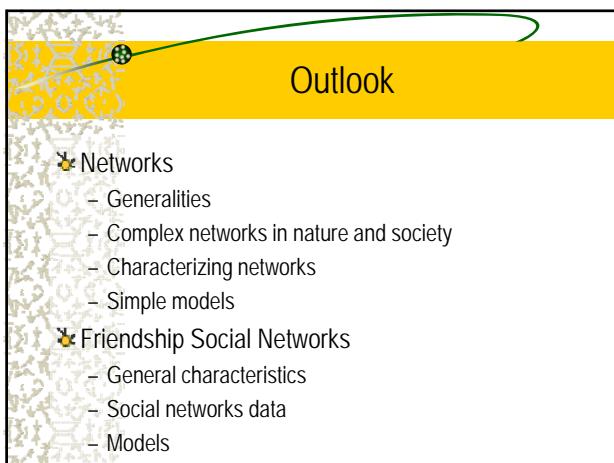


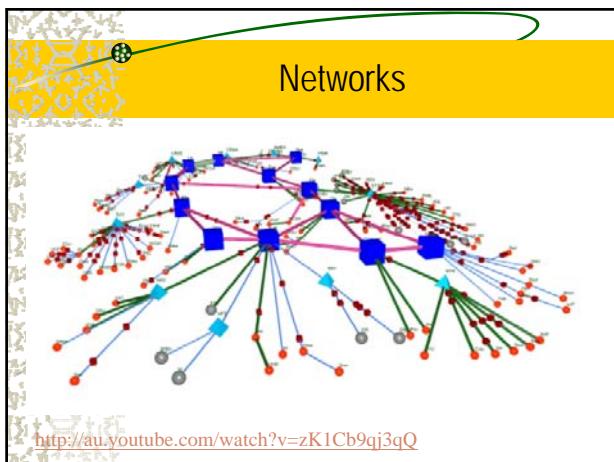
Friendship Network Formation

Albert Diaz-Guilera
<http://albert.diaz.guilera.googlepages.com>

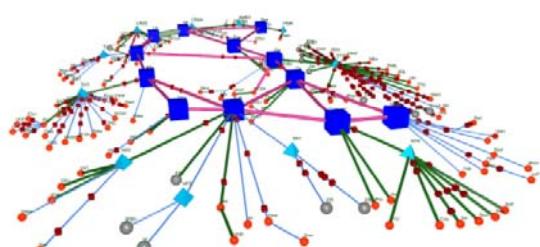


Outlook

- ✿ Networks
 - Generalities
 - Complex networks in nature and society
 - Characterizing networks
 - Simple models
- ✿ Friendship Social Networks
 - General characteristics
 - Social networks data
 - Models



Networks



<http://au.youtube.com/watch?v=zK1Cb9qj3qQ>

Why is network anatomy important

- Structure always affects function
- The topology of social networks affects the spread of information
- Internet
 - + access to the information
 - electronic viruses

Representations

- Network (graph): nodes (vertices) + links (edges)



- Network (matrix): adjacency matrix

A diagram showing three nodes labeled A, B, and C. Node A is at the top, B is at the bottom left, and C is at the bottom right. There is a horizontal edge between A and C. To the right of the nodes is a matrix labeled A =

$$A = \begin{pmatrix} 0 & 1 & 1 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}$$

How are the networks we find?

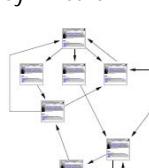
- Nodes and links: binary
- Bipartite: two types of nodes
- Directed: nonsymmetric adjacency matrices
- Weighted: the links can have strength

Common network

- Neural network of C. Elegans (a worm)
- Nodes = neurons
- Links = synapses

- Bipartite graph
- Coauthorship network

- ❖ Sometimes relational data has a direction
- ❖ The adjacency matrix is not symmetric
- ❖ Examples:
 - links to web pages
 - information
 - cash flow
 - trust networks



The diagram illustrates a directed graph structure representing the World Wide Web. It consists of several rectangular nodes, each containing a small icon and some text. Directed edges, shown as arrows, connect these nodes. Some nodes have multiple outgoing arrows pointing to other nodes, while others have incoming arrows pointing from other nodes. One node at the bottom is explicitly labeled "World Wide Web".

Weighted networks

- * Nodes and/or links can have different strengths
- * Population of cities (nodes)
- * Number of transported passengers (links)

Complex networks in nature and society

- * NOT regular lattices
- * NOT random graphs "simple" mathematical analysis
- * Huge databases and computer power

Social I: Networks of collaboration

- * Through collaboration acts
- * Examples:
 - movie actor
 - board of directors
 - musicians
 - scientific collaboration networks
(MEDLINE, Mathematical, neuroscience, e-archives,..)

Social networks (II)

- Email
- Trust networks
- Sexual contacts

Economical systems

- Relations between stocks
- Relations between economic agents
- World Trade Web

Directly measurable connectivities

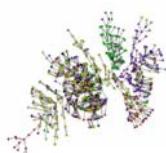
- World Trade Web (manufactured metal)

Indirect measures

- Correlations between dynamical properties:
returns, GDP,...

