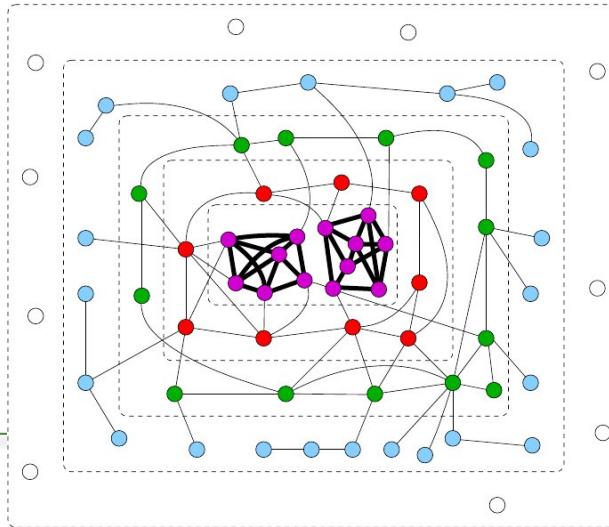


# Graph Mining



**Michalis Vazirgiannis**

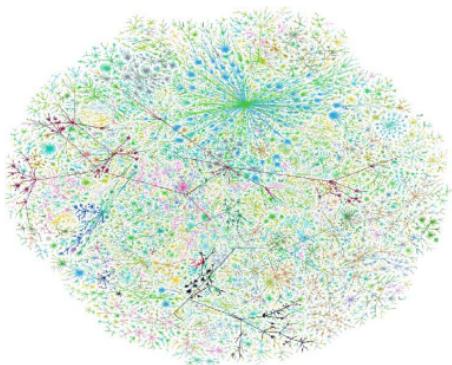
LIX @ Ecole Polytechnique

DSSP6

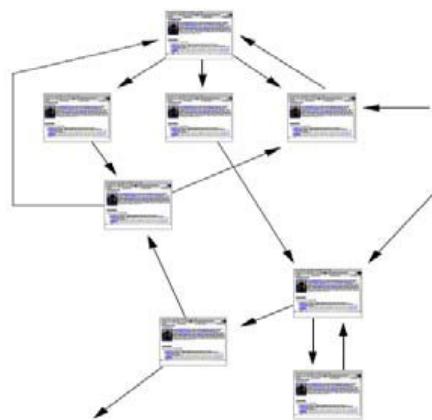
May 2018

# Networks are Everywhere

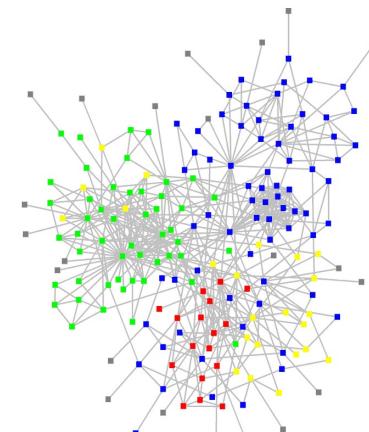
---



(a) Internet



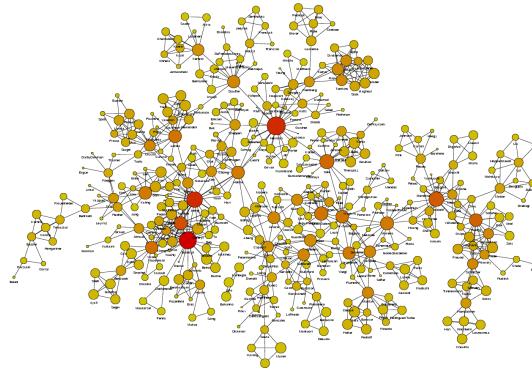
(b) World Wide Web



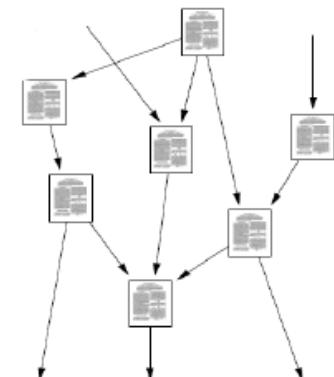
(c) Email network



(d) Social network



(e) Collaboration network



(f) Citation network

# Graphs are ubiquitous!

---

## ■ Technological networks:

- Internet
- Telephone networks
- Power grid
- Road, airline and rail networks

## ■ Information networks:

- World Wide Web
- Blog networks
- Citation networks

## ■ Social networks:

- Collaboration networks
- Organizational networks
- Communication networks

## ■ Biological networks:

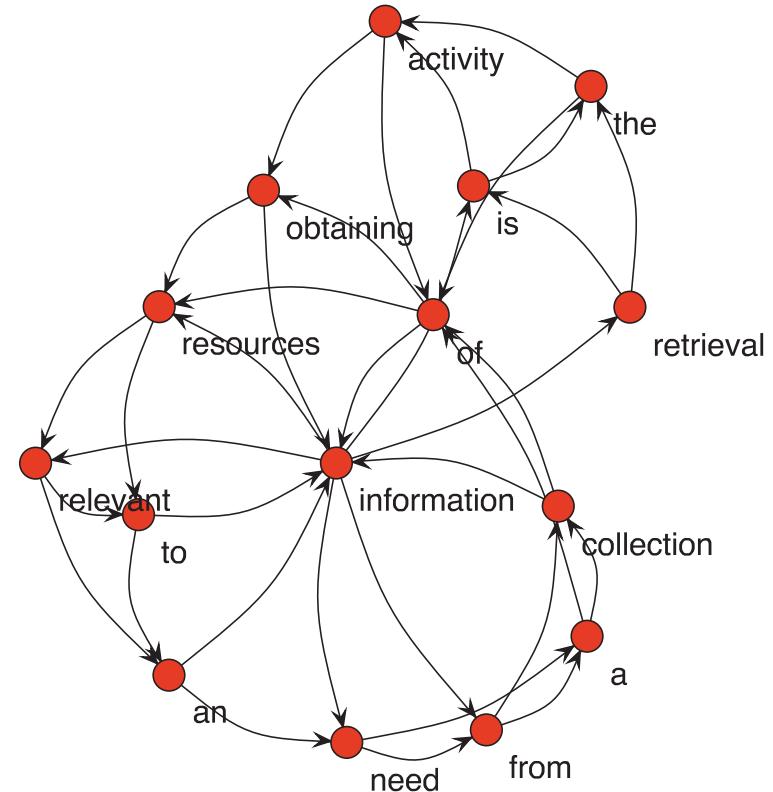
- Networks from Neuroscience
- Protein-protein interaction networks
- Gene regulatory networks
- Food webs

## ■ Software networks:

- Call graphs
- Software module/component interaction networks

# Even representing text - Graph-of-word

information retrieval is the activity of obtaining  
information resources relevant to an information need  
from a collection of information resources



"Graph of word approach for ad-hoc information retrieval", F. Rousseau, M. Vazirgiannis,  
Best paper mention award ACM CIKM 2013

# Elements of Learning from Graph data

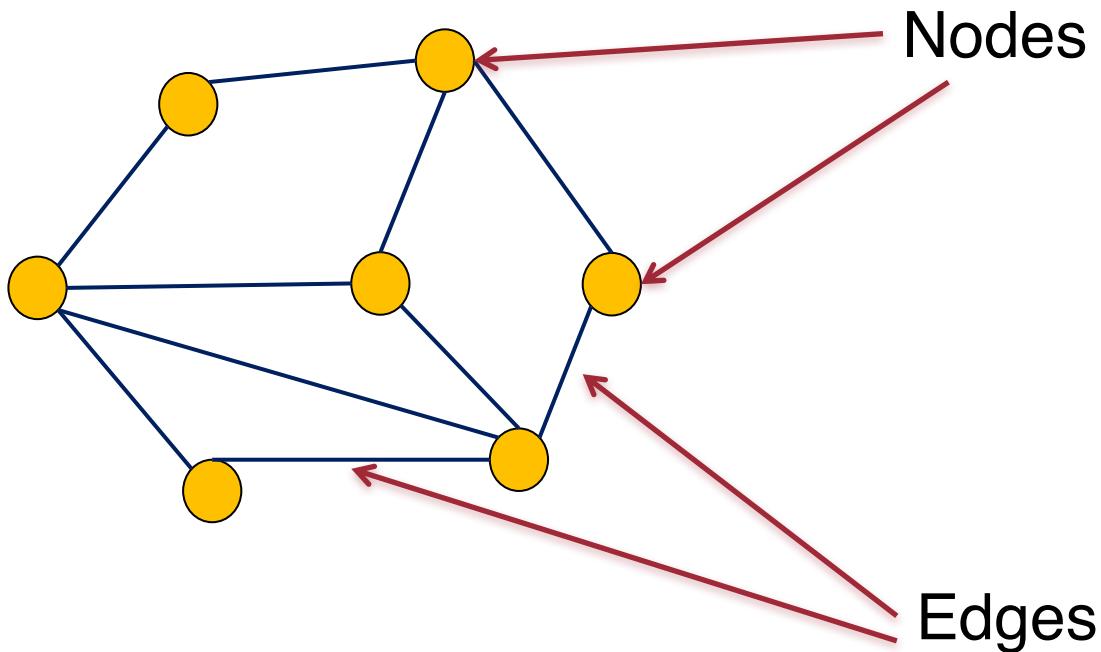
---

- **Graph models/ graph generators** graph generators  
(erdos reyni, preferential attachment, kronecker graphs)
- **Node base metrics:** - Ranking algorithms (Pagerank),  
Ranking evaluation measures (Kendal Tau, NDCG),
- **Graph exploration/preprocessing:** degree distributions,  
visualization
- **Supervised learning for graphs:** link prediction, graph  
kernels, graph classification
- **Unsupervised learning:** clustering, community mining,  
degeneracy.
- **Learning theory in graphs:** model ensembling/selection...

