

Social Network Analysis: Preferential Attachment

(Dynamic networks and the non-equilibrium modeling approach)

Donglei Du
(ddu@unb.ca)

Faculty of Business Administration, University of New Brunswick, NB Canada Fredericton
E3B 9Y2

Table of contents

- 1 Introduction
- 2 Preferential attachment models
 - Price's model for directed network (Price, 1976)
 - The vertex copying model (Aiello et al., 2000; Kumar et al., 1999; Vázquez et al., 2002)
 - Barabási-Albert's model for undirected network (Barabási and Albert, 1999)
- 3 A case study: How to become famous in the academic world? (Newman, 2009, 2013)
 - The dynamic of the degree distribution of Price's model
- 4 Appendices
 - Appendix A: The dynamical/non-equilibrium approach to networks
 - Appendix B: How to fit the Yule-Simon indegree distribution in R ?
 - Appendix C: The rank model (Fortunato et al., 2006)

- The material is adopted from (Newman, 2010).

The Preferential attachment model

- "Yule process"
- "cumulative advantage"
- "the rich get richer"
- "Matthew effect"
- "Gibrat's law"

Matthew effect (Matthew 25:29, New International Version)

- "For everyone who has will be given more, and he will have an abundance. Whoever does not have, even what he has will be taken from him."

Life is unfair...

- People are more likely to give credit to the famous than to the little known.
 - The classic example of the Matthew effect is a scientific discovery made simultaneously by two different people, one well known and the other little known.
 - It is claimed that under these circumstances people tend more often to credit the discovery to the well-known scientist.

A brief history of preferential attachment

- (Yule, 1925): explain the power-law distribution of the number of species per genus of flowering plants
- (Simon, 1955): the distribution of sizes of cities and other phenomena and the master equation method
- (Price, 1976): citation network
- (Barabási and Albert, 1999): World Wide Web

Price's preferential attachment model for directed network

- The initial network N_0 is immaterial for our conclusion as we focus on the stationary state in the long run limit
- Non-equilibrium systems may still have a stationary state in which the probability distribution is time-independent (See Appendix C for more details).

Price's preferential attachment model for directed network

- Node $t + 1$ is connected to on average $c \leq t$ existing nodes with a probability that is proportional to the indegree $k_i(t)$ plus a .
- Formally, the probability $p_i(t)$ that t is connected to node $i \in N_t$ is

$$\frac{k_i(t) + a}{\sum_{j=1}^t (k_j(t) + a)} = \frac{k_i(t) + a}{t(c + a)}, \quad (1)$$

where we used the fact that $\sum_{j=1}^t k_j(t) = tc$.

Degree distribution process in Price's preferential attachment model

- The preferential attachment above generates a stochastic process and we can study many quantities that are defined on this process.
- Let us first look at the degree distribution process $\{D(t)\}$, where $D(t)$ is the random variable denoting the degree distribution on network N_t for each time $t \geq 1$.

Degree distribution process in Price's preferential attachment model

- Let $p_q(t) = P(D(t) = q)$ be the fraction of vertices with in-degree q in N_t .
- The following formula will be used shortly: the expected number of new in-links to all vertices with degree q in $N(t)$ is

$$\boxed{c \frac{q+a}{t(c+a)} t p_q(t) = \frac{c(q+a)}{c+a} p_q(t)} \quad (2)$$

- The new node $t+1$ points to c others in $N(t)$ on average, so the expected number of new links to vertex i in $N(t)$ is c times (1).
- There are $t p_q(t)$ vertices with degree q in $N(t)$.

The evolution of the dynamic system: the master equation for the stochastic process

- In the new network $N(t+1)$ after the addition of the new node $t+1$, the number of nodes with indegree q changes from $tp_q(t)$ in N_t to $(t+1)p_q(t+1)$ in $N(t+1)$:

$$(t+1)p_q(t+1) - tp_q(t)$$

- The above net change is due to two factors:
 - the number of vertices with in-degree q increases by one for every vertex previously of in-degree $q-1$ in N_t that receives a new citation, thereby becoming a vertex of in-degree q . From (2), the expected number of such vertices is

$$\frac{c(q-1+a)}{c+a} p_{q-1}(t)$$

- Similarly, we lose one vertex of in-degree q every time such a vertex receives a new citation, thereby becoming a vertex of in-degree $q+1$. From (2), the expected number of such vertices receiving citations is

$$\frac{c(q+a)}{c+a} p_q(t)$$

The evolution of the dynamic system: the master equation for the stochastic process