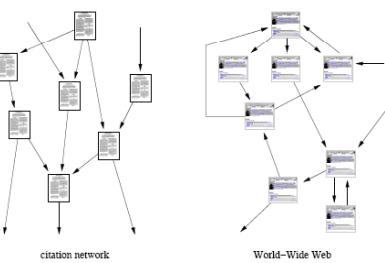
$$\rho_{ij}(\Delta t) = \frac{\langle p_i p_j \rangle - \langle p_i \rangle \langle p_j \rangle}{\sqrt{(\langle p_i^2 \rangle - \langle p_i \rangle^2)(\langle p_j^2 \rangle - \langle p_j \rangle^2)}}.$$

$$d_{ij} = \sqrt{2(1 - \rho_{ij})}$$



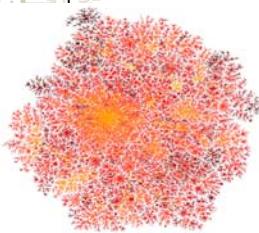
Information networks

- The World-Wide Web
 - Citation networks



Technological Networks

- Internet (not the same than WWW): hardwired
 - Power grid
 - Transportation networks



Characterizing networks: Topological properties

- Degree distribution
- Clustering
- Shortest paths
- Betweenness
- Correlations

Degree (microscopic scale)

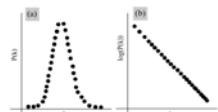
- Number of links that a node has
- It corresponds to the local centrality in social network analysis
- It measures how important is a node with respect to its nearest neighbors

Degree distribution (macroscopic scale)

- Gives an idea of the spread in the number of links the nodes have
- $P(k)$ is the probability that a randomly selected node has k links

What should we expect?

- In regular lattices all nodes are identical
- In random networks the majority of nodes have approximately the same degree



- Real-world networks: this distribution has a power tail

$$P(k) \approx k^{-\gamma} \text{ "scale-free" networks}$$

Clustering

- Cycles in social network analysis language
- Circles of friends in which every member knows each other

Clustering coefficient

- Clustering coefficient of a node

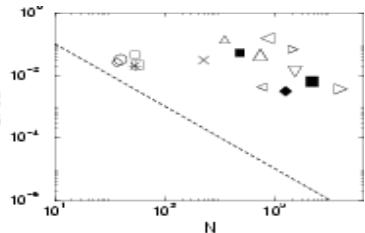
$$C_i = \frac{E_i}{k_i(k_i-1)/2}$$

- Clustering coefficient of the network

$$C = \frac{1}{N} \sum_{i=1}^N C_i$$

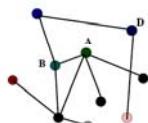
What happens in real networks?

- The clustering coefficient is much larger than it is in an equivalent random network



Distance between two nodes

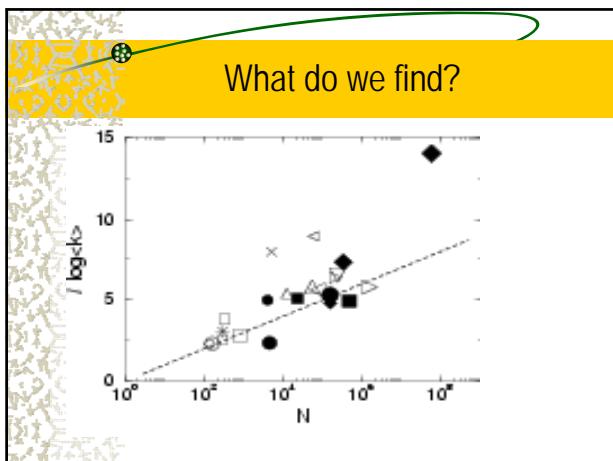
- Number of links that make up the shortest-path between two nodes



- Centrality: nodes that are "close" to many other nodes in the network.
- Global centrality: defined as the sum of minimum distances to any other nodes in the networks

Global centrality of the whole network?

Mean shortest path = average over all pairs of nodes in the network



Betweenness

- Measures the "intermediary" role in the network
- It is a set of matrices, one for each node

B_{ij}^k Ratio of shortest paths between i and j that go through k

$0 \leq B_{ij}^k \leq 1$ There can be more than one geodesic between i and j

$B_k = \sum_{ij} B_{ij}^k$ is a measure of the centrality, in terms of flow, of node k

Correlations

- Degree correlations: expected degree of the neighbors of a node as a function of its degree

$$k_{\text{nn}}(k) = \sum_{k'} k' P(k'|k)$$

Models of complex networks

- * The Erdős-Rényi (ER) random graph model
- * The Watts-Strogatz (WS) "small-world" model
- * The Barabasi-Albert (BA) "scale-free" model

Erdős-Renyi random model

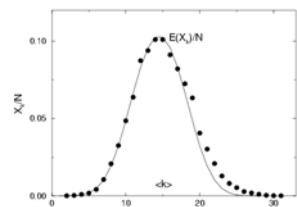
- * start with N nodes, every pair of nodes being connected with probability p
- * The total number of links, n , is a random variable
- $E(n)=pN(N-1)/2$
- * The degree of a node follows a binomial distribution

$$P(k_i = k) = \binom{N-1}{k} p^k (1-p)^{N-1-k}$$

Large N

- * Poisson distribution:

$$P(k_i = k) = \frac{\lambda^k}{k!} e^{-\lambda}$$



Mean shortest path

$$l_{\text{rand}} \approx \frac{\ln N}{\ln \langle k \rangle} \approx \frac{\ln N}{\ln pN}$$

Clustering coefficient

- * Probability that two nodes are connected (given that they are connected to a third)?



$$C_{\text{rand}} = p = \frac{\langle k \rangle}{N}$$

$$\frac{C_{\text{rand}}}{\langle k \rangle} \approx \frac{1}{N} \quad \text{while it is constant for real networks}$$

Watts- Strogatz small-world model

- * Small world: the average shortest path length in a real network is small
- * Six degrees of separation (Milgram, 1967)
- * Local neighborhood + long-range friends
- * A random graph is a small world

