

Social and Information Networks

University of Toronto CSC303
Winter/Spring 2023

Week 3: January 23-27

Chapter 4: The context of network formation

- In this chapter, we study social networks within their context, considering factors outside of the nodes and edges of the network that impact how the network structure evolves.
- The chapter introduces a very important (and often controversial) issue, namely the relative roles of selection (similarity) vs influence in social relations.
- As we have already noted, Easley and Kleinberg have already indicated that there is a limit to what one can understand just in terms of the network structure.

Word of caution from Chapter 3 repeated

Easley and Kleinberg (end of Section 3.3):

Given the size and complexity of the (who call whom) network, we cannot simply look at the structure. . . Indirect measures must generally be used and, because one knows relatively little about the meaning or significance of any particular node or edge, it remains an ongoing research challenge to draw richer and more detailed conclusions. . .

We should also add that we may know very little about the reasons for the formation (or disappearance) of an edge.

Unknown:

In theory there is no difference between theory and practice. In practice there is.

This week's agenda

- Homophily
 - ▶ Mutable and Immutable factors
 - ▶ Selection vs. Influence
 - ▶ Testing for Homophily
- Schelling Segregation model
- Social-Affiliation Networks
 - ▶ Triadic, Focal, and Membership closure
- Empirical Studies
 - ▶ Trends in closure probabilities
 - ▶ Selection and influence in Wikipedia editors

Homophily

- Why **triadic closure**? In Chapter 3: some network “**intrinsic**” reasons (opportunity, trust, incentive) for forming a friendship and now we consider “**contextual**” reasons
- **Homophily**: we tend to be similar to our friends.
- This observation is captured in various writings and proverbs perhaps most notably by “**Birds of a feather flock together**” suggesting that friendships (and membership in groups) are selectively formed due to similar interests.
- **Note**: But to what extent do we adopt similar interests based on friendship rather than conversely?

Characteristic factors

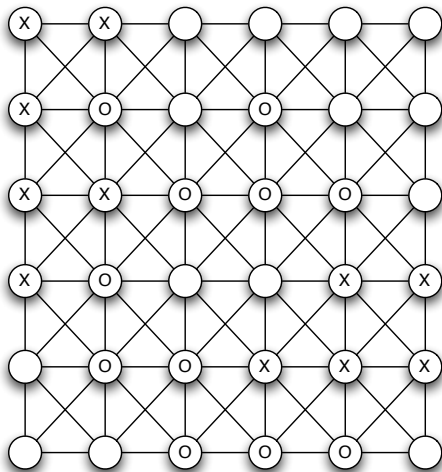
- **Factors** which help determine our friendships and relations can be **immutable** or **more transient**.
- Some (relatively) **immutable factors**: ethnicity, birth date, gender; religion, height. What other such (mainly permanent) factors exist?
- Some more **mutable (often related) factors**: membership in clubs or courses, educational level, recreational interests, professional interests, income level, residential neighbourhood
- Of course, immutable factors can and do **influence** mutable factors.

The Schelling Model

- A dramatic example of homophily is Schelling's Segregation model (end of Chapter 4)
 - ▶ Model motivated by racially segregated neighbourhoods
 - ▶ Highly simplified
- Schelling's model and his simulations led him to a fundamental observation:
 - ▶ Segregation can and will happen even if there is no explicit individual desire to avoid (say) people of a different race. All that is needed is some desire to be near enough similar people
- This observation isn't restricted to racial segregation
- The model also provides an interesting study of network dynamics, homophily driven by selection, and how local decisions can lead to global structure
- Importantly, Schelling's model does not preclude the presence of other economic and political factors, nor does it preclude explicit racism!

The Schelling model

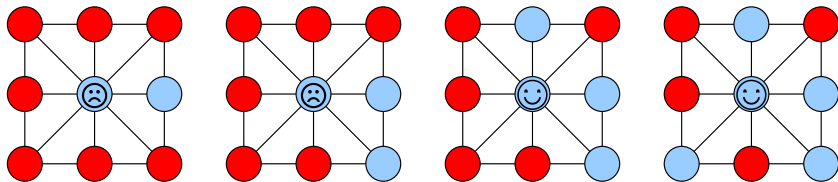
The model itself is quite simple but still hard to analyse analytically. In this model, we view two classes of individuals (X and O) living in a grid. More specifically, individuals occupy some subset of the nodes as depicted in figure 4.15 of the text.