- Now putting everything together, taking into the initial condition into consideration:

$$(t+1)p_q(t+1) - tp_q(t) = \begin{cases} 1 - \frac{c(q+a)}{c+a} p_q(t), & q = 0; \\ \frac{c(q-1+a)}{c+a} p_{q-1}(t) - \frac{c(q+a)}{c+a} p_q(t), & q \geq 1. \end{cases}$$

- Assume at the equilibrium, $\lim_{t \rightarrow \infty} p_q(t) = p_q$. Then we have, after rearranging,

$$p_q = \begin{cases} \frac{1+a/c}{a+1+a/c}, & q = 0; \\ \frac{q+a-1}{q+a+1+a/c} p_{q-1}, & q \geq 1. \end{cases}$$

- The above recursive relation can be solved as the **Yule-Simon** distribution

$$p_q = \frac{B(q+a, 2+a/c)}{B(a, 1+a/c)} \sim q^{-(2+a/c)}, \quad q \gg a.$$

- In the left, B is the Euler's beta function

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}.$$

- Γ is the gamma function (an extension of the factorial function) defined as

$$\Gamma(t) = \int_0^{\infty} x^{t-1} e^{-x} dx,$$

satisfying the functional equation:

$$\begin{cases} \Gamma(1) = 1 \\ \Gamma(t+1) = t\Gamma(t) \end{cases}.$$

- The Stirling approximation for the gamma function for large t :

$$\Gamma(t) \sim t^{t-\frac{1}{2}} e^{-t} \sqrt{2\pi}.$$

The Price's model: Illustration in R package: igraph

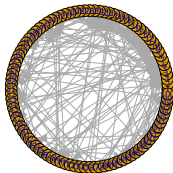
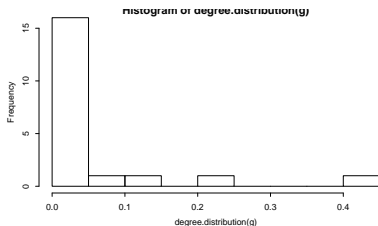
```
rm(list=ls()) # clear workspace
library(igraph) # call igraph
g <- barabasi.game(100, m=2, zero.appeal=5)
hist(degree.distribution(g))
plot(g,layout=layout.circle)
average.path.length(g)
```

```
## [1] 2.239456
```

```
transitivity(g, type="average")
```

```
## [1] 0.140033
```

```
#####
## The BA-model is a very simple stochastic algorithm for building a graph
#####
```



Alternative view of Price's model

- In Price's model, each new node is pointed to a vertex chosen in proportion to its in-degree plus a constant a . as in (1).
- This process can be viewed from an alternative point of view. When the next node $t + 1$ arrives:
 - With some probability ϕ , attach the edge to a vertex chosen strictly in proportion to its current in-degree, i.e., with probability

$$\frac{k_i(t)}{\sum_{j=1}^t k_j(t)} = \frac{k_i(t)}{tc}.$$

- With probability $1 - \phi$, attach to a vertex chosen uniformly at random from all t possibilities, i.e., with probability $1/t$.
- The total probability of attaching to vertex i in this process is

$$\phi \frac{k_i(t)}{tc} + (1 - \phi) \frac{1}{t}$$

Alternative view of Price's model

- When $\phi = c/(c + a)$, the above probability is exactly (1).
- So an alternative way of performing a step of Price's model is the following. When the next node $t + 1$ arrives:
 - With probability $c/(c + a)$ choose a vertex in strict proportion to in-degree.
 - With probability $a/(c + a)$ choose a vertex uniformly at random from the set of all vertices.

Why the alternative view is useful?

- Simulation of Price's model
- The copying model—related to Barabási-Albert's preferential attachment model.

