

Graph Limits & Subgraph Counts

**Econometric Methods for Social Spillovers and Networks,
University of St. Gallen, October 1st to 9th, 2018**

Bryan S. Graham

University of California - Berkeley

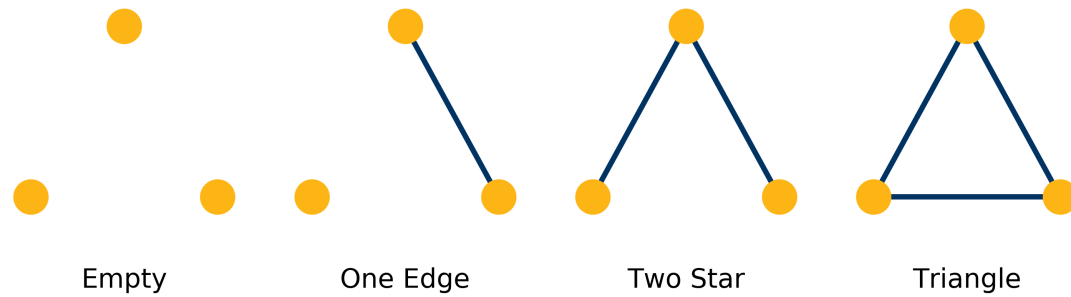
Introduction

In 1970 Paul Holland and Samuel Leinhardt (1970, *AJS*) introduced the *triad census*.

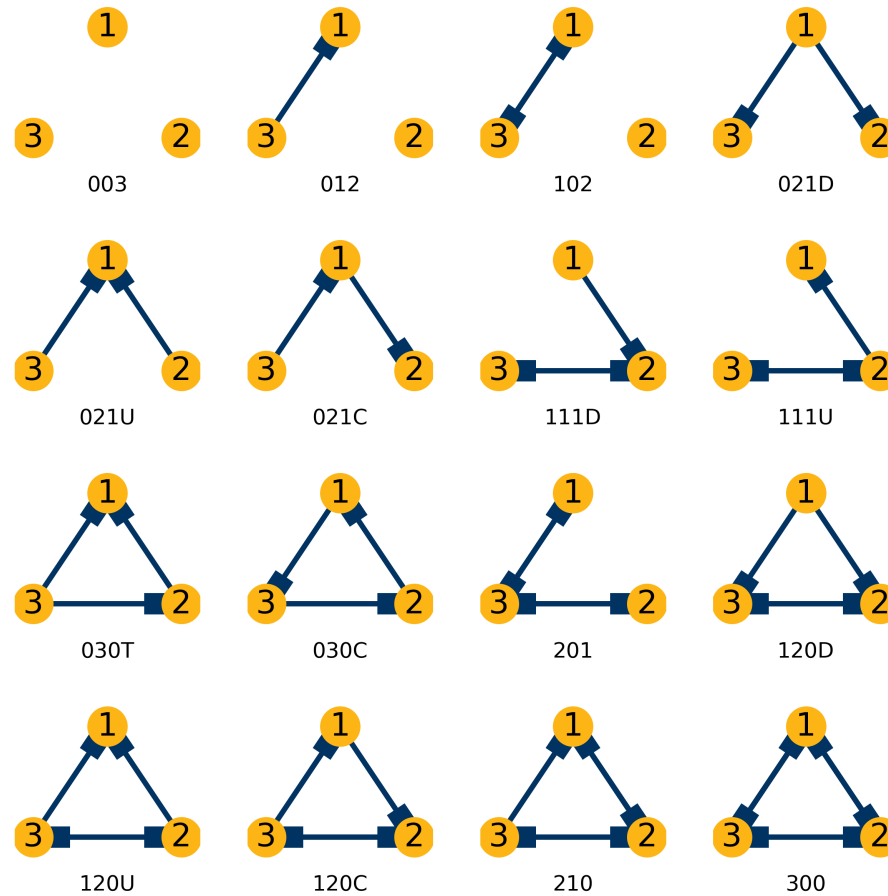
- counts of all 4 (16) unique triad isomorphisms in an undirected (directed) graph;
- can construct transitivity index (TI) from triad census...
- ...as well as the mean and variance of the degree sequence.