• Networks in nature (empirical observations)

$$l_{\text{network}} \approx \ln(N)$$

$$C_{\text{network}} \gg C_{\text{random graph}}$$

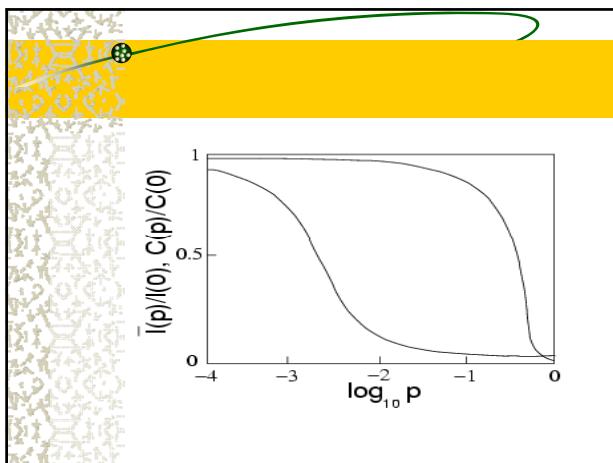
Model proposed

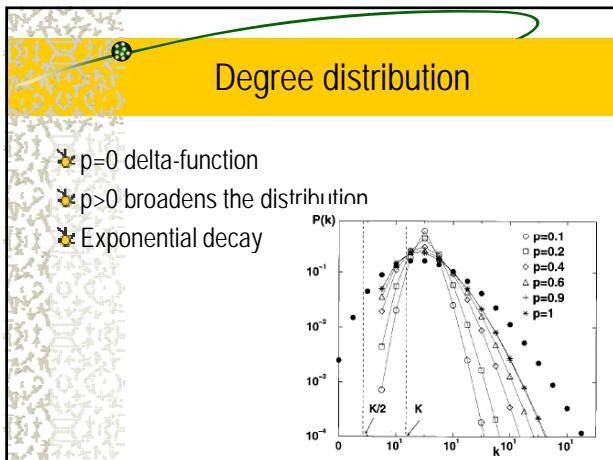
- Crossover from regular lattices to random graphs
- Tunable
- Small world network with (simultaneously):
 - Small average shortest path
 - Large clustering coefficient (not obeyed by ER-RG)

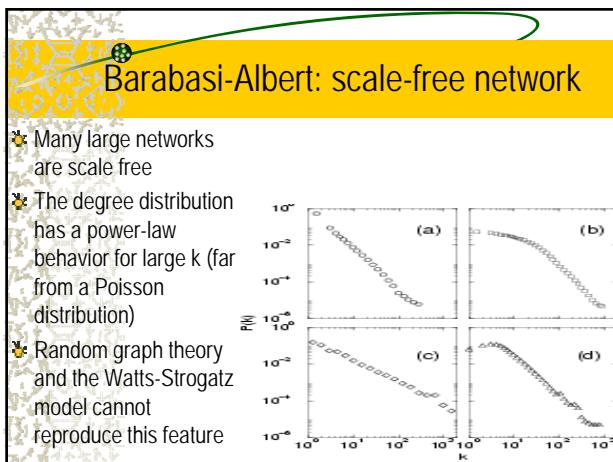
Two ways of constructing

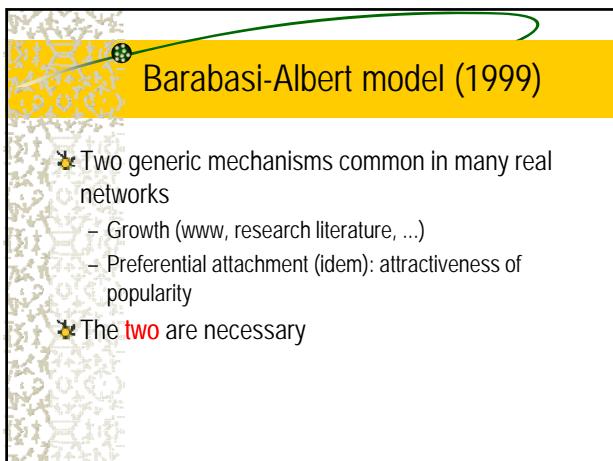
a)
rewiring of links

b)
addition of links

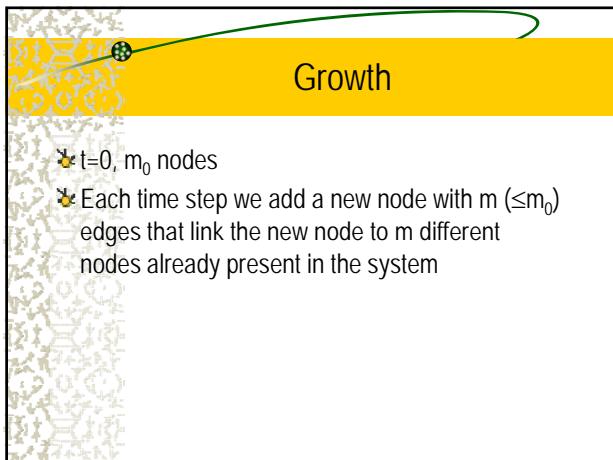




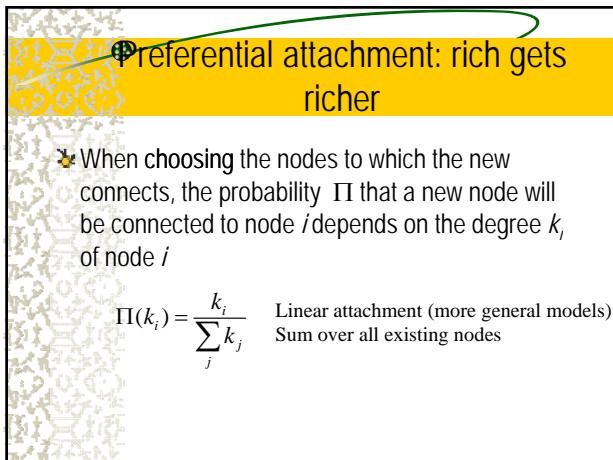


 Barabasi-Albert model (1999)