The copying model under Price's setting

- The copying model is local and no global knowledge (such as the degree of all the vertices) of the network is needed.
 - At time $t + 1$, a randomly chosen vertex $j \in \{1, \dots, t\}$ is duplicated as the new added node $(t + 1)$:
 - 1 For each of $(t + 1)$'s c out-going connections, with probability ϕ , keep the link;
 - 2 with probability $1 - \phi$, node $(t + 1)$ is rewired to a randomly chosen node.
- The end result is a set of outlinks for the new vertex in which, on average, ϕc of the entries are copied from the old vertex and the remainder are chosen at random.
- In effect, we made an imperfect copy of the old vertex in which the destinations of some fraction of the outgoing edges have been randomly reassigned.

The copying model under Price's setting

- The probability that vertex $i \in \{1, \dots, t\}$ receives a new incoming edge upon the addition of the new vertex $t + 1$ is

$$\phi \frac{k_i(t)}{t} + (1 - \phi) \frac{c}{t} \quad (3)$$

- Proof of the (3).
 - For i to receive a new edge, one of two things has to happen:
 - Case 1** . Either the newly added vertex happens to copy connections from a vertex that already points to vertex i , in which case with probability ϕ the connection to i will itself get copied;
 - Case 2** . or i could be one of the vertices chosen at random to receive a new edge.

Case 1 I

- The probability is $\phi k_i(t)/n$ as explained below:
 - The probability that node $(t + 1)$ copies own links from a randomly chosen j is $1/t$
 - Since i has $k_i(t)$ parents, the chance that any one of these $k_i(t)$ parents pointing to i is $k_i(t)/t$.
 - The chance of each of these links gets copied is ϕ .
- The above analysis actually implies that Step 1 in the copying model is equivalent to the following "preferential attachment" rule:
 - with probability ϕ , node $(t + 1)$ points to i with probability proportional to i 's indegree $k_i(t)$.

Case 2

- The probability is $(1 - \phi) \frac{c}{t}$ as explained below:
 - The average number of random links that node $(t + 1)$ randomly points to is $1 - \phi$ for each of its c outgoing edges, resulting a total of $(1 - \phi)c$ links.
 - The probability that vertex i is the target of one of these random links is $1/t$,

Results for the copying model

- Let

$$a = c \left(\frac{1}{\phi} - 1 \right)$$

- or equivalently

$$\phi = \frac{c}{a + c}$$

Barabási-Albert's preferential attachment model for undirected network (Barabási and Albert, 1999), (Barabási, 2009)

- The network begins with an initial connected network of m_0 nodes. New nodes are added to the network one at a time.
 - Each new node is connected to $c \leq m_0$ existing nodes with a probability that is proportional to the number of links that the existing nodes already have. Formally, the probability p_i that the new node is connected to node i is

$$p_i = \frac{k_i}{\sum_j k_j},$$

where k_i is the degree of node i and the sum is made over all pre-existing nodes j (i.e. the denominator results in the current number of edges in the network).

- Heavily linked nodes ("hubs") tend to quickly accumulate even more links, while nodes with only a few links are unlikely to be chosen as the destination for a new link.
- The new nodes have a "preference" to attach themselves to the already heavily linked nodes.

BA's model as a special case of Price's model

- Convert the undirected network into a directed one by orienting the more recent nodes to less recent nodes such that each vertex has out-degree exactly c .
- The total degree k_i in the original undirected network is the sum of the vertex's in-degree and out-degree in the directed network, $k_i = q_i + c$ where q_i is the in-degree.
- The probability of an edge attaching to a vertex is proportional to $k_i = q_i + c$, which is the same as in Price's model for $a = c$.
- Thus the distribution of in-degrees in this directed network is the same as for Price's model with $a = c$:

$$p_q = \frac{B(q + a, 3)}{B(a, 2)} \sim q^{-3}, \quad q \gg a.$$

The Preferential Attachment model: Illustration in Netlogo

- <http://ccl.northwestern.edu/netlogo/>
- Go to File/Model Library/Networks/Preferential Attachment
- Ref:
 - Barabási and Frangos (2002), Barabási and Albert (1999)