# Graphs and Networks

---

- Graphs allow for modeling dependencies



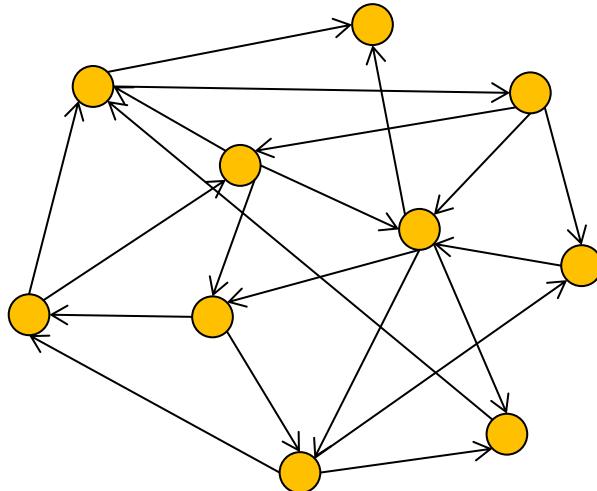
Nodes

Edges

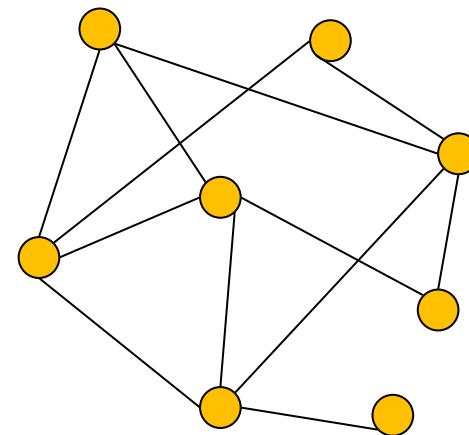
# Basic Graph Definitions

---

- A graph  $G=(V, E)$  consists of a set of **nodes**  $V$ ,  $|V|=n$  and a set of **edges**  $E$ ,  $|E|=m$
- Graphs can be **undirected** or **directed**



Directed



Undirected

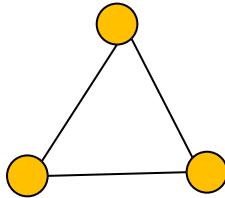
In-degree:  $d_{in}(i) = |\{j \mid (j,i) \text{ is edge}\}|$   
Out-deg:  $d_{out}(i) = |\{j \mid (i,j) \text{ is edge}\}|$

Degree:  $d(i) = d_{in}(i) = d_{out}(i)$

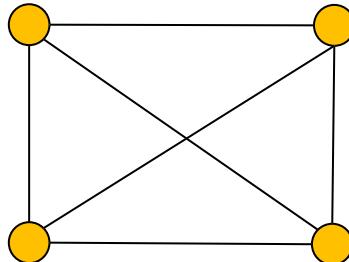
# Complete Graph

---

- **Definition:** A graph  $G=(V, E)$  is called complete  $K_n$  if every pair of nodes is connected by an edge



**Complete graph  
with 3 nodes:  
triangle**



**Complete graph  
with 4 nodes**

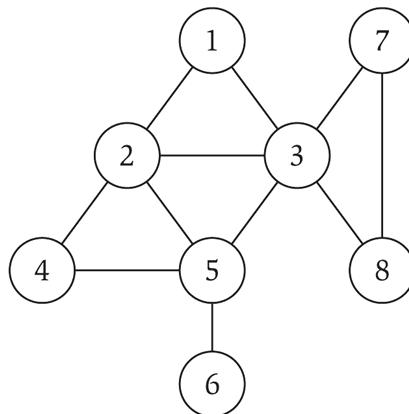
- What is the number of edges of a complete graph with  $n$  nodes?

- Note that, the notion of complete graphs is of particular importance for the problem of community detection
  - **Communities correspond to well-connected subgraphs**

# Graph Representation: Adjacency Matrix

---

- A graph can be represented by the adjacency matrix  $\mathbf{W}$ 
  - Matrix of size  $n \times n$ , where  $n$  is the number of nodes
  - $W_{ij} > 0$ , if  $i$  and  $j$  are connected
  - $W_{ij} = 0$ , if  $i$  and  $j$  are not connected
  - In case of unweighted graphs,  $W_{ij} = 1$ , if  $(i, j)$  is an edge of the graph
  - Space proportional to  $n^2$



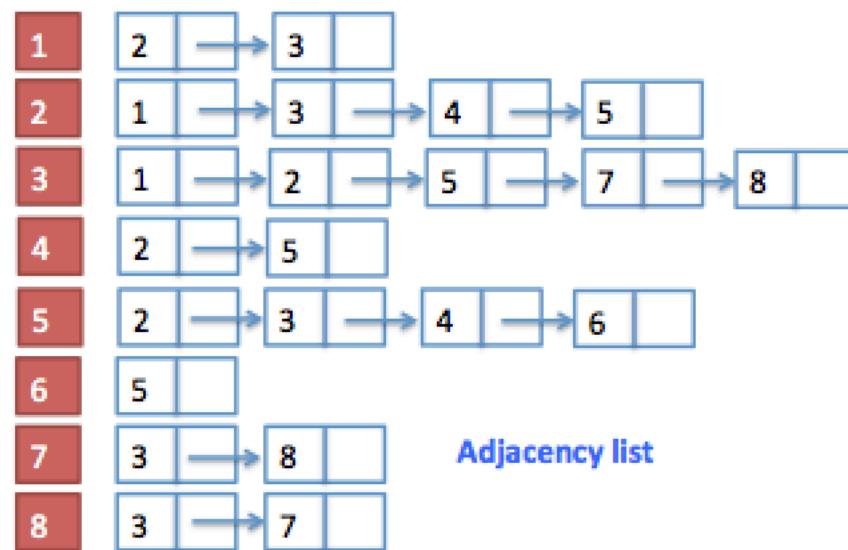
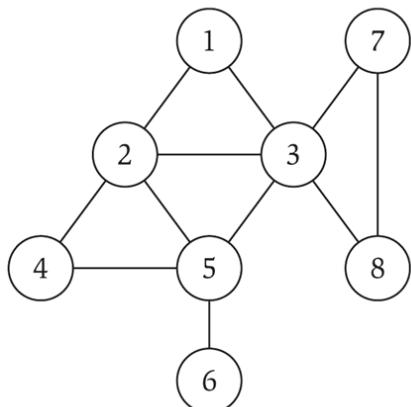
Undirected graph

0	1	1	0	0	0	0	0
1	0	1	1	1	0	0	0
1	1	0	0	1	0	1	1
0	1	0	0	1	0	0	0
0	1	1	1	0	1	0	0
0	0	0	0	1	0	0	0
0	0	1	0	0	0	0	1
0	0	1	0	0	0	1	0

Adjacency matrix

# Graph Representation: Adjacency Lists

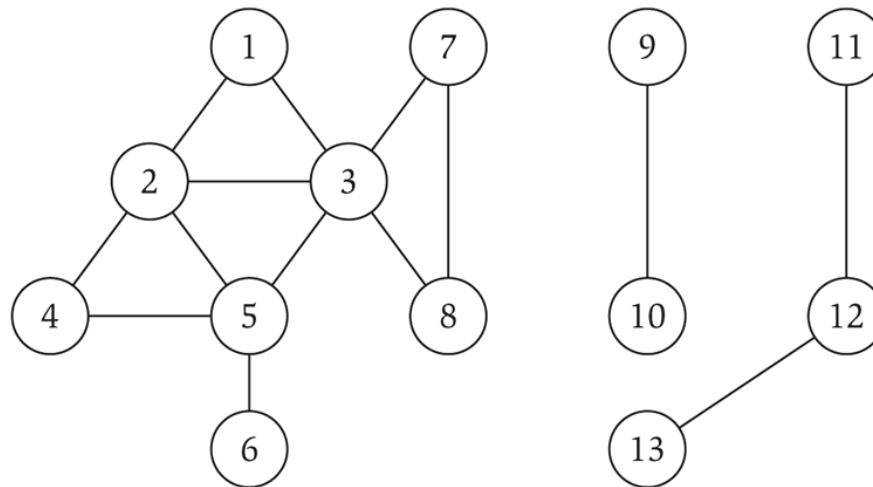
- Adjacency lists
  - Representation of a graph with  $n$  nodes using an array of  $n$  lists of nodes
  - List  $i$  contains node  $j$  if there is an edge  $(i, j)$
  - A weighted graph can be represented with a list of node/weight pairs
  - Space proportional to  $\Theta(m+n)$
  - Checking if  $(i, j)$  is an edge takes  $O(d_i)$  time



