

High-dimensional random geometric graphs

Gábor Lugosi

ICREA and Pompeu Fabra University, Barcelona

joint work with

Luc Devroye (McGill University, Montréal)

András Gyöngy (SZTAKI, Budapest)

Frederic Udina (Pompeu Fabra University, Barcelona),

motivation: alien detection

- satellites search for signs of extraterrestrial invasion.

motivation: alien detection

n satellites search for signs of extraterrestrial invasion.

Satellite i receives $\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,d})$.

motivation: alien detection

n satellites search for signs of extraterrestrial invasion.

Satellite i receives $\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,d})$.

If there is no signal, all $Z_{i,t}$ are i.i.d. standard normal.

motivation: alien detection

n satellites search for signs of extraterrestrial invasion.

Satellite i receives $\mathbf{Z}_i = (\mathbf{Z}_{i,1}, \dots, \mathbf{Z}_{i,d})$.

If there is no signal, all $\mathbf{Z}_{i,t}$ are i.i.d. standard normal.

Alternatively, a small subset $\mathbf{S} \subset \{1, \dots, n\}$ of satellites receives a common signal embedded in noise:

$$\mathbf{Z}_{i,t} = \begin{cases} \mathbf{N}_{i,t} & \text{if } i \notin \mathbf{S} \\ (\mathbf{N}_{i,t} + \mathbf{Y}_t) / \sqrt{1 + \sigma^2} & \text{if } i \in \mathbf{S} \end{cases}$$

where the $\mathbf{N}_{i,t}$ are i.i.d. standard normal and the \mathbf{Y}_t are independent normal $(\mathbf{0}, \sigma^2)$.

random correlation graph

For the testing problem, it is natural to calculate pairwise correlations

$$\frac{(\mathbf{Z}_i, \mathbf{Z}_j)}{\|\mathbf{Z}_i\| \cdot \|\mathbf{Z}_j\|}$$

and define a graph by connecting i and j if the correlation is large enough.

Under the null hypothesis, the $\mathbf{X}_i = \mathbf{Z}_i / \|\mathbf{Z}_i\|$ are uniformly distributed on the sphere and we have a **random geometric graph** in \mathbb{R}^d .

random geometric graph

Given n i.i.d. points in \mathbb{R}^d , connect two with an edge if their distance is $\leq r$.

Well understood if d is fixed and $n \rightarrow \infty$.

We are interested in the behavior of the graph when the dimension is large.

random geometric graph

Model: Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be independent vectors, uniform on $\mathbf{S}_{d-1} = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1\}$.

For a given $\mathbf{p} \in (0, 1)$, we define the **random geometric graph** $\overline{\mathbf{G}}(n, d, \mathbf{p})$

Vertex set $\mathbf{V} = \{1, \dots, n\}$.

i and j are connected by an edge if and only if

$$(\mathbf{X}_i, \mathbf{X}_j) \geq t_{\mathbf{p},d}$$

where $t_{\mathbf{p},d}$ is such that

$$\mathbb{P}\{(\mathbf{X}_i, \mathbf{X}_j) \geq t_{\mathbf{p},d}\} = \mathbf{p}.$$

Equivalently, $i \sim j$ if and only if $\|\mathbf{X}_i - \mathbf{X}_j\| \leq \sqrt{2(1 - t_{\mathbf{p},d})}$.

edge probability

For $\mathbf{p} = \mathbf{1}/2$, $\mathbf{t}_{\mathbf{p},d} = \mathbf{0}$.

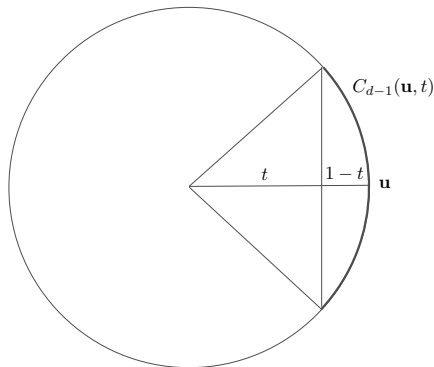
edge probability

For $\mathbf{p} = \mathbf{1}/2$, $t_{\mathbf{p},d} = 0$.

Let μ_{d-1} be the uniform probability measure over \mathbf{S}_{d-1} .

