

Combinatorial testing problems

the problem

A hypothesis testing problem.

One observes $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$.

Null hypothesis: the components of \mathbf{X} are i.i.d. standard normal.

Alternative hypothesis: We have a class $\mathcal{C} = \{\mathbf{S}_1, \dots, \mathbf{S}_N\}$ of sets $\mathbf{S}_k \subset \{1, \dots, n\}$.

Under H_1 , there exists an $\mathbf{S} \in \mathcal{C}$ such that

$$\mathbf{X}_i \text{ has distribution } \begin{cases} \mathcal{N}(0, 1) & \text{if } i \notin \mathbf{S} \\ \mathcal{N}(\mu, 1) & \text{if } i \in \mathbf{S} \end{cases}$$

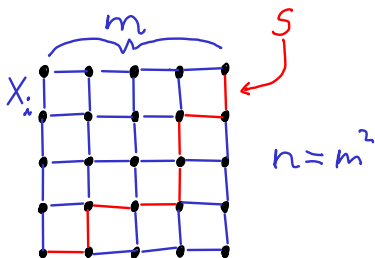
The components of \mathbf{X} are independent.

The distribution of \mathbf{X} is denoted by $\mathbb{P}_{\mathbf{S}}$.

Assume that $|\mathbf{S}| = k$ for every $\mathbf{S} \in \mathcal{C}$.

examples

- * the set of all subsets $S \subset \{1, \dots, n\}$ of size k (classical multiple testing, e.g., Ingster (1999), Baraud (2002), Donoho and Jin (2004))
- * set of all paths in a square grid (Arias-Castro et al., 2008)



examples

- ✱ set of all clusters of a certain structure on a grid
(Arias-Castro, Candès, and Durand, 2009)
- ✱ the set of all cliques of a given size in a complete graph;
- ✱ the set of all bicliques of a given size in a complete bipartite graph; (exploratory analysis of microarray data, see Shabalin, Weigman, Perou, and Nobel, 2009)
- ✱ the set of all spanning trees of a complete graph;
- ✱ the set of all perfect matchings in a complete bipartite graph;

the risk

A **test** is a function $f : \mathbb{R}^n \rightarrow \{0, 1\}$.

If $f(\mathbf{X}) = 0$ then the test accepts the null hypothesis.

The **risk** of f is

$$R(f) = \mathbb{P}_0\{f(\mathbf{X}) = 1\} + \frac{1}{N} \sum_{S \in \mathcal{C}} \mathbb{P}_S\{f(\mathbf{X}) = 0\}.$$

The optimal test is easy to determine: if $\mathbf{x}_S = \sum_{i \in S} \mathbf{x}_i$,

$$f^*(\mathbf{x}) = \begin{cases} 0 & \text{if } \frac{1}{N} \sum_{S \in \mathcal{C}} e^{\mu \mathbf{x}_S - k\mu^2/2} \leq 1 \\ 1 & \text{otherwise.} \end{cases}$$

The optimal risk is

$$R^* = R_c^*(\mu) = 1 - \frac{1}{2} \mathbb{E}_0 |L(\mathbf{X}) - 1|$$

where $L(\mathbf{X}) = \frac{1}{N} \sum_{S \in \mathcal{C}} e^{\mu \mathbf{X}_S - K\mu^2/2}$ is the likelihood ratio where $\mathbf{X}_S = \sum_{i \in S} \mathbf{X}_i$.

averaging test

Compute the average and threshold:

$$f(\mathbf{x}) = \mathbb{1}_{\{\sum_{i=1}^n x_i > \mu k/2\}} \cdot$$

Then for any $\delta > 0$, $R(f) \leq \delta$ whenever

$$\mu \geq \sqrt{\frac{8n}{k^2} \log \frac{2}{\delta}} \cdot$$

Surprisingly, this simple test is nearly optimal in many cases.

scan statistic

Another simple test is

$$f(\mathbf{x}) = 1 \quad \text{if and only if} \quad \max_{S \in \mathcal{C}} X_S \geq \frac{\mu k + \mathbb{E}_0 \max_{S \in \mathcal{C}} X_S}{2} .$$

