# Uniform estimation of some random graph parameters

Ismaël Castillo

LPMA – Université Paris VI

joint work with Peter Orbanz (Columbia)

Mathematical Methods of Statistics July 10th, 2017

Graph samples with n = 30 nodes



・ロト ・ 日本・ ・ ヨト・・

Graph samples with n = 30 nodes



How does the *information* grow with *n*?

Ismaël Castillo (LPMA – Paris VI)

Luminy - MMS Conference

July 10th, 2017 2 / 24

・ロト ・日子・ ・ヨト・

• Given a probability model for the random graph [next slides]

$$\mathcal{P} = \{P_{\eta}^{(n)}, \ \eta \in H\}$$

Collection of possible distributions for  $(X_{ij})_{1 \le i < j \le n}$ , where

▶  $X_{ij} \in \{0,1\}$  tells whether an edge is present or not between nodes *i* and *j* 

- One may be interested in estimation of
  - the 'parameters'  $\eta$
  - and/or of functionals  $\psi(P_{\eta})$  of those
- Example: edge density  $E_{P_{\eta}}[X_{ij}]$ 
  - Possible estimator

$$\frac{1}{\binom{n}{2}}\sum_{i< j}X_{ij}$$

イロト 不得下 イヨト イヨト

### Random graph models: SBM

The stochastic block model SBM with K classes

- Parameters
  - $p = (p_1, \ldots, p_K)$ •  $Q = (Q_{kl})$  symmetric  $K \times K$  matrix prob. of connection between classes

prob. of classes

- Notation: labels
  - $\varphi: \{1,\ldots,n\} \rightarrow \{1,\ldots,K\}$

assigns a class to each vertex

• Observations:  $(X_{ii})_{i < i}$ 

[the labelling map  $\varphi$  is not observed]

イロト イポト イヨト イヨト

Let  $\pi = p_1 \delta_1 + \ldots + p_K \delta_K$ . The data distribution is

$$egin{array}{lll} (arphi(1),\ldots,arphi(n)) &\sim & \pi^{\otimes n}, \ (X_{ij})_{i < j} \mid arphi &\sim & \bigotimes_{i < j} {\sf Be}(Q_{arphi(i)arphi(j)}) \end{array}$$

# Random graph models: SBM



- *K* = 5
- $p = (p_1, \dots, p_5)$ probabilities of classes
- Q a 5 × 5 matrix of connectivities

イロト イポト イヨト イヨト

Random graph models: graphon model

The graphon model

- Parameter
  - $f: [0,1]^2 \rightarrow [0,1]$  measurable symmetric
- Observations:  $(X_{ij})_{i < j}$  [the design variables  $U_i$  are not observed]

$$(U_i)_i \sim \operatorname{Unif}[0,1]^{\otimes n}$$
  
 $(X_{ij})_{i < j} \mid (U_i)_i \sim \bigotimes_{i < j} \operatorname{Be}(f(U_i, U_j))$ 

• The SBM model as a graphon model



$\frac{1}{2} +  heta$	$\frac{1}{2} -  heta$
$rac{1}{2} -  heta$	$\frac{1}{2} +  heta$

・ロト ・回ト ・ヨト

• The SBM model as a graphon model



• Matrix  $Q = (Q_{ij})$  of SBM model can be read as heights of histogram with partition on [0, 1] given by proportions vector p

[Bickel et al 2013] [here in case  $\mathbb{P}[X_{ij} = 1] =: \rho \sim \text{cst.}$ ]

・ロト ・回ト ・ヨト ・

[Bickel et al 2013] [here in case  $\mathbb{P}[X_{ij} = 1] =: \rho \sim \text{cst.}$ ]

Assume that

- $K = K_0$  is known and fixed
- no two lines of matrix  $Q_0$  are the same and  $\min_k p_{0,k} \neq 0$
- First, what happens if labels  $\varphi(\cdot)$  would be observed?