For $\mathbf{u} \in \mathbf{S}_{d-1}$ and $0 \leq t \leq 1$, a spherical cap of height $1 - t$ around \mathbf{u} is

$$C_{d-1}(\mathbf{u}, t) = \{\mathbf{x} \in \mathbb{R}^d : \mathbf{x} \in \mathbf{S}_{d-1}, (\mathbf{x}, \mathbf{u}) \geq t\}$$



edge probability

$\mathbf{p} = \mu_{d-1}(\mathbf{C}_{d-1}(\mathbf{e}, \mathbf{t}_{\mathbf{p},d}))$ is the normalized surface area of a spherical cap of height $\mathbf{1} - \mathbf{t}_{\mathbf{p},d}$.

edge probability

$\mathbf{p} = \mu_{d-1}(\mathbf{C}_{d-1}(\mathbf{e}, t_{\mathbf{p},d}))$ is the normalized surface area of a spherical cap of height $1 - t_{\mathbf{p},d}$.

It is useful to represent

$$\mathbf{x} = \frac{\mathbf{z}}{\|\mathbf{z}\|}$$

with $\mathbf{z} \in \mathbb{R}^d$ standard normal.

edge probability

$\mathbf{p} = \mu_{d-1}(\mathbf{C}_{d-1}(\mathbf{e}, t_{\mathbf{p},d}))$ is the normalized surface area of a spherical cap of height $1 - t_{\mathbf{p},d}$.

It is useful to represent

$$\mathbf{x} = \frac{\mathbf{z}}{\|\mathbf{z}\|}$$

with $\mathbf{z} \in \mathbb{R}^d$ standard normal.

Clearly, $\mathbb{E}\|\mathbf{z}\|^2 = d$.

edge probability

$\mathbf{p} = \mu_{d-1}(\mathbf{C}_{d-1}(\mathbf{e}, t_{\mathbf{p},d}))$ is the normalized surface area of a spherical cap of height $1 - t_{\mathbf{p},d}$.

It is useful to represent

$$\mathbf{x} = \frac{\mathbf{z}}{\|\mathbf{z}\|}$$

with $\mathbf{z} \in \mathbb{R}^d$ standard normal.

Clearly, $\mathbb{E}\|\mathbf{z}\|^2 = d$.

Since $\|\mathbf{z}\|$ is a Lipschitz function of \mathbf{z} , $\text{var}(\|\mathbf{z}\|) \leq 1$.

edge probability

$\mathbf{p} = \mu_{d-1}(\mathbf{C}_{d-1}(\mathbf{e}, \mathbf{t}_{\mathbf{p},d}))$ is the normalized surface area of a spherical cap of height $\mathbf{1} - \mathbf{t}_{\mathbf{p},d}$.

It is useful to represent

$$\mathbf{x} = \frac{\mathbf{z}}{\|\mathbf{z}\|}$$

with $\mathbf{z} \in \mathbb{R}^d$ standard normal.

Clearly, $\mathbb{E}\|\mathbf{z}\|^2 = d$.

Since $\|\mathbf{z}\|$ is a Lipschitz function of \mathbf{z} , $\text{var}(\|\mathbf{z}\|) \leq 1$.

In particular, $\|\mathbf{z}\|/\sqrt{d} \rightarrow \mathbf{1}$ in probability.

This implies $\mathbf{x}_1\sqrt{d}$ is approximately standard normal.

edge probability

Consequence: for any $\mathbf{s} > \mathbf{0}$,

$$\mu_{\mathbf{d}-1}(\mathbf{C}_{\mathbf{d}-1}(\mathbf{e}, \mathbf{s}/\sqrt{\mathbf{d}})) = \mathbb{P}\{\mathbf{X}_1 > \mathbf{s}/\sqrt{\mathbf{d}}\} \rightarrow 1 - \Phi(\mathbf{s})$$

as $\mathbf{d} \rightarrow \infty$.

edge probability

Consequence: for any $\mathbf{s} > \mathbf{0}$,

$$\mu_{\mathbf{d}-1}(\mathbf{C}_{\mathbf{d}-1}(\mathbf{e}, \mathbf{s}/\sqrt{\mathbf{d}})) = \mathbb{P}\{\mathbf{X}_1 > \mathbf{s}/\sqrt{\mathbf{d}}\} \rightarrow \mathbf{1} - \Phi(\mathbf{s})$$

as $\mathbf{d} \rightarrow \infty$.