Then $R(f) \leq \delta$ whenever

$$\mu \geq \frac{\mathbb{E}_0 \max_{S \in \mathcal{C}} X_S}{k} + 2\sqrt{\frac{2}{k} \log \frac{2}{\delta}} .$$

This is an easy consequence of Tsirelson's Gaussian concentration inequality.

lower bounds

Let \mathbf{S} and \mathbf{S}' be drawn independently, uniformly, at random from \mathcal{C} and let $Z = |\mathbf{S} \cap \mathbf{S}'|$. Then

$$R^* \geq 1 - \frac{1}{2} \sqrt{\mathbb{E} e^{\mu^2 Z} - 1}.$$

Proof.

$$R^* = 1 - \frac{1}{2} \mathbb{E}_0 |L(\mathbf{X}) - 1| \geq 1 - \frac{1}{2} \sqrt{\mathbb{E}_0 |L(\mathbf{X}) - 1|^2}$$

Since $\mathbb{E}_0 L(\mathbf{X}) = 1$,

$$\mathbb{E}_0 |L(\mathbf{X}) - 1|^2 = \text{Var}_0(L(\mathbf{X})) = \mathbb{E}_0 [L(\mathbf{X})^2] - 1.$$

a lower bound for gaussian processes

Let μ_C be such that

$$\mathbb{E} \exp(\mu_C^2 Z) = 2.$$

Then for all $\mu \leq \mu_C$,

$$R^* \geq \frac{1}{2}.$$

But recall from the scan statistic test that $R^* < 1/2$ whenever

$$\mu \geq \frac{\mathbb{E}_0 \max_{S \in \mathcal{C}} X_S}{k} + \frac{4}{\sqrt{k}}.$$

This implies that

$$\mathbb{E}_0 \max_{S \in \mathcal{C}} X_S \geq k\mu_C - 4\sqrt{k}.$$

negative association

Assume that the class \mathcal{C} is symmetric. Suppose that $S' = \{1, 2, \dots, k\} \in \mathcal{C}$ and let S be a randomly chosen from \mathcal{C} . If the random variables $\mathbb{1}_{\{1 \in S\}}, \dots, \mathbb{1}_{\{k \in S\}}$ are negatively associated then $R^* \geq 1/2$ whenever

$$\mu \leq \sqrt{\log \left(1 + \frac{n \log 2}{k^2} \right)}.$$

These tools allow us to determine the order of the critical μ in various applications.

examples: k -sets

Consider the class \mathcal{C} of all sets of size k . Thus, $N = \binom{n}{k}$.

This example is well-understood, see, e.g., Ingster (1999), Baraud (2002), Donoho and Jin (2004).

We have negative association, so $R^* \geq 1/2$ whenever

$$\mu \leq \sqrt{\log \left(1 + \frac{n \log 2}{k^2} \right)}.$$

Various regimes. Since

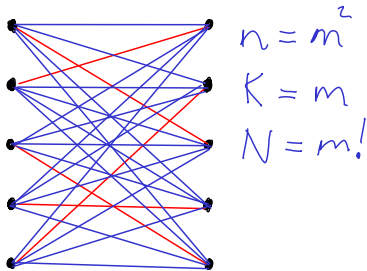
$$\frac{\mathbb{E}_0 \max_{S \in \mathcal{C}} X_S}{k} \leq \frac{\sqrt{2k \log \binom{n}{k}}}{k} \leq \sqrt{2 \log \left(\frac{ne}{k} \right)},$$

the scan test shows that for $k = O(n^{(1-\epsilon)/2})$ for some $\epsilon > 0$, then the threshold value is of the order of $\sqrt{\log n}$.

When k^2/n , is bounded away from zero, then the averaging test provides a matching upper bound.

perfect matchings

Let \mathcal{C} be the set of all perfect matchings of the complete bipartite graph $K_{m,m}$.



By the the averaging test, $R(f) \leq 1/2$ whenever $\mu \geq \sqrt{8 \log 4}$.

perfect matchings

$Z = |\mathcal{S} \cap \mathcal{S}'|$ is the number of fixed points in a random permutation.