イロト 不得下 イヨト イヨト

[Bickel et al 2013] [here in case  $\mathbb{P}[X_{ij} = 1] =: \rho \sim \text{cst.}$ ]

Assume that

- $K = K_0$  is known and fixed
- ▶ no two lines of matrix  $Q_0$  are the same and min<sub>k</sub>  $p_{0,k} \neq 0$
- First, what happens if labels  $\varphi(\cdot)$  would be observed?
  - the MLE  $(\tilde{p}^{ML}, \tilde{Q}^{ML})$  is asymptotically normal

$$egin{aligned} &\sqrt{n}( ilde{p}^{ML}-p_0) \mid arphi 
ightarrow \mathcal{N}(0,T_1) \ &n( ilde{Q}^{ML}-Q_0) \mid arphi 
ightarrow \mathcal{N}(0,T_2) \end{aligned}$$

$$\tilde{p}_{a}^{ML} = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}_{\varphi(i)=a}, \qquad \tilde{Q}_{ab}^{ML} = \frac{\sum_{i < j} X_{ij} \mathbb{1}_{\varphi(i)=a, \varphi(j)=b}}{\sum_{i < j} \mathbb{1}_{\varphi(i)=a, \varphi(j)=b}}$$

[Bickel et al 2013]

- Theorem. If labels are unknown, under the previous assumptions
  - the MLE  $(\hat{p}^{ML}, \hat{Q}^{ML})$  is asymptotically normal

$$\begin{aligned} &\sqrt{n}(\hat{p}^{ML*}-p_0) \rightarrow \mathcal{N}(0,T_1) \\ &n(\hat{Q}^{ML*}-Q_0) \rightarrow \mathcal{N}(0,T_2) \end{aligned}$$

where  $(p^*, Q^*)$  denotes a label switched-version of (p, Q)

• • • • • • • • • • • •

[Bickel et al 2013]

- Theorem. If labels are unknown, under the previous assumptions
  - the MLE  $(\hat{p}^{ML}, \hat{Q}^{ML})$  is asymptotically normal

$$\sqrt{n}(\hat{p}^{ML*} - p_0) \rightarrow \mathcal{N}(0, T_1)$$
$$n(\hat{Q}^{ML*} - Q_0) \rightarrow \mathcal{N}(0, T_2)$$

where  $(p^*, Q^*)$  denotes a label switched-version of (p, Q)

- Pointwise inference is asymptotically equivalent to given  $\varphi$  case
  - ▶ Idea : ML 'profiles out' the unknown  $\varphi$ → proof based on  $\hat{\varphi}$  that consistently estimates  $\varphi$  asymptotically

[Bickel et al 2013]

- Theorem. If labels are unknown, under the previous assumptions
  - the MLE  $(\hat{p}^{ML}, \hat{Q}^{ML})$  is asymptotically normal

$$\sqrt{n}(\hat{p}^{ML*} - p_0) \rightarrow \mathcal{N}(0, T_1)$$
  
 $n(\hat{Q}^{ML*} - Q_0) \rightarrow \mathcal{N}(0, T_2)$ 

where  $(p^*, Q^*)$  denotes a label switched-version of (p, Q)

- Pointwise inference is asymptotically equivalent to given  $\varphi$  case
  - Idea : ML 'profiles out' the unknown φ → proof based on φ̂ that consistently estimates φ asymptotically
- Note that
  - ▶  $p_0$  estimated at 'slow' rate  $\frac{1}{\sqrt{n}}$  →  $\frac{1}{n}$  in terms of quadratic risk
  - $Q_0$  estimated at 'fast' rate  $\frac{1}{n} \rightarrow \frac{1}{n^2}$  in terms of quadratic risk

イロト 不得下 イヨト イヨト

#### Some questions

The previous results are *asymptotic* and *pointwise* at  $(p_0, Q_0)$ 

Some questions

- *uniform* estimation of parameters?
  - $\rightarrow$  say of the connectivity parameters = the matrix Q
- in practice *n* and *K* are free
  - $\rightarrow$  non-asymptotic results where K is possibly 'large'?

Framework : SBM (to start with)

• • • • • • • • • • •

# Related topics

- Testing and 'community' detection
   [Arias-Castro, Candès and Durand 2011]
   [Butucea and Ingster 2013]
   [Arias-Castro and Verzelen 2014-15]
- Estimation of the graphon function f [wrt squared L<sup>2</sup>-type risk]
   [Olhede and Wolfe 2013]
   [Gao, Lu and Zhou 2015]
   [Klopp, Tsybakov and Verzelen 2015]
- Other random graph models 'sparsified' graphon model, preferential attachment, graphex model, ...