- * Two generic mechanisms common in many real networks
 - Growth (www, research literature, ...)
 - Preferential attachment (idem): attractiveness of popularity
- * The **two** are necessary

 Growth

- * $t=0$, m_0 nodes
- * Each time step we add a new node with m ($\leq m_0$) edges that link the new node to m different nodes already present in the system

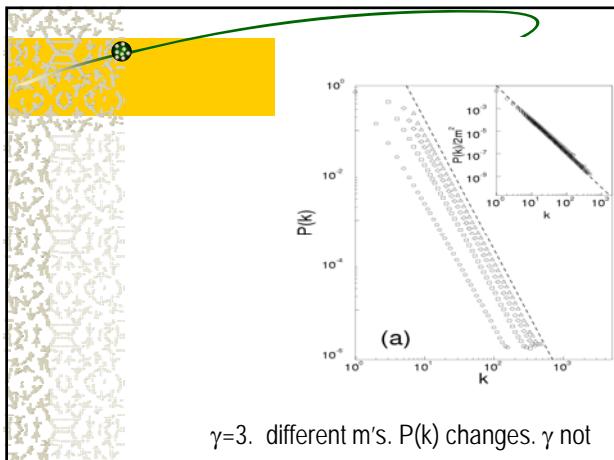
 Preferential attachment: rich gets richer

- * When choosing the nodes to which the new connects, the probability Π that a new node will be connected to node i depends on the degree k_i of node i

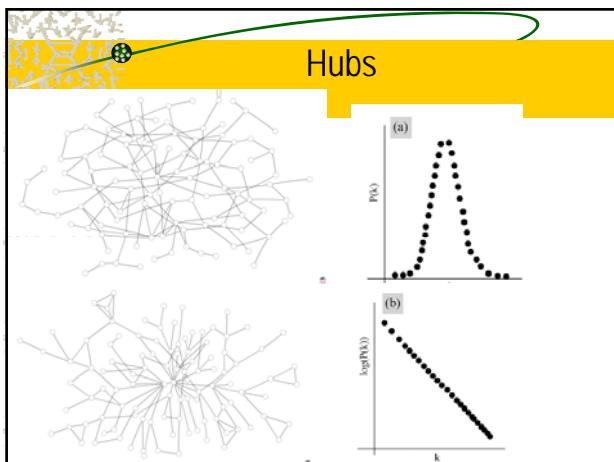
$$\Pi(k_i) = \frac{k_i}{\sum_j k_j} \quad \begin{array}{l} \text{Linear attachment (more general models)} \\ \text{Sum over all existing nodes} \end{array}$$

Numerical simulations

- Power-law $P(k) \approx k^{-\gamma}$ $\gamma_{SF} = 3$
- The exponent does not depend on m (the only parameter of the model)
- Calculation: continuity equation

$$P_k = \frac{2m(m+1)}{k(k+1)(k+2)}$$


$\gamma=3$. different m 's. $P(k)$ changes. γ not



The mesoscale: communities

- * Existence of communities or modules in networks
- * Technical issue: finding the best partition
- * Management issue: finding meaningful partitions

Communities

- * Social networks are formed by communities
- * Functional modules in biological networks
- * Many different origins: political, religious, function, ...

Technical issue

- * We have to identify the communities
- * How many possible partitions into communities?
- * NP problem to find the best one

Communities: intuitive picture

- Definition: subsets of nodes that are more densely linked, when compared with the rest of the network

Partition

- A partition is a division of the network into groups, communities or clusters
- The question is: Which of all possible partitions is the best?
- NP problem
- Community detection:
 - From computer scientists
 - To statistical physicists (Girvan-Newman, PNAS 99, 7821, 2002)

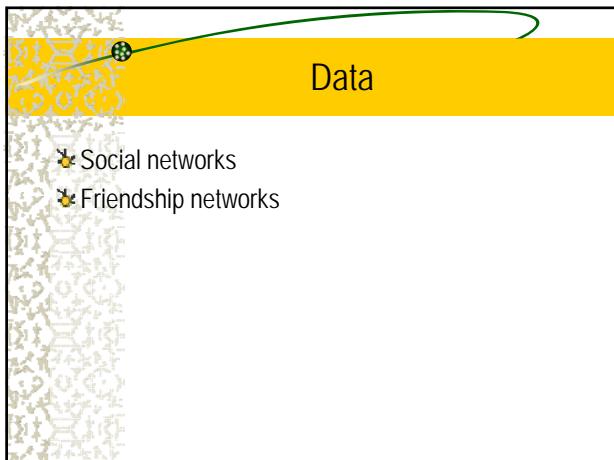
Quantifying a partition

- Modularity:
$$Q = \sum_i (e_{ii} - a_i^2)$$
- e_{ij} : fraction of total links starting at a node in partition i and ending at a node in partition j
- a_i : fraction of links connected to i
- a_i^2 : number of intracommunity links



