

An Incomplete Overview of Strategic Models of Network Formation

**Econometric Methods for Networks,
GCEP, May 8th & 9th, 2017**

Bryan S. Graham

University of California - Berkeley

Overview

- Why study models of network formation?
 - equilibrium \mathbf{D} may be inefficient.
 - planner may have preferences over \mathbf{D} and hence is interested in policies which influence it.
 - we might view network manipulations as a mechanism for influencing other outcomes.
 - correct for “endogenous network formation bias” (cf., Auerbach, 2016)

Existing Approaches

- “Applied theory approach”: posit generative models of network formation that match “stylized facts” (Albert and Barabasi, 2002).
- ERGM: directly write down likelihoods for $\Pr(\mathbf{D} = \mathbf{d})$ and try to maximize them (cf., Shalizi and Rinaldo, 2013, *Annals of Statistics*).
 - generally no micro-foundations...
 - ...but see Mele (2017, *Econometrica*)

Existing Approaches

- Random Utility Models (RUM): Sheng (2012), Christakis *et al.* (2012), Imbens and Goldsmith-Pinkham (2013), Graham (2013, 2016, 2017), de Paula *et al.* (2015)
- Specialized structural models: Banerjee *et al.* (2012).

Random Utility Approach

- This approach is both natural, and familiar, to economists
- Provides a framework for modeling the effect on link surplus (i.e., utility) of
 - observed agent/dyad covariates
 - unobserved agent attributes (heterogeneity)
 - preference interdependencies (e.g., a taste for transitivity)

A Simple Model of Network Formation

- Consider a network of three agents: i, j, k
 - Link formation: $D_{ij} = 1 \left(\alpha + \beta D_{ik} D_{jk} - U_{ij} \geq 0 \right)$ with $\beta \geq 0$ (returns to transitivity).
 - Three “types” of U_{ij} draws: $\mathbb{U}_L = (-\infty, \alpha]$, $\mathbb{U}_M = (\alpha, \alpha + \beta]$ or $\mathbb{U}_H = (\alpha + \beta, \infty)$.
 - Positive measure on the subset of the support of $\mathbf{U} = (U_{ij}, U_{ik}, U_{jk})'$ with multiple NE networks.
- The model is *incomplete* (cf., Bresnahan and Reiss, 1991; Tamer, 2003).

A Simple Model of Network Formation (continued)

- There are $3^3 = 27$ “configurations” of \mathbb{U} ...
- ...but only $\frac{(3+3-1)!}{3!(3-1)!} = 10$ non-isomorphic ones
- Two of these configurations admit multiple NE networks
 - if $U_{ij} \in \mathbb{U}_L$, $U_{ik} \in \mathbb{U}_M$ and $U_{jk} \in \mathbb{U}_M$
 - if $U_{ij} \in \mathbb{U}_M$, $U_{ik} \in \mathbb{U}_M$ and $U_{jk} \in \mathbb{U}_M$

A Simple Model of Network Formation (continued)

- Only one realization (out of 4 possible realizations of \mathbf{D}) is uniquely predicted
 - if $U_{ij} \in \mathbb{U}_L$, $U_{ik} \in \mathbb{U}_L$ and $U_{jk} \in \mathbb{U}_H$, then links $\{i, j\}$ and $\{i, k\}$ form and $\{j, k\}$ does not
- cf., Ciliberto and Tamer (2009), Sheng (2012), de Paula et al. (2015) provide methods for analyzing incomplete models of network formation
- serious challenges to implementation at scale

Christakis, Fowler, Imbens and Kalyanaraman (2010)

- Dyads form links *sequentially* and *myopically*.
- If the linking order is ij , ik and jk we have
 - $D_{ij} = \mathbf{1}(\alpha - U_{ij} \geq 0)$
 - $D_{ik} = \mathbf{1}(\alpha - U_{ik} \geq 0)$
 - $D_{jk} = \mathbf{1}(\alpha + \beta \mathbf{1}(\alpha - U_{ij} \geq 0) \mathbf{1}(\alpha - U_{ik} \geq 0) - U_{jk} \geq 0)$
- Conditional on the ij , ik and jk the realization of \mathbf{U} delivers a unique prediction of \mathbf{D} .

