sum_{r=1}^Q \pi_{qr} \alpha_r\right) \mathbf{1}_{l_q}(x)}{\sum_{q=1}^Q \sum_{r=1}^Q \pi_{qr} \alpha_q \alpha_r} dx.$$
(1)

## Limit of explored subgraph

## **Proposition**<sup>1</sup>

$$\lim_{n\to\infty} d_{sub}(\mathbf{G}_n,\kappa_{\Gamma^{-1}})=0,\qquad \text{a.s.}$$

where  $\Gamma$  is distribution function of m and  $\Gamma^{-1}$  is the generalized inverse of  $\Gamma$  and

$$\kappa_{\Gamma^{-1}}(x,y) = \kappa(\Gamma^{-1}(x),\Gamma^{-1}(y)).$$



•  $\Gamma^{-1}$  is not known if  $X^{(n)}$  are not observed.

 $\blacktriangleright$  How can we estimate  $\kappa$  from the subgraph  $G_n$ ?.

<sup>1.</sup> Dense graph limits under Respondent Driven Sampling; Athreya and Röllin. Annals of Applied Probability (2016).



**★** Suppose that  $X^{(n)}, Z, Y$  are observed:

 $N_n^q =$  number of nodes of type q $N_n^{q \leftrightarrow r} =$  number of edges of type qr.

For the SBM without bais:

$$\mathcal{L}(Z_{i}, Y_{ij}; i, j \in X^{(n)}; \theta) = \frac{n!}{N_{n}^{1}! \cdots N_{n}^{Q}!} \prod_{q=1}^{Q} \alpha_{q}^{N_{n}^{q}} \times \prod_{\substack{1 \leq i, j \leq n \\ i \neq j}} \pi_{Z_{i}Z_{j}}^{Y_{i,j}} (1 - \pi_{Z_{i}Z_{j}})^{(1 - Y_{i,j})}$$

★Without biases, the classical MLE:

$$\widehat{\alpha}_q^{\text{class}} = \frac{N_n^q}{n}, \qquad \widehat{\pi}_{qr}^{\text{class}} = \frac{N_n^{q \leftrightarrow r}}{N_n^q N_n^r}, \qquad \widehat{\pi}_{qq}^{\text{class}} = \frac{2N_n^{q \leftrightarrow q}}{N_n^q (N_n^q - 1)}.$$

 $\star$  With the biases:

$$\mathcal{L}(Z_{i}, Y_{ij}; i, j \in \boldsymbol{X}^{(n)}; \theta) = \frac{\prod_{i=1}^{n} \alpha_{Z_{i}}}{\prod_{i=1}^{n-1} \sum_{q=1}^{Q} \pi_{Z_{i}q} \alpha_{q}} \times \prod_{\substack{1 \le i, j \le n \\ i \ne j}} \pi_{Z_{i}Z_{j}}^{Y_{ij}} (1 - \pi_{Z_{i}Z_{j}})^{1 - Y_{ij}},$$
(2)

#### Proposition

The ML estimator  $\hat{\theta} = (\hat{\pi}, \hat{\alpha})$  is solution of:

$$\frac{N_n^q}{\widehat{\alpha}_q} - \sum_{p=1}^Q \frac{(N_n^p - \mathbf{1}_{Z_n = p})\widehat{\pi}_{pq}}{\sum_{q'=1}^Q \widehat{\pi}_{pq'}\widehat{\alpha}_{q'}} = \frac{N_n^r}{\widehat{\alpha}_r} - \sum_{p=1}^Q \frac{(N_n^p - \mathbf{1}_{Z_n = p})\widehat{\pi}_{pr}}{\sum_{q'=1}^Q \widehat{\pi}_{pq'}\widehat{\alpha}_{q'}};$$

$$\frac{N_n^{q \leftrightarrow q}}{\widehat{\pi}_{qq}} - \frac{N_n^{q \leftrightarrow q}}{1 - \widehat{\pi}_{qq}} - \frac{(N_n^q - \mathbf{1}_{Z_n = q})\widehat{\alpha}_q}{\sum_{q'=1}^Q \widehat{\pi}_{qq'}\widehat{\alpha}_{q'}} = 0;$$

$$\frac{N_n^{q \leftrightarrow r}}{\widehat{\pi}_{qr}} - \frac{N_n^{q \leftrightarrow r}}{1 - \widehat{\pi}_{qr}} - \frac{(N_n^q - \mathbf{1}_{Z_n = q})\widehat{\alpha}_r}{\sum_{q'=1}^Q \widehat{\pi}_{qq'}\widehat{\alpha}_{q'}} - \frac{(N_n^r - \mathbf{1}_{Z_n = r})\widehat{\alpha}_q}{\sum_{q'=1}^Q \widehat{\pi}_{rq'}\widehat{\alpha}_{q'}} = 0 \quad \text{if } q \neq r.$$

