

MU5IN075

Network Analysis and Mining

6. SNA and spreading processes

Esteban Bautista-Ruiz, Lionel Tabourier

LIP6 – CNRS and Sorbonne Université

first_name.last_name@lip6.fr

October 19, 2021

1/18

Outline

- 1 Applying our tools to Social Network Analysis
 - The homophily phenomenon
 - Local density and community structure
 - Impact on social contagion
- 2 Epidemic spreading models on graphs
 - Compartmental models in epidemiology
 - What networks bring to the models

2/18

About this section

From Part A (Complex Networks Analysis tools)
→ to Part B (Graph mining)

- use our knowledge and tools to explore real networks
- adapt them to specific problems

→ **problem-oriented view**

Illustration: a few Social Network Analysis concepts

3/18

About this section

From Part A (Complex Networks Analysis tools)
→ to Part B (Graph mining)

- use our knowledge and tools to explore real networks
- adapt them to specific problems

→ **problem-oriented view**

Illustration: a few Social Network Analysis concepts

3/18

What is homophily?

From Greek, *homo*: same, similar and *philos*: friend of, to like
→ “birds of a feather flock together”

Observed for a long time in sociology:

- smoking habits, food habits
- residential segregation
- voting behavior
- ...

How to observe this phenomenon through SNA?

What is homophily?

From Greek, *homo*: same, similar and *philos*: friend of, to like
→ “birds of a feather flock together”

Observed for a long time in sociology:

- smoking habits, food habits
- residential segregation
- voting behavior
- ...

How to observe this phenomenon through SNA?

Measuring homophily in a network

Consider a quantitative property related to a node:
age, revenue, number of children, number of connections ...

- q_k probability distribution of this property in the network
- e_{jk} joint probability of this property considering edges

Correlation coefficient across ties r

$$r = \frac{\sum_{j,k} j \cdot k \cdot (e_{jk} - q_j q_k)}{\sigma_q^2}$$

$$\sigma_q^2 = \frac{1}{n} \sum_k q_k (k - \bar{k})^2 \text{ is the variance of } q_k$$

- r is a normalized quantity ($\in [-1 : 1]$)
- $r \simeq 1$: strong correlation, $r \simeq -1$: strong anti-correlation
- $r \simeq 0$: no correlation

Measuring homophily in a network

Consider a quantitative property related to a node:
age, revenue, number of children, number of connections ...

- q_k probability distribution of this property in the network
- e_{jk} joint probability of this property considering edges

Correlation coefficient across ties r

$$r = \frac{\sum_{j,k} j \cdot k \cdot (e_{jk} - q_j q_k)}{\sigma_q^2}$$

$$\sigma_q^2 = \frac{1}{n} \sum_k q_k (k - \bar{k})^2 \text{ is the variance of } q_k$$

- r is a normalized quantity ($\in [-1 : 1]$)
- $r \simeq 1$: strong correlation, $r \simeq -1$: strong anti-correlation
- $r \simeq 0$: no correlation

Structural homophily, based on degree

Do high-degree nodes connect to high-degree nodes?

Assortativity (degree-based homophily)

q_k probability distribution of the *remaining degree*

We measure r called here **assortativity coefficient**:

$$r = \frac{\sum_{j,k} j \cdot k \cdot (e_{jk} - q_j q_k)}{\sigma_q^2} \text{ with } \sigma_q^2 = \frac{1}{n} \sum_k q_k (k - \bar{k})^2$$

r is the ratio between covariance across ties of the rdd and variance of the rdd

Newman - *Phys. Rev. E*, 2003

6/18

Structural homophily, based on degree

Do high-degree nodes connect to high-degree nodes?

Typical degree-assortativity on social networks:

- astrophysics coauthorship: $r = 0.235$ (Georgia Tech data)
- actor collaboration: $r = 0.227$ (Notre-Dame Univ data)
- friendship network: $r = 0.039$ (Livejournal data)

→ **social networks are usually degree-assortative**

Note that **all complex networks are not degree-assortative**:

- Internet AS level: $r = -0.215$ (UCLA data)
- human protein network: $r = -0.126$ (Vidal data)
- US power grid network: $r = 0.003$ (Tore Opsahl data)

6/18

Structural homophily, based on degree

Do high-degree nodes connect to high-degree nodes?

Typical degree-assortativity on social networks:

- astrophysics coauthorship: $r = 0.235$ (Georgia Tech data)
- actor collaboration: $r = 0.227$ (Notre-Dame Univ data)
- friendship network: $r = 0.039$ (Livejournal data)

→ **social networks are usually degree-assortative**

Note that **all complex networks are not degree-assortative**:

- Internet AS level: $r = -0.215$ (UCLA data)
- human protein network: $r = -0.126$ (Vidal data)
- US power grid network: $r = 0.003$ (Tore Opsahl data)

6/18

Transitivity and clustering

Reminder: social networks have a high average clustering

Triadic closure phenomenon

Old concept in sociology (Simmel, 1908).