Christakis, Fowler, Imbens and Kalyanaraman (2010)

- Since we don't observe the order of link formation we
 - assign a (prior) distribution to it and
 - work with an integrated likelihood.
- With three dyads there are $3! = 6$ possible link orderings. Let $O \in \mathbb{O} = \{1, 2, 3, 4, 5, 6\}$ be the possible orderings. Our integrated likelihood is

$$\Pr(\mathbf{D} = \mathbf{d}) = \sum_{o \in \mathbb{O}} \Pr(\mathbf{D} = \mathbf{d} | O = o) \Pr(O = o)$$

Christakis, Fowler, Imbens and Kalyanaraman (2010)

- Christakis et al. (2010) use Bayesian MCMC methods
 - provides a method of inference as well
 - (no large sample theory for their estimator)
- Simulation methods, and assumptions about the timing of link formation, are also central to work by Goldsmith-Imbens and Pinkham (2013) and Mele (2017)

Tetrad Logit Redux: Graham (2014/2017)

- Simple model of dyadic link formation with *unobserved* agent-level degree heterogeneity
- Related to the β -model introduced earlier
- Natural generalization of modeling approaches currently used in empirical work
- Model does *not* allow for network interdependencies

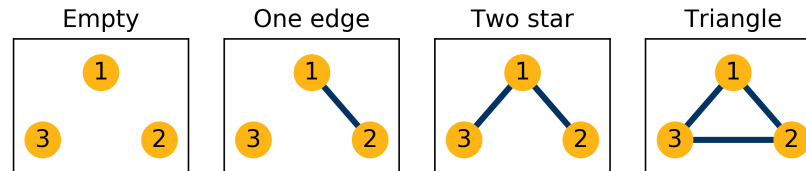
Detailed Overview of Graham (2016)

- Focuses of discriminating between *homophilous sorting* on unobserved attributes and a structural taste for *transitivity*
 - both effects generated clustered links
 - heterogeneity and interdependencies
 - timing used to side-step issues of incompleteness
- Uses availability of multiple network observations over time

Setup

- Large (sparse) network consisting of $i = 1, \dots, N$ potentially connected agents
- Observe all ties in each of $t = 0, 1, 2, 3$ periods
- \mathbf{D}_t denotes the period t adjacency matrix
 - $D_{ijt} = 1$ if agents i and j are connected in period t and zero otherwise
 - Ties are undirected: $D_{ijt} = D_{jit}$
 - No self-ties: $D_{iit} = 0$

Fact: Links are clustered



- Real world networks exhibit substantial clustering/transitivity in ties
- Transitivity indices often substantially exceed network densities

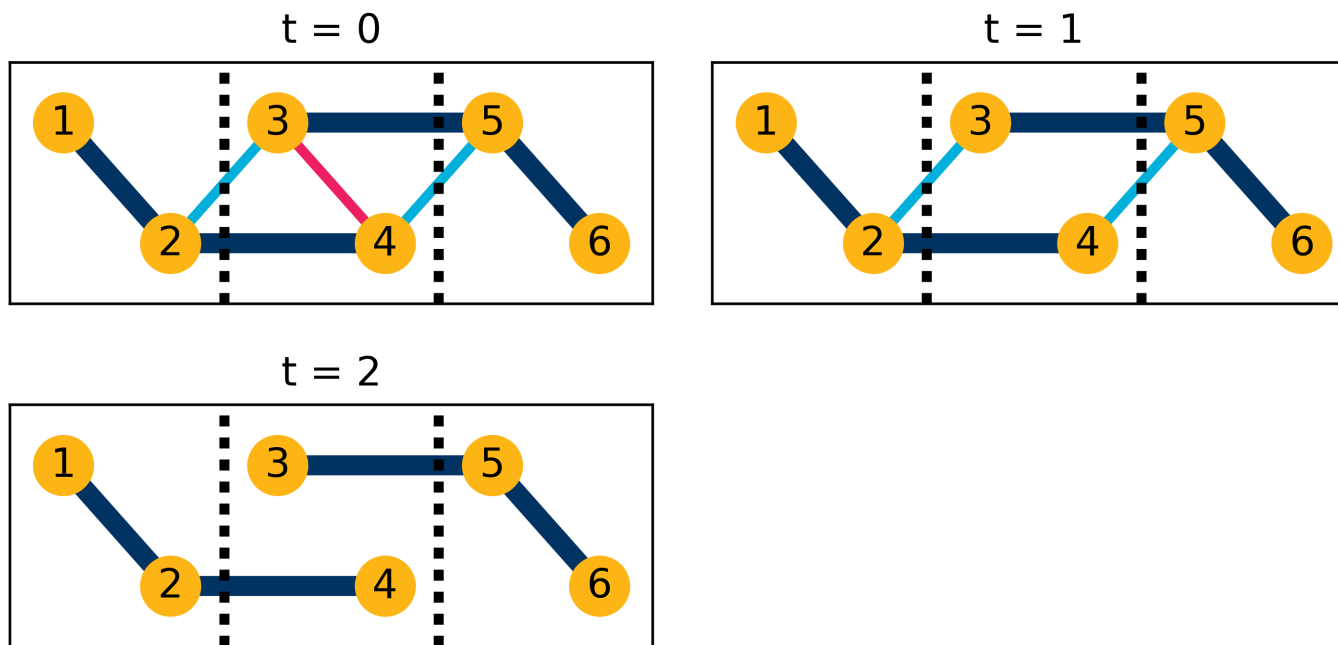
$$\begin{aligned}\rho_{CC} &= \Pr(D_{ij} = 1 \mid D_{ik} = 1, D_{jk} = 1) \\ &> \Pr(D_{ij} = 1) = \rho_D\end{aligned}$$