The dynamics of the Schelling model

- Schelling hypothesizes that every individual wants to have at least some threshold t of their immediate neighbours be of the same class

Example: Assume threshold, $t = 3$:



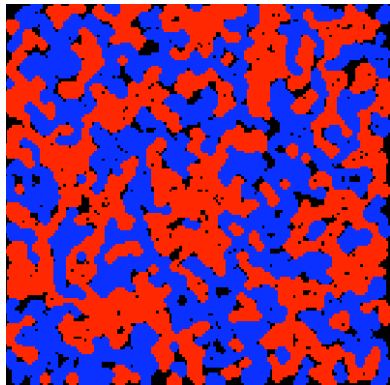
- When a individual's threshold is not met, they move
- Variants of the model alter the order in which individuals move, and where they randomly move to in order to satisfy the desire for similarity
 - ▶ The claim is that the results do not change qualitatively

The dynamics of the Schelling model

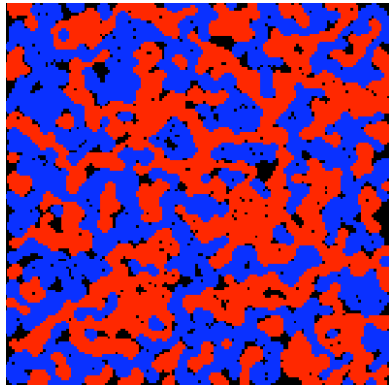
To observe the dynamics, simulations of the network are conducted for different threshold values. What is very apparent is the segregated structure of the network as it evolves.

The following figures show the results for thresholds $t = 3$ (i.e. an individual desires less than a majority of his/her neighbours to be similar) and $t = 4$. The grid is a 150 by 150 (i.e., 12,500 cells, with 10,000 cells occupied), and both groups are equally represented.

Simulations for $t = 3$



(a) *A simulation with threshold 3.*



(b) *Another simulation with threshold 3.*

Figure 4.17: Two runs of a simulation of the Schelling model with a threshold t of 3, on a 150-by-150 grid with 10,000 agents of each type. Each cell of the grid is colored red if it is occupied by an agent of the first type, blue if it is occupied by an agent of the second type, and black if it is empty (not occupied by any agent).

Simulation for $t = 4$

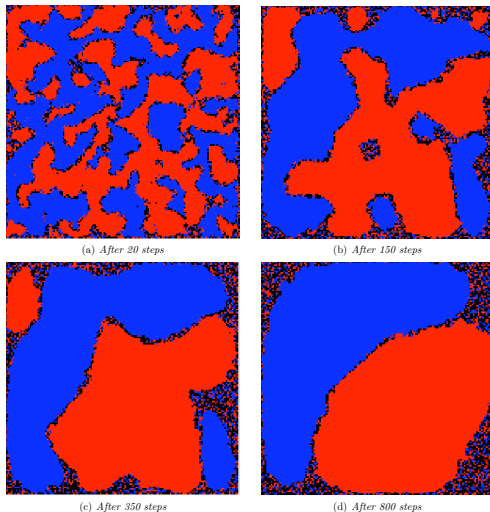


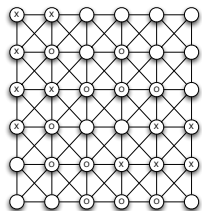
Figure 4.19: Four intermediate points in a simulation of the Schelling model with a threshold t of 4, on a 150-by-150 grid with 10,000 agents of each type. As the rounds of movement progress, large homogeneous regions on the grid grow at the expense of smaller, narrower regions.