Holland and Leinhardt (1976, *SM*) provided variance expressions for these counts (brute force).

Triads: Undirected Case



Triads: Directed Case



Introduction (continued)

In early work normality of these counts was assumed (w/o proof).

Nowicki (1989, 1991) showed asymptotic normality of counts for homogenous random graphs.

Bickel, Chen & Levina (2011, AS) demonstrated asymptotic normality in the “general” case under specific conditions.

Introduction (continued)

Subgraph counts, called *network moments* by Bickel, Chen and Levina (2011), summarize average local properties of a network.

Large literature in sociology which uses triad counts to “test” various hypotheses

- see Holland and Leinhardt (1976, SM) and Wasserman and Faust (1994)
- cf., computational biology (e.g., Milo et al., 2002)

Asymptotic distribution theory puts these tests on firmer ground.

Introduction (continued)

Subgraph frequencies might be used to (partially) identify structural models of network formation (e.g., de Paula et al., 2018).

indirect inference approach:

1. use structural model to simulate networks...and count subgraphs;
2. compare simulated counts with actual counts;
3. estimate structural parameters by minimum distance.

Setup

Let $G(\mathcal{V}, \mathcal{E})$ be a finite undirected random graph with



- agents/vertices $\mathcal{V} = \{1, \dots, N\}$,
- links/edges $\mathcal{E} = \{\{i, j\}, \{k, l\}, \dots\}$, and
- adjacency matrix $\mathbf{D} = [D_{ij}]$ with

$$D_{ij} = \begin{cases} 1 & \text{if } \{i, j\} \in \mathcal{E} \\ 0 & \text{otherwise} \end{cases}$$

Subgraphs

- (Partial Subgraph) Let $\mathcal{V}(S) \subseteq \mathcal{V}(G)$ be any subset of the vertices of G and $\mathcal{E}(S) \subseteq \mathcal{E}(G) \cap \mathcal{V}(S) \times \mathcal{V}(S)$, then $S = (\mathcal{V}(S), \mathcal{E}(S))$ is an *partial subgraph* of G .
- (Induced Subgraph) Let $\mathcal{V}(S) \subseteq \mathcal{V}(G)$ be any subset of the vertices of G and $\mathcal{E}(S) = \mathcal{E}(G) \cap \mathcal{V}(S) \times \mathcal{V}(S)$, then $S = (\mathcal{V}(S), \mathcal{E}(S))$ is an *induced subgraph* of G .

Subgraphs (continued)

- The induced subgraph S includes *all* edges in G connecting any two agents in $\mathcal{V}(S)$
 - a (partial) subgraph may include only a subset of such edges
 - $S =$  is a partial subgraph of $G =$ , but not an induced subgraph

Graph Isomorphism

- Consider two graphs, R and S , of the same order.
- Let $\varphi : \mathcal{V}(R) \rightarrow \mathcal{V}(S)$ be a bijection from the nodes of R to those of S .
- The bijection $\varphi : \mathcal{V}(R) \rightarrow \mathcal{V}(S)$
 - *maintains adjacency* if for every dyad $i, j \in \mathcal{V}(R)$ if $\{i, j\} \in \mathcal{E}(R)$, then $\{\varphi(i), \varphi(j)\} \in \mathcal{E}(S)$;
 - *maintains non-adjacency* if for every dyad $i, j \in \mathcal{V}(R)$ if $\{i, j\} \notin \mathcal{E}(R)$, then $\{\varphi(i), \varphi(j)\} \notin \mathcal{E}(S)$.

Graph Isomorphism (continued)

- If the bijection maintains both adjacency and non-adjacency we say it *maintains structure*.
- (Graph Isomorphism) The graphs R and S are *isomorphic* if there exists a structure-maintaining bijection $\varphi : \mathcal{V}(R) \rightarrow \mathcal{V}(S)$.
- Notation: $R \cong S$ means “ R is isomorphic to S .”