Homophily versus Transitivity

Two explanations for clustering:

- Homophily – ‘*birds of a feather flock together*’ (assortative mixing, community structure)
 - sorting may be on both observed and unobserved agent attributes
- (Structural) taste for transitivity (‘triadic closure’) – ‘*a friend of a friend is also my friend*’

Homophily versus Transitivity: Policy implications



Link formation model

- Agents i and j form a link in periods $t = 1, \dots, 3$ according to the rule

$$D_{ijt} = 1 \left(\beta D_{ijt-1} + \gamma R_{ijt-1} + A_{ij} - U_{ijt} > 0 \right)$$

- $R_{ijt} = \sum_{k=1}^N D_{ikt} D_{jkt}$ equals the number of period t friends i and j have in common
- $A_{ij} = A_{ji}$ is dyad-specific unobserved heterogeneity
- U_{ijt} is iid across links and over time with distribution function $F(u)$

Comments on model

Model captures three key features of link formation

1. State dependence – β
2. Structural taste for transitivity or ‘triadic closure’ – γ
3. (Time invariant) dyad-specific heterogeneity, A_{ij}
 - (a) Degree heterogeneity (van Duijn et al., 2004; Graham, *forthcoming*)
 - (b) Homophily (Assortative Mixing on *unobservables*)

Comments on model (continued)

Dyad-specific heterogeneity, A_{ij} , admits many specifications (cf., Krivitsky, Handcock, Raftery and Hoff, 2009; Zhao, Levina, Zhu, 2012).

Example #1

$$A_{ij} = v_i + v_j - g(\xi_i, \xi_j)$$

where

- v_i induces degree heterogeneity,
- $g(\xi_i, \xi_j)$ measures distance in ξ_i attribute space (assortative matching on ξ_i)

Comments on model (continued)

Example #2

$$A_{ij} = v_i + v_j + C_i' P C_j$$

- C_i a $K \times 1$ vector with a 1 in k^{th} row if i belongs to community k and zeros elsewhere (and P a $K \times K$ real symmetric matrix)

In what follows $\mathbf{A} = (A_{12}, \dots, A_{N-1N})'$ is left unrestricted

Comments on model (continued)

In each period agents take initial structure of the network as fixed when deciding whether to form, maintain or dissolve links.