Extensions to the Schelling model

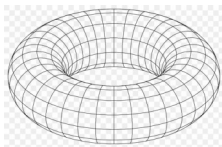
- Rogers & McKane published a work accessible here <https://arxiv.org/abs/1104.1971> defining a more general model, analyzing a simplified model, and demonstrating similar behaviour to multiple other variants
- They assumed individuals lived on an arbitrary undirected graph with nodes $V = \{1 \dots n\}$
- $\sigma \in \{-1, 0, 1\}^n$ defined the state of their world
 - ▶ -1 signifies occupied by individual of type A
 - ▶ 0 signifies empty
 - ▶ 1 signifies occupied by individual of type B
- $\sigma^{(ij)}$ is σ with the entries i and j swapped
- They assumed that each timestep a pair of nodes is chosen at random and their occupants swapped – this is done with probability T_{ij} , subject to some constraints
- $P(\sigma, t + 1) = \sum_{ij} P(\sigma^{(ij)}, t) T_{ij}(\sigma^{(ij)})$

Extensions to the Schelling model

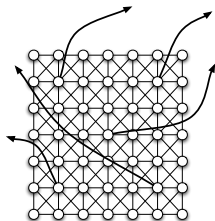
- They compared their analysis to various pre-existing models that could be expressed in this general form
- The graphs considered were:



Grid; E&K Fig 4.15



Toroidal Grid. Image
by Adith George
1840427, CC BY-SA
4.0, via Wikimedia
Commons

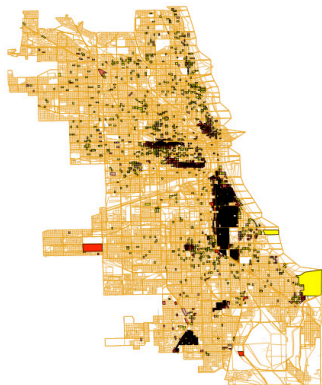


Small world network.
E&K Fig 20.3

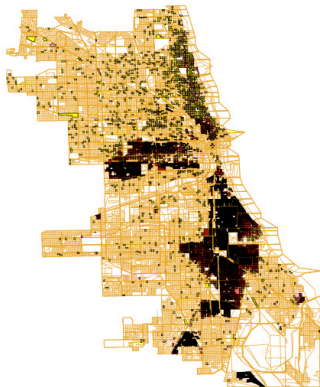
Some concluding comments on the Schelling study

- The model is not constructed so as to build in segregation. More specifically, the model allows for stable configurations that are well integrated.
- However, given a random starting configuration, the simulations show that people will gravitate to a segregated structure.
- There is a compounding effect of the model dynamics. Namely, when one person leaves, it can result in other neighbours falling below their threshold and hence a new desire to leave the current location has been created.
- The word segregation is a term with a very negative connotation due to the use of the term with respect to racial (e.g., Jim Crow legislation in the US) and religious segregation (e.g., ghettos in Europe) which was forced by governments.
- What if we used the word "clustered" instead of segregated? Do we think that neighbourhoods that are concentrated along say ethnic lines is a bad thing? Be careful to distinguish between the model's assumptions and reality.

The reality of neighbourhood segregation in Chicago (1970s)



(a) *Chicago, 1940*



(b) *Chicago, 1960*

And how integrated or segregated is Toronto?

It seems much easier to talk about Chicago (as we all know about racial segregation in the US) but perhaps more difficult to talk about Canada and Toronto.

At some level (i.e., Metro Toronto), Toronto may be the most ethnically diverse city as is claimed. But at a more detailed level, many neighbourhoods are far from being “integrated”. I am posting a September 2018 newspaper article and a February 2019 talk by David Hulchanski (UT Faculty of Social Work) that describes the changes in income levels and neighbourhoods in Toronto. The title of the article more or less summarizes his major conclusion: “Toronto is segregated by race and income. And the numbers are ugly”.

I call attention to this (and other similar studies) to indicate that while homophily is a factor (especially with regard to ethnicity), there are clearly many other factors that are prevalent and arguably dominant.

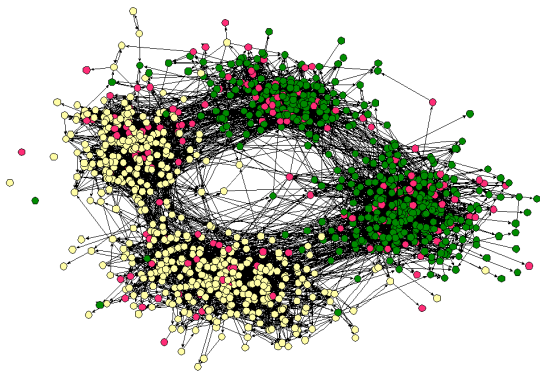
The influence vs selection issue

- Thusfar, we've seen that immutable factors can and do **influence** mutable factors.
- Furthermore, one's friendships can and do **influence** mutable factors such as say recreational interests.
- So the selection vs influence issue can be seen as the relative extent to which our friendships are formed selectively due to similarity vs friendships influencing our interests and other traits.
- Homophily refers to the **correlation** between friendships and similarity; should it be attributed to friendships causing similarity, or vice versa?
 - ▶ For immutable factors: more easily attributed (directly or indirectly) to similarity leading to friendships
 - ▶ For mutable factors: much less clear, may be quite controversial!