イロト イポト イヨト イヨト

For simplicity assume equiproportions  $p_1 = p_2 = \frac{1}{2}$ 

• Consider the SBM submodel  $p = [\frac{1}{2}, \frac{1}{2}], Q = Q^{\theta}, K = 2$ , with

$$Q^{ heta} = egin{bmatrix} rac{1}{2}+ heta & rac{1}{2}- heta\ rac{1}{2}- heta & rac{1}{2}+ heta\end{bmatrix}$$

For simplicity assume equiproportions  $p_1 = p_2 = \frac{1}{2}$ 

• Consider the SBM submodel  $p = [\frac{1}{2}, \frac{1}{2}], Q = Q^{\theta}, K = 2$ , with

$$Q^{\theta} = \begin{bmatrix} \frac{1}{2} + \theta & \frac{1}{2} - \theta \\ \frac{1}{2} - \theta & \frac{1}{2} + \theta \end{bmatrix}$$

• Theorem 1 [minimax lower bound]. For some  $c_0 > 0$ ,  $\inf_{\hat{\theta}} \sup_{|\theta| \le \frac{1}{2}} E_{\theta} \left[ (\hat{\theta} - \theta)^2 \right] \ge \frac{c_0}{n}.$ 

For simplicity assume equiproportions  $p_1 = p_2 = \frac{1}{2}$ 

• Consider the SBM submodel  $p = [\frac{1}{2}, \frac{1}{2}], Q = Q^{\theta}, K = 2$ , with

$$Q^{ heta} = egin{bmatrix} rac{1}{2}+ heta & rac{1}{2}- heta\ rac{1}{2}- heta & rac{1}{2}+ heta\end{bmatrix}$$

• Theorem 1 [minimax lower bound]. For some  $c_0 > 0$ ,

$$\inf_{\hat{\theta}} \sup_{|\theta| < \frac{\delta_0}{n}} E_{\theta} \left[ (\hat{\theta} - \theta)^2 \right] \geq \frac{c_0}{n}.$$

For simplicity assume equiproportions  $p_1 = p_2 = \frac{1}{2}$ 

• Consider the SBM submodel  $p = [\frac{1}{2}, \frac{1}{2}], Q = Q^{\theta}, K = 2$ , with

$$Q^{\theta} = \begin{bmatrix} \frac{1}{2} + \theta & \frac{1}{2} - \theta \\ \frac{1}{2} - \theta & \frac{1}{2} + \theta \end{bmatrix}$$

• Theorem 1 [minimax lower bound]. For some  $c_0 > 0$ ,

$$\inf_{\hat{\theta}} \sup_{|\theta| < \frac{\delta_0}{n}} E_{\theta} \left[ (\hat{\theta} - \theta)^2 \right] \geq \frac{c_0}{n}.$$

• minimax local lower bound, relevant not only at  $\theta = 0$ , but also locally

can show that if 
$$t = \frac{\delta_0}{2n}$$
,  $\inf_{\hat{\theta}} \sup_{|\theta - t| < \frac{\delta_0}{2n}} E_{\theta} \left[ (\hat{\theta} - \theta)^2 \right] \ge \frac{c_0}{n}$ .

• Consider the submodel

$$Q^{ heta} = egin{bmatrix} rac{1}{2} + heta & rac{1}{2} - heta \ rac{1}{2} - heta & rac{1}{2} + heta \end{bmatrix}$$

One observes SBM-data with  $p = [\frac{1}{2}, \frac{1}{2}]$ ,  $Q = Q^{\theta}$ , K = 2

- $\bullet~ {\rm Let}~ \hat{\theta}^{\rm ML}$  be the MLE in the above submodel
- Theorem 1 (bis) [minimax upper bound]. For some  $d_0 > 0$ ,

$$\sup_{|\theta| \le \frac{1}{2}} E_{\theta} \left[ (\hat{\theta}^{ML} - \theta)^2 \right] \le \frac{d_0}{n}$$

Proof: profile (pseudo)-maximum likelihood

Main result for K = k classes, motivation

An example with n = 30, K = 5 and a 'difficult' matrix Q

$$\begin{bmatrix} .55 & .45 & .4 & .1 & .7 \\ .45 & .55 & .4 & .1 & .7 \\ .4 & .4 & .6 & .2 & .1 \\ .1 & .1 & .2 & .1 & .4 \\ .7 & .7 & .4 & .4 & .2 \end{bmatrix}$$

A B > A B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B >
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 B
 A
 A