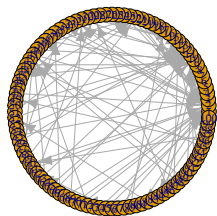
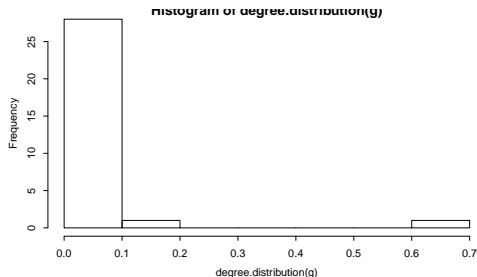
The Preferential Attachment model: Illustration in R package: igraph

```
rm(list=ls()) # clear workspace
library(igraph) # call igraph
g <- barabasi.game(100)
hist(degree.distribution(g))
plot(g,layout=layout.circle)
average.path.length(g)
```

```
## [1] 1.889796
```

```
transitivity(g, type="average")
```

```
## [1] 0
```



How to become famous? A case study based on Price's model Newman (2009), Newman (2013)

—MEJ Newman

"Were we wearing our cynical hat today, we might say that the scientist who wants to become famous is better off—by a wide margin—writing a modest paper in next year's hottest field than an outstanding paper in this year's."

- If the creation times are known (instead of a snapshot of the network, where the ERGM later will be useful), Price's model (parameters c , a) can be estimated directly from the empirical degree distribution by fitting the predicted form to the empirical data via, for example, maximum likelihood.

Empirical tests of Price's model

- Objective: test the first-mover effect predicted from the Price model.
- Data:
 - 2009: A citation network of 2407 papers on network science theory, spanning a ten-year period from June 1998 to June 2008, including five early and well-cited papers in the field along with all papers that cite them, but excluding review articles and restricted to papers in physics and related areas.
 - 2013: updated to 6976 papers between 1998-2013 (first figure on the right).
- the Yule-Simon distribution in (last formula on Slide 15) was fitted to the degree distribution via maximum likelihood to recover estimates of a and c (second figure on the right for the 2009 data).

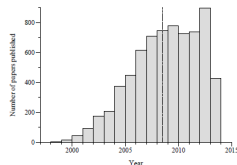


Figure: The historical number of papers between 1998-2013

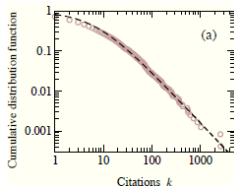


Figure: Empirical measurements in brown; theoretical predictions in black. The ccdf (complement cumulative density function) of the degree distribution. The best fit is achieved for $\alpha = 2 + a/c = 2.28$ and $a = 6.38$.

Observations

- In general, Price's model agrees with the empirical data.
- This suggests that pure citation numbers may not be a good indicator of paper impact, since much of their variation can be predicted from publication dates, without reference to paper content.
- Instead, they proposed an alternative measure of impact.
 - Rather than looking for papers with high total citation counts, we should look for papers with counts higher than expected given their date of publication.
 - The appropriate calculation is to count the citations a paper has received and compare that figure to the counts for other papers on the same topic that were published around the same time.
- A simple z-score is good for this purpose:

$$z = \frac{\text{number of citations of a paper} - \text{mean number of citations in a window close to the date of a paper of interest}}{\text{standard deviation number of citations in a window close to the date of a paper of interest}}$$

- Papers with high z-scores are conjectured to be of particular interest within the field.

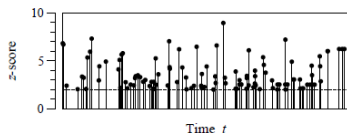


Figure: Only z-scores greater than two (the dashed line) are plotted. The 'time' in the horizontal axis is the publication date of each paper.

Degree distribution as a function of the time of creation in Price's model

- We extend the analysis in Slides 14-15 to include the dynamics.
- Let $p_q(i, t)$ be the average fraction of vertices in N_t that were created at time i and have in-degree q .
- Following analogous analysis, we have

$$(t+1)p_q(i, t+1) - tp_q(i, t) = \begin{cases} \delta_{it} - \frac{c(q+a)}{c+a} p_q(i, t), & q = 0; \\ \frac{c(q-1+a)}{c+a} p_{i, q-1}(t) - \frac{c(q+a)}{c+a} p_{i, q}(t), & q \geq 1. \end{cases} \quad (4)$$

where the Kronecker delta $\delta_{it} = 1$ if $i = t$ and 0 otherwise.