The influence vs selection issue

- Further complicating matters, the “environment” of various (perhaps unobserved) external events or hidden influences can also impact one’s friendships and/or interests and affiliations.
- For example, Alejandro and Betty are not friends nor have any interest in political issues. Then a popular entertainer is performing in a rally for a political candidate. Alejandro and Betty meet at the event and become friends as well as becoming more politically involved.
 - ▶ Therefore, this homophily is caused by... neither selection nor influence; the rally is a confounding variable

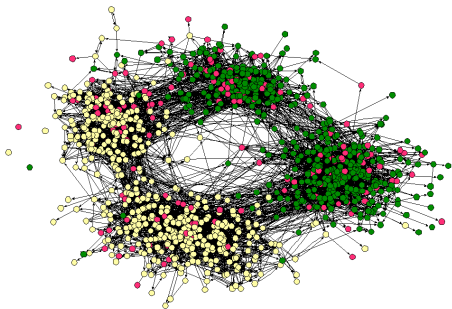
Graphic visualization of homophily



[Fig. 4.1, textbook]

- Homophily can divide a social network into **densely-connected, homogeneous parts that are weakly connected to each other.**
- In this social network from a town's middle school and high school, **two divisions** are apparent: one based on **race** (students of different races drawn as differently-colored circles), and the other based on friendships in the **middle and high schools.**

Comments on figure 4.1



[Fig. 4.1, textbook]

- Such a visualization is not at a scale that one can see most of the individual relations. The visualization clearly shows homophily based on race and the junior/senior high split (both immutable factors).
- We can measure the extent of homophily (as we will next see) but observing any such phenomena (even for immutable factors) is just the **starting point** in truly understanding the phenomena.
- The figure does show some detailed information; i.e. individuals without any friends (isolated nodes) or with few friends (low degree).

Measuring homophily

- As mentioned before, when networks are large (and/or when homophily is less dramatic) it is difficult if not impossible to visualize various aspects of a network and so one needs a **measure of homophily** (whatever the cause or the consequence of the network).
- **How might we measure homophily?**
- **Think Big!**: Lets think in terms of large social networks where the presence or absence of a given individual will not have any noticeable impact
- Suppose we wish to study the **likelihood of friendships** according to some factor (with say two values). For instance, high school vs. middle school (i.e., junior vs. senior students).

Thought experiment

- What would it mean to say that a social network does or does not exhibit homophily according to some factor such as school?
- For our network, let the fraction (i.e. probability) of middle school students be p and the fraction of high school students be q .
 - ▶ Consider a given edge (u, v) in the network.
 - ▶ If school has no correlation with relations, then the probability that the schools of u and v are different is $2pq$. Why?
- This leads to a **homophily test**: If the actual fraction of cross-school edges is “significantly less than” $2pq$ then there is evidence for homophily.
What would this say about same school (middle-middle) or (high-high) edges?
- Clearly the meaning of an edge is an essential aspect of any study; e.g. consider the difference between an edge representing collaboration in a course project vs an edge meaning a romantic relationship.

Reviewing selection vs social influence

- With **immutable factors** (such as race and for the most part gender), when we observe evidence of homophily, we often attribute increased friendships to **selection**, which is the tendency to form friendships with others who are like you in some way(s). (But note that race often correlates with neighbourhoods or academic programs.)
- But when considering more **mutable factors**, there is a feedback between similar characteristics and social links.
 - ▶ To what extent does behaviour get modified by our social network?
 - ▶ That is, to what extent is **social influence** determining interests and behaviour?
- Of course, both selection and social influence can be interacting in the same social network. How does one understand the relative interplay?

