

Asymptotic normality of global clustering coefficient of uniform random intersection graph

Mindaugas Bloznelis*
joint work with Jerzy Jaworski

*Vilnius University, Lithuania

<http://www.mif.vu.lt/~bloznelis>

May 16, 2018

Global clustering coefficient

The global clustering coefficient of a finite graph \mathcal{G} is the ratio

$$C_{\mathcal{G}} = 3N_{\Delta}/N_V,$$

N_{Δ} is the number of triangles, N_V is the number of 2-paths.

$C_{\mathcal{G}}$ represents the probability that a randomly selected path of length 2 induces triangle in \mathcal{G} .

$C_{\mathcal{G}}$ is a commonly used network characteristic, assessing the strength of the statistical association between neighboring adjacency relations.

For example, in a social network the tendency of linking actors which have a common neighbor is reflected by a non-negligible value of the global clustering coefficient.

Random intersection graph, RIG

Clustering in a social network can be explained by an auxiliary bipartite structure: each actor is prescribed a collection of attributes and any two actors sharing a common attribute have high chances of being adjacent, cf. [Newman et al. *Random graphs with arbitrary degree distributions and their applications*. Phys.Rev.E, 2002].

The respective random intersection graph (RIG) on the vertex set $V = \{v_1, \dots, v_n\}$ and with the auxiliary attribute set $W = \{w_1, \dots, w_m\}$ defines adjacency relations with the help of a random bipartite graph H linking actors (=vertices) to attributes: two actors are adjacent in RIG if they have a common neighbour in H .

For H drawn uniformly at random from the class of bipartite graphs where each actor $v_i \in V$ has exactly r neighbours in W we obtain the uniform RIG denoted $G = G(n, m, r)$.

Alternative definition: vertices v_1, \dots, v_n of G are represented by iid random subsets $S_1, \dots, S_n \subset W$, each of size r . Vertices v_i, v_j are adjacent whenever $S_i \cap S_j \neq \emptyset$.

Uniform random intersection graph $G(n, m, r)$

Vertices v_1, \dots, v_n of $G = G(n, m, r)$ are represented by iid random subsets $S_1, \dots, S_n \subset W = \{w_1, \dots, w_m\}$, each of size r . Vertices v_i, v_j are adjacent whenever $S_i \cap S_j \neq \emptyset$.

The graph has been widely studied in the literature mainly as a model of secure wireless sensor network that uses random predistribution of keys [Eschenauer and Gligor (2002)]: each sensor v_i is prescribed a collection S_i of “secret” keys; two sensors establish a secure communication link if they share a secret key.

We consider large random intersection graphs, where $m, n \rightarrow \infty$. For $r^2 = o(m)$ the edge probability is

$$p_e = \mathbf{P}(v_i \sim v_j) = r^2 m^{-1} + O(r^4 m^{-2}).$$

We are interested in sparse graphs with bounded average degree

$$\eta := (n-1)p_e \approx (n-1)r^2 m^{-1}.$$

Paveiksliaukas !

Asymptotic normality of $C_G = 3N_\Delta/N_V$ (I)

Let \mathcal{G} be an instance of $G(n, m, r)$. Denote

$$\mu = \frac{3\mathbf{E}N_\Delta}{\mathbf{E}N_V}, \quad N_\Delta^* = \frac{N_\Delta - \mathbf{E}N_\Delta}{\mathbf{E}N_\Delta}, \quad N_V^* = \frac{N_V - \mathbf{E}N_V}{\mathbf{E}N_V}.$$

We approximate C_G by μ . The approximation error

$$C_G - \mu = \mu \left(\frac{N_\Delta^* + 1}{N_V^* + 1} - 1 \right) = \mu(N_\Delta^* - N_V^*) - R \quad (1)$$

has the leading term $\mu(N_\Delta^* - N_V^*)$ and (negligible) remainder

$$R = \mu(N_\Delta^* - N_V^*) \frac{N_V^*}{N_V^* + 1}.$$

Assume that $\mathbf{E}N_\Delta, \mathbf{E}N_V \rightarrow +\infty$ and properly standardized vector (N_Δ^*, N_V^*) is asymptotically (bivariate) normal. Then (1) implies the asymptotic normality of $C_G - \mu$.

Asymptotic normality of $C_G = 3N_\Delta/N_V$ (I)

Let \mathcal{G} be an instance of $G(n, m, r)$. Denote

$$\mu = \frac{3\mathbf{E}N_\Delta}{\mathbf{E}N_V}, \quad N_\Delta^* = \frac{N_\Delta - \mathbf{E}N_\Delta}{\mathbf{E}N_\Delta}, \quad N_V^* = \frac{N_V - \mathbf{E}N_V}{\mathbf{E}N_V}.$$

We approximate C_G by μ . The approximation error

$$C_G - \mu = \mu \left(\frac{N_\Delta^* + 1}{N_V^* + 1} - 1 \right) = \mu(N_\Delta^* - N_V^*) - R \quad (1)$$

has the leading term $\mu(N_\Delta^* - N_V^*)$ and (negligible) remainder

$$R = \mu(N_\Delta^* - N_V^*) \frac{N_V^*}{N_V^* + 1}.$$

Assume that $\mathbf{E}N_\Delta, \mathbf{E}N_V \rightarrow +\infty$ and properly standardized vector (N_Δ^*, N_V^*) is asymptotically (bivariate) normal. Then (1) implies the asymptotic normality of $C_G - \mu$.

Remark. We have

$\mu = \mathbf{P}(v_2 \sim v_3 | v_1 \sim v_2, v_1 \sim v_3) = \mathbf{P}(v_2^* \sim v_3^* | v_1^* \sim v_2^*, v_1^* \sim v_3^*),$
where the vertex triple (v_1^*, v_2^*, v_3^*) is drawn at random.

