Veronica Ciocanel

Static Networks I - Reading Group 02/05/2015

I  Structural Properties of Networks

- characterise any real-world networks, not just random graphs.

0. Intro    $G = (N, L)$ where $N$ = nodes
                        $L$ = ordered pairs of elements in $N$

$|N| = N \Rightarrow$ can have 0 to $\frac{N(N-1)}{2}$ edges.

Component of a graph : maximally connected induced subgraph
Giant component : component with size $O(N)$.

Matricial Representation :

- Adjacency matrix $A$ : $N \times N$ square matrix w/
$$a_{ij} = \begin{cases} 1 & \text{when } l_{ij} \in L \iff (i,j) \in L \\ 0 & \text{else.} \end{cases}$$
   (symmetric for undirected graphs)

1. Degree Distributions
node $i \longrightarrow$ degree $k_i = \sum_{j \in N} a_{ij}$  (from adjacency matrix)
  For directed graphs, outgoing links $k_i^{out} = \sum_j a_{ij}$
                    ingoing  —''—  $k_i^{in} = \sum_j a_{ji}$
Degree dist'n    $p(k)$ = prob. that a node chosen uniformly at random
                    has degree $k$
                = fraction of nodes in the graph having degree $k$.
Moments of dist'n
  $m$-moment of $p(k)$ :  $\langle k^m \rangle = \sum_k k^m p(k)$.
   $\langle k \rangle$ = mean degree of $G$
   $\langle k^2 \rangle$ = fluctuations of the degree distribution.

Exe: Exponential dist'n : $p_k \sim e^{-k/\kappa}$

Power law : $p_k \sim k^{-\alpha} \rightarrow$ scale-free-network

($1^{st}$ exe: Price's network of citations between scientific papers, w/ $\alpha = 3.04$)

## 2. Shortest path, Diameter

matrix $D$ : $d_{ij}$ = geodesic (shortest/optimal path) from node $i$ to node $j$

$Diam(G)$ = diameter of graph $G = \max_{i,j \in N} \{ d_{ij} \}$

Typical measure : average shortest path length / characteristic path length
= mean of geodesic lengths over all couples of nodes.

$$L = \frac{1}{N(N-1)} \sum_{\substack{i,j \in N \\ i \neq j}} d_{ij} \,.$$

Issue : $L$ diverges if there are disconnected components in the graph.

Alternative : harmonic mean of geodesic lengths (efficiency of $G$)

$$E = \frac{1}{N(N-1)} \sum_{\substack{i,j \in N \\ i \neq j}} \frac{1}{d_{ij}} \,.$$

## 3. Clustering / Transitivity

$\rightarrow$ real-world network property
— clear deviation from behavior of r.graph

— vertex A connected to vertex B, B with C $\rightarrow$ higher prob A connected to C.

— heightened # of $\triangle$'s in the network.

a) Clustering Coefficient $C = \dfrac{3 \times \text{\# of } \triangle\text{'s in the network}}{\text{\# of connected triples of vertices}}$

Exe: $\qquad C = 3 \cdot \frac{1}{8} = \frac{3}{8} \,.$

— measures how clique-like the friendship network is.

b) Local Clustering Coefficient $c_i = \dfrac{2 e_i}{k_i(k_i-1)} = \dfrac{\sum_{j,m} a_{ij}\, a_{jm}\, a_{mi}}{k_i(k_i-1)}$

where $e_i$ = # of edges in $G_i$ (subgraph of neighbors of node $i$)

Then $C = \langle c \rangle = \dfrac{1}{N} \sum_{i \in N} c_i \,.$

C: easier to compute via a); c: in b): use numerical methods; efficient algorithms are an area of research.

✦   It is suspected that for many types of networks the probability that a friend of your friend is also a friend should $\to$ a nonzero limit as network gets large.

i.e. $c = O(1)$ as $n \to \infty$.

But for random graphs (we'll see): $c = O(\frac{1}{n})$.

Clustering coefficient can be generalised to density of $k$-loops, etc.