# Paths and Connectivity in Graphs

---

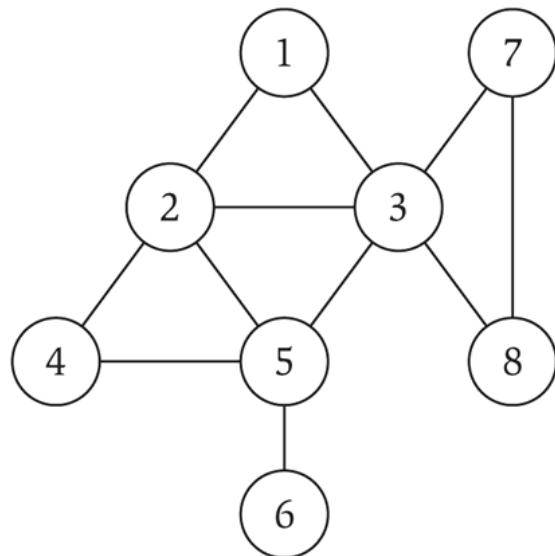
- **Definition:** A path in an undirected graph  $G=(V,E)$  is a sequence of nodes  $v_1, v_2, \dots, v_k$  with the property that each consecutive pair  $v_{i-1}, v_i$  is joined by an edge in  $E$
- **Definition:** An undirected graph is connected if for every pair of nodes  $u$  and  $v$ , there is a path between  $u$  and  $v$



# Cycles in Graphs

---

- **Definition:** A cycle is a path  $v_1, v_2, \dots, v_k$  in which  $v_1 = v_k$ ,  $k > 2$  and the first  $k-1$  nodes are all distinct



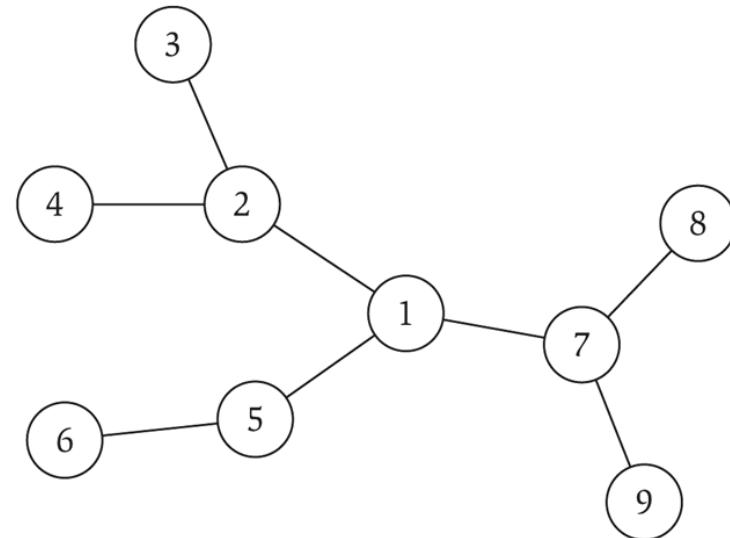
**Cycle**  $C = 1 - 2 - 4 - 5 - 3 - 1$

# Trees

---

- **Definition:** An undirected graph is a tree if it is connected and does not contain a cycle
- **Theorem:** Let **G** be an undirected graph with **n** nodes. Then, any two of the following statements imply the third:

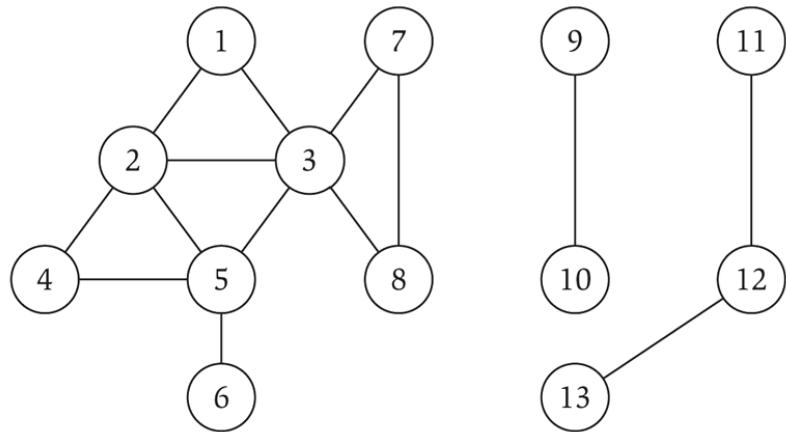
- **G** is connected
- **G** does not contain a cycle
- **G** has **n-1** edges



# Connected Components

---

- A **connected component** is a maximal connected subgraph of a graph **G** (there is a path between any pair of nodes)



Connected component containing node 1:  
 $\{1, 2, 3, 4, 5, 6, 7, 8\}$

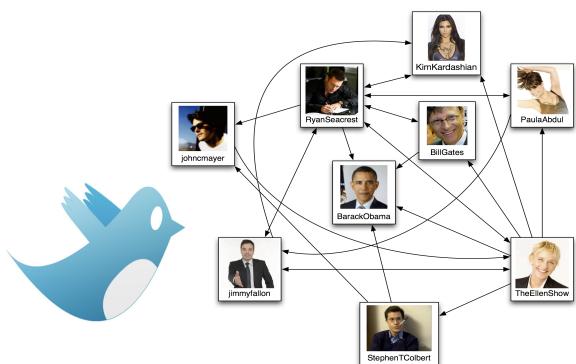
Graph with 3 connected components

**Question:** How can we compute the connected components of a graph?

**A:** Apply BFS

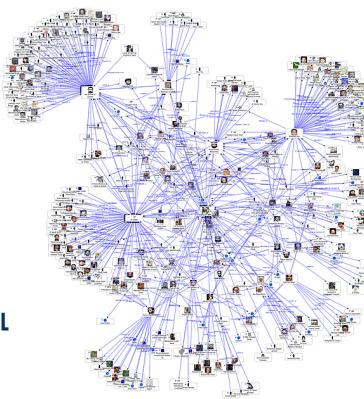
# Connectivity in Directed Graphs

- A plethora of network data from several applications is from their nature **directed**



Twitter

[Image: <http://sites.davidson.edu/mathmovement/>]

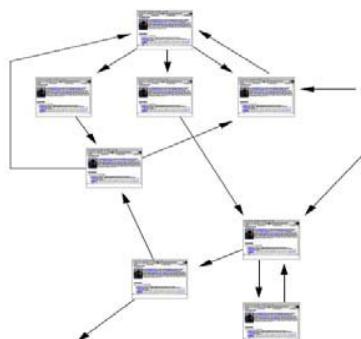


flickr



LIVEJOURNAL

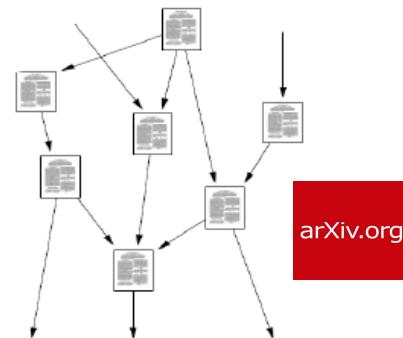
Online Social Networks



Web Graph



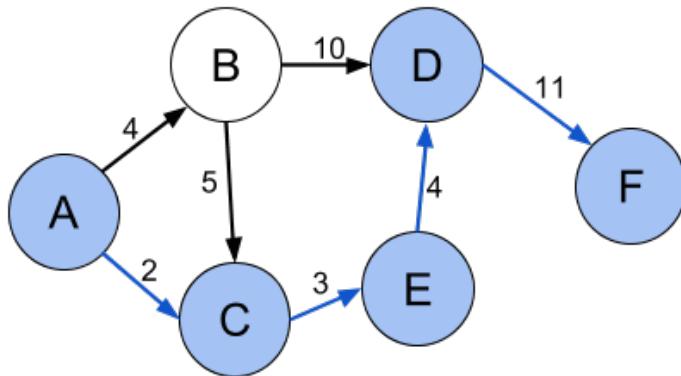
Wikipedia



Citation Graph

# Shortest Paths

- **Definition:** find a path between two nodes in a graph, in such a way that the sum of the weights of its constituent edges is minimized
  - Many applications (e.g., road networks)
  - **Single-source** shortest path problem
  - **Single-destination** shortest path problem
  - **All-pairs** shortest path problem



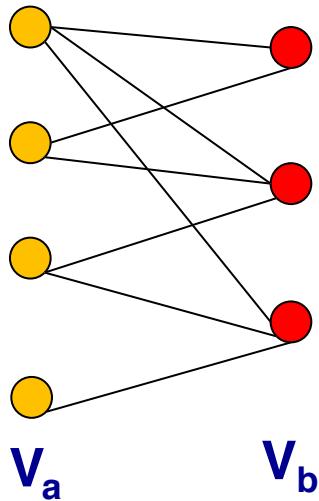
Many algorithms:  
• Dijkstra  
• Bellman-Ford

Shortest path (A, C, E, D, F) between vertices A and F in the weighted directed graph

# Bipartite Graphs

---

■ **Definition:** A graph  $G=(V,E)$  is called **bipartite** if the node set  $V$  can be partitioned into two disjoint sets  $V_a, V_b$  and every edge  $(u,v)$  connects a node of  $V_a$  to a node of  $V_b$



- Strong modeling capabilities and many real-world applications
- E.g., **Collaborative filtering** in recommender systems
  - Model the customer-product space using a bipartite graph (who-purchased-what)
  - If a user A has purchased the same product with a user B, then it is more likely to purchase another product as B did, than of a person selected randomly

# Properties of Real-World Graphs

---

■ Networks arising from **real-world** applications obey fascinating properties

## ■ **Static networks**

- Heavy-tailed degree distribution
- Small diameter
- Giant connected component (GCC)
- Triangle Power Law
- Community structure
- ...