- Note that $p_q(i, t) \rightarrow 0$, $t \rightarrow \infty$, because only one vertex is created at any particular i .
- Introduce the rescaled time of creation $\tau = i/t$, and the density function $\pi_q(\tau, t)$ such that

$$\pi_q(\tau, t) = tp_q(i, t). \quad (5)$$

- Explanation: $\pi_q(\tau, t)d\tau$ is the fraction of vertices that have in-degree q and fall in the interval $[\tau, \tau + d\tau]$;
- The number of vertices in the interval $d\tau$ is $td\tau$.
- So $\pi_q(\tau, t)d\tau = p_q(i, t)td\tau$.

Master equation in terms of τ , t and π

- Due to (5), the master equation in (4) becomes

$$\pi_q \left(\frac{t\tau}{t+1}, t+1 \right) - \pi_q(\tau, t) = \begin{cases} \delta_\tau - \frac{c(q+a)}{c+a} \frac{\pi_q(\tau, t)}{t}, & q = 0; \\ \frac{c(q-1+a)}{c+a} \frac{\pi_{q-1}(\tau, t)}{t} - \frac{c(q+a)}{c+a} \frac{\pi_q(\tau, t)}{t}, & q \geq 1. \end{cases} \quad (6)$$

- Let $\epsilon = 1/t$ and $\lim_{t \rightarrow \infty} \pi_q(\tau, t) = \pi_q(\tau)$. Dropping terms of order ϵ^2 , the above becomes

$$\frac{\pi_q(\tau - \epsilon\tau) - \pi_q(\tau)}{\epsilon} = \begin{cases} \delta_\tau - \frac{c(q+a)}{c+a} \pi_q(\tau), & q = 0; \\ \frac{c(q-1+a)}{c+a} \pi_{q-1}(\tau) - \frac{c(q+a)}{c+a} \pi_q(\tau), & q \geq 1. \end{cases} \quad (7)$$

- Equivalently, master equation becomes a differential equation

$$\tau \frac{d\pi_q(\tau)}{d\tau} = \begin{cases} \frac{c(q+a)}{c+a} \pi_q(\tau), & q = 0; \\ \frac{c(q-1+a)}{c+a} \pi_{q-1}(\tau) - \frac{c(q+a)}{c+a} \pi_q(\tau), & q \geq 1. \end{cases} \quad (8)$$

along with the initial conditions $\pi_0(1) = 1$ and $\pi_k(1) = 0$ for $q \geq 1$.

Solution to the master equation (8)

- For $q = 0$

$$\pi_0(\tau) = \tau^{\frac{ca}{c+a}}$$

- General term:

$$\begin{aligned}\pi_q(\tau) &= \frac{\Gamma(q+a)}{\Gamma(q+1)\Gamma(a)} \tau^{\frac{ca}{c+a}} \left(1 - \tau^{\frac{c}{c+a}}\right)^q \\ &= \frac{1}{qB(q,a)} \tau^{\frac{ca}{c+a}} \left(1 - \tau^{\frac{c}{c+a}}\right)^q \\ &\sim q^{a-1} \left(1 - \tau^{\frac{c}{c+a}}\right)^q\end{aligned}$$

The first-mover effect

- **The first-mover effect:** papers published early on should on average receive far more citations than those published later.
- One can calculate the average number of citations a paper receives as a function of its time of publication:

$$\gamma(\tau) = \sum_{q=0}^{\infty} q\pi_q(\tau) = a \left(\tau^{-\frac{c}{a+c}} - 1 \right) \quad (9)$$

The first-mover effect: illustration of (9)

- The out-degree parameter c was in each case $c = 2a$, so that the exponent of the power-law degree distribution $\alpha = 2 + a/c$ (Slide 14) is 2.5 for all curves, which is a typical value for real-world networks.

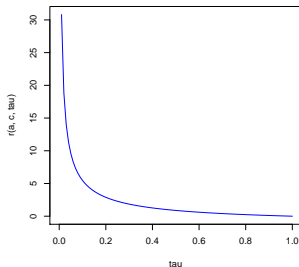


Figure: $\gamma(\tau)$: The average in-degree of vertices in Price's network model decreases a function of the rescaled entering time τ at which they were added to the network

How the expected in-degree of a vertex varies with its age after it enters the network?