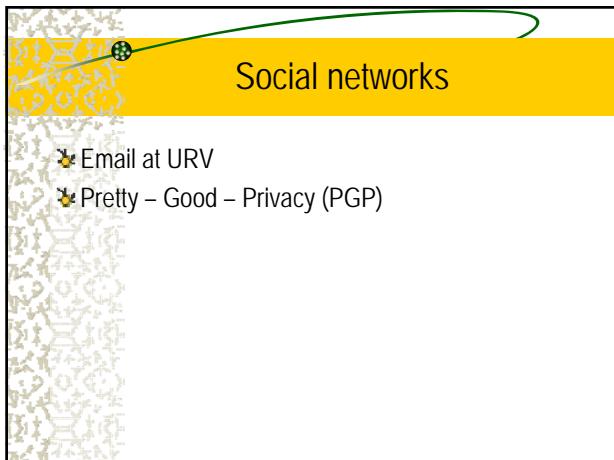
Identifying communities

- Identifying what communities are
- How a group is created, organized, grown, ...



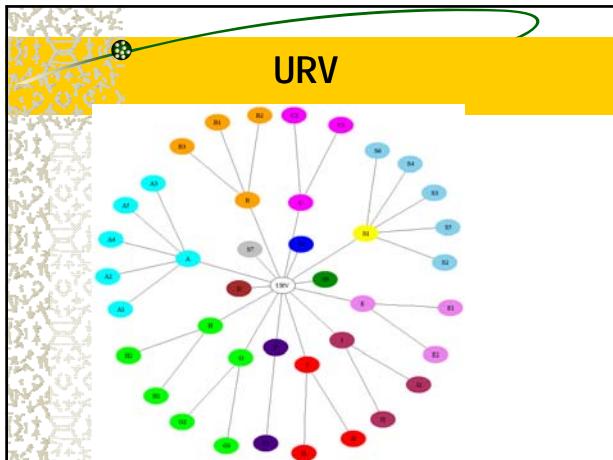
Data

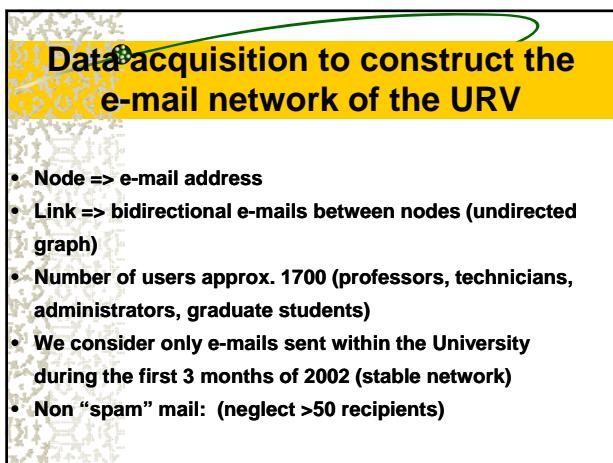
- Social networks
- Friendship networks

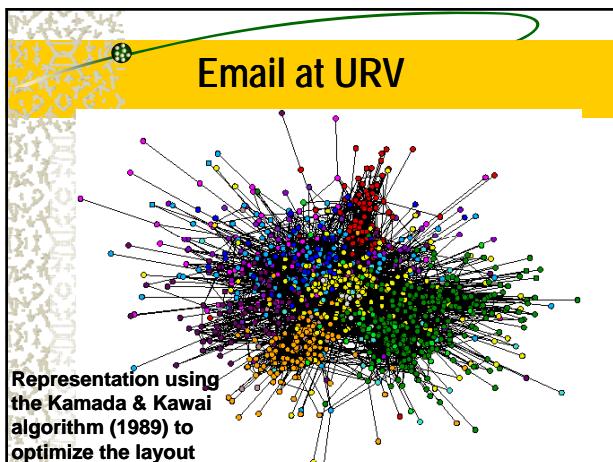


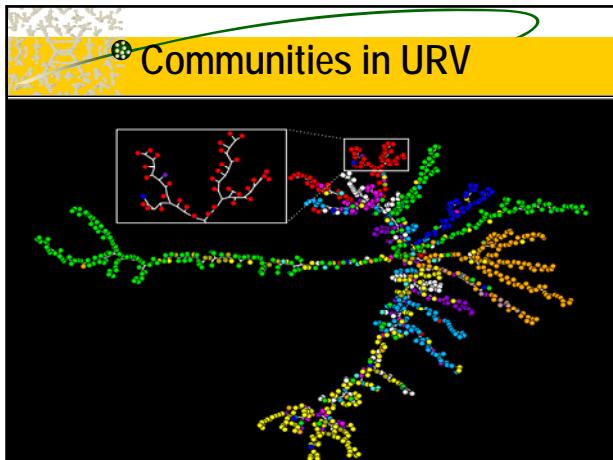
Social networks

- Email at URV
- Pretty – Good – Privacy (PGP)