- (myopic) Best-reply type dynamics (e.g., Jackson & Wolinsky, 1996)
- No completeness/coherence problems
- Measurement challenges (cf. Chamberlain, 1985; Snijders, 2011)

Comments on model (continued)

A link forms if its net surplus is positive; *utility is transferrable*

R_{ijt-1} measures opportunities to engineer 'triadic closure' or the number of triangles an agent (myopically forecasts) a period t ij link will create

If agents have a structural taste for transitivity the network will evolve in a way that fills these so-called 'structural holes'

Comments on model (continued)

The link rule specified above applies only to periods $t = 1, \dots, 3$. The *initial condition* is unspecified. Assume

$$(\mathbf{D}_0, \mathbf{A}) \sim \Pi_0$$

with \mathbf{A} denoting the $\frac{1}{2}N(N-1)$ vector of dyad-specific heterogeneity terms

Comments on model (continued)

Π_0 is unrestricted

- D_0 and A may covary
- Elements of A may also be dependent

In a single cross-section any network configuration can be generated by an appropriately chosen draw of A (graphon).

Likelihood

The joint probability density at $\mathbf{D}_0^T = \mathbf{d}_0^T$ and $\mathbf{A} = \mathbf{a}$ is:

$$\begin{aligned} p(\mathbf{d}_0^T, \mathbf{a}, \theta) &= \pi(\mathbf{d}_0, \mathbf{a}) \\ &\times \prod_{i < j} \prod_{t=1}^T F(\beta d_{ijt-1} + \gamma r_{ijt-1} + a_{ij})^{d_{ijt}} \\ &\times \left[1 - F(\beta d_{ijt-1} + \gamma r_{ijt-1} + a_{ij})\right]^{1-d_{ijt}} \end{aligned}$$

$\pi(\mathbf{d}_0, \mathbf{a})$ is the density of the ‘initial network condition’ (high dimensional nuisance parameter)

Comments on likelihood

Since \mathbf{A} is unobserved, the econometrician has three options:

- **random effects**: specify a distribution for \mathbf{A} given \mathbf{D}_0 and base inference on the corresponding integrated likelihood; also specify distribution of U_{ij}
- **joint fixed effects**: treat the $\binom{N}{2}$ components of \mathbf{A} as additional (incidental) parameters to be estimated; also specify distribution of U_{ij}
- **conditional fixed effects**: find an (identifying) implication of the model that is invariant to \mathbf{A} ; distribution of U_{ij} may or may not be specified

Comments on likelihood (continued)

First option (random effects) is difficult conceptually and computationally. Second option (joint fixed effects) will have poor statistical properties in the present setting. Third option (conditional fixed effects) is pursued here.

Comments on likelihood (continued)

- Can we learn anything about β and γ without imposing (strong) restrictions on $\pi(d_0, a)$ and/or $F(\bullet)$?
- Need an (identifying) implication of the model that is invariant to A
 - This is a high-dimensional object
 - Initial condition is also high dimensional

Comments on likelihood (continued)

- If we change the value of a single link (i, j) from, say, zero to one, many components of the likelihood may change
- Dyad-specific decisions today may alter the incentives for link formation across many other pairs in subsequent periods

Stable neighborhoods

- Idea: we can learn about the β and γ by comparing the frequency of different link histories for a given pair (i, j) holding other (local) features of the network fixed
- Problem: Changing the link history of a single (i, j) pair has effects which cascade throughout the likelihood
- Solution: Look for pairs embedded in 'stable neighborhoods'

Stable neighborhoods (continued)

The pair (i, j) are embedded in a stable neighborhood if

- all their links, except possibly those with each other, are stable across periods 1, 2, 3
- the links belonging to their friends are stable in periods 1, 2

Let $Z_{ij} = 1$ if (i, j) is a *stable dyad*: embedded in a stable neighborhood *and* $D_{ij1} \neq D_{ij2}$ and zero otherwise

Let $\mathcal{D}_s = \{\mathbf{i} \mid Z_{i_1 i_2} = 1\}$ denote the set of all stable dyads

Conditioning Set

Consider the set of network sequences

$$\begin{aligned} \mathbb{V}^s = \{ & \mathbf{v}_0^3 = (\mathbf{v}_0, \mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3) \mid \mathbf{v}_t \in \mathbb{D} \text{ for } t = 0, \dots, 3, \\ & \mathbf{v}_0 = \mathbf{d}_0, \mathbf{v}_1 + \mathbf{v}_2 = \mathbf{d}_1 + \mathbf{d}_2, \mathbf{v}_3 = \mathbf{d}_3, \\ & v_{ij1} = d_{ij1} \ \& \ v_{ij2} = d_{ij2} \\ & \text{if } z_{ij} = 0, \text{ for } i, j = 1, \dots, N \} \end{aligned}$$

\mathbb{V}^s contains all network sequences constructed by permutating the period 1 and 2 link decisions of the $\mathbf{m}_N \stackrel{def}{=} |\mathcal{D}_s|$ stable dyads.