## 4. Graph Spectra

$A$ = adjacency matrix $\to$ eigenvalues form the spectrum of the graph.

$G$: undirected $\Rightarrow$ $A$ real and symmetric $\to$ real eigvals $\mu_1 \leq \mu_2 \leq \dots \leq \mu_N$ and eigenvectors corresponding to distinct eigvals are orthogonal.

Perron-Frobenius: $\exists$ real $\mu_N$ s.t. $|\mu| \leq \mu_N$ $\forall$ eigvals $\mu$ of $A$

$-\mu_N$ = spectral radius of $A := \varsigma(A) = \|A\|$.

where $\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}$.

✗Why important? Spectral eigvals and eigenvectors are closely related to topological features such as diameter, # of cycles, connectivity ...

Ex: • Thm $\varsigma(A) \leq \text{Diam}(G)$ = diameter = $\max \{d_{ij}\}$

• $(i,j)^{th}$ entry of $A^k$: # of walks of length $k$ from node $i$ to node $j$.

• Eigvals sum to 0 since $\text{Tr}(A) = 0$.

etc

→ Diam(G) < # of distinct eigvals in a generic graph G.

Other important info on connectivity properties of G:

normed matrix $N = D^{-1} A$ , $D$ = diag. matrix; $D_{ii} = \sum_j a_{ij} = k_i$.

Laplacian matrix $\Lambda = D - A$ → symmetric positive semi-def. matrix.

(Kirchhoff matrix) → all $\lambda$'s of $\Lambda$ are real & non-neg., full set of $N$ real, orthogonal eigenvectors.

→ all rows of $\Lambda$ sum to 0 → $\Lambda$ admits

the lowest eigvalue $\lambda_1 = 0$, w/ eigvector $(1, 1, \dots 1)$

Corollaire: • multiplicity of $\lambda_1 = 0$ is # of comps of G.

• Theorem for $\lambda_2$ → the larger it is, the more difficult to cut G into pieces.


5. Small - World Effect.

Milgram experiment : degree b connectivity on average between any 2 nodes.

Can define small-world networks as :

• networks whose $L$ = average shortest path length scales as $\log(n)$

Recall $L = \frac{1}{N(N-1)} \sum_{\substack{i,j \in N \\ i \neq j}} d_{ij}$ (mean of optimal paths)

$E = \frac{1}{N(N-1)} \sum_{\substack{i,j \in N \\ i \neq j}} \frac{1}{d_{ij}}$

or

• networks that have small value of $L$, like r. graphs $(\log n)$ and a high clustering coefficient $c$.

I.  Random graphs    (Particularly Erdos-Renyi r. graph).

   ↳ initially by Erdos & Renyi in 1959

E-R r. graph :   $G_{m,p}$ :— $m$ nodes & probability $p$ of connecting each

   pair of nodes.

   — graphs w/ $m$ edges appear w/ probability:
   $$p^m (1-p)^{M-m} \; ; \; M = \frac{m(m-1)}{2} = \text{max } \# \text{ of poss. edges.}$$

· So here have set of vertices $V_{m,p} = \{1, 2, \ldots m\}$

   · & introduce $\{\xi_{xy}\}_{1 \leq x < y \leq m}$ : iid r.v's, bernoulli(p)
   $$\begin{cases} P(\xi_{xy} = 1) = p. \\ P(\xi_{xy} = 0) = 1-p \end{cases}$$

   If $\xi_{xy} = 1$ : ∃ edge between $x$ & $y$.

   · consider undirected graphs :   $\xi_{xy} = \xi_{yx}$.

· # of neighbors ~ $\text{Bin}(m-1, p) \Rightarrow E(\text{\# neighbors}) = (m-1) \cdot p.$

   $\underline{\qquad} = \text{average degree} = <k> \text{ from before.}$

   Large $m \Rightarrow p_k \longrightarrow \text{Poisson}(\lambda) \; ; \; \lambda = mp$  is more convenient

   to consider than $(m-1) \cdot p.$

   This is why  E-R random graphs are sometimes called Poisson

   random graphs.