For any fixed $\mathbf{p} \in (0, 1)$,

$$\lim_{\mathbf{d} \rightarrow \infty} t_{\mathbf{p}, \mathbf{d}} \sqrt{\mathbf{d}} = \Phi^{-1}(1 - \mathbf{p}) .$$

very large dimension

$\mathbf{G}(\mathbf{n}, \mathbf{p})$ denotes the Erdős-Rényi random graph. (\mathbf{n} vertices, edges are present independently, with probability \mathbf{p} .)

Total variation distance between two random graphs \mathbf{G} and \mathbf{G}' :

$$d_{\text{TV}}(\mathbf{G}, \mathbf{G}') = \max_{\mathcal{G}} |\mathbb{P}\{\mathbf{G} \in \mathcal{G}\} - \mathbb{P}\{\mathbf{G}' \in \mathcal{G}\}|$$

where the maximum is over all $2^{\binom{\mathbf{n}}{2}}$ sets of graphs over \mathbf{n} vertices.

very large dimension

$\mathbf{G}(n, p)$ denotes the Erdős-Rényi random graph. (n vertices, edges are present independently, with probability p .)

Total variation distance between two random graphs \mathbf{G} and \mathbf{G}' :

$$d_{TV}(\mathbf{G}, \mathbf{G}') = \max_{\mathcal{G}} |\mathbb{P}\{\mathbf{G} \in \mathcal{G}\} - \mathbb{P}\{\mathbf{G}' \in \mathcal{G}\}|$$

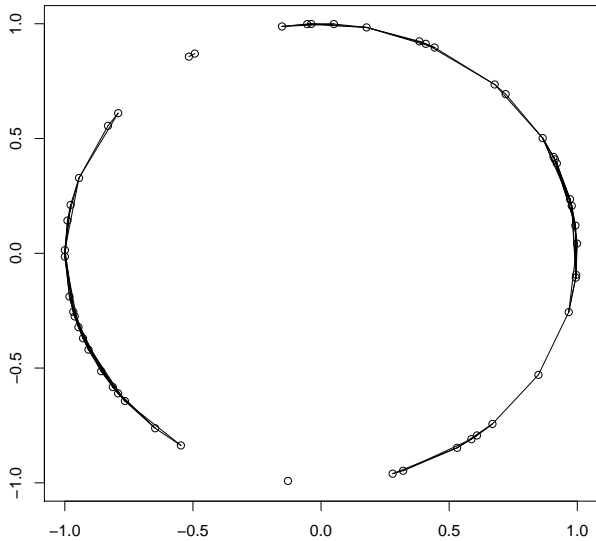
where the maximum is over all $2^{\binom{n}{2}}$ sets of graphs over n vertices.

THEOREM. Fix n and p . Then

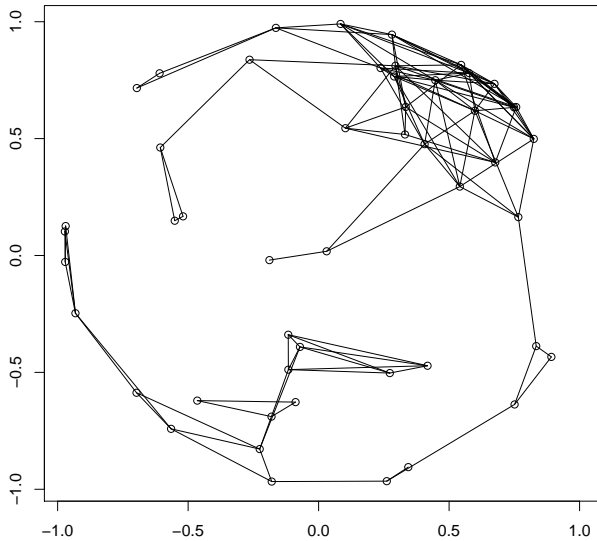
$$\lim_{d \rightarrow \infty} d_{TV}(\overline{\mathbf{G}}(n, d, p), \mathbf{G}(n, p)) = 0 .$$

Follows from a multivariate central limit theorem.