All other link decisions are held fixed at their observed values.

The set \mathbb{V}^s contains $2^{|\mathcal{D}_s|} = 2^{\mathbf{m}_N}$ elements.

Stable neighborhoods

Permutation Lemma: For all $l \neq i, j$ let (R_{il1}^*, R_{il2}^*) denote the values of (R_{il1}, R_{il2}) after permuting D_{ij1} and D_{ij2} . If the pair (i, j) is a stable dyad, then $(R_{il1}^*, R_{il2}^*) = (R_{il2}, R_{il1})$.

- Permuting D_{ij1} and D_{ij2} does alter period 2 and 3 link incentives for other agents to which i and j are linked, but in a controlled way
- Neighborhood stability implies that $D_{il1} = D_{il2}$, so the change of incentives is entirely via transitivity effects

Stable neighborhoods (continued)

- Consider the period 2 and 3 likelihood contributions of an (i, l) pair that is linked in both periods.
- After permutation

$$\begin{aligned} & F(\beta d_{il1} + \gamma r_{il1}^* + a_{il}) F(\beta d_{il2} + \gamma r_{il2}^* + a_{il}) \\ = & F(\beta d_{il1} + \gamma r_{il2} + a_{il}) F(\beta d_{il2} + \gamma r_{il1} + a_{il}) \\ = & F(\beta d_{il2} + \gamma r_{il2} + a_{il}) F(\beta d_{il1} + \gamma r_{il1} + a_{il}) \\ = & F(\beta d_{il1} + \gamma r_{il1} + a_{il}) F(\beta d_{il2} + \gamma r_{il2} + a_{il}) \end{aligned}$$

- This coincides with the pre-permutation contribution!

Stable neighborhoods (continued)

- If i and j are embedded in a stable neighborhood, then permuting the D_{ij1} and D_{ij2} leaves
 - initial condition unaffected
 - all period 1 likelihood contributions, except those associated with (i, j) , are unaffected
 - (net) period 2 and 3 contributions from (i, l) and (j, l) dyads are unaffected (use permutation lemma)
 - period 2 and 3 contributions from all (k, l) dyads are unaffected (D_{ij1} and D_{ij2} do not enter the likelihood contributions of these pairs)

Main result: Notation

Let $S_{ij} \stackrel{def}{=} D_{ij2} - D_{ij1}$, $Q_{ij} \stackrel{def}{=} (D_{ij0}, D_{ij3}, R_{ij0}, R_{ij1})'$ and

$$b_{ij}^{01}(q_{ij}, a_{ij}, \theta) \stackrel{def}{=} \frac{1 - F(\beta d_{ij0} + \gamma r_{ij0} + a_{ij})}{F(\beta d_{ij0} + \gamma r_{ij0} + a_{ij})} \frac{F(\beta d_{ij3} + \gamma r_{ij1} + a_{ij})}{1 - F(\beta d_{ij3} + \gamma r_{ij1} + a_{ij})}$$
$$b_{ij}^{10}(q_{ij}, a_{ij}, \theta) \stackrel{def}{=} \frac{F(\beta d_{ij0} + \gamma r_{ij0} + a_{ij})}{1 - F(\beta d_{ij0} + \gamma r_{ij0} + a_{ij})} \frac{1 - F(\beta d_{ij3} + \gamma r_{ij1} + a_{ij})}{F(\beta d_{ij3} + \gamma r_{ij1} + a_{ij})}.$$

Main Result (continued)