Note    Many properties of  E-R random graphs come from the

   limit of large graph size $m \to \infty$, but while keeping the

   mean degree  $<k> = \lambda$ constant.

Reed-Frost epidemic on an E-R random graph (SIR)

at $t=0$; $\begin{cases} S_0 = \langle 2, \ldots, m \rangle . \\ J_0 = \langle 1 \rangle . \\ R_0 = \phi . \end{cases}$

Update Rule: $R_{t+1} = R_t \cup J_t$     (infected → recovered in one timestep)

$J_{t+1} = \{ y \in S_t / \xi_{xy} = 1 \text{ for some } x \in J_t \}$

$S_{t+1} = S_t \setminus J_{t+1}$   ( $\xi_{xy}$ Bernoulli in E-R)

$v \in V_{mp}$:   $C(v)$ = connected component in $G_{m,p}$ containing $v$.

(people that get infected by an infection starting at $v$).

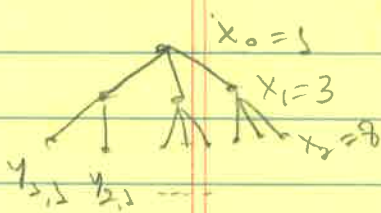Interested in:   asymptotic size of $C(v)$ as $m \to \infty$ ?

Setting above: $C(1)$ of interest. Note $C(1) = \bigcup_{t \geq 0} J_t$.

$\boxed{\lambda = mp \text{ constant}}$ One approach: Construct a branching process (BP) approximation to $J_t$.

Key: Identify a BP $\{Z_t\}$ s.t. $|J_t| \leq Z_t$ and $\sum_{t=1}^{\infty} E(Z_t - |J_t|) \leq \frac{c}{m}$ if $\lambda < 1$.

What is a BP? a seq. of R.V's $\{X_k\}_{k \geq 0}$ s.t. $X_{m+1} = \sum_{i=1}^{X_m} Y_{m,i}$   ①

$X_0 = 1$    $Y_{m,i}$ indep & dist'd according to $\langle P_k \rangle_{k \geq 0}$.

$X_1 = 3$    $X_m =$ # of nodes in gen' $m$.

$X_2 = 8$    $Y_{m,i}$ = # of offspring of node $i$ in generation $m$.

$Y_{2,1}$ $Y_{2,1}$ ...

$\boxed{BP}$ - constructing $Z_t$ is a bit technical, but it depends on the $\xi_{xy}$'s & is chosen s.t. ① is satisfied.

Note: $\lambda < 1$: $E(Z_t - |J_t|) \leq \frac{c}{m}$.

$\lambda > 1$: $E(Z_t - |J_t|) \leq \frac{c}{m} \cdot \lambda^{2t+2}$ → good approx. for initial times.

This approx'n using the BP $Z_t$ is important because it is one way in which important connectivity properties of the E-R graph are proven.

→ Important because <u>giant/largest component</u> comes up in applications of any real networks, not just E-R r. graphs.

Thm  Case 1  <u>Subcritical regime</u>  $\boxed{\lambda < 1}$.

∃ $\rho = \rho(\lambda) = \rho(\infty p) > 0$, t.

$$\lim_{m \to \infty} P\left( |C_1| \leq \rho \log m \right) = 1.$$

i.e.  Largest connected component is at most size $O(\log m)$, smaller than the whole population.

Case 2  <u>Supercritical regime</u>  $\boxed{\lambda > 1}$.

$P_{ext}$ : extinction prob of BP $Z$ w/ offspring distribution Poisson ($\lambda$).

$$0 < P_{ext} < 1!$$

Then ∃ $\rho = \rho(\lambda) > 0$, $\forall \, \xi > 0$ :

$$\lim_{m \to \infty} P\left( \left| \frac{|C_1|}{m} - (1 - P_{ext}) \right| < \xi, \; |C_2| \leq \rho \log m \right) = 1 \; \forall \xi.$$

i.e.  Largest component has size a fixed fraction of $m$, all others are pockets of size $O(\log m)$.