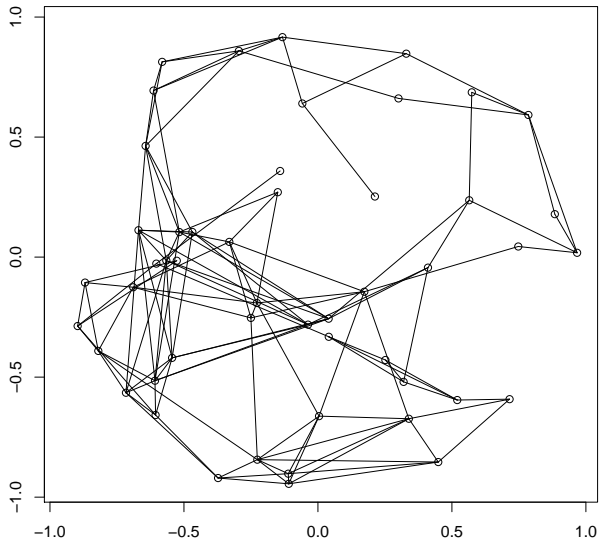
n= 50 p= 0.1 d= 2



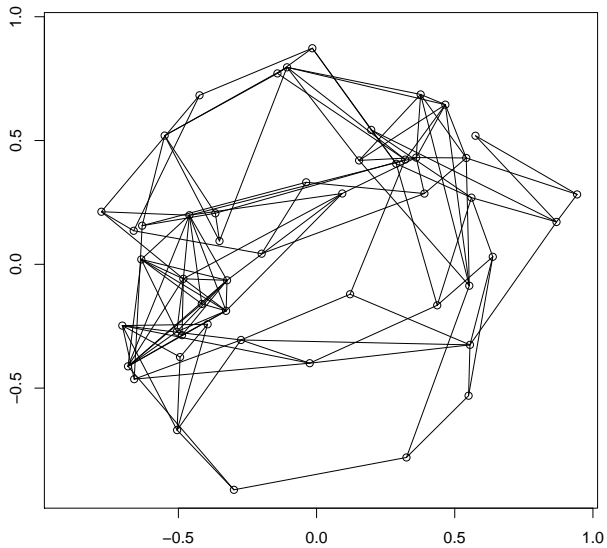
n= 50 p= 0.1 d= 3



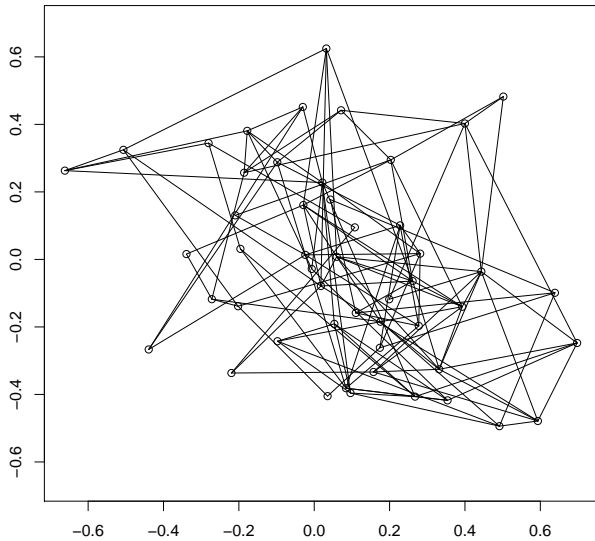
n= 50 p= 0.1 d= 4



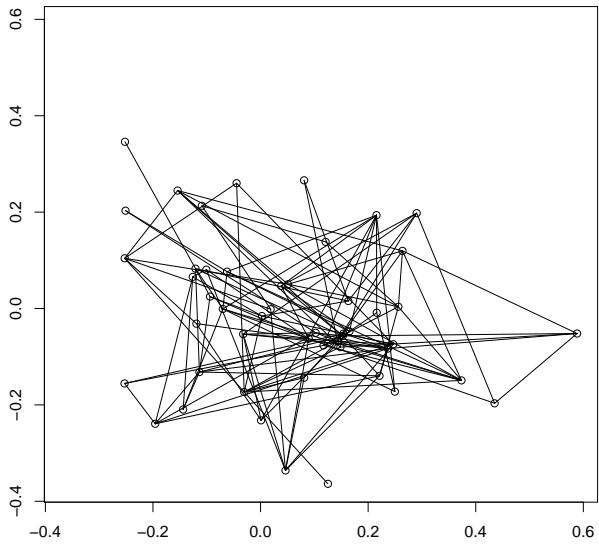
n= 50 p= 0.1 d= 5



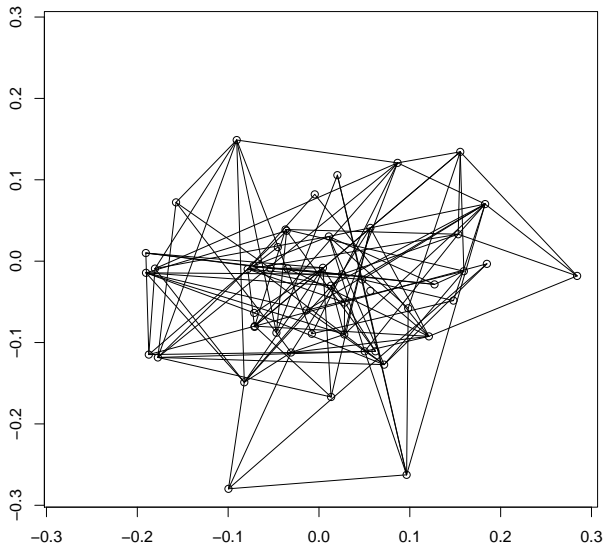
$n=50$ $p=0.1$ $d=10$



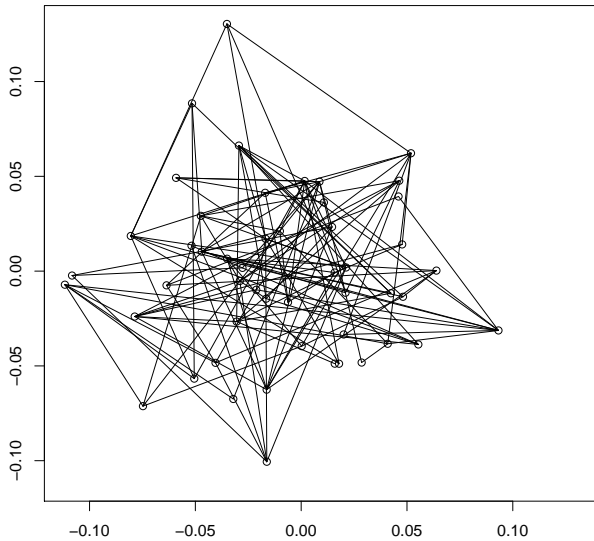
$n=50$ $p=0.1$ $d=30$



n= 50 p= 0.1 d= 100



n= 50 p= 0.1 d= 500



clique number of $\overline{G}(n, d, p)$

For fixed d and p , the clique number $\omega(n, d, p)$ grows linearly with n .

For $d = \infty$ the behavior is very different:

$$\omega(n, \infty, p) = 2 \log_{1/p} n - 2 \log_{1/p} \log_{1/p} n + O(1).$$

How fast does $\omega(n, d, p)$ approach the clique number of $G(n, p)$?

How large does d need to be for similar behavior?

clique number bounds

p is fixed, n grows.

if $d \sim \text{const.}$, then $\omega(n, d, p) = \Omega_p(n)$

if $d \rightarrow \infty$, then $\omega(n, d, p) = o_p(n)$

if $d = o(\log n)$, then $\omega(n, d, p) \geq n^{1-o_p(1)}$

if $d \sim \log^2 n$, then $\omega(n, d, p) = O_p(\log^3 n)$

if $d \gg \log^3 n$, then $\omega(n, d, p) = (2 + o_p(1)) \log_{1/p} n$

if $d \sim \log^5 n$, then

$\omega(n, d, p) = 2 \log_{1/p} n - 2 \log_{1/p} \log_{1/p} n + O_p(1)$

proof ideas

first three statements are easy (from area estimate of a spherical cap)

proof ideas

first three statements are easy (from area estimate of a spherical cap)

third follows from Jung's theorem and Vapnik-Chervonenkis inequality

proof ideas

first three statements are easy (from area estimate of a spherical cap)

third follows from Jung's theorem and Vapnik-Chervonenkis inequality

Jung's theorem (1901): For every set $A \subset \mathbb{R}^d$ of diameter at most **1** there exists a closed ball of radius $\sqrt{d/(2(d+1))}$ that includes **A**.