P-Cycles

A p -cycle is p^{th} order graphlet with nodes labeled (or relabeled) such that its edges form a cycle:

$$\mathcal{E}(S) = \{(i_1, i_2), (i_2, i_3), \dots, (i_p, i_1)\}.$$

A p -cycle is a connected graphlet with p edges on p nodes.

Examples: triangles ($S = \triangle$) and 4-cycles ($S = \square$).

Trees



A *tree* is a connected graph with no cycles.

The number of edges on a p^{th} order tree is $p - 1$; a feature which will prove highly convenient.

Examples: p -star graphlets, such as two-stars ($S = \text{🔺}$) and three-stars ($S = \text{🔻}$).

Also called connected acyclic graphs.

Induced Subgraph Density

- S is a p^{th} -order graphlet of interest (e.g., $S =$  or $S =$ )
- G_N is the network/graph under study
- $\mathbf{i}_p \subseteq \{1, 2, \dots, N\}$ is a set of p integers with $i_1 < i_2 < \dots < i_p$
 - $\mathcal{C}_{p,N}$ is set of all $\binom{N}{p}$ such integer sets
 - $G[\mathbf{i}_p]$ is the induced subgraph of G associated with vertex set \mathbf{i}_p

Induced Subgraph Density (continued)

- The *induced subgraph density* of S in G_N , denoted by $t_{\text{ind}}(S, G_N)$ or $P_N(S)$ equals the probability that $G_N[\mathbf{i}_p]$, for \mathbf{i}_p chosen uniformly at random from $C_{p,N}$, is isomorphic to S :

$$\begin{aligned} t_{\text{ind}}(S, G_N) &= \binom{N}{p}^{-1} \sum_{\mathbf{i}_p \in C_{p,N}} \mathbf{1}(S \cong G_N[\mathbf{i}_p]) \\ &= \Pr(S = G_N[\mathbf{i}_p]) \\ &= P_N(S) \end{aligned}$$

- Slightly different definition used in some of the technical literature...(see *Handbook* chapter)

Induced Subgraph Density (Examples)

- $t_{\text{ind}}(\text{triangle}, \text{square}) = \frac{2}{4}$, $t_{\text{ind}}(\text{V}, \text{square}) = \frac{2}{4}$

and $t_{\text{ind}}(\text{edge}, \text{square}) = \frac{0}{4}$

- $t_{\text{ind}}(\text{triangle}, \text{K}_4) = \frac{1}{4}$, $t_{\text{ind}}(\text{V}, \text{K}_4) = \frac{2}{4}$

and $t_{\text{ind}}(\text{edge}, \text{K}_4) = \frac{1}{4}$

Goal

We would like a result of the form...

$$\sqrt{N} \left(\begin{pmatrix} \hat{P}_N \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) \\ \hat{P}_N \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) \end{pmatrix} - \begin{pmatrix} P \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) \\ P \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) \end{pmatrix} \right) \xrightarrow{D} \mathcal{N}(0, \Sigma)$$

...under conditions we can understand

...with a covariance Σ we can estimate

...and interpretable limit values $P \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right)$ and $P \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right)$

Goal (continued)

With this result we can conduct inference on *transitivity*...

$$\text{TI} = \frac{3P\left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array}\right)}{P\left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array}\right) + 3P\left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array}\right)}$$

Is $\text{TI} > P\left(\begin{array}{c} \bullet \\ \text{---} \\ \bullet \end{array}\right)$ (see Jackson et al. (2012) for some motivation)?

cf., Blitzstein and Diaconis (2011)

Induced Subgraph Density: Graphon Case

Let $h(U_i, U_j)$ be a valid graphon.

Let $\text{iso}(S)$ be the group of isomorphisms of S , and $|\text{iso}(S)|$ its cardinality.