Longitudinal studies may make it possible to see the behavioral changes that occur after changes in an individual's network connections, as opposed to the changes to the network that occur after an individual changes his or her behavior.

A study of similarity and interaction

We point ahead to a study by Crandall et al [2008] that suggests that in certain settings, it may be possible to gain some insight into the selection vs influence issue. We return to this study later in lecture.

Using Wikipedia data, the text presents one study that speaks to the manner in which selection and influence combine to result in observed homophily. The nodes are Wikipedia editors, and edges correspond to communication via a user-talk page for a Wikipedia page. So we know what the graph means and can observe the emergence of edges over time.

The study defines a numerical similarity measure between two users A and B as a small variation on the following ratio which is analogous to the way neighbourhood overlap was defined:

$$\frac{\text{number of articles edited by both } A \text{ and } B}{\text{number of articles edited at least one of } A \text{ or } B}$$

Fortunately, every action on Wikipedia is recorded and time-stamped so it is possible to conduct a meaningful longitudinal study by looking at each “time step” defined by an “action” by either of the editors where an action is either an article edit, or a communication.

Average level of similarity before and after the first Wikipedia communication

The figure below plots the level of similarity as a function of the number of edits before and after the first communication. Time 0 is defined to be the time of the first interaction between a pair (A, B) of editors. This is then averaged over all the (A, B) plots.

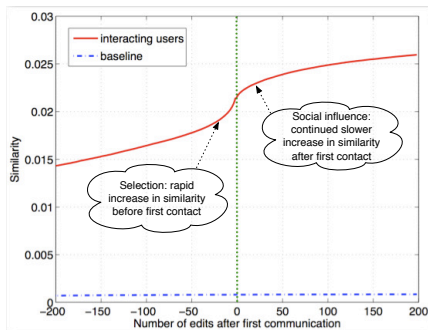


Figure: [E&K, Fig 4.13]

Two interesting and opposing longitudinal studies

- In academic success (or drug usage) in teenage friendship networks, Cohen (1977) and Kandel (1978) claim that peer pressure (i.e. **social influence**) is less a factor here than previously believed. We can speculate that (for example) similar family environments is a significant determining factor for such behaviour amongst friends.
- In contrast to the above example, in a controversial report on obesity patterns of 32,000 people observed over a 32 year period, Christakis and Fowler (2007) claim: **obesity** or keeping fit is (perhaps surprisingly) to some extent **a contagious disease spread within a social network**. “You don’t necessarily catch it from your friends the way you catch the flu, but it nonetheless can spread through the underlying social network via the mechanism of social influence.” (Later in the course we will discuss models for the spread of influence in a network.)

Why the obesity homophily?

- Three possibilities identified by Christakis and Fowler:
 - ① [1] selection
 - ② [2] homophily being driven by other factors that correlate with obesity (e.g. poverty)
 - ③ [3] the social influence of peer pressure say as in the case of drug use or academic performance or fitness.

- Christakis and Fowler conclude that even accounting for [1] and [2], social influence is a significant factor.
Aside: I am not sure as to the extent that they consider the relative role of genetics vs diet.

- Once again, we caution that observing homophily is clearly only a starting point.

Why do we care?

- Why do we care about the relative interplay (selection vs. social influence) and how could we model this?
- If indeed social influence is a significant factor, then targeting key individuals and trying to modify undesirable behaviour (or promote positive behaviour) can be effective since we are then viewing such behaviour as a process of influence spread.
- If not, focusing on a few individuals will at best change the behaviour of a few individuals.

Recap

- Homophily
 - ▶ Mutable and Immutable factors
 - ▶ Selection vs. Influence
 - ▶ Testing for Homophily
- Schelling Segregation model

Social-affiliation networks: incorporating context into the network

- Up to now we have viewed contextual (mutable and immutable) factors that affect the formation of links to be outside of the social network being considered.
- Section 4.3 discusses how to include context in the network so as to have a common framework for studying the interplay between the extent of (social) triadic closure (common friendships induce new friendships), homophily determined by selection, and mutual activity determined by social influence.
- Let's consider the (mutable) context of **affiliation** in a group/participation in an activity. Such an activity is referred to as a **foci**, a focal point for social interaction.
- We incorporate such foci into social networks by considering a focus to be a different type of node, distinct from a node representing an individual. We first consider a pure **affiliation network**, an example being of which we have already seen in a bipartite graph with individuals and corporate boards.