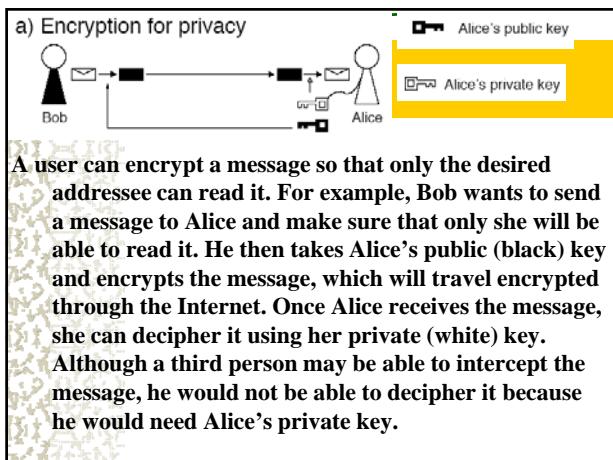












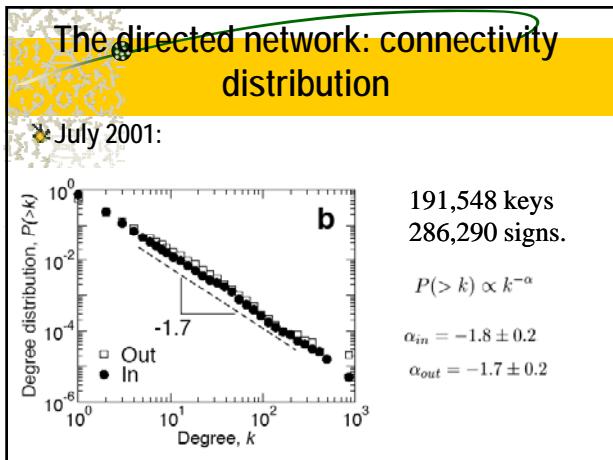
b) Encryption for authentication

Alice's public key
Alice's private key

A user can digitally authenticate a message so that the recipient of the information can verify the authenticity of the information's origin. For example, Alice wants to send a message to Bob and wants to make sure that he knows with absolute certainty that she is the sender of the message. To this end, she encrypts the message with her private key. If Bob is able to decipher it using Alice's public key, then the message must have originated with Alice's private key and, moreover, the message must be intact.

The web of PGP signatures

- Use keys of users you know directly, over the phone, for instance
- A central authority that certifies signatures?
- PGP web of trust users sign public keys certifying a user is who it claims to be
- This creates a trust relation
- A directed network of trust



Friendship networks data

- 友谊 of graduate students in Japan
- Instant messaging
- Erasmus and Eurovision
- Students in US

Japanese graduate students

Table 1. Observed data

Organization	School of Knowledge Science, Japan Advanced Institute of Science and Technology (JAIST).
Participants	Graduate students who have just entered the school.
Period	Every week from April 2005 to May 2005. 8 weeks for group A and 7 weeks for group B.
Group	People who are attending a lecture form a group. Two groups (two lectures) are investigated. The number of members in group A and group B are 26 and 16, respectively. The groups do not overlap.
Method	Questionnaire that enquires about the friendship among the people in a particular group.

Figure 1: Start-up network (left) and the network at the end of the survey (right) in group A.

Figure 2: Start-up network (left) and the network at the end of the survey (right) in group B.

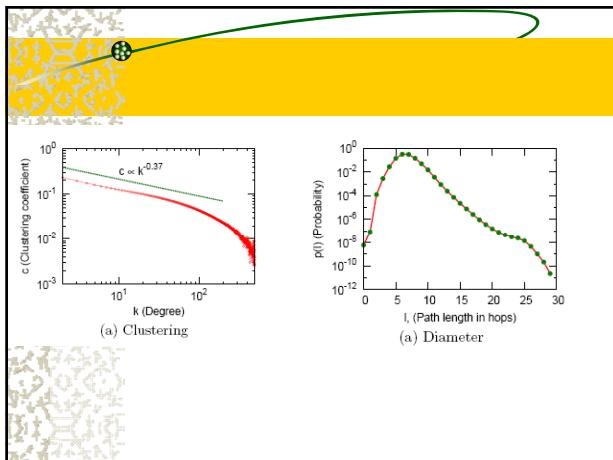
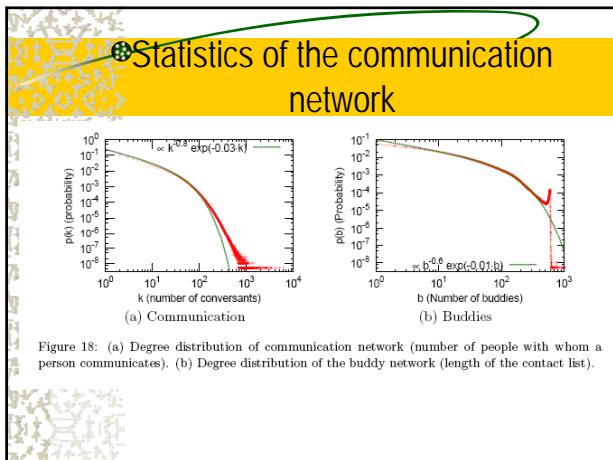
Biases

- Reciprocity and sibling: links created based on the network structure. Important for the model

Figure 8: Ratio of types of biased links

Instant messaging: a world-wide networks

- Presence events:** These include login, logout, first ever login, add, remove and block a buddy, add unregistered buddy (invite new user), change of status (busy, away, be-right-back, idle, etc.). Events are user and time stamped.
- Communication:** For each user participating in the session, the log contains the following tuple: session id, user id, time joined the session, time left the session, number of messages sent, number of messages received.
- User data:** For each user, the following self-reported information is stored: age, gender, location (country, ZIP), language, and IP address. We use the IP address to decode the geographical coordinates, which we then use to position users on the globe and to calculate distances.