## Method 2: Complete observation + de-biased graphon(1/2)

★ By Athreya & Röllin: 
$$G_n \longrightarrow \kappa_{\Gamma^{-1}}$$
, where  $\kappa_{\Gamma^{-1}} =: \kappa_{\tilde{\theta}}$  and  $\tilde{\theta} := (\tilde{\alpha}, \pi)$ 

The classical estimator for  $\tilde{\alpha}, \pi$  (neglecting the biases):

$$\begin{split} \widehat{\lambda}_{q}^{n} &:= \frac{N_{n}^{q}}{n}; \\ \widehat{\pi}_{qr}^{n} &:= \frac{N_{n}^{q \leftrightarrow r}}{N_{n}^{q} N_{n}^{r}} \quad \text{for} \quad q \neq r \quad \text{and} \quad \widehat{\pi}_{qq}^{n} &:= \frac{2N_{n}^{q \leftrightarrow q}}{N_{n}^{q} (N_{n}^{q} - 1)}. \end{split}$$

 $\star \widehat{\chi}_n(x, y)$  the graphon associated to  $(\widehat{\lambda}^n, \widehat{\pi}^n)$ .

## Proposition

(i) When  $n \to +\infty$ ,

$$\lim_{n \to +\infty} d_{sub}(G_n, \widehat{\chi}_n) = 0.$$
(3)

(ii) The limit  $\hat{\chi}_n$  is then the biased graphon  $\kappa_{\Gamma^{-1}}$ .

$$\lim_{n \to +\infty} d_{\rm sub}(\widehat{\chi}_n, \kappa_{\Gamma^{-1}}) = 0.$$
(4)

★ The 2-stage estimation:

**1st step:** Estimate  $\tilde{\theta} = (\tilde{\alpha}, \pi)$ :

•  $\widehat{\pi}^n$  is a consistent estimator of  $\pi$ :

$$\lim_{n\to+\infty}\widehat{\pi}^n=\pi_{qr},$$

• and  $\widehat{\lambda}_q^n$  is a consistent estimator of  $\widetilde{\alpha}$ :

$$\lim_{n \to +\infty} \widehat{\lambda}_q^n = \Gamma(\sum_{r=1}^q \alpha_r) - \Gamma(\sum_{r=1}^{q-1} \alpha_r) = \widetilde{\alpha}_q.$$

**2nd step:** Correct the estimator  $\widetilde{\theta}$  to obtain  $\theta$ 

A consistent estimator of  $\alpha_q$  is

$$\widehat{\alpha}_{q}^{n} = \Gamma_{n}^{-1} \Big( \sum_{r=1}^{q} \widehat{\lambda}_{r}^{n} \Big) - \Gamma_{n}^{-1} \Big( \sum_{r=1}^{q-1} \widehat{\lambda}_{r}^{n} \Big).$$
(5)

In the case Q = 2, an estimator for  $\alpha_1$  is  $\widehat{\alpha}_1^n = \Gamma_n^{-1}(\widehat{\lambda_1^n})$ .

Suppose that we observe only  $Y_{ij}$  and  $Z_i$  are unknown.

★ The incomplete likelihood:

$$\mathcal{L}(\mathbf{Y}_{ij}; i, j \in \llbracket 1, n \rrbracket; \theta) = \sum_{q_1, \dots, q_n=1}^{Q} \Big[ \prod_{i=1}^{n} \mathbf{1}_{Z_i = q_i} \frac{\prod_{i=1}^{n} \alpha_{q_i}}{\prod_{i=1}^{n-1} \sum_{q=1}^{Q} \pi_{q_i q} \alpha_q} \\ \times \prod_{\substack{1 \le i, j \le n \\ i \neq j}} b(\mathbf{Y}_{ij}, \pi_{q_i q_j}) \Big],$$

➡ The sum of q ∈ {1,.., Q} is not tractable.
➡ Use the SAEM approach the MLE numerically.

## Method 3 (2/2): Incomplete observations + SAEM

Given  $\theta^{(k-1)} = (\alpha^{(k-1)}, \pi^{(k-1)})$ , at the iteration  $k^{\text{eme}}$ :  $\bigstar$  Step 1: Choose the appropriate proposal Z; We follow the variational approach of Daudin et al.<sup>1</sup>: choose  $Z_i$  by a multinomal distribution of parameter  $\tau_{iq}$ ,

$$\tau_{iq} \propto \frac{\alpha_q}{\sum_{\ell=1}^{Q} \pi_{q\ell} \alpha_{\ell}} \prod_{i \neq j} \prod_{\ell=1}^{Q} b(Y_{ij}, \pi_{q\ell})^{\tau_{j\ell}}.$$
(6)

★ Step 2: Stochastic approximation, update the quantity:

$$\mathcal{Q}^{(k)}( heta) = \mathcal{Q}^{(k-1)}( heta) + extsf{s}_k \left(\log \mathcal{L}(Z^{(k)}_i, Y_{ij}, heta) - \mathcal{Q}^{(k-1)}( heta)
ight);$$

★ Step 3: Maximization:

$$\theta^{(k)} := rg\max_{\theta} \mathcal{Q}^{(k)}(\theta).$$

Coupling a stochastic approximation version of EM with an MCMC procedure; Kuhn and Lavielle. ESAIM:ps (2004).