proof ideas

first three statements are easy (from area estimate of a spherical cap)

third follows from Jung's theorem and Vapnik-Chervonenkis inequality

Jung's theorem (1901): For every set $\mathbf{A} \subset \mathbb{R}^d$ of diameter at most **1** there exists a closed ball of radius $\sqrt{d/(2(d+1))}$ that includes **A**.

last two statements are the main result.

upper bound for $p = 1/2$

N_k is the number of cliques of size k . For $G(n, p)$,

$$\mathbb{E}N_k = \binom{n}{k} 2^{-\binom{k}{2}}$$

Let $\delta > 0$ and $K > 2$. If

$$d \geq \frac{K^3}{\delta^2},$$

then, for $1 \leq k \leq K$,

$$\mathbb{E}N_k(n, d, 1/2) \leq \binom{n}{k} \Phi(\delta)^{\frac{(k-1)(k-2)}{2}}.$$

upper bound for $p = 1/2$

N_k is the number of cliques of size k . For $G(n, p)$,

$$\mathbb{E}N_k = \binom{n}{k} 2^{-\binom{k}{2}}$$

Let $\delta > 0$ and $K > 2$. If

$$d \geq \frac{K^3}{\delta^2},$$

then, for $1 \leq k \leq K$,

$$\mathbb{E}N_k(n, d, 1/2) \leq \binom{n}{k} \Phi(\delta)^{\frac{(k-1)(k-2)}{2}}.$$

Follows from an inductive argument, using approximate orthogonality of $\mathbf{X}_1, \dots, \mathbf{X}_n$.

clique number estimates

The upper bounds for $\omega(\mathbf{n}, \mathbf{d}, \mathbf{p})$ follow from the **first moment method**:

$$\mathbb{P}\{\omega(\mathbf{n}, \mathbf{d}, 1/2) \geq k\} = \mathbb{P}\{\mathbf{N}_k \geq 1\} \leq \mathbb{E}\mathbf{N}_k ,$$

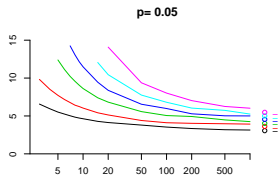
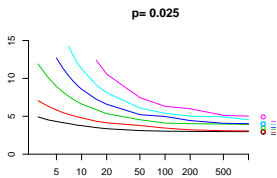
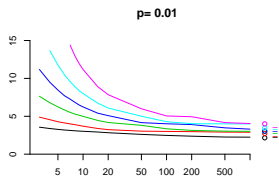
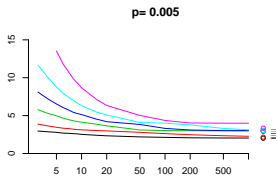
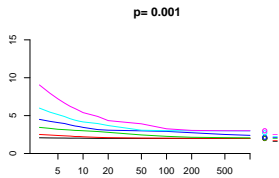
Lower bounds for $\omega(\mathbf{n}, \mathbf{d}, \mathbf{p})$ follow from the **second moment method**.

First we prove a similar lower bound for $\mathbb{E}\mathbf{N}_k$ and then show

$$\frac{\text{var}(\mathbf{N}_k)}{(\mathbb{E}\mathbf{N}_k)^2} \rightarrow 0$$

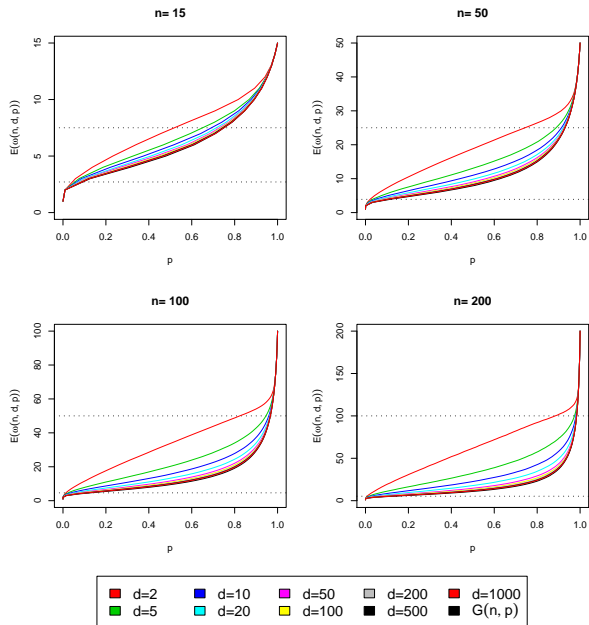
for the relevant values of k .