Then $R^* \geq 1/2$ whenever

$$\mu \leq \sqrt{\log(1 + \log 2)} .$$

The optimal test f^* can be computed efficiently: computing

$$\frac{1}{N} \sum_{S \in \mathcal{C}} e^{\mu X_S} = \frac{1}{m!} \sum_{\sigma} \prod_{j=1}^m e^{\mu X_{(j, \sigma(j))}}$$

is equivalent to computing the permanent of a non-negative $m \times m$ matrix. [Jerrum, Sinclair, and Vigoda \(2004\)](#) show that this may be done by a polynomial-time randomized approximation.

spanning trees

Let \mathcal{C} be the class of all spanning trees of the complete graph K_m .

$$n = \binom{m}{2} \quad k = m - 1 \quad N = m^{m-2}.$$

Negative association holds by a result of Feder and Mihail (1992).

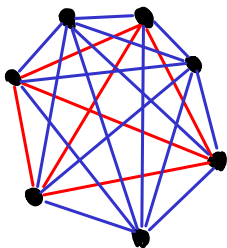
$R^* \geq 1/2$ whenever

$$\mu \leq \sqrt{\log \left(1 + \frac{\log 2}{2} \right)}.$$

The likelihood ratio $(1/N) \sum_{s \in \mathcal{C}} e^{\mu X_s}$ may be computed by the Propp-Wilson (1998) algorithm.

cliques

Let \mathcal{C} contain all cliques of size k of the complete graph K_m .



$$n = \binom{m}{2}$$

$$K = \binom{m}{k}$$

$$N = \binom{n}{k}$$

Negative association doesn't hold anymore.

cliques

Assume $k \leq \sqrt{m(\log 2)/e}$. Then

(i) $R^* \leq 1/2$ whenever

$$\mu \geq 2\sqrt{\frac{1}{k-1} \log\left(\frac{me}{k}\right)} + 4\sqrt{\frac{\log 4}{k(k-1)}},$$

(ii) $R^* \geq 1/2$ whenever

$$\mu \leq \sqrt{\frac{1}{k} \log\left(\frac{m}{2k}\right)}.$$

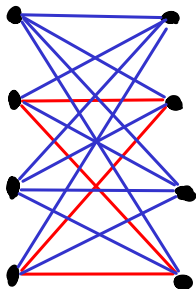
Upper bound holds by scan test. Both scan test and optimal test are difficult to compute.

bi-cliques

A similar model, relevant in the exploratory analysis of microarray data (Shabalin, Weigman, Perou, and Nobel, 2009):

\mathcal{C} contains all bi-cliques of size k of the complete bipartite graph $K_{m,m}$.

Then $n = m^2$, $K = k^2$, and $N = \binom{m}{k}$.



The analysis and results are similar to the case of cliques.

research problem: sharp threshold?

The optimal risk $R^* = R_C^*(\mu)$ is a decreasing function of μ .

For what classes is there a sharp threshold?

If μ_a is such that $R_C^*(\mu) = a$, is it true that

$$\mu_\epsilon - \mu_{1-\epsilon} = o(\mu_{1/2}) \quad \text{as } k \rightarrow \infty?$$

detection of correlations

n sensors receive noisy signals.

Sensor i receives $(X_{i,1}, \dots, X_{i,d})$.

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

In the presence of an event, a small subset $S \subset \{1, \dots, n\}$ of sensors receives a common signal covered by noise:

$$X_{i,t} = \begin{cases} N_{i,t} & \text{if } i \notin S \\ (N_{i,t} + Y_t)/\sqrt{1 + \rho} & \text{if } i \in S \end{cases}$$

where $N_{i,t}$ are standard normal and Y_t are normal $(0, \rho)$ independent of $(\rho \ll 1)$.

hypothesis testing problem

One observes d i.i.d vectors $\mathbf{X}_1, \dots, \mathbf{X}_d \in \mathbb{R}^n$.

Under the **null hypothesis**, the components of $\mathbf{X}_1, \dots, \mathbf{X}_d$ are independent.