Under the “Aldous-Hoover DGP” the *ex ante* probability that an induced p-subgraph is isomorphic to S is given by

$$\begin{aligned} t_{\text{ind}}(S, h) &= |\text{iso}(S)| \\ &\times \mathbb{E} \left[\prod_{\{i,j\} \in \mathcal{E}(S)} h(U_i, U_j) \prod_{\{i,j\} \in \mathcal{E}(\bar{S})} [1 - h(U_i, U_j)] \right] \\ &= P(S). \end{aligned}$$

Graph Limits

Let $\{G_N\}_{N=1}^{\infty}$ be a sequence of networks. If

$$\lim_{N \rightarrow \infty} t_{\text{ind}}(S, G_N) = t_{\text{ind}}(S, h)$$

for some graphon $h(\cdot, \cdot)$ and *all* fixed subgraphs S , then we say that G_N converges to $h(\cdot, \cdot)$.

- Lovász (2012) for complete development.
- Diaconis and Janson (2008) for connections with Aldous-Hoover Theorem.
- Result establishes a connection between subgraph counts and the graphon.

(Injective) Homomorphism Density

The homomorphism density gives the probability that S is (isomorphic to) a subgraph of a randomly selected induced subgraph of G_N of order $p = |\mathcal{V}(S)|$

Alternatively the homomorphism density equals fraction of injective mappings $\varphi : \mathcal{V}(S) \rightarrow \mathcal{V}(G_N)$ that preserve edge adjacency

$$\begin{aligned} t_{\text{hom}}(S, G_N) &= \frac{1}{\binom{N}{p} |\text{iso}(S)|} \sum_{R \subseteq K_N, R \cong S} \mathbf{1}(R \subseteq G_N) \\ &= \frac{1}{\binom{N}{p} |\text{iso}(S)|} \sum_{R \subseteq K_N, |V(R)|=p} \mathbf{1}(R \cong S) \prod_{\{i,j\} \in \mathcal{E}(R)} D_{ij} \\ &= Q_N(S) \end{aligned}$$

Homomorphism Density (continued)

Summation in $t_{\text{hom}}(S, G_N) = Q_N(S)$ is over the $\binom{N}{3} \left| \text{iso}(\triangle) \right| = \frac{3}{6}N(N-1)(N-2)$ (partial) subgraphs of K_N (the complete graph) which are isomorphic to $S = \triangle$.

We count the number of these subgraphs which are also *partial* subgraphs of G_N

Homomorphism Density (continued)

The expected value of $Q_N(S)$ is:

$$\begin{aligned}
 \mathbb{E}[Q_N(S)] &= \frac{1}{\binom{N}{p} |\text{iso}(S)|} \sum_{R \subseteq K_N, |V(R)|=p} \{ \mathbf{1}(R \cong S) \\
 &\quad \times \mathbb{E} \left[\mathbb{E} \left[\prod_{\{i,j\} \in \mathcal{E}(R)} D_{ij} \middle| U_1, \dots, U_N \right] \right] \} \\
 &= \mathbb{E} \left[\prod_{\{i,j\} \in \mathcal{E}(S)} h(U_i, U_j) \right] \\
 &= Q(S) \stackrel{\text{def}}{=} t_{\text{hom}}(S, h)
 \end{aligned}$$

Can also use $t_{\text{hom}}(S, G_N)$ to define graph convergence.

Recap

Induced subgraph density, $P_N(S)$: probability that $G_N[\mathbf{i}_p]$, for \mathbf{i}_p chosen uniformly at random from $C_{p,N}$, is isomorphic to S .

Homomorphism density, $Q_N(S)$: probability that a (partial) subgraph of $G_N[\mathbf{i}_p]$, for \mathbf{i}_p chosen uniformly at random from $C_{p,N}$, is isomorphic to S .

If $\lim_{N \rightarrow \infty} P_N(S) = t_{\text{ind}}(S, h)$ for some graphon $h(\cdot, \cdot)$ and all fixed subgraphs S , then we say that G_N converges to $h(\cdot, \cdot)$.

Computation


Useful to reformulate definition of $\hat{P}_N(S)$.



Let $\mathbf{D}_{[\mathbf{i}_p, \mathbf{i}_p]}$ be the $p \times p$ sub-adjacency matrix constructed by removing all rows and columns of \mathbf{D} except those in $\mathbf{i}_p = \{i_1, \dots, i_p\}$.