The conditional likelihood of $D_0^3 = \mathbf{d}_0^3$ given $\mathbf{d}_0^3 \in \mathbb{V}^s$,

$$l^c(\mathbf{d}_0^3, \mathbf{a}, \theta) = \frac{p(\mathbf{d}_0^3, \mathbf{a}, \theta)}{\sum_{\mathbf{v} \in \mathbb{V}^s} p(\mathbf{v}_0^3, \mathbf{a}, \theta)}, \quad (1)$$

equals

$$l^c(\mathbf{d}_0^3, \mathbf{a}, \theta) = \prod_{\mathbf{i} \in \mathcal{D}_s} \left[\frac{1}{1 + b_{i_1 i_2}^{01}(q_{ij}, a_{ij}, \theta)} \right]^{1(s_{i_1 i_2} = 1)} \\ \times \left[\frac{1}{1 + b_{i_1 i_2}^{10}(q_{ij}, a_{ij}, \theta)} \right]^{1(s_{i_1 i_2} = -1)}$$

- Denominator in (1) is a summation over $2^{\mathbf{m}_N}$ elements

Main Result (continued)

- ...surprisingly this sum is not intractable (“binomial theorem”).
- The ratio (1) can be expressed as a product of just m_N terms!

Main Result (comments)

An unexpected byproduct of conditioning is (conditional) independence.

Link histories of stable dyads are conditionally independent!

Distribution of U_{ij} unspecified \Rightarrow maximum score approach to estimation (Manski, 1975, 1987; Honore and Kyriazidou, 2000)

If U_{ij} is logistically distributed, then \mathbf{A} doesn't enter the conditional likelihood; criterion function takes familiar logit form

Nonparametric case

Under the data generating process specified above

$$\Pr(D_{ij1} = 0, D_{ij2} = 1 \mid Q_{ij} = q, Z_{ij} = 1) - \Pr(D_{ij1} = 1, D_{ij2} = 0 \mid Q = q, Z_{ij} = 1) \begin{matrix} \leq \\ \geq \end{matrix} 0$$

according to whether

$$\beta(d_3 - d_0) + \gamma(r_1 - r_0) \begin{matrix} \leq \\ \geq \end{matrix} 0.$$

cf. Manski (1987); suggests the following estimator:

$$\sup_{\theta: \|\theta'\theta\|=1} \binom{N}{2}^{-1} \sum_{i=1}^N \sum_{j < i} Z_{ij} (D_{ij2} - D_{ij1}) \operatorname{sgn} \{X'_{ij}\theta\} \quad (2)$$

for $x = (d_3 - d_0, r_1 - r_0)'$.

Logit case

When the idiosyncratic component of surplus U_{ijt} is logistic

$$\Pr \left(D_{ij1} = d_1, D_{ij2} = d_2 \mid Q_{ij} = q, Z_{ij} = 1 \right) = \left(\frac{\exp(x'\theta)}{1 + \exp(x'\theta)} \right)^{1(s=1)} \left(\frac{\exp(x'\theta)}{1 + \exp(x'\theta)} \right)^{1(s=-1)}$$

Note: A_{ij} does not enter to the right of the equality (\Rightarrow point identification up to scale)

Logit case (continued)

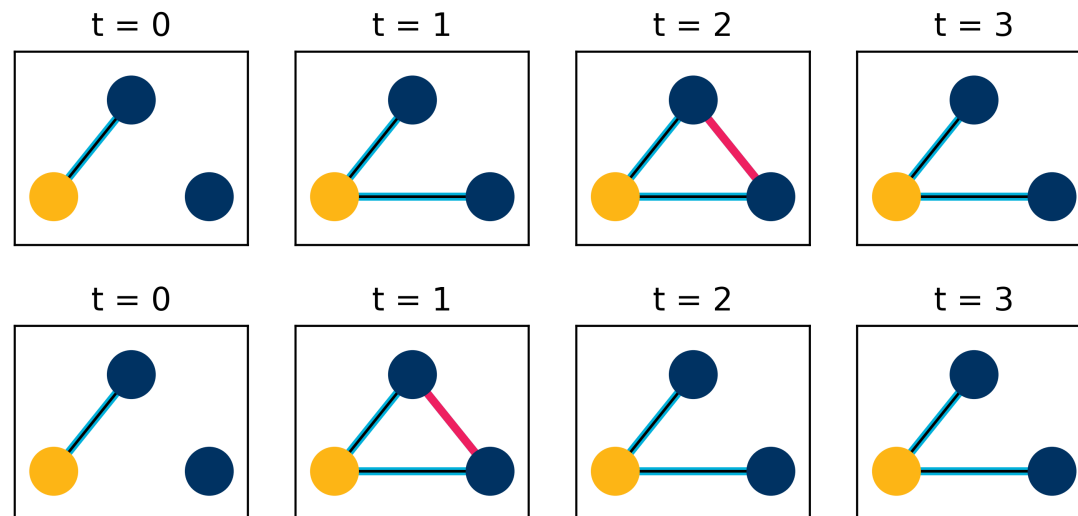
The *stable neighborhood logit* estimate of θ_0 is the maximizer of