$\omega(n, d, p)$ as a function of d



$n=5000$
 $n=2000$
 $n=1000$
 $n=500$
 $n=200$
 $n=100$

$\omega(n, d, p)$ as a function of p



testing hidden dependencies

$$\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,d}), i = 1, \dots, n.$$

Null hypothesis: all $\mathbf{Z}_{i,t}$ are i.i.d. standard normal.

testing hidden dependencies

$$\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,d}), i = 1, \dots, n.$$

Null hypothesis: all $\mathbf{Z}_{i,t}$ are i.i.d. standard normal.

Alternative hypothesis: $\exists \mathbf{S} \subset \{1, \dots, n\}$ with $|\mathbf{S}| \geq m$ such that

$$\mathbf{Z}_{i,t} = \begin{cases} \mathbf{N}_{i,t} & \text{if } i \notin \mathbf{S} \\ (\mathbf{N}_{i,t} + \mathbf{Y}_t) / \sqrt{1 + \sigma^2} & \text{if } i \in \mathbf{S} \end{cases}$$

where the $\mathbf{N}_{i,t}$ are i.i.d. standard normal and the \mathbf{Y}_t are independent normal $(\mathbf{0}, \sigma^2)$.

test

Define $\mathbf{X}_i = \mathbf{Z}_i / \|\mathbf{Z}_i\|$ and form the graph $\overline{\mathbf{G}}(\mathbf{n}, \mathbf{d}, 1/2)$.

accept the null hypothesis if and only if $\omega(\mathbf{n}, \mathbf{d}, 1/2) \leq 3 \log_2 \mathbf{n}$.

test

Define $\mathbf{X}_i = \mathbf{Z}_i / \|\mathbf{Z}_i\|$ and form the graph $\overline{\mathbf{G}}(\mathbf{n}, \mathbf{d}, 1/2)$.

accept the null hypothesis if and only if $\omega(\mathbf{n}, \mathbf{d}, 1/2) \leq 3 \log_2 \mathbf{n}$.

$\exists \mathbf{C}, \epsilon_n \rightarrow 0$ such that if

$$\mathbf{d} \geq \mathbf{C} \max \left(\frac{\ln m}{\sigma^4}, \log_2^3 \mathbf{n} \right) \quad \text{and} \quad m > 3 \log_2 \mathbf{n}$$

then the test errs with probability $< \epsilon_n$ under both the null and alternative hypotheses.

test

Define $\mathbf{X}_i = \mathbf{Z}_i / \|\mathbf{Z}_i\|$ and form the graph $\overline{\mathbf{G}}(\mathbf{n}, \mathbf{d}, 1/2)$.

accept the null hypothesis if and only if $\omega(\mathbf{n}, \mathbf{d}, 1/2) \leq 3 \log_2 \mathbf{n}$.

$\exists \mathbf{C}, \epsilon_n \rightarrow 0$ such that if

$$\mathbf{d} \geq \mathbf{C} \max \left(\frac{\ln m}{\sigma^4}, \log_2^3 \mathbf{n} \right) \quad \text{and} \quad m > 3 \log_2 \mathbf{n}$$

then the test errs with probability $< \epsilon_n$ under both the null and alternative hypotheses.

This test is computationally very expensive.

questions

sharper bounds for the value of \mathbf{d} ?

conjecture: $\mathbb{E}\omega(\mathbf{n}, \mathbf{d}, \mathbf{p})$ is nonincreasing in \mathbf{d} for fixed \mathbf{n}, \mathbf{p} .

when does two-point concentration kick in?

connectivity threshold? giant component?

a computationally efficient test? (related to hidden clique problem of Alon, Krivelevich, and Sudakov).