Let S be a graphlet of interest.

We can check for whether S is an isomorphism of $G[\mathbf{i}_p]$ by inspecting the elements of the $\mathbf{D}_{[\mathbf{i}_p, \mathbf{i}_p]}$ sub-adjacency matrix.


Computation (continued)

Consider the two star triad $S =$ , we can express $\mathbf{1}(S \cong G_N[\mathbf{i}_p])$ in terms of $\mathbf{D}_{[\mathbf{i}_p, \mathbf{i}_p]}$ as

$$\begin{aligned} \mathbf{1} \left(\text{ \cong G_N[\mathbf{i}_3] \right) &= D_{i_1 i_2} D_{i_1 i_3} (1 - D_{i_2 i_3}) + D_{i_1 i_2} (1 - D_{i_1 i_3}) D_{i_2 i_3} \\ &\quad + (1 - D_{i_1 i_2}) D_{i_1 i_3} D_{i_2 i_3} \\ &\stackrel{\text{def}}{=} V \\ &\quad \text{ , \mathbf{i}_3 \end{aligned}$$


Computation (continued)

Let $\text{iso}(S)$ be the group of isomorphisms of S , and $|\text{iso}(S)|$ its cardinality (i.e., number of subgraphs of K_p that are isomorphic to S).

We have $|\text{iso}(\text{triangle})| = 3$; three terms to the right of the (first) equality are indicators for three isomorphisms of  on $\{i_1, i_2, i_3\}$.

Computation (continued)

In general $\mathbf{1}(S \cong G_N[\mathbf{i}_p])$ may be defined in terms of $\mathbf{D}_{[\mathbf{i}_p, \mathbf{i}_p]}$ with number of components equal to the number of possible isomorphisms of S .

There is only one isomorphism of the  configuration, yielding a second example of

$$\mathbf{1} \left(\begin{array}{c} \text{triangle} \\ \cong G_N[\mathbf{i}_3] \end{array} \right) = D_{i_1 i_2} D_{i_1 i_3} D_{i_2 i_3} \underset{\substack{\text{def} \\ \equiv V}}{\text{triangle}} , \mathbf{i}_3$$

Unbiasedness

Two star configuration; iterated expectations and conditional independence of edges given $\mathbf{U} = (U_1, \dots, U_N)'$ yields

$$\begin{aligned}\mathbb{E} \left[D_{i_1 i_2} D_{i_1 i_3} (1 - D_{i_2 i_3}) \right] &= \mathbb{E} \left[\mathbb{E} \left[D_{i_1 i_2} D_{i_1 i_3} (1 - D_{i_2 i_3}) \mid \mathbf{U} \right] \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[D_{i_1 i_2} D_{i_1 i_3} (1 - D_{i_2 i_3}) \mid U_{i_1}, U_{i_2}, U_{i_3} \right] \right] \\ &= \mathbb{E} \left[h(U_{i_1}, U_{i_2}) h(U_{i_1}, U_{i_3}) [1 - h(U_{i_2}, U_{i_3})] \right]\end{aligned}$$

Unbiasedness (continued)

Value of $\mathbb{E} \left[D_{i_1 i_2} D_{i_1 i_3} (1 - D_{i_2 i_3}) \right]$ is invariant to permutations of its indices.

Recalling that $\left| \text{iso} \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) \right| = 3$ we have

$$\mathbb{E} \left[\mathbf{1} \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \cong G_N [\mathbf{i}_p] \right) \right] = 3 \cdot \int \int \int h(t, u) h(t, v) [1 - h(u, v)] dt du dv$$

$$\stackrel{\text{def}}{=} P \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right)$$

Large Sample Properties

Our estimator is

$$\begin{pmatrix} \hat{P}_N \left(\begin{array}{c} \text{triangle} \end{array} \right) \\ \hat{P}_N \left(\begin{array}{c} \text{triangle} \end{array} \right) \end{pmatrix} = \binom{N}{3}^{-1} \sum_{i_1 < i_2 < i_3} \begin{pmatrix} V \begin{array}{c} \text{triangle} \\ V \begin{array}{c} \text{triangle} \end{array} \end{array} , i_3 \end{pmatrix}.$$

It is not a U-Statistics, but has many U-Statistic-like properties.