Under the **alternative hypothesis**, a subset \mathbf{S} of components ($|\mathbf{S}| = k \ll n$) has correlation $\rho > 0$, the rest are independent.

One does not know the (possibly) “contaminated” subset.

tests and their risk

A **test** is a function $f : \mathbb{R}^{nd} \rightarrow \{0, 1\}$.

$f(\mathbf{X}_1, \dots, \mathbf{X}_d) = 0$ means that the null is accepted.

The **risk** of a test is

$$R(f) = \mathbb{P}_0\{f(\mathbf{X}_1, \dots, \mathbf{X}_d) = 1\} + \frac{1}{\binom{n}{k}} \sum_{\mathbf{S}} \mathbb{P}_{\mathbf{S}}\{f(\mathbf{X}_1, \dots, \mathbf{X}_d) = 0\}$$

where \mathbb{P}_0 is the probability under the null.

$\mathbb{P}_{\mathbf{S}}$ is the probability if the set of contaminated components is \mathbf{S} .

Another measure of risk:

$$\bar{R}(f) = \mathbb{P}_0\{f(\mathbf{X}_1, \dots, \mathbf{X}_d) = 1\} + \max_{\mathbf{S}} \mathbb{P}_{\mathbf{S}}\{f(\mathbf{X}_1, \dots, \mathbf{X}_d) = 0\}$$

optimal test

There exists a test f^* , (likelihood ratio test) with optimal risk $R^* = R(f^*)$.

This test is difficult to compute. Requires calculating a sum of $\binom{n}{k}$ terms.

a lower bound

$$R^* \geq \frac{1}{2} - \frac{1}{4} \sqrt{\exp\left(\left(e^{\rho d} - 1\right) \frac{k^2}{n}\right) - 1}.$$

(Arias-Castro, Bubeck, Lugosi, 2015)

Proof idea: The bound

$$R^* \geq 1 - \frac{1}{2} \sqrt{\mathbb{E}_0[L(\mathbf{X})^2] - 1}$$

does not work. Instead, we condition on the value of the common “signal” and treat the problem as shifted means.

If the class of possible “contaminated sets” is not necessarily the class of all sets of size k , then again R^* maybe bounded from below in terms of the moment generating function of

$$Z = |\mathbf{S} \cap \mathbf{S}'|.$$

a lower bound

$$R^* \geq \frac{1}{2} - \frac{1}{4} \sqrt{\exp\left(\left(e^{\rho d} - 1\right) \frac{k^2}{n}\right) - 1}.$$

If $d \ll (1/\rho) \log n$, it is impossible to have $R^* \rightarrow 0$ unless $k = n^{1/2 - o(1)}$.

However, if $d \geq (9/\rho) \log n$, small risk is possible even if k is a constant.

This can be achieved by a simple test based on a “scan statistic”.

scan statistic

$$T_{n,k,d} = \max_S \sum_{t=1}^d \sum_{i,j \in S, i \neq j} X_{i,t} X_{j,t}.$$

The test f_n that accepts the null hypothesis if and only if

$$T_{n,k,d} \leq dk(k-1)\rho/4$$

has risk $R(f_n) \rightarrow 0$ if $(k-1)\rho > 8$ and $d \geq (9/\rho) \log n$.

The proof is based on the analysis of Gaussian quadratic forms.

random correlation graph

Let $\mathbf{X}_i = (X_{i,1}, \dots, X_{i,d})$.

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

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

and define a graph by connecting i and j if the correlation is large enough (e.g., > 0).

In the presence of a signal, one expects a large clique.

What is large?

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

random geometric graph

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

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 $\|\mathbf{X}_i - \mathbf{X}_j\| \leq \sqrt{2(1 - t_{\mathbf{p},d})}$.

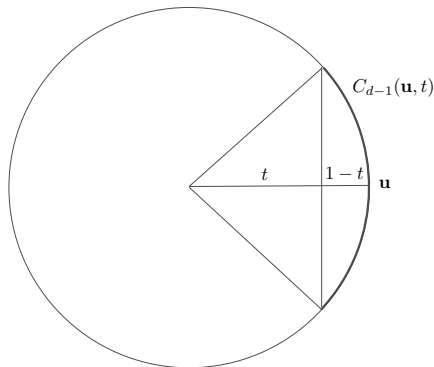
edge probability

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

Let μ_{d-1} be the uniform probability measure over S_{d-1} .