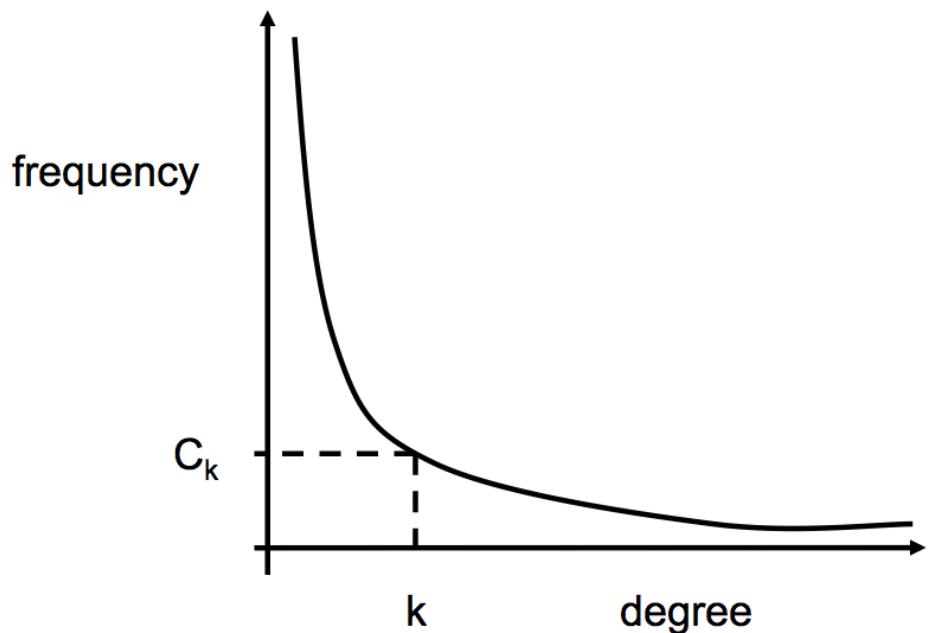
## ■ **Dynamic networks**

- Densification
- Small and shrinking diameter

# Degree Distribution

---

- The **probability distribution** of the degrees over the network



- Let  $C_k$  = number of nodes with degree  $k$
- **Problem:** find the probability distribution that **fits** best the **observed data**

# Power-law Degree Distribution

---

- Let  $C_k$  = number of nodes with degree  $k$

$$C_k = c k^{-\gamma}$$

with  $\gamma > 1$  and  $c$  a constant

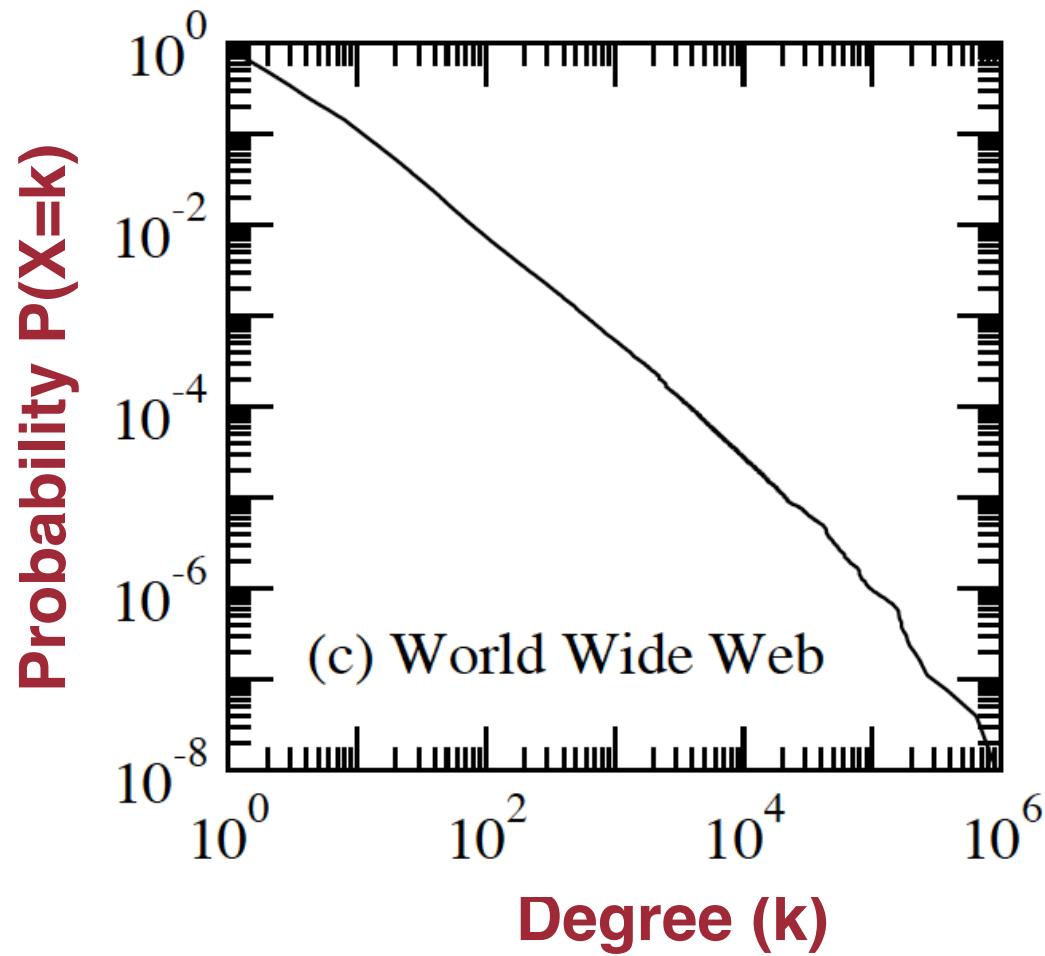
- How to recognize a power-law distribution?

$$\ln C_k = \ln c - \gamma \ln k$$

- Plotting  $\ln C_k$  versus  $\ln k$  gives a straight line with slope  $-\gamma \ln k$

## Power-law Degree Distribution in Real-Networks (1/2)

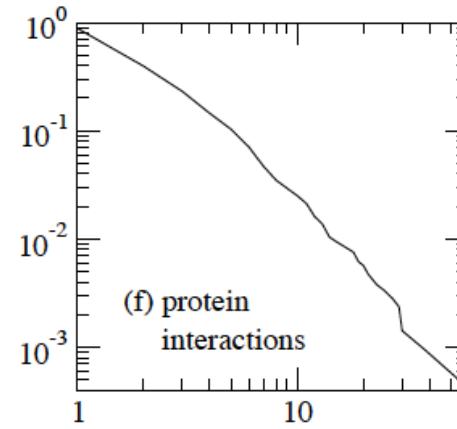
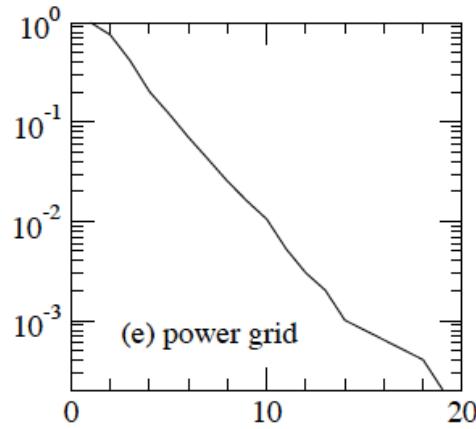
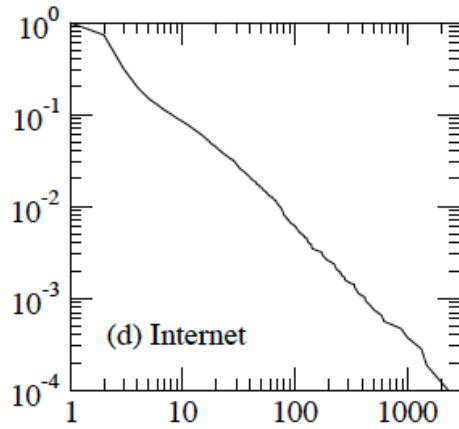
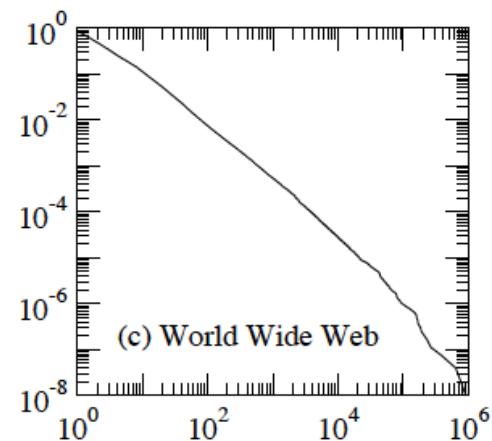
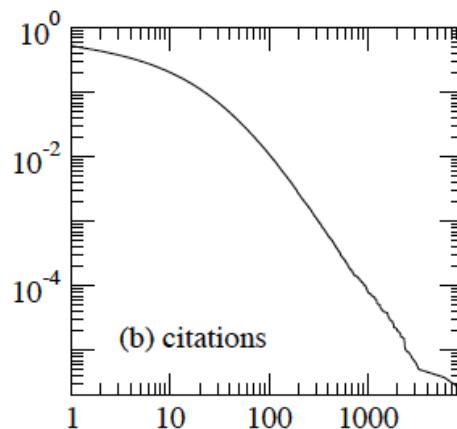
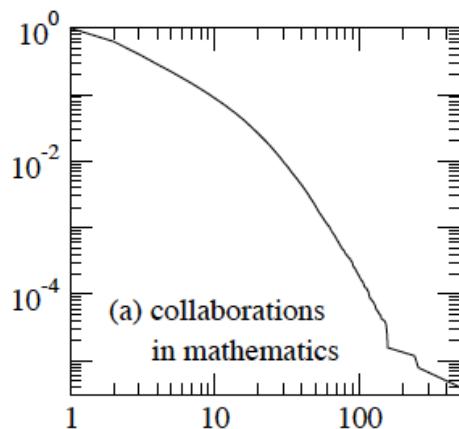
---