Case 3  <u>Critical regime</u>  $\boxed{\lambda = 1}$.

$$P\left( |C_1| = O\left( N^{2/3} \right) \right) = 1 \quad \text{a.s.}$$

> Note : v. similar to theory of phase transitions in material science.

Proof.

<u>For case 1</u>.  Introduce an important way of exploring nodes via r. walk.

Pick arbitrary node $v \in \{1, 2, ..., m\}$. $C(v)$ its connected component

$A_k$ = set of "active" nodes in $C(v)$

$B_k$ = set of "explored" nodes in $C(v)$.

<u>Initially</u> :  $\begin{cases} A_0 = \{v\} \\ B_0 = \phi. \end{cases}$
($t=0$)

<u>Iteration</u> : ① At step $k$, choose arbitrary $v_{k-1} \in A_{k-1}$.

② $D_k$ = neighbors of $v_{k-1}$

③ $A_k = A_{k-1} \cup D_k \setminus \{v_{k-1}\}$

④ $B_k = B_{k-1} \cup \{v_{k-1}\}$

$$|A_k| = |A_{k-1}| + z_k - 1.$$

$$T = \min\{k > 0 \mid |A_k| = 0\}. \qquad \text{(i.e. done exploring)}$$

$$|A_T| = 1 + \sum_{i=1}^{T} z_i - T \quad \Rightarrow \quad 0 = 1 + \sum_{i=1}^{T} z_i - T.$$

$$\Rightarrow T = 1 + \sum_{i=1}^{T} z_i.$$

$$T = |B_T| = C(v) \qquad \text{all nodes that have been explored.}$$

$$P(|C(v)| > k) = P(T > k) = P(|A_0| > 0, |A_1| > 0, \dots, |A_k| > 0)$$

$$\leq P(|A_k| > 0) = \qquad\qquad \text{(Bin. dist'n)}$$

$$= P\left(\text{Bin}\left(m-1, 1-(1-p)^k\right) \geq k\right) \leq$$

$$\leq P\left(\text{Bin}(m, kp) \geq k\right) \qquad\qquad \left(1-(1-p)^k \leq kp\right)$$

$$= P\left(e^{\theta \cdot \text{Bin}(m, kp)} \geq e^{k\theta}\right) \leq$$

$$\leq E\left[e^{\theta \cdot \text{Bin}(m, kp)}\right] \cdot e^{-k\theta} \qquad\qquad \text{(Markov inequality)}$$

$$P(|C(v)| > k) \leq \left(1 + kp(e^\theta - 1)\right)^m e^{-k\theta} \leq$$

$$\leq e^{mkp(e^\theta - 1)} \cdot e^{-k\theta} = \qquad\qquad (1 + x \leq e^x)$$

$$= e^{-k(\theta - \lambda(e^\theta - 1))} \qquad\qquad (mp = \lambda)$$

For $\lambda < 1$, small $\theta$; $\theta - \lambda(e^\theta - 1) > 0$.

$$\Rightarrow \text{some choice of } \theta: P(|C(v)| > k) \leq e^{-b k}; \; b > 0.$$

$$P(|C_1| > b^{-1} \cdot s \cdot \log m) \leq e^{-b \cdot b^{-1} s \log m} =$$

$$= m^{-s}.$$

Choose $s > 0$; $m \to \infty \;\Rightarrow\; p \to 0.$ $\qquad (\lambda = mp = \text{cst.})$

Let $\rho = \text{cst} = b^{-1} s > 0$.

$$\Rightarrow P(|C_1| \leq \rho \log m) = 1 \text{ as } m \to \infty.$$

- Connectivity in E-R r. graph $G(n,p)$ (continued)

$u \in V$; $\deg(u) = \sum_{v \in V} \xi_{uv}$; $\xi_{uv} = \begin{cases} 1 & \text{if } (u,v) \in E \\ 0 & \text{else} \end{cases}$

$u$ isol'd if $\deg(u) = 0$. Int'd in = # of isolated nodes.