Main result for K = k classes, motivation

An example with n = 30, K = 5 and a 'difficult' matrix Q

.55	.45	.4	.1	.7
.45	.55	.4	.1	.7
.4	.4	.6	.2	.1
.1	.1	.2	.1	.4
.7	.7	.4	.4	.2

Local uniform estimation rate of elements of this submatrix will be 'slow'

イロト イポト イヨト イ

Suppose equiproportions for simplicity [balanced proportions would work] Let us consider estimation along the submodel

$$\begin{bmatrix} a_0 & a_1 & \dots & a_{k-2} \\ a_1 & & & \\ \vdots & B & \\ a_{k-2} & & \end{bmatrix}$$

Suppose equiproportions for simplicity [balanced proportions would work] Let us consider estimation along the submodel

$$\begin{bmatrix} a_0 & a_1 & \dots & a_{k-2} \\ a_1 & & & \\ \vdots & B & & \\ a_{k-2} & & & \end{bmatrix} \rightarrow \begin{bmatrix} a_0 + \theta & a_0 - \theta & a_1 & \dots & a_{k-2} \\ a_0 - \theta & a_0 + \theta & a_1 & \dots & a_{k-2} \\ a_1 & a_1 & & & \\ \vdots & \vdots & B & & \\ a_{k-2} & a_{k-2} & & & \end{bmatrix} = Q^{\theta}$$

Set  $A = [a_0 \ a_1 \ \cdots \ a_{k-2}]$  and  $B = (b_{ij})$ 

Suppose proportions are equidistributed [extends to 'balanced proportions'] • main message of Theorem 2 below

Around *any* point in the interior of the K = k - 1 classes model,

there is a direction coming from the K = k classes model such that

rate of estimation of matrix parameter along this direction no better than  $\frac{k}{n}$ 

Suppose proportions are equidistributed [extends to 'balanced proportions'] • main message of Theorem 2 below

Around any point in the interior of the K = k - 1 classes model,

there is a direction coming from the K = k classes model such that

rate of estimation of matrix parameter along this direction no better than  $\frac{k}{n}$ 

• One observes SBM-data with law  $E_{\theta}$  specified by

$$K = k, \ p = [\frac{1}{K}, \dots, \frac{1}{K}], \ Q = Q^{\theta}$$
 the previous  $K \times K$  matrix

• Theorem 2 [minimax lower bound]. For some  $c_1 > 0$ , for any A and B,

$$\inf_{\hat{\theta}} \sup_{|\theta| \leq \frac{1}{2}} E_{\theta} \left[ (\hat{\theta} - \theta)^2 \right] \geq c_1 \frac{k}{n}.$$

Comments on the lower bound

- The bound is minimax *local*
- Idea of proof of lower bound

a 'mixture vs mixture' lower bound argument

more involved than before, as 'null hypothesis' is a mixture

Implications

- The 'boundary' of SBM model with K = k moves with k and n
- If one is interested in estimating (some of the) heights = elements of the Q matrix, one should take this moving boundary into account to determine precisely the accuracy of estimation
- The rate along constructed submodel deteriorates with *k*.

The lower bound is non-asymptotic.

For many k, n, the boundary area, of size at least  $\frac{k}{n}$ , can be 'large'

Upper-bound result

Recall 
$$Q^{\theta} = \begin{bmatrix} a_0 + \theta & a_0 - \theta & a_1 & \dots & a_{k-2} \\ a_0 - \theta & a_0 + \theta & a_1 & \dots & a_{k-2} \\ a_1 & a_1 & & & \\ \vdots & \vdots & B = (b_{ij}) \\ a_{k-2} & a_{k-2} & & \end{bmatrix}$$

• Require, for  $\mathcal{C}:=\{a_i,b_{ij},\ 1\leq i,j\leq k-2\}$ ,

$$\begin{split} \min_{c \in \mathcal{C}} \left\{ |c - a_0| \right\} &\geq \kappa > 0 \qquad (C) \\ k^3 \log k \lesssim \kappa^4 n \qquad (D) \end{split}$$

・ロト ・回ト ・ヨト ・

Upper-bound result

Recall 
$$Q^{\theta} = \begin{bmatrix} a_0 + \theta & a_0 - \theta & a_1 & \dots & a_{k-2} \\ a_0 - \theta & a_0 + \theta & a_1 & \dots & a_{k-2} \\ a_1 & a_1 & & & \\ \vdots & \vdots & B = (b_{ij}) \\ a_{k-2} & a_{k-2} & & & \end{bmatrix}$$