Hypothesis on the growth dynamics of a network:



Consequences:

- **high clustering**
- **large number of cliques** (complete subgraphs)
- **densification over time**

7/18

Transitivity and clustering

Reminder: social networks have a high average clustering

Triadic closure phenomenon

Old concept in sociology (Simmel, 1908).

Hypothesis on the growth dynamics of a network:



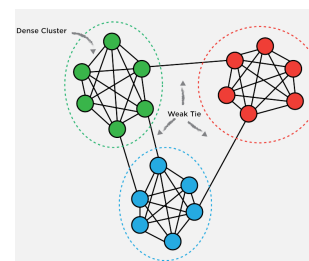
Consequences:

- high clustering
- large number of cliques (complete subgraphs)
- densification over time

7/18

Transitivity and clustering

⇒ schematic picture of social networks



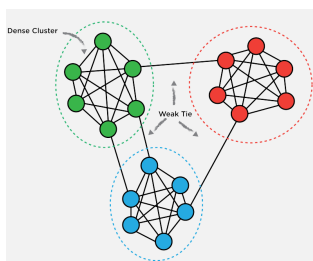
credits image: V.Gauthier

Warning: only a schematic representation
misses overlaps in clusters, groups hierarchy, core/periphery...

7/18

Transitivity and clustering

⇒ schematic picture of social networks



credits image: V.Gauthier

Warning: only a schematic representation
misses overlaps in clusters, groups hierarchy, core/periphery...

7/18

About weak ties...

Hypothesis of the strength of weak tie Granovetter - 1973

A "weak tie" is a link in a social network which represents a relation which is not frequently maintained

It is argued that weak ties play an essential role as they **ensure connections between groups**

What is considered a strong tie in social sciences?

- frequent contacts
- strong affinity (if measurable)
- structural criterion: many mutual neighbors

8/18

About weak ties. . .

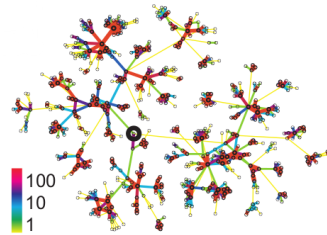
Onnela et al. - 2007

Experimental validation on a phonecall network

→ are weaker links between clusters?

- strength of a relationship = cumulative duration of calls
- link between groups measured with link betweenness

weight (color) = cumulative duration of calls



8/18

About weak ties. . .

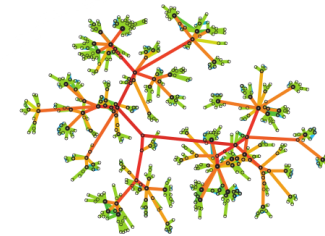
Onnela et al. - 2007

Experimental validation on a phonecall network

→ are weaker links between clusters?

- strength of a relationship = cumulative duration of calls
- link between groups measured with link betweenness

weight (color) = link betweenness



8/18

Effects on spreading in a social network

Examples: innovation spreading, rumor spreading, advertising. . .

What can we expect from the previous observations?

- fast spreading within a community
- use of weak links to spread from a group to another

In practice hard to measure experimentally:

- "contagion" hard to track and isolate
- spreading rarely reaches a large part of a network

→ very active field of research

9/18

Effects on spreading in a social network

Examples: innovation spreading, rumor spreading, advertising. . .

What can we expect from the previous observations?

- fast spreading within a community
- use of weak links to spread from a group to another

In practice hard to measure experimentally:

- "contagion" hard to track and isolate
- spreading rarely reaches a large part of a network

→ very active field of research

9/18

Effects on spreading in a social network

Examples: innovation spreading, rumor spreading, advertising. . .

What can we expect from the previous observations?

- fast spreading within a community
- use of weak links to spread from a group to another

In practice hard to measure experimentally:

- “contagion” hard to track and isolate
- spreading rarely reaches a large part of a network

→ very active field of research

Outline

- 1 Applying our tools to Social Network Analysis
 - The homophily phenomenon
 - Local density and community structure
 - Impact on social contagion
- 2 Epidemic spreading models on graphs
 - Compartmental models in epidemiology
 - What networks bring to the models

From social networks to epidemic spreading

In SNA, innovation spreading dates back to the 50s

In parallel, epidemic modeling developed

Late 90s, data availability ⇒ take into account the social network that supports the spreading

Networks and epidemic models - Keeling and Eames, 2005

Traditional (simplified) approach of epidemic model

Compartmental models

Basic assumption:

- random mixing hypothesis
each individual has an equal chance to come into contact with anyone
- ⇒ homogeneous description of individual behaviors

Infection and population complexity → compartments

- S: susceptible
- I: infected
- R: recovered
- E: exposed
- M: maternally-immune . . .