Example of a pure affiliation network

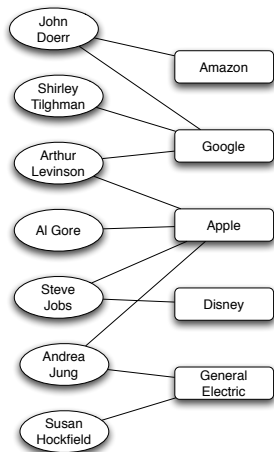


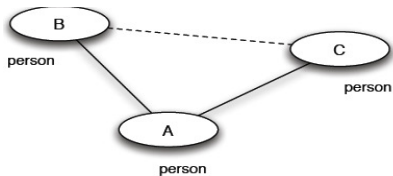
Figure: [E&K, Fig 4.4] One type of affiliation network that has been widely studied is the memberships of people on corporate boards of directors. A very small portion of this network (as of mid-2009) is shown here.

Social-affiliation networks continued

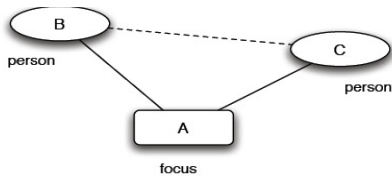
We can then combine the people-people edges of a social network with the people-focus edges of an affiliation network to form a **social-affiliation network**. Within such a combined network, we can discuss three types of graph triangle closures:

- **triadic closure** as introduced in chapter 3 where common friends of one or more individuals become friends
- **focal closure** where individuals become friends based on their common interest(s)
- **membership closure** where an individual joins an activity because a friend (or a group of friends) is (are) already in that activity

Three types of closure

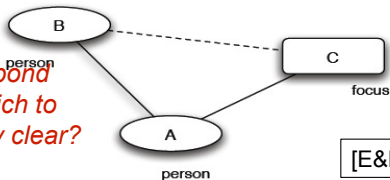


(a) *Triadic closure*



(b) *Focal closure*

Which of these correspond to social influence, which to selection? Is it still fully clear?



(c) *Membership closure*

[E&K, Ch.4, Fig. 4.6]

Figure: [E&K, Fig 4.6] Three types of closure

Toy example of a social-affiliation network

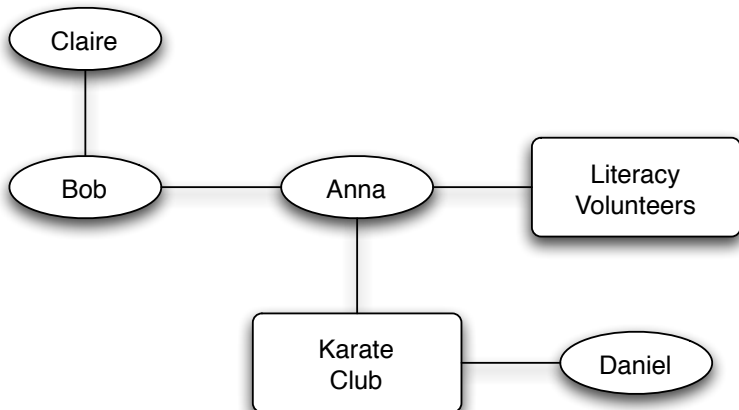


Figure: [E&K, Fig 4.5] In this social-affiliation network, the oval nodes are people and the rectangular nodes are activities. What kinds of triangular closures can occur?

Toy example showing three types of closure

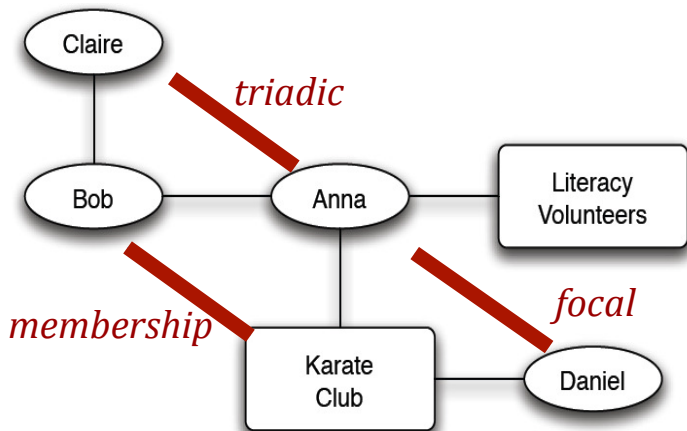


Figure: [E&K, Fig 4.7] We can observe the three types of triangular closures that have occurred in some time period.