- The expected indegree $\gamma_i(s)$ of the vertex added at time i , as a function of its age s :

$$\gamma_i(s) = a \left(\left(1 + \frac{s}{i} \right)^{\frac{c}{a+c}} - 1 \right) \quad (10)$$

- When a vertex is first added to the network and $s \ll i$, the in-degree of a vertex initially grows linearly with the age of the vertex, on average, but with a constant of proportionality that is smaller the later the vertex entered the network—again we see that there is a substantial advantage for vertices that enter early.

$$\gamma_i(s) \approx \frac{ca}{c+a} \frac{s}{i}$$

- As the vertex ages, there is a crossover to another regime around the point $s = i$, i.e., at the point where the vertex switches from being in the younger half of the population to being in the older. For $s \gg i$, the growth is slower than linear but still favors vertices that appear early:

$$\gamma_i(s) \approx a \left(\frac{s}{i} \right)^{\frac{c}{c+a}}$$

The first-mover effect: illustration of (10)

- The out-degree parameter c was in each case $c = 2a$, so that the exponent of the power-law degree distribution $\alpha = 2 + a/c$ (Slide 14) is 2.5 for all curves, which is a typical value for real-world networks.

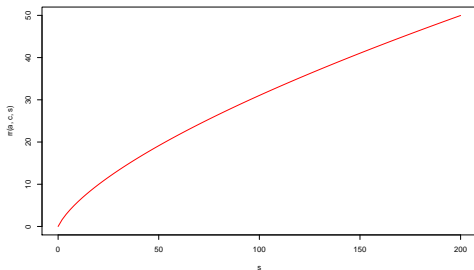
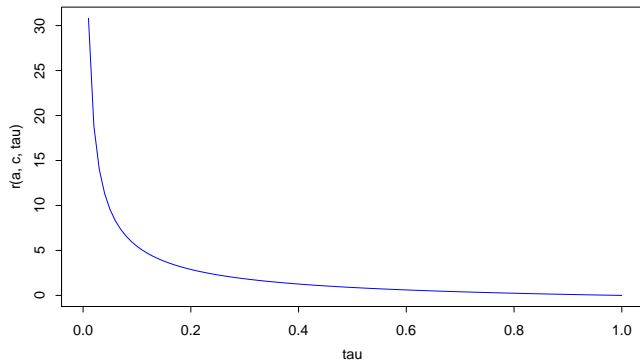


Figure: $\gamma_i(s)$: the horizontal axis is s/i . The average in-degree of a vertex initially grows linearly with the age and then slows down after the crossover point $s/i = 1$.

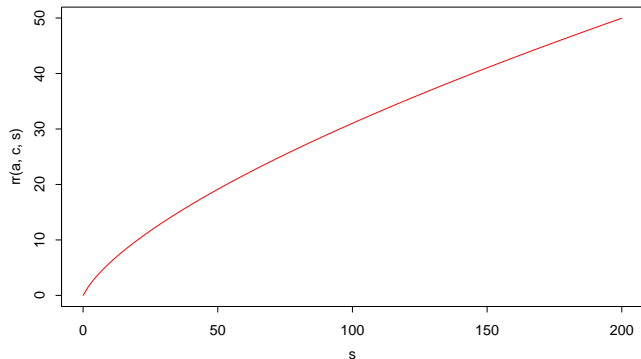
R Code

```
# define the average in-degree  $\bar{r}(t)$  for a vertex created at time  $t$ :  
r<-function(a,c,tau){  
  return(a*(tau^(-c/(a+c))-1))  
}  
a<-1.5  
c<-3  
# plot the curve  
curve(r(a,c,tau), 0, 1, xname = "tau", col = "blue")
```



R Code

```
# define the expected indegree  $\tau_i(s)$  of the vertex added at time  $i$ , as a function  
rr<-function(a,c,s){  
  return(a*((1+s)^(c/(a+c))-1))  
}  
a<-1.5  
c<-3  
# plot the curve  
curve(rr(a,c,s), 0, 200, xname = "s", col = "red")
```