[Newman, 2003]

---

## Power-law Degree Distribution in Real-Networks (2/2)



Cumulative degree distribution for six different networks [Newman 2003]

# Power-law Degree Exponents

---

## ■ Power law degree exponent is typically $2 < \gamma < 3$

- Web graph [Broder et al., 2000]
  - $\gamma_{\text{in}} = 2.1, \gamma_{\text{out}} = 2.4$
- Autonomous systems (Internet graph) [Faloutsos et al., 1999]
  - $\gamma = 2.4$
- Actor collaborations [Barabasi and Albert, 2000]
  - $\gamma_{\text{in}} = 2.3$
- Citation graphs [Redner, 1998]
  - $\gamma_{\text{in}} = 3$
- MSN messenger graph [Leskovec et al., 2007]
  - $\gamma_{\text{in}} = 2$

[Leskovec, ICML, 2009]

---

# Summary – Degrees in Real Networks

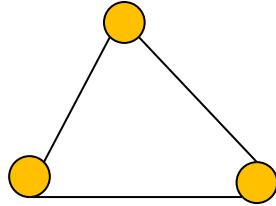
---

- The degree distribution is **heavily skewed**
  - Distribution is **heavy-tailed** (heavier tails compared to the exponential distribution)

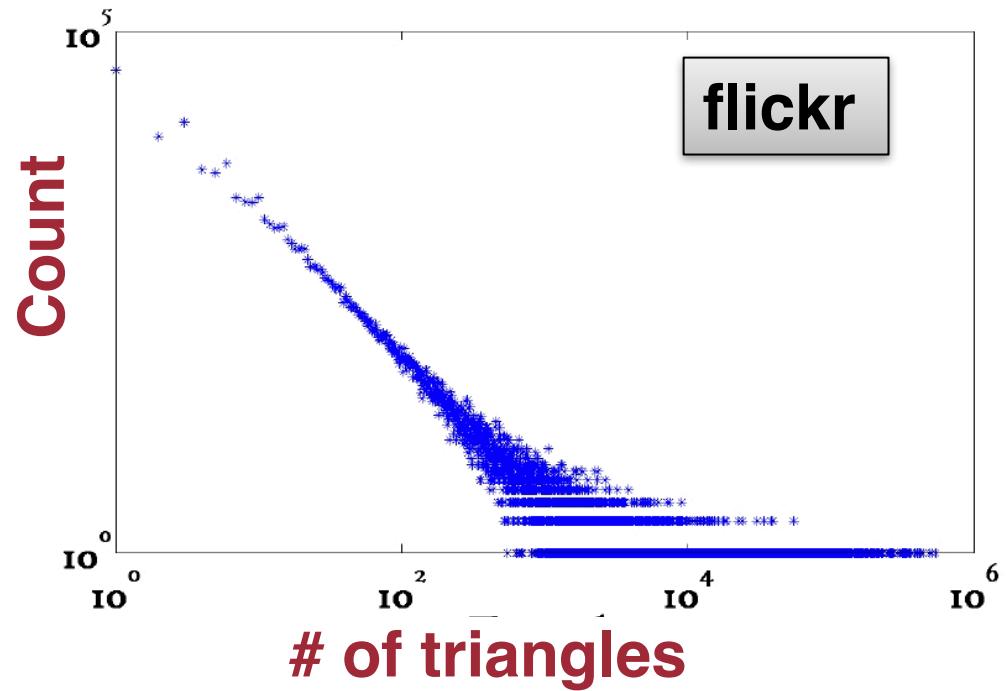
$$\lim_{x \rightarrow \infty} \frac{Pr(X > x)}{e^{-\epsilon x}} = \infty$$

- Various names and forms
  - Long tail, Zipf's law, Pareto distribution

# Triangle Participation Distribution



Complete graph  
with 3 nodes:  
triangle



- Number of nodes that participate in  $k$  triangles vs.  $k$  in log-log scale
- **Heavy-tailed** distribution

# Clustering Coefficient

---

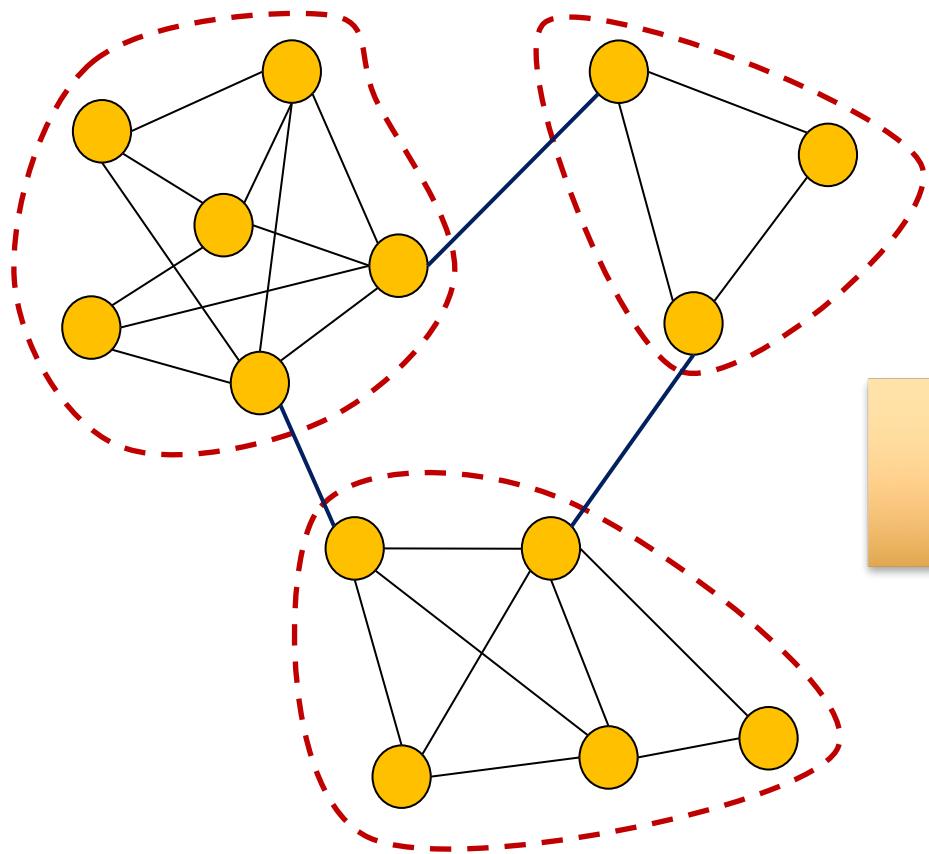
- Captures the tendency of the nodes of a graph to cluster together

$$T(G) = 3 \times \# \text{ of triangles in } G / \# \text{ of connected triplets}$$

- Captures the transitivity of clustering
  - If  $u$  is connected to  $v$  and  $v$  is connected to  $w$  ...
  - ... it is likely that  $u$  is also connected to  $w$
- Real-world networks tend to have high clustering coefficient
  - Connections to the existence of clustering and community structure property

# Community Structure

---

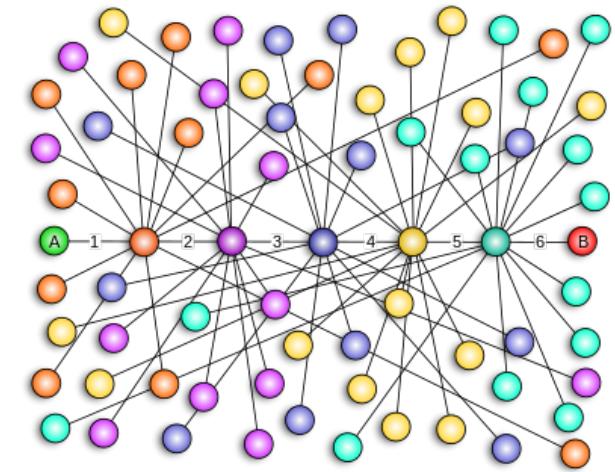
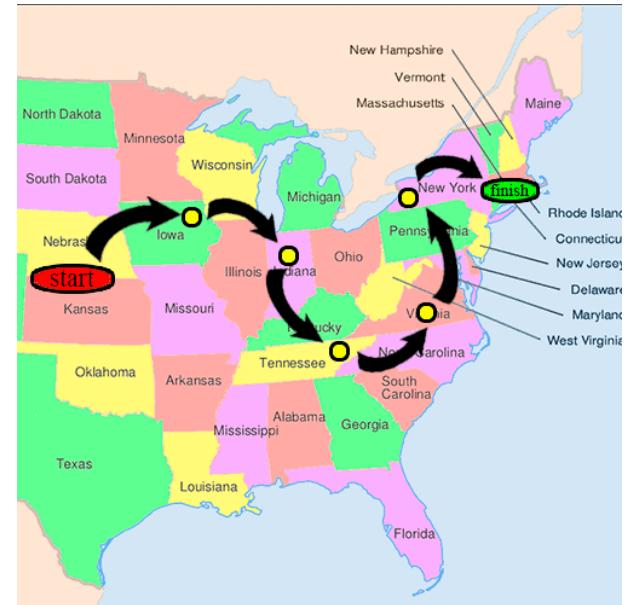