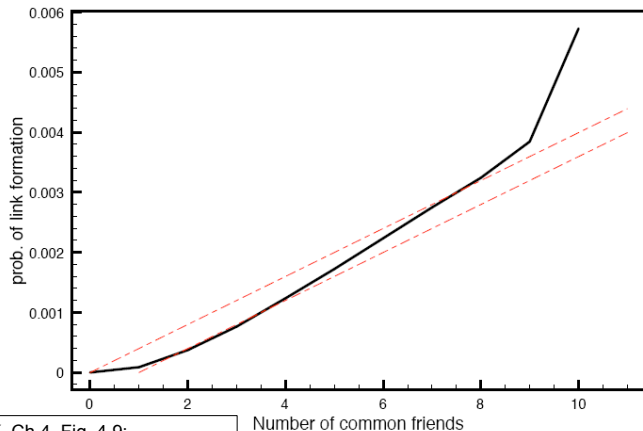
Empirically measuring these processes

- Closure is inherently **dynamic**
 - ▶ So we need to **take snapshots of the network at different times** to see how the relationships evolve and to what extent each form of closure occurs
 - ▶ If common friends or common interests are causing new links (i.e., closures) then **the more friends or interests in common, the more we should see this effect.**
- We briefly look at a couple studies stemming from online interactions, but realize the usual warning about limitations of such studies
 - ▶ As in all modeling we may be missing many factors
 - ▶ The timing of the snapshots may influence results
 - ▶ These particular studies look at link formation, but not link dissolution. What would the network look like if links formed but never dissolved?

Triadic closure: dependence on number mutual friends

- Email exchanges (over a year) by 45,000 students & staff in large US university [Kossinets, Watts 2006]
- “Friends” defined as two-way email communication (prev. 60 days)
- Measure probability $T(k)$ of a new friendship emerging between a pair of students as a function of the number k of mutual friends
- That is, the probability of it happening in any given day (averaging over many such pairs)
- Compare data (black) with baseline theoretical model (red) baseline: **assume** any single mutual friend will generate a new friendship with probability p and that this will happen *independently* for each common friend. Thus $T(k) = 1 - (1 - p)^k$ **Why?**
- For **small** p , $(1 - p)^k \approx 1 - pk$ so that $T(k) \approx pk$.

Probability (per-day) of triadic closure as a function of the number of common friends



[E&K, Ch.4, Fig. 4.9;
from Kossinets and Watts, 2006]

Figure: [E&K, Fig 4.9]

Observations

- Data does not show much more propensity for friendship when going from zero to one mutual friend.
 - ▶ The second dashed red line shifts the curve over by one friend so as to better compare the actual data and baseline model.
 - ▶ Why no major impact with one common friend?
- Increasing from 1 to 9 friends shows linear curve (greater slope than baseline)
- A sharp difference going beyond 9 friends
 - ▶ The theoretical model (and its assumption of independence) no longer supported.
 - ▶ Is there some threshold of mutual friends which escalates the pressure for triadic closure?

Exercise: translate per-day probability into per-month or per-year probability

Probability of focal closure as a function of the number of common classes

Kossinets and Watts also studied focal closure where a focus means a class in which a student is enrolled.

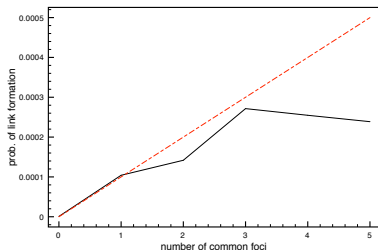


Figure: [E&K, Fig 4.10]

Clearly the theory and the actual data do not correspond especially when considering students going from 3 to 4 common classes. **Can you speculate on a reason?** If you haven't formed a friendship having attend 3 classes together, then perhaps there is a reason?