Suppose that  $(Z_1, ..., Z_n)$  are unobserved, but the positions  $(X_1, ..., X_n)$  are observed.

**Step 1**: Neglecting the sampling biases and using the variational EM algorithm (VEM):

- Using EM algorithm to estimate  $(\lambda, \pi)$ ;
- Choosing the types  $Z_i$  based on the information of  $X^{(n)}$ .

**Step 2**: Estimate the cumulative distribution function  $\Gamma_n$ , then deduce the estimator  $\hat{\alpha}^n$  of  $\alpha$  and thus the estimator of  $\kappa$ :

$$\widehat{\kappa}_{n}(x,y) := \sum_{q=1}^{Q} \sum_{r=1}^{Q} \widehat{\pi}_{qr}^{n} \mathbf{1}_{[\sum_{k=1}^{q-1} \widehat{\alpha}_{k}^{n}, \sum_{k=1}^{q} \widehat{\alpha}_{k}^{n}]}(x) \mathbf{1}_{[\sum_{k=1}^{r-1} \widehat{\alpha}_{k}^{n}, \sum_{k=1}^{r} \widehat{\alpha}_{k}^{n}]}(y).$$
(7)

Method 4b: Incomplete observations (Z is unobserved and  $X^{(n)}$  is observed) + graphon de-biasing

When 
$$Z = (Z_1, ..., Z_n)$$
 and  $X^{(n)} = (X_1, ..., X_n)$  are unobserved:

$$\widetilde{lpha}_q = rac{lpha_q \overline{\pi}_q}{\overline{\pi}}, \;\; ext{for all } q \in \{1, \dots Q\} \;\;\; \Leftrightarrow \widetilde{lpha} = rac{lpha \odot (\pi lpha)}{lpha^T \pi lpha},$$

Estimator  $\hat{\alpha}$  for the vector  $\alpha = (\alpha_1, \dots, \alpha_Q)$  can be obtained from solving the equation:

$$(\widehat{\boldsymbol{\alpha}}^T \widehat{\boldsymbol{\pi}} \widehat{\boldsymbol{\alpha}}) \widehat{\lambda} = \widehat{\boldsymbol{\alpha}} \odot (\widehat{\boldsymbol{\pi}} \widehat{\boldsymbol{\alpha}}).$$

It leads to solve the optimization problem

$$\min_{\mathbf{x}\in S} \| (\mathbf{x}^T \widehat{\pi} \mathbf{x}) \widehat{\lambda} - \mathbf{x} \odot (\widehat{\pi} \mathbf{x}) \|,$$

where  $S = \{ \mathbf{x} = (x_1, \cdots, x_Q) \in [0; 1]^Q : x_1 + ... + x_Q = 1 \}.$ 

## Simulations



**Figure 5:** Estimation by the complete data for a graph of n = 60 vertices with Q = 2 classes and parameters  $\alpha_1 = 2/3$ ,  $\pi_{11} = 0.7$ ,  $\pi_{12} = \pi_{21} = 0.4$  and  $\pi_{22} = 0.8$ . 500 such graphs are simulated and the empirical distributions of the estimators are represented here with the true parameters in red line. (a): estimator of  $\alpha$ , (b):estimator of  $\pi_{11}$ .

	Complete	SAEM	De-biased	De-biasing	De-biasing
Parameters	likelihood		graphon	& SAEM	& alg. eq.
$\pi_{11}$	$3.52 \ 10^{-4}$	$5.25 \ 10^{-3}$	$3.52 \ 10^{-4}$	$3.54 \ 10^{-4}$	$3.54 \ 10^{-4}$
$\pi_{12}$	$4.99 \ 10^{-4}$	$5.14 \ 10^{-3}$	$4.99 \ 10^{-4}$	$6.65 \ 10^{-4}$	$4.99 \ 10^{-4}$
$\pi_{22}$	$1.41 \ 10^{-3}$	$1.45 \ 10^{-2}$	$1.41 \ 10^{-3}$	$1.42 \ 10^{-3}$	$1.41 \ 10^{-3}$
$\alpha$	$7.01 \ 10^{-3}$	$3.80 \ 10^{-2}$	$6.80 \ 10^{-4}$	$5.31 \ 10^{-4}$	$4.51 \ 10^{-3}$

 Table 1: Mean square errors.

# Merci de votre attention !