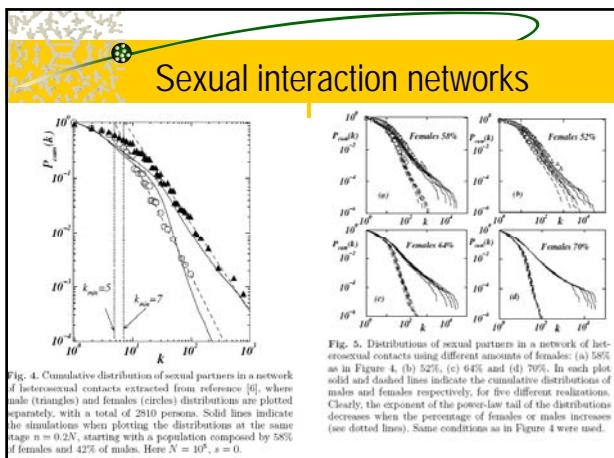


Fig. 4. Cumulative distribution of sexual partners in a network of heterosexual contacts extracted from reference [6], where male (triangles) and female (circles) distributions are plotted respectively. The distributions are fitted with power-law tails. We also show simulations when plotting the distributions at the same stage $n = 0.2N$, starting with a population composed by 58% of females and 42% of males. Here $N = 10^5$, $s = 0$.

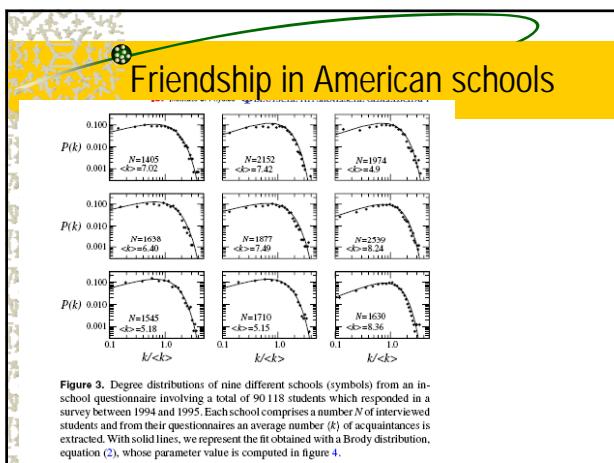
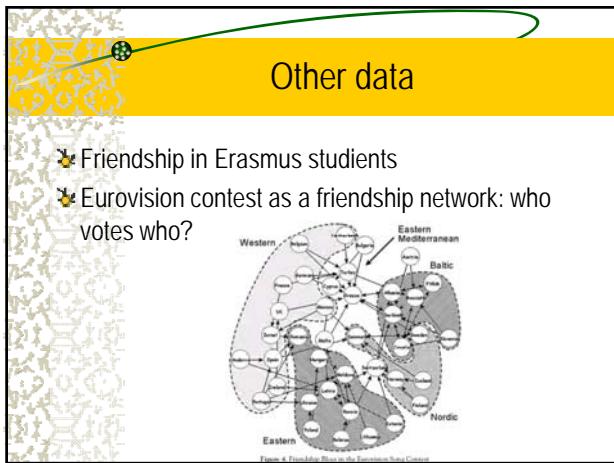


Figure 3. Degree distributions of nine different schools (symbols) from an in-school questionnaire involving a total of 90 118 students which responded in a survey between 1994 and 1995. Each school comprises a number N of interviewed students and from their questionnaires an average number $\langle k \rangle$ of acquaintances is extracted. With solid lines, we represent the fit obtained with a Brody distribution, whose parameter value is computed in figure 4.



Models

- Model for Japanese graduate schools
- Model based on collision
- Model based on social distance

Japanese model: biases

1. Add n_0 nodes to the network.
2. Add m links to the network based on the probability calculated by equation 6.
3. Add n_t nodes having no links to the network if necessarily.
4. Repeat steps 2 to 3 an arbitrary number of times.

$$\text{cad}(x \rightarrow y) = \theta_{\text{rep}} * \text{rep}(x \rightarrow y) + \theta_{\text{sib}} * \text{sib}(x \rightarrow y) + \theta_{\text{nob}} * \text{nobia}(x \rightarrow y),$$

$$\theta_{\text{rep}} \geq 0, \theta_{\text{sib}} \geq 0,$$

$$\theta_{\text{nob}} \geq 0,$$

$$\theta_{\text{rep}} + \theta_{\text{sib}} + \theta_{\text{nob}} = 1 \quad (5)$$

$$P(x \rightarrow y) = \text{cad}(x \rightarrow y) / \sum_{i \in X} \sum_{j \in Y} \text{cad}(i \rightarrow j) \quad (6)$$

where $\text{rep}(x \rightarrow y)$, $\text{sib}(x \rightarrow y)$, and $\text{nobia}(x \rightarrow y)$ denote the reciprocity bias, sibling bias, and absence of bias from nodes x to y , respectively. θ_{rep} , θ_{sib} , θ_{nob} are the simulation parameters.

Model for mobile agents

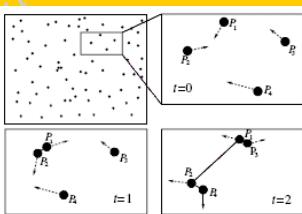


Figure 1. Illustration of the two-dimensional mobile agents system. Initially there are no connections between nodes and nodes move with some initial velocity v_0 in a randomly chosen direction (arrows). At $t = 1$ two nodes, P_1 and P_2 collide and a connection between them is introduced (solid line), velocities are updated increasing their magnitude and choosing a new random direction. At $t = 2$ two other collisions occur, between nodes P_2 and P_4 and between nodes P_1 and P_3 . In this way a network of nodes and connections between them emerges as a straightforward consequence of their motion (see text).