Large Sample Properties (continued)

It is unbiased for $\left(\begin{array}{c} P \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) \\ P \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \\ \hline \bullet \quad \bullet \end{array} \right) \end{array} \right)$ under joint exchangeability (iterated expectations).

Can use Hoeffding (1948) arguments to study variance-covariance (cf., Holland and Leinhardt, 1976).

Network moments: Large N behavior

Projecting $\hat{P}_N \left(\triangle \right)$ on $\mathbf{U} = (U_1, \dots, U_N)'$ gives:

$$\begin{aligned} \hat{P}_N \left(\triangle \right) = & \binom{N}{3}^{-1} \sum_{i_1 < i_2 < i_3} h(U_{i_1}, U_{i_2}) h(U_{i_1}, U_{i_3}) h(U_{i_2}, U_{i_3}) \\ & + \binom{N}{3}^{-1} \sum_{i_1 < i_2 < i_3} \left\{ D_{i_1 i_2} D_{i_1 i_3} D_{i_2 i_3} \right. \\ & \left. - h(U_{i_1}, U_{i_2}) h(U_{i_1}, U_{i_3}) h(U_{i_2}, U_{i_3}) \right\}. \end{aligned}$$

Second term is mean independent of first with conditionally independent summands.

First term is a 3^{rd} order U-Statistic (large sample properties well-understood).

Network moments: Large N behavior (continued)

Under some conditions (most important of which is that average degree grows with N) $\hat{P}_N \left(\text{triangle} \right)$ behaves like a U-Statistic s.t.

$$\sqrt{N} \left(\begin{pmatrix} \hat{P}_N \left(\text{triangle} \right) \\ \hat{P}_N \left(\text{triangle} \right) \end{pmatrix} - \begin{pmatrix} P \left(\text{triangle} \right) \\ P \left(\text{triangle} \right) \end{pmatrix} \right) \xrightarrow{D} \mathcal{N}(0, 9\Sigma_1)$$

...with Σ_1 estimable (analog estimate involves $O(N^5)$ operations!).

Use delta method to conduct inference on transitivity.

Intellectual history

Some basic ideas (e.g., use of Hoeffding-like variance decompositions) go back (at least) to Holland and Leinhardt (1976).

Subsequent work by Nowicki (1991), Picard et al. (2008) and others.

Big breakthrough by Bickel et al. (2011) – abstract (proof uses lots of “tricks”) and limiting variance is not characterized.

Bhattacharya and Bickel (2015) – explicit characterization of variance and an estimator (cf., Menzel, 2017).

Some (interesting and empirically-relevant) subtleties ignored today.

Intellectual history (continued)

My exposition (anchored in textbook U-Statistic theory) is based on basic approach of Graham (2017).

Challenge is finding a notation that can neatly handle all cases.

Some open questions regarding sparse graph sequences.

Second (Simple) Example Density

We estimate $\rho_N = \Pr(D_{ij} = 1)$ by

$$\hat{\rho}_N = \frac{2}{N(N-1)} \sum_{i < j} D_{ij}.$$

Projecting onto U_1, \dots, U_N yields the decomposition:

$$\begin{aligned} \hat{\rho}_N &= \underbrace{\frac{2}{N(N-1)} \sum_{i < j} h_N(U_i, U_j)}_{\text{U-Statistic}} + \underbrace{\frac{2}{N(N-1)} \sum_{i < j} (D_{ij} - h_N(U_i, U_j))}_{\text{"Poisson Binomial R.V."}} \\ &= U_N + T_N. \end{aligned}$$

Observe that T_N is mean independent of U_N .