Example graph with  
three communities

- Will be covered later on in detail
-

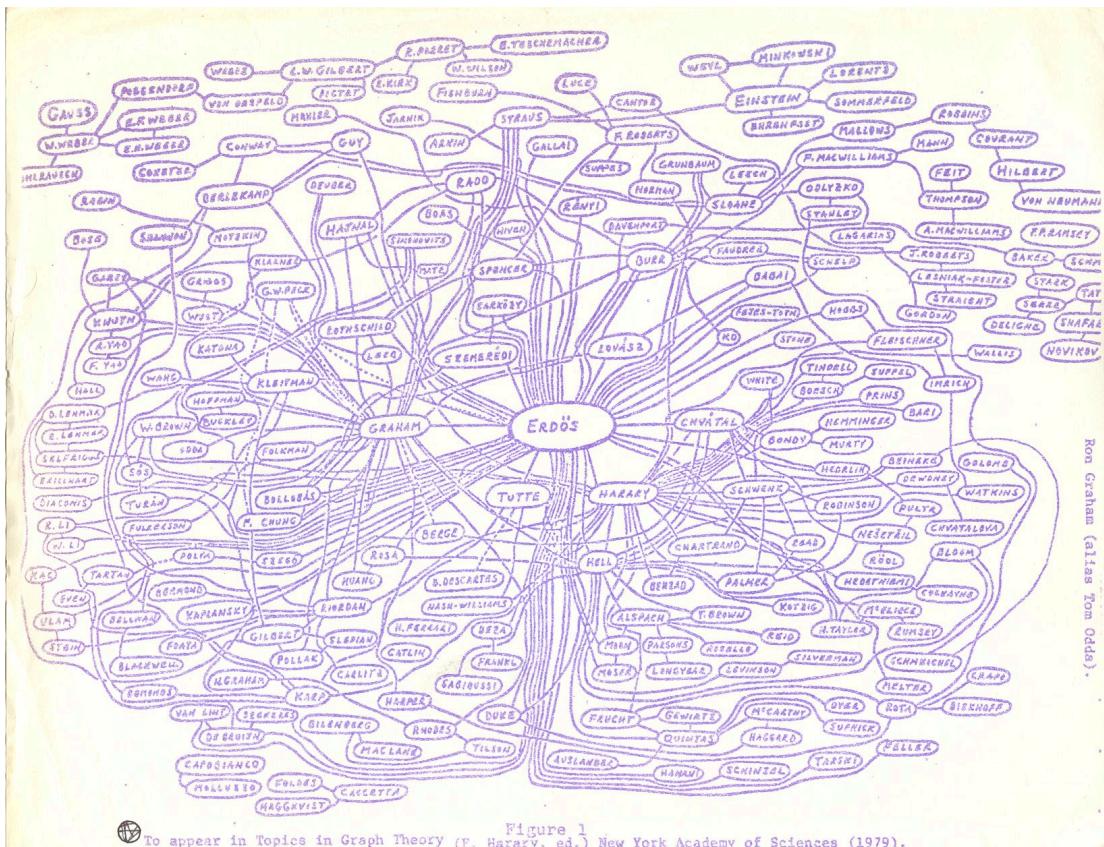
# Small-world Phenomenon (1/4)

- Six degrees of separation
  - Experiment done by sociologist Stanley Milgram (1960's)
  - Randomly selected people in Nebraska were asked to send letters to Boston, by contacting somebody with whom they had direct connection
    1. People either sent the letter directly to the recipient
    2. Or to somebody they believed had a high likelihood of knowing the target
- For those letters that reached their destination, the **average path length was 5.5 to 6**
  - Short paths are abundant in the networks
  - **Decentralized routing:** people are capable of discovering which links to follow to reach faster the target

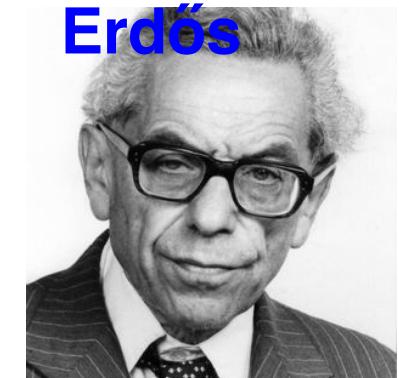


# Small-world Phenomenon (3/4)

- The small-world phenomenon appears in various network settings



# Paul Erdős



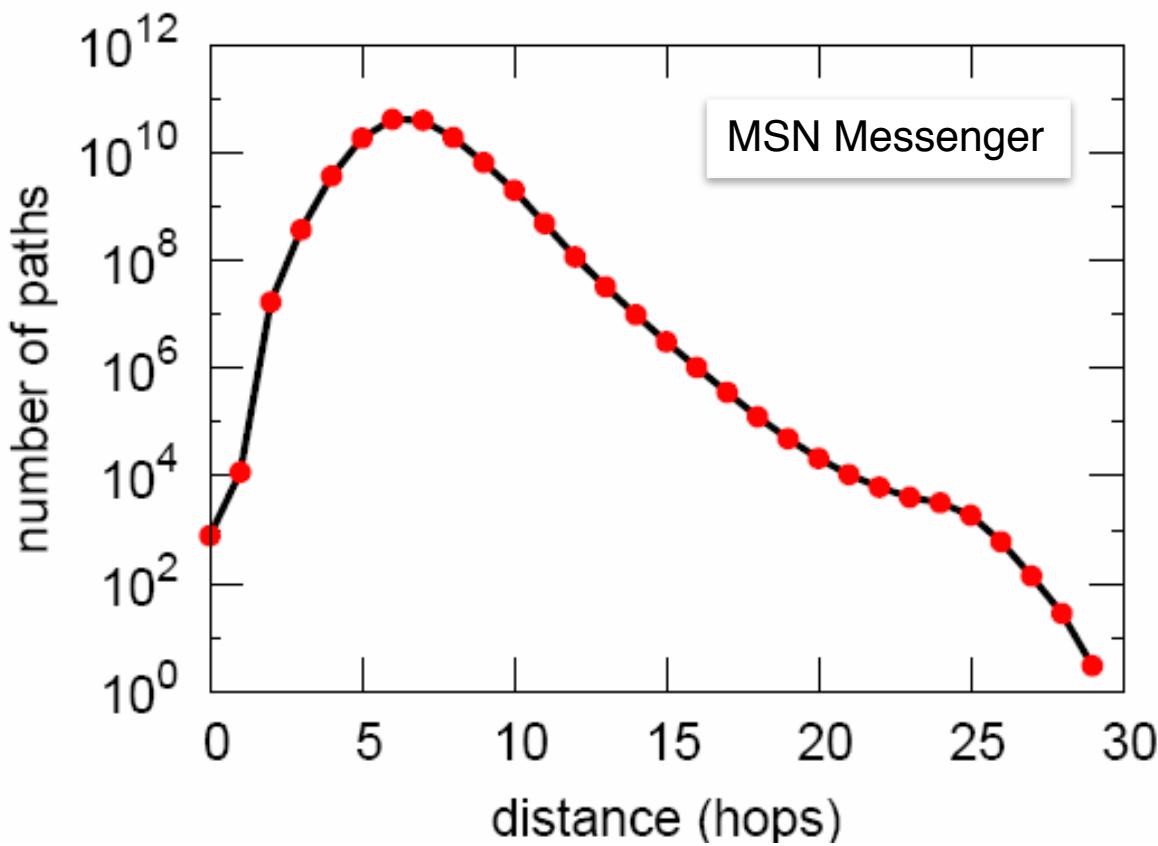
Source: UCSD

**Source:**  
[physicsbuzz.physicscentral.com](http://physicsbuzz.physicscentral.com)

**Erdős number:** # of hops needed to connect the author of a paper to Paul Erdős

# Small-world Phenomenon (4/4)

- The small-world phenomenon appears in various network settings

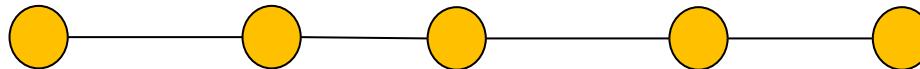


- Average path length is **6.6**
- 90% of the nodes are reachable in less than 8 steps
- Facebook** network:
  - Average distance is **4.7**
  - [Ugander et al., 2011]

# Small Diameter

---

- **Diameter** is the largest shortest path in the graph
  - Diameter is often sensitive to **chains** of nodes