• Require, for  $\mathcal{C}:=\{a_i,b_{ij},\ 1\leq i,j\leq k-2\}$ ,

$$\begin{split} \min_{c \in \mathcal{C}} \left\{ |c - a_0| \right\} &\geq \kappa > 0 \qquad (C) \\ k^3 \log k \lesssim \kappa^4 n \qquad (D) \end{split}$$

Theorem 3 [upper bound]. For some C<sub>1</sub> > 0, for any A, B such that conditions (C)–(D) hold, for θ̂ a profile-MLE estimate,

$$\sup_{|\theta|\leq\kappa} E_{\theta}\left[(\hat{\theta}-\theta)^2\right]\leq C_1\frac{k}{n}.$$

Upper-bound result – checking (C) and (D)

Conditions (C) and (D) are satisfied in the following settings

• Example 1 (well-separated block)

If  $\kappa$  is a given positive constant e.g.  $\kappa = 1/4$ , then (C) means that the coefficients of A, B are fairly different from  $a_0$ 

• Example 2 (randomly sampled matrices A, B)

If the vector defining A and the upper half of the matrix B are sampled iid  $\mathcal{U}[0,1]$ , then (C)-(D) is satisfied with high probability if

$$k \lesssim n^{1/7}$$

#### Dependence on the 'environment' A and B

Consider a 'least-favorable case"

$$Z^{\theta} = \begin{bmatrix} \frac{1}{2} + \theta & \frac{1}{2} - \theta & \frac{1}{2} & \cdots & \frac{1}{2} \\ \frac{1}{2} - \theta & \frac{1}{2} + \theta & \frac{1}{2} & \cdots & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \cdots & \frac{1}{2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \cdots & \frac{1}{2} \end{bmatrix}.$$

[or a perturbation thereof]

Observe data from the corresponding SBM model with equiproportions

• Theorem 4. There exists  $c_2 > 0$  such that

$$\inf_{\hat{\theta}} \sup_{\theta \in (-1/2, 1/2)} E_{\theta} (\hat{\theta} - \theta)^2 \geq c_2 \frac{k^2}{n}$$

### Functionals of graphon model

In the more general graphon model, consider the functional

$$\psi(\langle f \rangle) = \left[ \int_{[0,1]^2} \left( f(x,y) - \int_{[0,1]^2} f \right)^2 dx dy \right]^{1/2}$$

[Continuous 'graphon-analogue' of previous parameter  $\theta$ ]

• Theorem 5. For  $\mathcal{P}$  a  $\mathcal{C}^1(M)$ -class of graphons, for some  $c_3, c_4 > 0$ ,

$$\frac{c_3}{n} \leq \inf_{\hat{\psi}} \sup_{f \in \mathcal{P}} E_f\left[(\hat{\psi} - \psi(f))^2\right] \leq \frac{c_4}{n}.$$

Ismaël Castillo (LPMA – Paris VI)

• • • • • • • • • • • •

### Functionals of graphon model

• Lower bound also holds for other functionals such as

$$\psi(P_w) = \int_{[0,1]^2} \left| f(x,y) - \int_{[0,1]^2} |f(x,y)| dx dy \right| dx dy$$

• The quadratic 'slow rate' 1/n [or  $1/\sqrt{n}$  non-quadratic] appears to be quite universal for uniform estimation of many functionals

# Conclusion

Previously known results for SBMs (K, p, Q): given fixed K = k,

- $\rightarrow$  proportions p estimated at 'slow' rate  $\frac{1}{p}$  asymptotically
- $\rightarrow$  connectivities Q estimated at 'fast' rate  $\frac{1}{n^2}$  asymptotically, pointwise, in the interior of set of SBMs with k classes

Conclusions

• Fast rate for connectivities Q not achievable uniformly

Uniform rates are at most k/n [=generic lower bound if k not too large] and can be as slow as  $k^2/n$ , from a non-asymptotic perspective in n and k

• This phenomenom happens close to any k-1 classes model, not only around a 'least favorable one'  $\rightarrow$  local lower bound

Open questions

- SBM: 'Continuum' of rates inbetween k/n and  $k^2/n$  depending on A, B?
- More general estimation theory of graphon functionals?

イロト イポト イヨト イヨト