$$L_N(\theta) = \binom{N}{2}^{-1} \sum_{i=1}^N \sum_{j < i} l_{ij}(\theta)$$

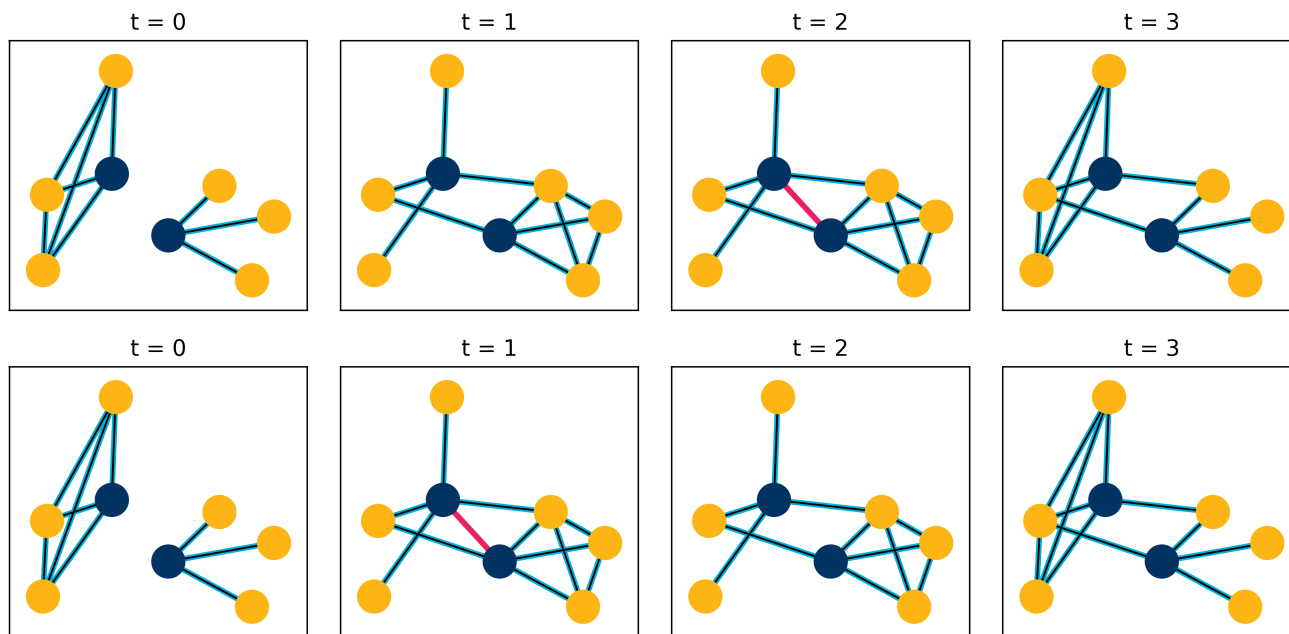
with

$$l_{ij}(\theta) = Z_{ij} \left\{ S_{ij} X'_{ij} \theta - \ln \left[1 + \exp \left(S_{ij} X'_{ij} \theta \right) \right] \right\}.$$