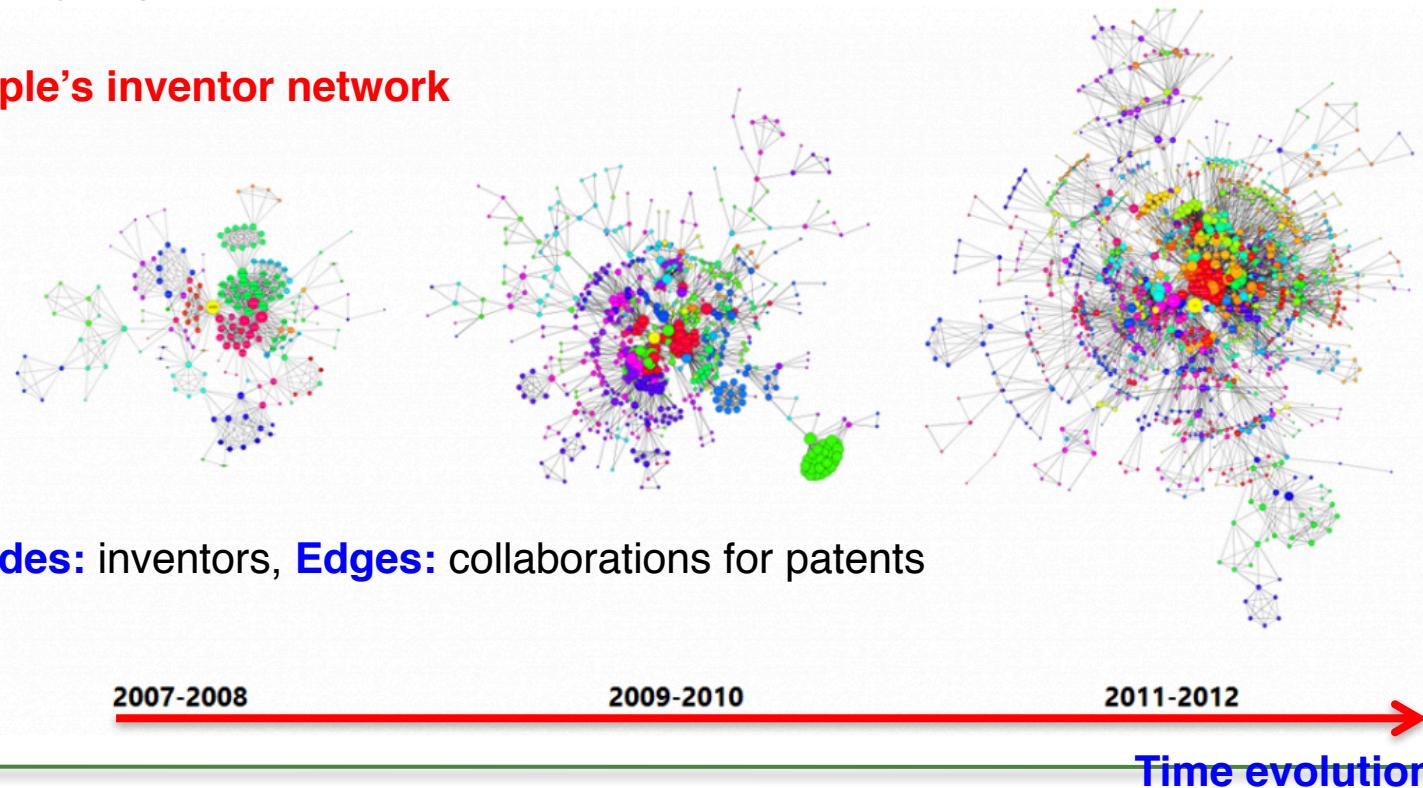


- In practice, we use the **effective diameter**
  - Upper bound of the shortest path over 90% of the pairs of nodes
- As an effect of the small-world phenomenon, real networks have **small diameter**

# Network Evolution

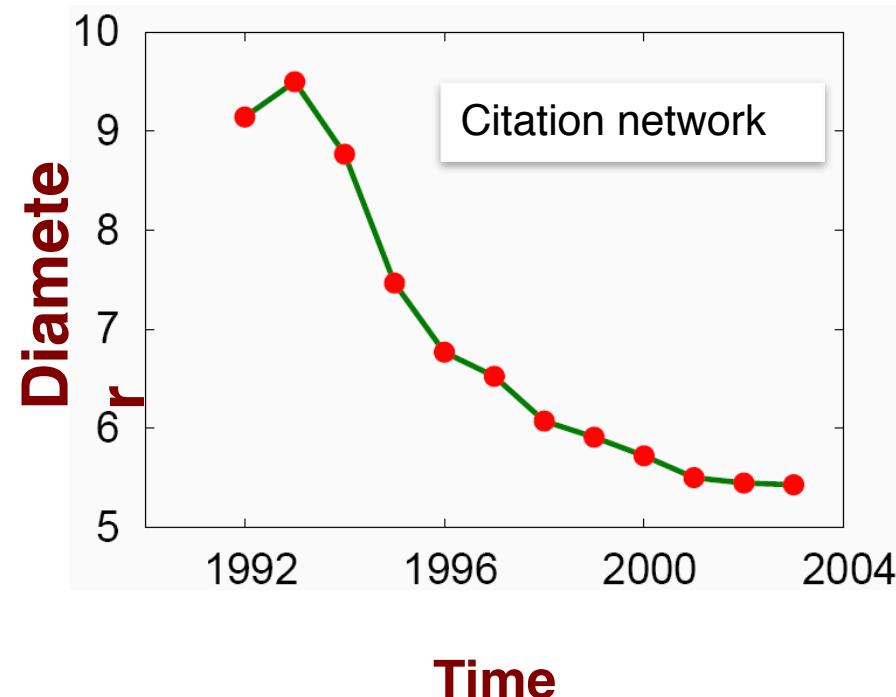
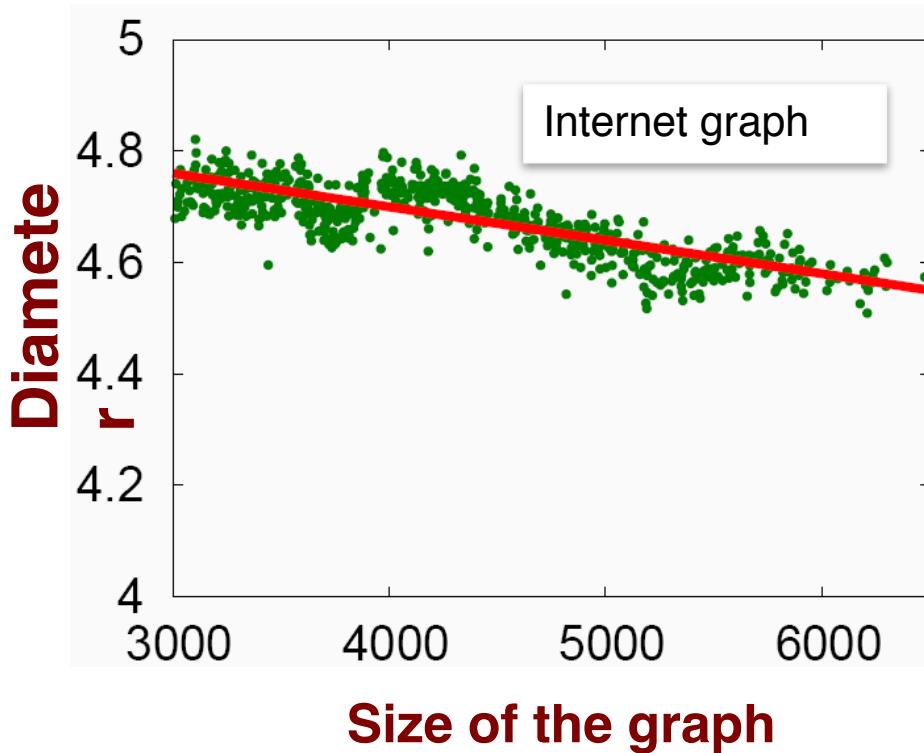
- Real-world networks are not static, but they evolve over time
  - New nodes/edges are added and/or deleted
  - We are interested in making predictions about the structure of the network

**Apple's inventor network**



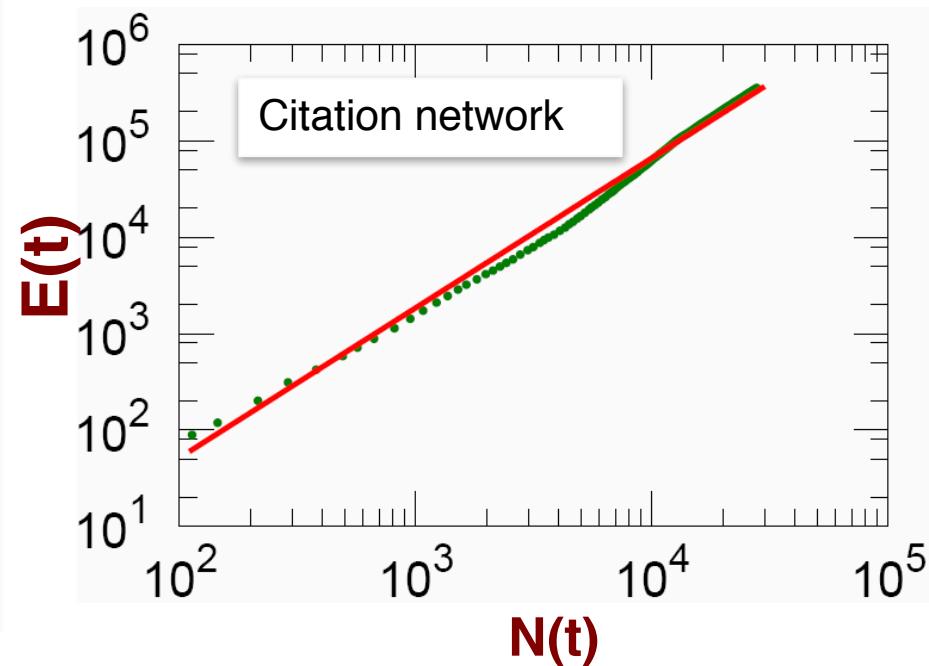
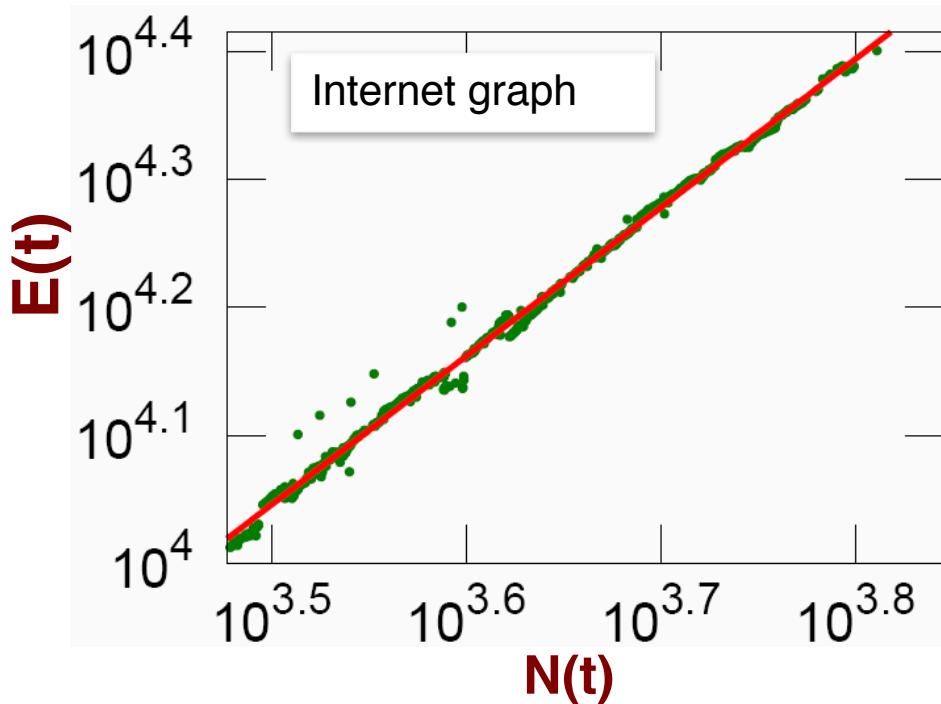
# Shrinking Diameter