Probability of membership closure as a function of the number of common friends

The text presents two studies of membership closure where there is data concerning both person-to-person interactions and person-foci affiliations. The first study shows the probability of joining a community in the blogging site LiveJournal where “friendship” is self-identified within a user’s profile.

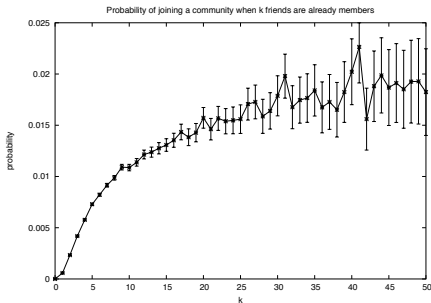


Figure: [E&K, Fig 4.11]

Second study of membership closure as a function of the number of common friends

The second study concerns Wikipedia editors and foci are specific Wikipedia pages. Here “friendship” is defined as having communicated together on a user-talk page and membership in a foci corresponds to having edited a Wikipedia page.

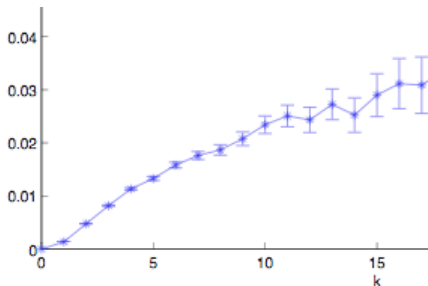


Figure: [E&K, Fig 4.12]

The interplay between selection and influence

Using the same Wikipedia data as in the previous focal closure example, The text presents one study that speaks to the manner in which selection and influence combine to result in observed homophily. Once again, the nodes are Wikipedia editors, the foci are articles, and edges correspond to communication via a user-talk page.

In addition, the study defines a numerical similarity measure between two users A and B as a small variation on the following ratio which is analogous to the way neighbourhood overlap was defined:

$$\frac{\text{number of articles edited by both } A \text{ and } B}{\text{number of articles edited at least one of } A \text{ or } B}$$

Fortunately, every action on Wikipedia is recorded and time-stamped so it is possible to conduct a meaningful longitudinal study by looking at each “time step” defined by an “action” of an editor where an action is either an article edit, or a communication.

Average level of similarity before and after the first Wikipedia communication

The figure below plots the level of similarity as a function of the number of edits before and after the first communication. Time 0 is defined to be the time of the first interaction between a pair (A, B) of editors. This is then averaged over all the (A, B) plots.

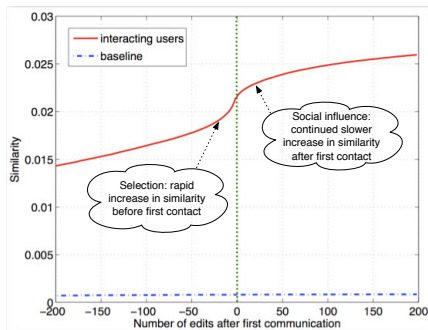


Figure: [E&K, Fig 4.13]

Observations on similarity vs. interactions (Figure 4.13)

There are a number of interesting observations and caveats regarding Figure 4.13. First some notable observations.

- The level of similarity is increasing over “time” before and after the first interaction.
- The steepest increase in similarity occurs just before the first interaction suggesting that selection is playing a pronounced role in forming this “friendship link” in the networks that are being dynamically created.
- The bottom dashed line indicates the level of similarity for those who never communicate. Clearly those who eventually interact are more similar.

Some caveats

- Like any averaging of individual data, we cannot say why any particular pair of editors have decided to communicate.
- Because the defined time 0 corresponds to different moments in “real time” for each pair, we cannot understand to what extent real time events may also be a factor leading communication.
- In this study, links are never eliminated. Other “fully dynamic” network settings would have node and/or links that are not permanent.
- The biggest question about such a study is the extent to which any observations may or may not extend to different settings. In what settings do we have the same kind of detailed time stamping of events?

Recap

- Homophily
 - ▶ Mutable and Immutable factors
 - ▶ Selection vs. Influence
 - ▶ Testing for Homophily
- Schelling Segregation model
- Social-Affiliation Networks
 - ▶ Triadic, Focal, and Membership closure
- Empirical Studies
 - ▶ Trends in closure probabilities
 - ▶ Selection and influence in Wikipedia editors