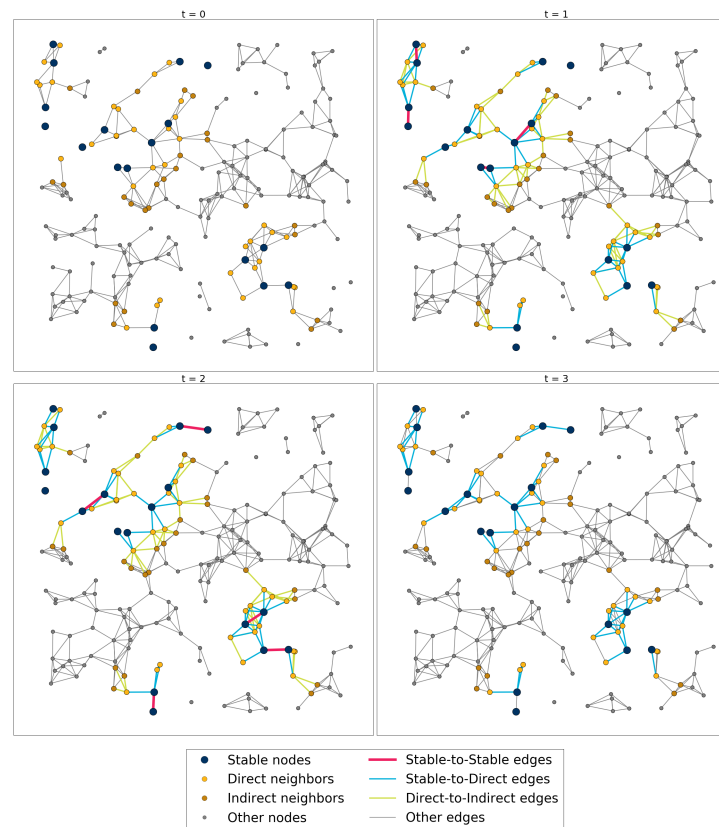
Stable neighborhood example



Stable neighborhood example



Stable neighborhoods in large network



Monte Carlo

Agents are scattered uniformly on the two-dimensional plane

$$\left[0, \sqrt{N}\right] \times \left[0, \sqrt{N}\right].$$

Initial network is generated according to

$$D_{ij0} = \mathbf{1} \left(A_{ij} - U_{ij0} \geq 0 \right),$$

with U_{ij0} logistic and A_{ij0} taking one of two values.

Monte Carlo (continued)

1. If the Euclidean distance between i and j is less than or equal to r , then $A_{ij0} = \ln\left(\frac{0.75}{1-0.75}\right)$, otherwise $A_{ij0} = -\infty$
2. Agents less than r apart link with probability 0.75, while those greater than r apart link with probability zero.

Network in $t = 1, 2, 3$ generated using link rule with $\beta = \gamma = 1$ and U_{ijt} logistic.

Properties of simulated networks

| Asymptotic Degree | 4 | | |
|-------------------|------------------------------|------|------|
| Period | $(N - 1) \mathbb{E}[D_{it}]$ | T | GC |
| $t = 0$ | 3.94 | 0.44 | 0.58 |
| $t = 1$ | 4.98 | 0.58 | 0.83 |
| $t = 2$ | 5.12 | 0.59 | 0.84 |
| $t = 3$ | 5.14 | 0.59 | 0.85 |

Notes: The table reports period-specific network summary statistics across the $B = 1,000$ Monte Carlo simulations for each design ($N = 5,000$). See paper for other design details. The $(N - 1) \mathbb{E}[D_{it}]$ column gives the average degree, T the global clustering coefficient or transitivity index and GC the fraction of agents that are part of the largest giant component.

Sampling properties of SN logit

| Asymptotic Degree | 4 | |
|------------------------|---------|----------|
| $N = 5,000$ | β | γ |
| Mean | 1.0438 | 1.0456 |
| Median | 1.0410 | 1.0133 |
| Std. Dev. | 0.4575 | 0.2976 |
| Mean Std. Err. | 0.4493 | 0.2917 |
| Coverage | 0.9620 | 0.9650 |
| Avg. # of Stable Dyads | 110.6 | |
| # of cvg. failures | 1 | |

Final Thoughts

- The availability of multiple observations of a network over time is potentially very informative
- Fruitful to compare the relative frequency of certain sequences of link formation for a given pair, holding the link history of other pairs fixed
- Consistent estimation using a single (large sparse) network is possible (primitive conditions for $\binom{N}{2}\alpha_N \rightarrow \infty$ with $\alpha_N = \Pr(Z_{ij} = 1)$).

Final Thoughts (continued)

- 'Fixed effect' identification analysis can also help formulate more realistic random effects models (cf., Goldsmith-Pinkham and Imbens, 2013)
- Computational challenge: efficient algorithm to find all stable dyads
- Directed graphs, covariates, efficiency bound, empirical application...