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

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



edge probability

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

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, $\|\mathbf{Z}\|/\sqrt{d} \rightarrow 1$ in probability.

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

edge probability

Consequence: for any $s > 0$,

$$\mu_{d-1}(\mathbf{C}_{d-1}(\mathbf{e}, s/\sqrt{d})) = \mathbb{P}\{X_1 > s/\sqrt{d}\} \rightarrow 1 - \Phi(s)$$

as $d \rightarrow \infty$.

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

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

very large dimension

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

Total variation distance between two random graphs G and G' :

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

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

THEOREM. Fix n and p . Then

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

very large dimension

For the proof, write

$$\begin{aligned}\mathbb{P}\{\bar{\mathbf{G}}(n, d, p) = \mathbf{g}\} &= \mathbb{P}\left\{\bigcap_{1 \leq i < j \leq n} \mathbb{1}_{\{(\mathbf{X}_i, \mathbf{X}_j) > t_{p,d}\}} = \mathbf{g}_{i,j}\right\} \\ &= \mathbb{P}\left\{\bigcap_{1 \leq i < j \leq n} \mathbb{1}_{\{\sum_t \sqrt{d} \mathbf{X}_{i,t} \mathbf{X}_{j,t} > \sqrt{d} t_{p,d}\}} = \mathbf{g}_{i,j}\right\}\end{aligned}$$

and

$$\mathbb{P}\{\mathbf{G}(n, p) = \mathbf{g}\} = \mathbb{P}\left\{\bigcap_{1 \leq i < j \leq n} \mathbb{1}_{\{N_{i,j} > \Phi^{-1}(1-p)\}} = \mathbf{g}_{i,j}\right\}$$

where the $N_{i,j}$ are independent standard normal. The proof follows from a multivariate central limit theorem.

clique number of $\overline{G}(n, d, 1/2)$

For fixed d , the clique number $\omega(n, d, 1/2)$ grows linearly with n .
For $d = \infty$ the behavior is very different:

$$\omega(n, \infty, 1/2) = 2 \log_2 n - 2 \log_2 \log_2 n + O(1).$$

How fast does $\omega(n, d, 1/2)$ approach the clique number of

$G(n, 1/2)$?

How large does d need to be for similar behavior?

clique number bounds

if $d \sim \text{const.}$, then $\omega(n, d, 1/2) = \Omega_p(n)$

if $d \rightarrow \infty$, then $\omega(n, d, 1/2) = o_p(n)$

if $d = o(\log n)$, then $\omega(n, d, 1/2) \geq n^{1-o_p(1)}$

if $d \sim \log^2 n$, then $\omega(n, d, p) = O_p(\log^3 n)$ when $p < 1/2$.

if $d \gg \log^3 n$, then $\omega(n, d, 1/2) = (2 + o_p(1)) \log_2 n$

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

$$\omega(n, d, 1/2) = 2 \log_2 n - 2 \log_2 \log_2 n + O_p(1)$$

if $d \geq \log((2 \log 2)n)$,

$$\omega(n, d, 1/2) \geq \exp\left(\log^2(n)/(20d)\right) - 1.$$

proof ideas

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

fourth 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 .

clique number bounds

the points \mathbf{X}_i in a clique \mathbf{K} form a set of diameter at most $\sqrt{2(1 - t_{p,d})}$.

by Jung's theorem, \mathbf{K} is contained in a spherical cap $\mathbf{C}_{d-1}(\mathbf{u}, s_p)$ for some $\mathbf{u} \in \mathbf{S}_{d-1}$ and $s_p = \sqrt{(dt_{p,d} + 1)/(d + 1)}$.

$$\omega(d, n, p) \leq n \sup_{\mathbf{C}=\mathbf{C}_{d-1}(\mathbf{u}, s_p)} \mu_{d-1,n}(\mathbf{C}),$$

where $\mu_{d-1,n}$ is the empirical measure.