Bivariate asymptotic normality of (N_Δ, N_V) .

Denote $\sigma_\Delta^2 = \mathbf{Var}N_\Delta$, $\sigma_V^2 = \mathbf{Var}N_V$, $\sigma_{\Delta V} = \mathbf{Cov}(N_\Delta, N_V)$.

Proposition. Let $m, n \rightarrow +\infty$. Assume that $r \geq 2$ and $r = O(1)$. Assume that $\sigma_\Delta^2 \rightarrow +\infty$, $\sigma_V^2 \rightarrow +\infty$. Suppose that the ratio $\sigma_{\Delta V}/(\sigma_\Delta \sigma_V)$ converges to a limit. We denote the limit ρ . Then the random vector $(\sigma_\Delta^{-1}(N_\Delta - \mathbf{E}N_\Delta), \sigma_V^{-1}(N_V - \mathbf{E}N_V))$ converges in distribution to a Gaussian random vector (ξ_Δ, ξ_V) , where $\mathbf{E}\xi_j = 0$, $\mathbf{E}\xi_j^2 = 1$, $j = \Delta, V$, and $\mathbf{E}\xi_\Delta \xi_V = \rho$.

Remark about the scaling of σ_Δ and σ_V . We have for $r^3 = o(m)$

$$\sigma_\Delta^2 \approx \frac{1}{2r^2}m\eta^3 + \frac{1}{6r}m\eta^2, \quad \sigma_V^2 \approx m\eta^3\left(2 + \frac{2}{r} + \frac{1}{2r^2}\right) + m\eta^2\left(\frac{1}{2} + \frac{1}{r}\right),$$
$$\sigma_{\Delta V} \approx m\eta^3\left(\frac{1}{2r^2} + \frac{1}{r}\right) + m\eta^2\frac{1}{2r}.$$

Recall that $\eta = (n-1)p_e \approx nr^2m^{-1}$ denotes the average degree. For fixed r and bounded $\eta \approx \eta_0 > 0$ we have

$$\rho \approx \frac{\eta_0(1+2r) + r}{\sqrt{(\eta_0 + 3^{-1}r)(\eta_0(4r^2 + 4r + 1) + r^2 + 2r)}}.$$

Asymptotic normality of $C_G = 3N_\Delta/N_V$ (II)

The bivariate asymptotic normality (Proposition)

$$\left(\frac{N_\Delta - \mathbf{E}N_\Delta}{\sigma_\Delta^{-1}}, \frac{N_V - \mathbf{E}N_V}{\sigma_V^{-1}} \right) \rightarrow (\xi_\Delta, \xi_V)$$

implies the approximation (in distribution)

$$C_G \approx \mu \left(\frac{N_\Delta - \mathbf{E}N_\Delta}{\mathbf{E}N_\Delta} - \frac{N_V - \mathbf{E}N_V}{\mathbf{E}N_V} \right) \approx \mu \left(\frac{\sigma_\Delta}{\mathbf{E}N_\Delta} \xi_\Delta - \frac{\sigma_V}{\mathbf{E}N_V} \xi_V \right) \quad (*)$$

Hence the asymptotic normality of C_G .

Comment on the contribution of N_Δ and N_V to (*):

- for $r = O(1)$ the coefficients $\sigma_\Delta/\mathbf{E}N_\Delta$ and $\sigma_V/\mathbf{E}N_V$ are of the same order of magnitude, thus, N_Δ and N_V enter on equal terms;
- for $r \rightarrow +\infty$, $\eta = o(r)$ we have $\sigma_V/\mathbf{E}N_V = o(\sigma_\Delta/\mathbf{E}N_\Delta)$, that is, the triangle count N_Δ outplays N_V .

These observations follow easily from the relations

$$\mathbf{Var}\xi_\Delta = \mathbf{Var}\xi_V = 1, \quad \mu = 3\mathbf{E}N_\Delta/\mathbf{E}N_V \approx r^{-1},$$

$$\sigma_\Delta^2 \approx \frac{1}{2r^2} m\eta^3 + \frac{1}{6r} m\eta^2, \quad \sigma_V^2 \approx m\eta^3 \left(2 + \frac{2}{r} + \frac{1}{2r^2} \right) + m\eta^2 \left(\frac{1}{2} + \frac{1}{r} \right).$$

Approach to the asymptotic normality of (N_Δ, N_V) .

It suffices to show that for any reals a, b the random variables

$$X_{n,m} = a \frac{N_\Delta - \mathbf{E}N_\Delta}{\sigma_\Delta} + b \frac{N_V - \mathbf{E}N_V}{\sigma_V}$$

converge in distribution to $a\xi_\Delta + b\xi_V$ as $n, m \rightarrow +\infty$.

To this aim we use Stein's method. In our case this approach leads to awkward calculations. To make calculations feasible we “preprocess” the subgraph counts

$$N_\Delta = N_\Delta(S_1, \dots, S_n) = \sum_{\{i,j,k\} \subset [n]} \Delta_{i,j,k}, \quad (1)$$

$$N_V = N_V(S_1, \dots, S_n) = \sum_{\{i,j,k\} \subset [n]} (\vee_{ijk} + \vee_{jki} + \vee_{kij}).$$

Namely, we decompose them into series of uncorrelated “polynomials” of increasing order, called Hoeffding's decomposition.

The decomposition helps to tackle statistical dependencies. In particular, we obtain reasonably simple expressions of $\mathbf{Var}(N_\Delta)$, $\mathbf{Var}(N_V)$ and $\mathbf{Cov}(N_\Delta, N_V)$.