Preferential

Assuming that large numbers of acquaintances tend to favour the occurrence of new contacts, the velocity should increase with degree k , namely

$$\vec{v}(k_i) = (\bar{v}k_i^\alpha + v_0)\vec{\omega}, \quad (1)$$

- ✿ Correctly parameterized it reproduces results from sexual interactions, American schools, ...

A model based on social distance

- ✿ Relations are established between individuals that feel "close"
- ✿ Assign a position to each individual = a set of values in a number of social dimensions

1 social dimension

- ✿ Position of individual i in social space h_i
- ✿ Define a density $\rho(h)$
- ✿ Probability of establishing a relation decreases with social distance

$$r_n(h_i^n, h_j^n) = \frac{1}{1 + [b_n^{-1} d_n(h_i^n, h_j^n)]^{\alpha_n}}$$

- ✿ α measures the degree of homophily

Some results

- ★ Homogeneous density $\rho(h) = 1/h_{\max}$
- ★ Average degree

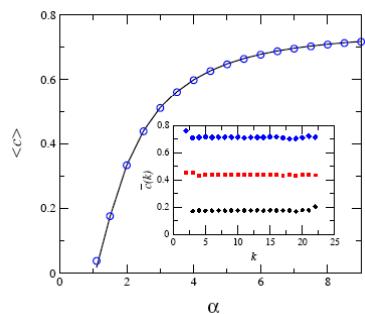
$$\langle k \rangle = \frac{2\delta b\pi}{\alpha \sin \pi/\alpha} \quad \delta = N/h_{\max}$$

- ★ Clustering coefficient

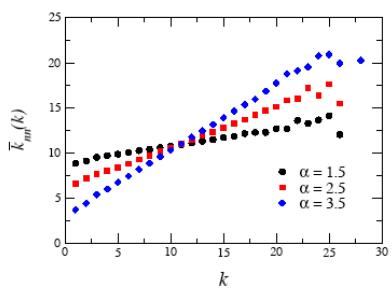
$$\langle c \rangle = \frac{\alpha^2}{4\pi^2} f(\alpha) \sin^2 \frac{\pi}{\alpha}$$

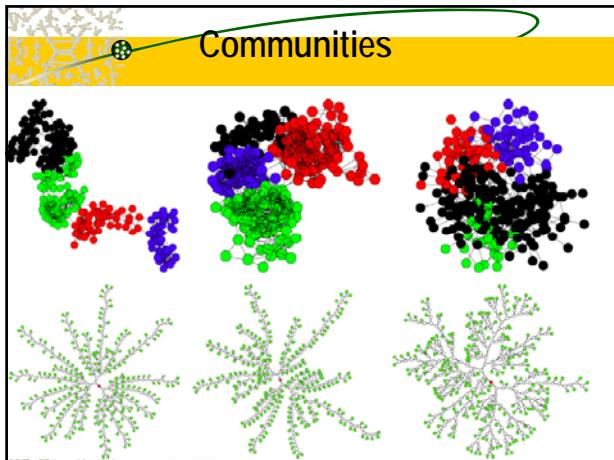
$$f(\alpha) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{dxdy}{(1+|x|^{\alpha})(1+|x-y|^{\alpha})(1+|y|^{\alpha})}$$

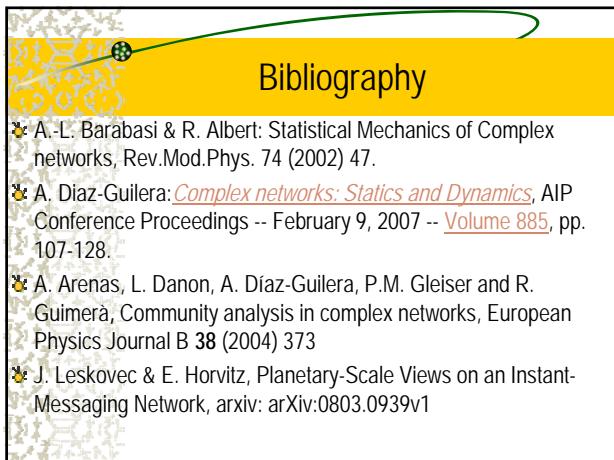
Clustering coefficient

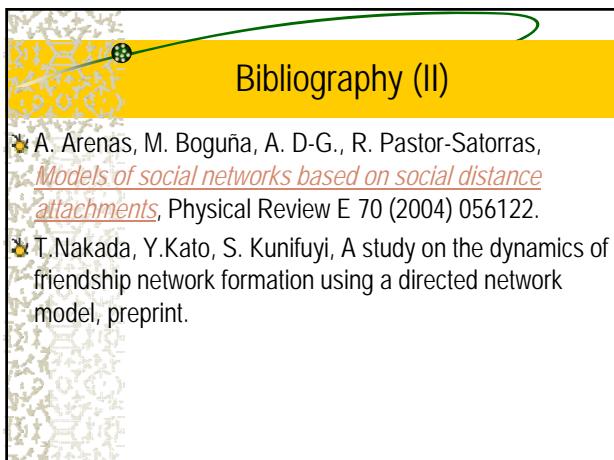


Assortativity: degree correlations









Bibliography (III)

- ★ A. de Federico de la Rúa, A., L'analyss longitudinal de reseaux sociaux totaux avec SIENA- Methode, discussion et application. *Bulletin de Méthodologie Sociologique*. N° 84, octobre, pp. 5-39.
- ★ A. Dekker, The Eurovision Song Contest as a « Friendship» Network, *Connections* 28 (2008) 59
- ★ P.J.Lind and H.J. Herrmann, *New approaches to model and study social networks*, *New.J.Phys.* 9 (2007) 228.