by Vapnik and Chervonenkis, with probability at least $1 - \eta$, for all $\mathbf{C} = \mathbf{C}_{d-1}(\mathbf{u}, s_p)$,

$$\mu_{d-1,n}(\mathbf{C}) \leq 2\mu_{d-1}(\mathbf{C}) + \frac{4(d+1)}{n} \ln \frac{2ne}{d+1} + \frac{4}{n} \ln \frac{4}{\eta}.$$

and

$$\mu_{d-1,n}(\mathbf{C}) \leq \frac{1}{s_p \sqrt{d}} e^{-s_p^2(d-1)/2} + \frac{4(d+1)}{n} \ln \frac{2ne}{d+1} + \frac{4}{n} \ln \frac{4}{\eta}.$$

clique number bounds

if $d \sim \text{const.}$, then $\omega(n, d, 1/2) = \Omega_p(n)$

if $d \rightarrow \infty$, then $\omega(n, d, 1/2) = o_p(n)$

if $d = o(\log n)$, then $\omega(n, d, 1/2) \geq n^{1-o_p(1)}$

if $d \sim \log^2 n$, then $\omega(n, d, p) = O_p(\log^3 n)$ when $p < 1/2$.

if $d \gg \log^3 n$, then $\omega(n, d, 1/2) = (2 + o_p(1)) \log_2 n$

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

$$\omega(n, d, 1/2) = 2 \log_2 n - 2 \log_2 \log_2 n + O_p(1)$$

if $d \geq \log((2 \log 2)n)$,

$$\omega(n, d, 1/2) \geq \exp\left(\log^2(n)/(20d)\right) - 1.$$

upper bound for $d \gg \log^3 n$

N_k is the number of cliques of size k . For $G(n, 1/2)$,

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

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 $\omega(n, d, p)$ follow from the **first moment method**:

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

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

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

$$\frac{\text{var}(N_k)}{(\mathbb{E}N_k)^2} \rightarrow 0$$

for the relevant values of k .

clique number bounds

if $d \sim \text{const.}$, then $\omega(n, d, 1/2) = \Omega_p(n)$

if $d \rightarrow \infty$, then $\omega(n, d, 1/2) = o_p(n)$

if $d = o(\log n)$, then $\omega(n, d, 1/2) \geq n^{1-o_p(1)}$

if $d \sim \log^2 n$, then $\omega(n, d, p) = O_p(\log^3 n)$ when $p < 1/2$.

if $d \gg \log^3 n$, then $\omega(n, d, 1/2) = (2 + o_p(1)) \log_2 n$

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

$$\omega(n, d, 1/2) = 2 \log_2 n - 2 \log_2 \log_2 n + O_p(1)$$

if $d \geq \log((2 \log 2)n)$,

$$\omega(n, d, 1/2) \geq \exp\left(\log^2(n)/(20d)\right) - 1.$$

lower bound

The lower bound follows from the lower bound for the risk of any test.

If the largest clique was small then there would exist a test with small risk.

lower bound

Suppose $0 < \rho \leq 1/2$. Under the alternative hypothesis, with probability at least $1 - \delta$, the graph contains a clique of size k whenever

$$\binom{k}{2} \leq \delta \exp\left(\frac{d\sigma^4}{10}\right),$$

where $\sigma^2 = \rho/(1 - \rho)$.

references

G. Lugosi. Lectures on Combinatorial Statistics.

.../SaintFlour.pdf

L. Addario-Berry, N. Broutin, L. Devroye, and G. Lugosi (2010). Combinatorial testing problems. *Annals of Statistics*.

L. Devroye, A. György, G. Lugosi, and F. Udina (2011). High-dimensional random geometric graphs and their clique number. *Electronic Journal of Probability*.

E. Arias-Castro, S. Bubeck, and G. Lugosi, (2012). Detection of correlations. *Annals of Statistics*.

E. Arias-Castro, S. Bubeck, and G. Lugosi, (2015). Detecting positive correlations in a multivariate sample. *Bernoulli*.