References I

- Aiello, W., Chung, F., and Lu, L. (2000). A random graph model for massive graphs. In *Proceedings of the thirty-second annual ACM symposium on Theory of computing*, pages 171–180. Acm.
- Barabási, A.-L. (2009). Scale-free networks: a decade and beyond. *Science*, 325(5939):412–413.
- Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks. *science*, 286(5439):509–512.
- Barabási, A.-L. and Frangos, J. (2002). *Linked: The New Science Of Networks Science Of Networks*. Basic Books.
- Erdős, P. and Rényi, A. (1959). On random graphs. *Publicationes Mathematicae Debrecen*, 6:290–297.
- Fortunato, S., Flammini, A., and Menczer, F. (2006). Scale-free network growth by ranking. *Physical review letters*, 96(21):218701.

References II

- Holland, P. W. and Leinhardt, S. (1981). An exponential family of probability distributions for directed graphs. *Journal of the American Statistical Association*, 76(373):33–50.
- Kumar, R., Raghavan, P., Rajagopalan, S., and Tomkins, A. (1999). Extracting large-scale knowledge bases from the web. In *VLDB*, volume 99, pages 639–650. Citeseer.
- Newman, M. (2009). The first-mover advantage in scientific publication. *EPL (Europhysics Letters)*, 86(6):68001.
- Newman, M. (2010). *Networks: an introduction*. Oxford University Press.
- Newman, M. (2013). Prediction of highly cited papers. *arXiv preprint arXiv:1310.8220*.

References III

- Price, D. d. S. (1976). A general theory of bibliometric and other cumulative advantage processes. *Journal of the American Society for Information Science*, 27(5):292–306.
- Simon, H. A. (1955). On a class of skew distribution functions. *Biometrika*, 42(314):425–440.
- Vázquez, A., Flammini, A., Maritan, A., and Vespignani, A. (2002). Modeling of protein interaction networks. *Complexus*, 1(1):38–44.
- Watts, D. J. and Strogatz, S. H. (1998). Collective dynamics of “small-world” networks. *nature*, 393(6684):440–442.
- Yule, G. U. (1925). A mathematical theory of evolution, based on the conclusions of dr. jc willis, frs. *Philosophical Transactions of the Royal Society of London. Series B, Containing Papers of a Biological Character*, 213(402-410):21–87.

Static vs dynamic approaches

- Static modeling approach
 - The static modeling approach focuses on the stationary properties of the network, assuming the system is in the equilibrium state.
 - The emphasis is on the statistical behavior of the system.
 - Examples: (Erdős and Rényi, 1959; Watts and Strogatz, 1998; Holland and Leinhardt, 1981)
- Dynamical modeling approach
 - The dynamical modeling approach tries to capture the emergent behavior and complex features of networks by studying the collective dynamics of its constituents.
 - The emphasis is on the evolutionary micro-mechanisms that generate the macro-topological properties, which become a by-product of the system's dynamics.
 - Examples: (Barabási and Albert, 1999)

The master equation approach for the dynamic modeling approaches

- $P(X, t)$: the probability of a particular network realization X at time t .
- The master equation is a linear differential equation for the probability that any network, owing to the microscopic dynamics, is in the configuration X .
- Assuming a Markovian process, then

$$\frac{\partial P(X, t)}{\partial t} = \sum_{Y \neq X} [P(Y, t)r_{Y \rightarrow X} - P(X, t)r_{X \rightarrow Y}]$$

How to fit the Yule-Simon indegree distribution in R?

- Assume the indegree follows the Yule-Simon distribution

$$p_q = \frac{B(q + a, 2 + a/c)}{B(a, 1 + a/c)}, \quad q \geq 1.$$

- Given a directed network with observed indegrees $d = (d_1, \dots, d_n)$, we want to estimate the parameters via the maximum likelihood method.