$v \in V$; $\mathfrak{z}_v = \begin{cases} 1 & \text{if } \deg(v) = 0 \\ 0 & \text{else} \end{cases}$

$$\mathfrak{z}_v = \prod_{v \neq u} \mathbb{1}\{\xi_{uv} = 0\} \stackrel{=}{} \prod_{v \neq u} (1 - \xi_{v,u})$$

Note: If one node is not isolated, the other node is automatically not isolated either, so r.v's not iid.

$X = \#$ of isolated nodes $= \sum_v \mathfrak{z}_v \longrightarrow$ not quite Poisson, but v. close

[Thm] Scaling: $np = \log n + c$ (constant)

Then $d_{TV}(X, \text{Poi}(e^{-c})) \longrightarrow 0$ as $n \to \infty$.

Proof: Uses Stein-Chen method.

So, $P(\text{no isol'd nodes}) = P(X=0) = e^{-e^{-c}}$ as $n \to \infty$.

In fact, can show that $P(\exists \text{ conn'd comp of size } 2, \ldots, k) \xrightarrow[n \to \infty]{} 0$.

So $P(G_{n,p} \text{ connected}) = P(\text{no isol'd nodes}) = e^{-e^{-c}}$.

Note: for $\lambda = np = $ cst scaling, $P(G_{n,p} \text{ connected})$ goes to 0 (Thm p.7).

Diameter of E-R r. graph: it can be proven (technical) that the diameter has values in a small range of values around $\text{Diam} = \dfrac{\ln N}{\ln(pN)} = \dfrac{\ln N}{\ln \langle k \rangle}$ ☆

$\longrightarrow$ Same for the average shortest path $L \sim O\left(\dfrac{\ln N}{\ln \langle k \rangle}\right)$.

Why? Average # of neigbors a distance $\ell$ away is $\lambda^\ell$ where
$\lambda = (n-1) \cdot p \simeq np$ for $n \to \infty$.
To get to the whole network, $\lambda^L = N \Rightarrow L = \dfrac{\ln N}{\ln \lambda} = \dfrac{\ln N}{\ln \langle k \rangle}$.

<u>Note</u> : Since $L \sim O\left(\frac{\ln N}{\ln <k>}\right)$, slower than $\log(N)$ $\implies$ the E-R r-graph reproduces the small-world scenario.

- <u>Clustering coefficient</u> $\longrightarrow$ edges among neighbors of a node

$$C = p = \frac{<k>}{m} \quad ; \quad \text{Why?} \quad c_i = \frac{p k_i(k_i-1)/2}{k_i(k_i-1)/2} = p$$

$$\& \quad C = \frac{\sum c_i}{m} = \frac{p \cdot m}{m} = p \cdot \checkmark$$

2) $C = p = \frac{\lambda}{m} \rightarrow 0$ as $m \rightarrow \infty$ (Not realistic!)

- <u>Degree distribution</u> : Poisson $\rightarrow$ unrealistic, no correlation between degs of adjacent vertices, no community structure $\rightarrow$ inadequate to describe most observed dist'n ($\sim$ power laws).

  There exist extensions of the E-R r-graph.

  [But]: important main model b/c ideas of <u>giant component</u>, <u>phase transitions</u> are present in all the more sophisticated models.

<u>Summary</u> E-R r-graph

- Poisson degree distribution: not realistic
- clustering : $O\left(\frac{1}{m}\right) \underset{m\to\infty}{\rightarrow} 0$ : not realistic
- no community structure
- characteristic path length : $O(\log m)$ $\rightarrow$ reproduces small-world phenomenon.

III  Generalized random graphs

- make E-R more realistic
- easiest property to change: non-Poisson degree distribution

1) The Configuration Model

Def'd in the following way: specify degree dist'n $p_k$ ; $p_k = $ fraction of vertices w/ degree $k$.
- Choose a degree sequence: $n$ vals of the degrees $k_i$ of vertices $i = 1, n$, from this dist'n.
- Give each vertex $i$ in our graph $k_i$ "stubs" or "spokes" sticking out of it → i.e. ends of edges-to-be.
- Then choose pairs of stubs at random from the network & connect them together   — gives every top. of a graph w/ the given seq w/ equal prob'ty

Main results — results on size of giant component can be proven here via powerful formalism of a generating function.