Traditional (simplified) approach of epidemic model

Compartmental models

Basic assumption:

- **random mixing** hypothesis
each individual has an equal chance to come into contact with anyone
- \Rightarrow homogeneous description of individual behaviors

Infection and population complexity \rightarrow **compartments**

- **S**: susceptible
- **I**: infected
- **R**: recovered
- **E**: exposed
- **M**: maternally-immune ...

12/18

Traditional (simplified) approach of epidemic model

Equations

Model: **set of transition rules from a compartment to another**

- when S encounters I , uniform probability of contamination
- after being contaminated, I has uniform probability of recovery
- etc.

Example: SIR model

$$\begin{cases} \frac{dS}{dt} = -\beta \frac{IS}{N} \\ \frac{dI}{dt} = +\beta \frac{IS}{N} - \gamma I \\ \frac{dR}{dt} = +\gamma I \end{cases}$$

12/18

Traditional (simplified) approach of epidemic model

Example: SIR model

$$\begin{cases} \frac{dS}{dt} = -\beta \frac{IS}{N} \\ \frac{dI}{dt} = +\beta \frac{IS}{N} - \gamma I \\ \frac{dR}{dt} = +\gamma I \end{cases}$$

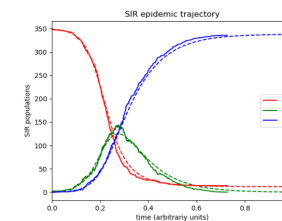
- $N = S + I + R$: size of the population
- β : infection rate (parameter)
- γ : recovery rate (parameter)

12/18

Traditional (simplified) approach of epidemic model

Example: SIR model

$$\begin{cases} \frac{dS}{dt} = -\beta \frac{IS}{N} \\ \frac{dI}{dt} = +\beta \frac{IS}{N} - \gamma I \\ \frac{dR}{dt} = +\gamma I \end{cases}$$



12/18

Classic models

- **SIR:**
 - 3 compartments S, I, R
 - used for disease with lifelong immunity
ex: *measles (rougeole)*, *whooping cough (coqueluche)*
- **SIS:**
 - 2 compartments S, I
 - used for disease with possible reinfections
ex: *STD such as chlamydia*
 - equations

$$\begin{cases} \frac{dS}{dt} = -\beta \frac{IS}{N} + \gamma I \\ \frac{dI}{dt} = +\beta \frac{IS}{N} - \gamma I \end{cases}$$
- SEIS, SEIR, SEIRS, MSEIR, ...

13/18

A few useful concepts in epidemiology

- **Basic reproductive number R_0 :** expected number of new infections from a single infection if everyone is susceptible
ex: for *SIR with random mixing*, $R_0 = \frac{\beta}{\gamma}$
some estimated values (without intervention):
 - *measles*: 12–18
 - *seasonal flu*: 1–2
 - *covid-19*: 3.3–5.7
- **k value:** shape parameter, dispersion parameter, related to the inverse of the dispersion
random mixing \Rightarrow *homogeneous behavior* \Rightarrow *high values of k*
some estimated values (without intervention):
 - *measles*: 0.22
 - *seasonal flu*: 2 – 50
 - *covid-19*: 0.16

14/18

A few useful concepts in epidemiology

- **Basic reproductive number R_0 :** expected number of new infections from a single infection if everyone is susceptible
ex: for *SIR with random mixing*, $R_0 = \frac{\beta}{\gamma}$
some estimated values (without intervention):
 - *measles*: 12–18
 - *seasonal flu*: 1–2
 - *covid-19*: 3.3–5.7
- **k value:** shape parameter, dispersion parameter, related to the inverse of the dispersion
random mixing \Rightarrow *homogeneous behavior* \Rightarrow *high values of k*
some estimated values (without intervention):
 - *measles*: 0.22
 - *seasonal flu*: 2 – 50
 - *covid-19*: 0.16

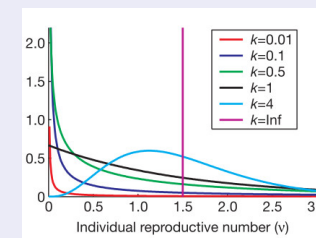
14/18

More about the k value

Superspreading and the effect of individual variation on disease emergence -
Lloyd-Smith et al., 2005

Origin

Model for the individual reproductive number distribution:
negative binomial distribution