Density: Variance Calculation

We have

$$\begin{aligned}\mathbb{V}(\hat{\rho}_N) &= \mathbb{V}(U_N) + \mathbb{V}(T_N) + 2\mathbb{C}(U_N, T_N) \\ &= \mathbb{V}(U_N) + \mathbb{V}(T_N).\end{aligned}$$

A Hoeffding (1948) variance decomposition gives

$$\mathbb{V}(U_N) = \binom{N}{2}^{-2} \sum_{q=1}^2 \binom{N}{2} \binom{2}{q} \binom{N-2}{2-q} \Omega_q$$

for

$$\Omega_q = \mathbb{C}\left(h_N(U_{i_1}, U_{i_2}), h_N(U_{j_1}, U_{j_2})\right)$$

with $\{i_1, i_2\}$ and $\{j_1, j_2\}$ sharing $q = 1, 2$ indices in common.

Density: Variance Calculation (continued)

Evaluating Ω_1 yields

$$\begin{aligned}\Omega_1 &= \mathbb{E}[h_N(U_1, U_2) h_N(U_1, U_3)] - \mathbb{E}[h_N(U_1, U_2)] \mathbb{E}[h_N(U_1, U_3)] \\ &= Q\left(\begin{array}{c} \bullet \\ / \quad \backslash \\ \bullet \quad \bullet \end{array}\right) - P\left(\begin{array}{c} \bullet \text{---} \bullet \end{array}\right) P\left(\begin{array}{c} \bullet \text{---} \bullet \end{array}\right).\end{aligned}$$

Evaluating Ω_2 yields

$$\begin{aligned}\Omega_2 &= \mathbb{E}[h_N(U_1, U_2)^2] - \mathbb{E}[h_N(U_1, U_2)] \mathbb{E}[h_N(U_1, U_2)] \\ &= \mathbb{V}(\mathbb{E}[D_{12} | \mathbf{U}]).\end{aligned}$$

Density: Variance Calculation (continued)

Evaluating the variance of $\mathbb{V}(T_N)$ we get

$$\begin{aligned}\mathbb{V}(T_N) &= \mathbb{V}(\mathbb{E}[T_N | \mathbf{U}]) + \mathbb{E}[\mathbb{V}(T_N | \mathbf{U})] \\ &= 0 + \left(\frac{2}{N(N-1)}\right)^2 \mathbb{E}\left[\mathbb{V}\left(\sum_{i < j} (D_{ij} - h_N(U_i, U_j)) \middle| \mathbf{U}\right)\right] \\ &= \left(\frac{2}{N(N-1)}\right)^2 \mathbb{E}\left[\sum_{i < j} \mathbb{V}(D_{ij} - h_N(U_i, U_j) | \mathbf{U})\right] \\ &= \frac{2}{N(N-1)} \mathbb{E}[\mathbb{V}(D_{12} | \mathbf{U})].\end{aligned}$$

Density: Variance Calculation (continued)

Collecting terms we have:

$$\begin{aligned}
 \mathbb{V}(\hat{\rho}_N) &= \frac{4(N-2)}{N(N-1)} \left[Q \left(\begin{array}{c} \bullet \\ / \quad \backslash \\ \bullet \quad \bullet \end{array} \right) - P \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) P \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \right] \\
 &\quad + \frac{2}{N(N-1)} \mathbb{V}(\mathbb{E}[D_{12} | \mathbf{U}]) + \frac{2}{N(N-1)} \mathbb{E}[\mathbb{V}(D_{12} | \mathbf{U})] \\
 &= \frac{4(N-2)}{N(N-1)} \left[Q \left(\begin{array}{c} \bullet \\ / \quad \backslash \\ \bullet \quad \bullet \end{array} \right) - P \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) P \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \right] \\
 &\quad + \frac{2}{N(N-1)} P \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \left(1 - P \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \right).
 \end{aligned}$$

Density: Variance Calculation (continued)

To allow for graph sequences where $\rho_N \rightarrow 0$ as $N \rightarrow \infty$ we normalize''

- Let $\tilde{Q} \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) = \frac{Q \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right)}{\rho_N^2}$ and $\tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) = \frac{P \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right)}{\rho_N}$.
- Recall that $\lambda_N = (N - 1) \rho_N$.