Probability generating fcn of $\overline{X}$ (takes values $k$ w/ probability $p(k)$) is:
$$G(z) = E(z^X) = \sum_{k=0}^{\infty} p(k) z^k.$$
Note : $G'(z) = \sum_{k=1}^{\infty} k \, p(k) \, z^{k-1}$
$$G'(1) = \sum_{k=1}^{\infty} k \, p(k) = E(X)$$

Back to config model: Degree of a vertex that we reach by following a randomly chosen edge is not $p_k$.

∃ $k$ edges that arrive at a vertex of deg $k$ ⇒ $k$ times as likely to arrive at that vertex than at one of degree 1.

⇒ Deg. dist'n of the vertex @ the end of a randomly chosen edge is $\underline{k \cdot p_k}$. Many times int'd in the # of edges that leave a vertex (excess degree)

→ dist'n $\quad q_k = \dfrac{(k+1) \cdot p_{k+1}}{\sum_k k \, p_k} = \dfrac{(k+1) \, p_{k+1}}{\langle k \rangle} = \dfrac{(k+1) \, p_{k+1}}{z}$

Recall: $g_k = \dfrac{(k+1) p_{k+1}}{\sum k p_k}$

Define 2 gen'ing fcns for dist's $p_k$ & $g_k$:

$$G_0(x) = \sum_{k=0}^{\infty} p_k x^k, \qquad G_1(x) = \sum_{k=0}^{\infty} g_k x^k.$$

Note: $G_1(x) = \dfrac{G_0'(x)}{z}$ . $\qquad , z = \langle k \rangle = \sum k p_k .$

• Gen'ing fcn $H_1(x)$ for the total # of vertices reachable by following an edge :

$$\boxed{H_1(x) = x \, G_1(H_1(x))}$$

Won't prove, but here's the $\boxed{\text{intuition}}$ :

— when following an edge, we find at least a vertex at the other end $(x)$ ; + some other clusters of vertices (repr'd by $H_1$) reachable by following other edges attached to that one vertex .

$\underset{\text{excess degree} \sim g_k \Rightarrow G_1(x)}{}$

• Gen'ing fcn $H_0(x)$ = total # of vertices reachable from a randomly chosen vertex :

$$\boxed{H_0(x) = x \, G_0(H_1(x))}$$

$\hookrightarrow$ idea of giant comp

• Mean component size in the region of no giant component is :

$$\langle s \rangle = \underset{\text{exp value}}{H_0'(1)} = \; + \frac{G_0'(1)}{1 - G_1'(1)} = 1 + \frac{z_1^2}{z_1 - z_2} \quad \bigstar$$

where $z_1 = z = \langle k \rangle = G_0'(1)$.

$z_2 = \langle k^2 \rangle - \langle k \rangle = G_0''(1) - G_1'(1)$.

• Divergence in $\bigstar$ when $z_1 = z_2$, i.e. when $G_1'(1) = 1$, i.e

when $\boxed{\sum_k k (k-2) p_k = 0.}$ $\longleftarrow$ critical cond'n

$\longleftarrow$ phase transition at which a giant comp. appears.

i.e. $\sum > 0 \Rightarrow$ giant comp. a.s. (occupies a fraction of the graph).

$\sum < 0 \Rightarrow$ largest comp. is $O(\log N)$.

more extensions: — directed graphs, bipartite graphs w/ 2 types of nodes.

— still use generating fcn framework.

References used :

1. M. E. J. Newman, "The Structure and Function of Complex Networks", SIAM REVIEW, 45, 2003 (167-256)

2. S. Boccaletti et al, "Complex networks: Structure and dynamics", Physic Reports, 424, 2006 (175-308)

3. M. Draief and L. Massoulie, "Epidemics and Rumours in Complex Networks", Cambridge University Press, 2010.

4. D. J. Watts, "Small Worlds: The Dynamics of Networks between order and randomness", Princeton University Press, 2003.