R Code

```
rm(list = ls())
library(igraph)
# generate a network using the Price's model
g <- barabasi.game(100, m=5, zero.appeal=5) #m=c and zero.appeal=c
deg<-degree(g, mode="in")
deg_dist<-degree.distribution(g)
```

R Code

```
library(distr) # user-defined distribution package
# define the maximum likelihood function
fn <- function(coef) {
  -sum ( log(beta(deg+coef[1], 2+coef[1]/coef[2]))
}
# find the MLE
output<-nlm(fn, coef<-c(1,1), hessian=TRUE)
```

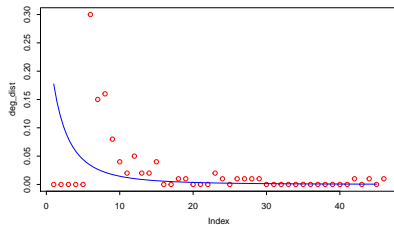
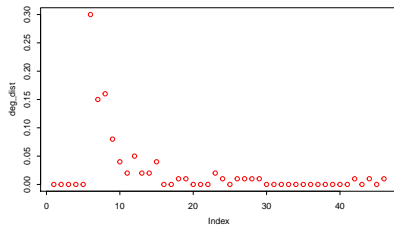
```
## Warning in beta(deg + coef[1], 2 + coef[1]/coef[2]): NaNs
produced
## Warning in beta(coef[1], 1 + coef[1]/coef[2]): NaNs produced
## Warning in nlm(fn, coef <- c(1, 1), hessian = TRUE): NA/Inf
replaced by maximum positive value
```

R Code

```
a<-output[[2]][1]
c<-output[[2]][2]
# the density function for Yule-Simon distribution
yule_density<-function(q){
  density <- beta(q+a, 2+a/c)/beta(a,1+a/c)
  return(density)
}
# the complement cumulative density function for Yule-S
yule_ccdf<-function(q){
  ccdf <- 1-sum(yule_density(1:q))
  return(ccdf)
}
```

R Code

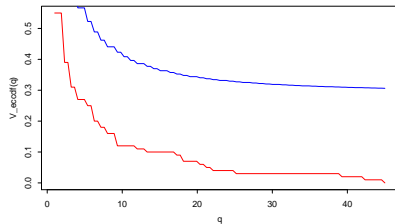
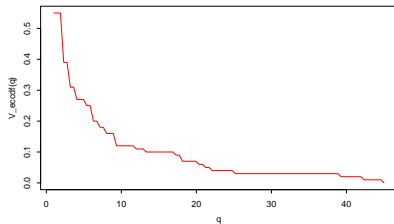
```
# compare the fitted and empirical pmf
plot(deg_dist,col = "red")
V_density<-Vectorize(yule_density, vectorize_args=c("q"))
curve(V_density, 1, max(deg), xname = "q", add=T,col = "blue" )
```



R Code

```
# compare the fitted and empirical ccmf
cf<-ecdf(deg)
eccdf <- function( q){
  1- cf(q)
}
V_eccdf<-Vectorize(eccdf, vectorize.args=c("q"))
curve(V_eccdf, 1, max(deg), xname = "q",col = "red") # empirical ccdf

V_ccdf<-Vectorize(yule_ccdf, vectorize.args=c("q"))
curve(V_ccdf, 1, max(deg), xname = "q", add=T,col = "blue")# theoretical ccdf
```



The rank model (Fortunato et al., 2006)

- First, a prestige ranking criterion is selected, and the previous t nodes are ranked according to prestige.
- At the $t + 1$ th iteration, the new node $t + 1$ is created and new links are set from it to m preexisting nodes with the linking probability that node $t + 1$ be connected to node j depends only on the rank:

$$P(t + 1 \rightarrow j) = \frac{r_j^{-\alpha}}{\sum_{k=1}^t r_k^{-\alpha}},$$

where $\alpha > 0$ is a real-valued parameter. The linking probability clearly decreases with increasing rank.

- **Open question:** how are about social distance instead of rank?

An example

- If $r_t = t$, at each iteration, a constant number m of links are created between the new node and the older ones. The expected total number of links $k^N(r)$ that the r th node has attracted after N nodes have been created is

$$k^N(r) = \sum_{t=r+1}^N \frac{mr^{-\alpha}}{\sum_{j=1}^t j^{-\alpha}}$$

- The ratio $k^N(r)/N$ in the limit of large N yields the probability $p(k)$ that a node of the network has degree k :

$$P(k) \sim k^{-(1+1/\alpha)}.$$

- Preferential attachment mechanisms in growing networks lead to scale-free networks, even in the absence of a complete knowledge of the values of the nodes' attributes.