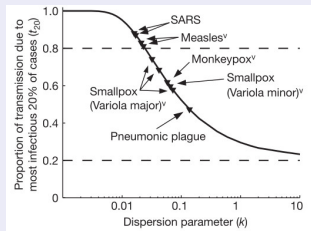
Average $\mu = R_0$, variance $\mu + \frac{\mu^2}{k}$

15/18

More about the k value

Superspreading and the effect of individual variation on disease emergence -
Lloyd-Smith et al., 2005

Another interpretation



* error of a factor 10 on the X-axis range

More about the importance of k for Covid-19 pandemic:

<https://www.theatlantic.com/health/archive/2020/09/k-overlooked-variable-driving-pandemic/616548/>

What networks bring

Random mixing assumption implies uniform contact patterns
we know it is not true

Network-based models

- keep the compartments (S, I, R, E, \dots)
- spreading occurs on the contact network
- at each step, nodes may change compartment:
 - I may contaminate S (or E) neighbors
 - I may recover and turns R
 - ...

What networks bring

Random mixing assumption implies uniform contact patterns
we know it is not true

Network-based models

- keep the compartments (S, I, R, E, \dots)
- spreading occurs on the contact network
- at each step, nodes may change compartment:
 - I may contaminate S (or E) neighbors
 - I may recover and turns R
 - ...

What networks bring

Random mixing assumption implies uniform contact patterns
we know it is not true

Network-based models

- keep the compartments (S, I, R, E, \dots)
- spreading occurs on the contact network
- at each step, nodes may change compartment:
 - I may contaminate S (or E) neighbors
 - I may recover and turns R
 - ...

Data collection and issues

Network-based epidemiology develops **because of data**
but **how available and reliable** are the data?

What is the meaning of an edge?

- **potentially infectious contact**
disease specific: relatively clear for some diseases (STDs),
but airborne diseases?
ex: TousAntiCovid definition of a contact
- ⇒ some degree of arbitrariness

17/18

Data collection and issues

Network-based epidemiology develops **because of data**
but **how available and reliable** are the data?

What is the meaning of an edge?

- **potentially infectious contact**
disease specific: relatively clear for some diseases (STDs),
but airborne diseases?
ex: TousAntiCovid definition of a contact
- ⇒ some degree of arbitrariness

17/18

Data collection and issues

Network-based epidemiology develops **because of data**
but **how available and reliable** are the data?

What are the data collection methods?

- **infection tracing**: look for infectious individuals that have transmitted the disease
- **contact tracing**: interview individuals to collect their potentially infectious contacts
- **diary-based**: based on day-to-day collection by individuals
ex: cattle breeders in Europe since the 90s



credits image: M.Keeling, K.Eames

17/18

Data collection and issues

Network-based epidemiology develops **because of data**
but **how available and reliable** are the data?

Issues with data collection methods

- **infection tracing**: focus on infectious contacts not all contacts, costly
- **contact tracing**: individual bias, sensitive data, subjective evaluation of danger, heterogeneity of data, costly
- **diary-based**: individual bias, sensitive data, disconnected network, heterogeneity of data



credits image: M.Keeling, K.Eames

17/18

Modeling on artificial networks

Difficult data collection
⇒ large use of **network models**

A few key-results in this field:

- “shortcuts” (e.g. air connections) break the locality of spreading and geographic wave-like patterns
- hubs (**superspreaders**) have a dramatic effect on an epidemic, can re-ignite the spreading
- **large variability** of spreading simulations: heterogeneity of networks ⇒ fluctuations

18/18

Modeling on artificial networks

Difficult data collection
⇒ large use of **network models**

A few key-results in this field:

- “shortcuts” (e.g. air connections) break the locality of spreading and geographic wave-like patterns
- hubs (**superspreaders**) have a dramatic effect on an epidemic, can re-ignite the spreading
- **large variability** of spreading simulations: heterogeneity of networks ⇒ fluctuations

18/18

Modeling on artificial networks

Difficult data collection
⇒ large use of **network models**

A few key-results in this field:

- “shortcuts” (e.g. air connections) break the locality of spreading and geographic wave-like patterns
- hubs (**superspreaders**) have a dramatic effect on an epidemic, can re-ignite the spreading
- **large variability** of spreading simulations: heterogeneity of networks ⇒ fluctuations

18/18

Modeling on artificial networks

Difficult data collection
⇒ large use of **network models**

A few key-results in this field:

- “shortcuts” (e.g. air connections) break the locality of spreading and geographic wave-like patterns
- hubs (**superspreaders**) have a dramatic effect on an epidemic, can re-ignite the spreading
- **large variability** of spreading simulations: heterogeneity of networks ⇒ fluctuations

18/18