Density: Variance Calculation (continued)

After normalization:

$$\begin{aligned} \mathbb{V} \left(\frac{\hat{\rho}_N}{\rho_N} \right) &= \frac{4(N-2)}{N(N-1)} \left[\tilde{Q} \left(\begin{array}{c} \bullet \\ \diagup \quad \diagdown \\ \bullet \quad \bullet \end{array} \right) - \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \right] \\ &\quad + \frac{2}{N\lambda_N} \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) - \frac{2}{N(N-1)} \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right)^2 \\ &= O \left(\frac{1}{N} \right) + O \left(\frac{1}{N\lambda_N} \right) + O \left(\frac{1}{N^2} \right). \end{aligned}$$

- If $\lambda_N \rightarrow \infty$ first term dominates.
- If $\lambda_N \rightarrow \lambda_0 > 0$, first two terms dominate.

Asymptotic Inference

Asymptotic theory for U-Statistics gives, for $\lambda_N \rightarrow \infty$ as $N \rightarrow \infty$

$$\sqrt{N} \left(\frac{\hat{\rho}_N}{\rho_N} - 1 \right) \xrightarrow{D} \mathcal{N} \left(0, 4 \left[\tilde{Q} \left(\begin{array}{c} \bullet \\ / \quad \backslash \\ \bullet \quad \bullet \end{array} \right) - \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \right] \right).$$

Result (in high level form) due to Bickel, Chen and Levina (2011, *Annals of Statistics*).

Comment: Under Erdos-Renyi $\tilde{Q} \left(\begin{array}{c} \bullet \\ / \quad \backslash \\ \bullet \quad \bullet \end{array} \right) = \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right) \tilde{P} \left(\begin{array}{c} \bullet \text{---} \bullet \end{array} \right).$

Variance Estimation

We can estimate the asymptotic variance using the analog estimators:

$$\begin{aligned}\hat{Q} \left(\text{triangle} \right) &= \binom{N}{3}^{-1} \sum_{i < j < k} \frac{1}{3} \left\{ D_{ij} D_{ik} + D_{ij} D_{jk} + D_{ik} D_{jk} \right\} \\ &= \binom{N}{3}^{-1} \frac{1}{3} [T_{TS} + 3T_T]\end{aligned}$$

and

$$\hat{P} \left(\text{edge} \right) = \binom{N}{2}^{-1} \sum_{i < j} D_{ij}$$

Nyakatoke



Variance Estimation for $\hat{P} \left(\text{---} \right)$: Nyakatoke

For Nyakatoke we have

$$\hat{Q} \left(\text{^} \right) \cong 0.006105$$

and

$$\hat{P} \left(\text{---} \right) \simeq 0.0698$$

which gives

$$\begin{matrix} \hat{\rho}_N \\ \text{(a.s.e)} \end{matrix} = \begin{matrix} 0.0698 \\ (0.0072) \end{matrix}, \quad \begin{matrix} \hat{\lambda}_N \\ \text{(a.s.e)} \end{matrix} = \begin{matrix} 8.2364 \\ (0.8459) \end{matrix}$$

Note: Estimate above includes first two terms.

Standard Error Estimation for $\hat{T}I$: Nyakatoke

In Nyakatoke there are $\binom{119}{3} = 273,819$ triad configurations to count and a total of $\binom{119}{5} = 182,637,273$ pentads that need to be inspected in order to calculate variances.

Direct calculation gives

$$P_N(\triangle) = \frac{0.00115}{(0.00030)}, \quad P_N(\triangle) = \frac{0.00496}{(0.00100)}$$

Standard Error Estimation for \hat{T}_I : Nyakatoke (continued)

Applying the delta method we get

$$\hat{T}_I = \begin{matrix} 0.188 \\ (0.011) \end{matrix}$$

which suggests that transitivity is greater than what we would expect to observe under the Erdős-Renyi random graph null.

Wrapping Up

In large graphs subgraph counting is computationally challenging

- implications for feasibility of both estimation and inference.
- see Bhattacharya and Bickel (2015) for a subsampling approach.

Very little (i.e., essentially none) empirical work using these results.

Tremendous scope for using these methods in empirical analysis; but not easy!