- **Q:** How does the **diameter** change, while the graph evolves with the addition of nodes and edges?
  - **Intuition:** the diameter should slowly grow (e.g., **log N**, **log log N**)
- **Diameter shrinks over time**



# Densification Power Law

- **Q:** What is the relation between the number of nodes and edges over time?
- Networks become **denser** over time
  - $\alpha$  is the densification exponent ( $1 \leq \alpha \leq 2$ )  $E(t) \propto N(t)^\alpha$



# Graph Generators - Network Evolution

---

**Goal: Characterize, model and understand** the structure of real networks

- How do real-world networks look like?
  1. **Empirical: statistical properties of networks** (e.g., degree distribution, diameter) **[Previous part]**
  2. **Generative models of network structure** **[Current part]**
    - Mechanisms that reproduce the underlying generative processes

# Why do we Care?

---

- Creating models for real-world graphs is important for several reasons
  - Help us to **understand** and **reason** about the observed properties
  - Create **artificial data** for simulation purposes
  - **Predict** the evolution of networks
  - **Privacy preservation:** release the parameters of the generative model, instead of the network itself

# What is a Network Model?

---

- Informally, it is a process (randomized or deterministic) for generating a graph
- Models of **static** graphs
  - **Input:** a set of parameter  $\Pi$  and the size of the graph  $n$
  - **Output:** a graph  $G(\Pi, n)$
- Models of **evolving** graphs
  - **Input:** a set of parameter  $\Pi$  and an initial graph  $G_0$
  - **Output:** a graph  $G_t$  for each time step  $t$

# Erdős–Rényi Random Graph Model

---

- Suppose that we want to generate a network with  $n$  nodes
- The  $G_{n,p}$  model:
  - Graph with  $n$  nodes and edge probability  $p$
  - For each pair of nodes  $(u, v)$ , add the edge  $(u, v)$  **independently** with probability  $p$
  - Family of graphs, in which a graph with  $m$  edges appears with probability
- The  $G_{n,m}$  model:  $p^m(1-p)^{\binom{n}{2}-m}$ 
  - Select  $m$  edges uniformly at random

# Degree Distribution of the ER Model (1/2)

---

- **Q:** Do Erdős–Rényi graphs look **realistic**?
- The degree distribution is **Binomial**
  - Let  $C_k$  denote the number of nodes with degree  $k$
- What if  $n \rightarrow \text{infinity}$  and we fix the expected degree =  $c$ ?

If  $n \rightarrow \infty$  and  $np \rightarrow c$  (with  $c > 0$ ) then

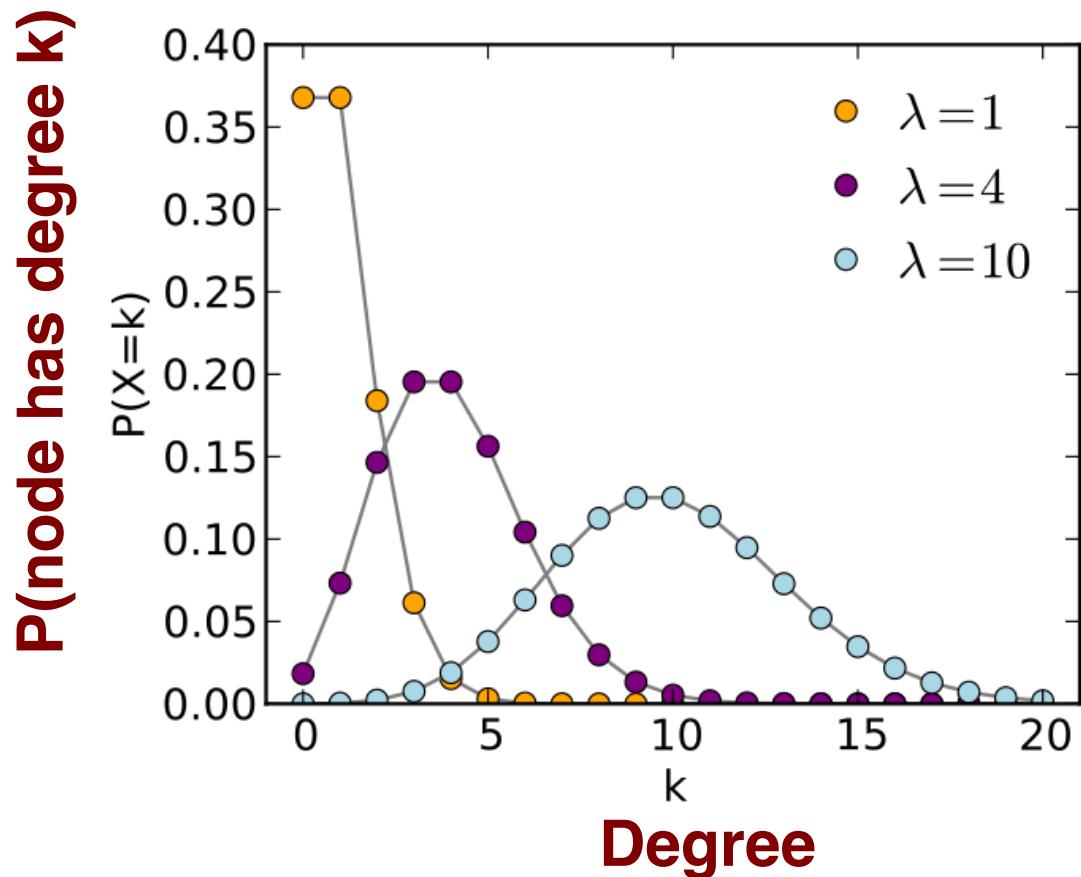
$$\frac{n!}{(n-k)!k!} p^k (1-p)^{n-k} \rightarrow e^{-c} \frac{e^c}{k!}$$

Poisson distribution

# Degree Distribution of the ER Model (2/2)

Poisson distribution

$$\frac{\lambda^k e^{-\lambda}}{k!}$$



The degree distribution of ER random graph model is  
**not realistic** for real-world graphs

# Preferential Attachment Model – General Idea

---

- Recall that real-world networks tend to have **power-law** (or in general heavy-tailed) degree distribution
  - **Barabasi-Albert** (BA) model
    - Based on the idea of preferential attachment
  - Intuition
    - Design a graph generating model that produces a small number of high degree nodes (hubs) and ...
    - ... also captures the long-tail (nodes with small degree)
- Idea:** Consider nodes that are more likely to connect to high-degree nodes
-

# Barabasi-Albert Model (1/2)

---

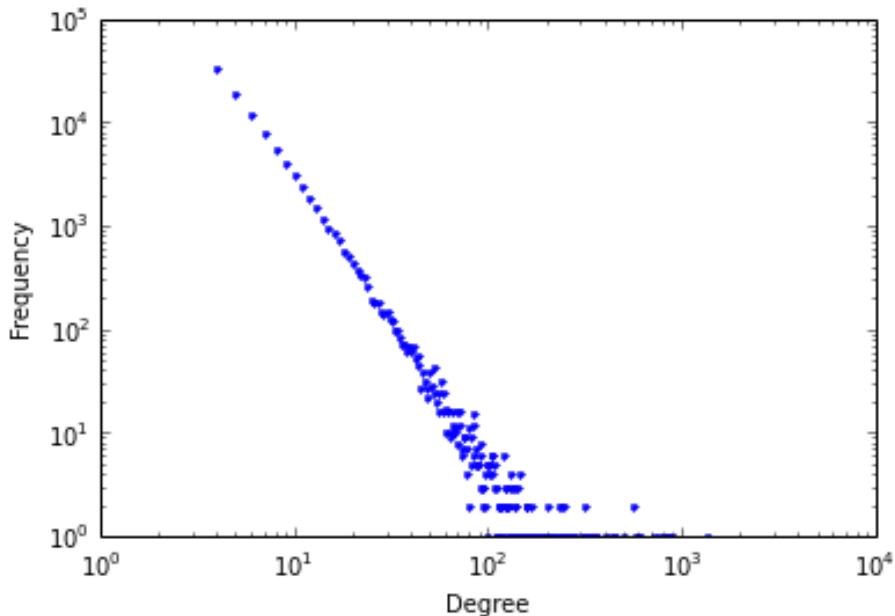
- The **Barabasi-Albert** model:
  - **Input:** some initial subgraph  $G_0$  and a parameter  $m$  that corresponds to the number of edges per new node
  - The process:
    - The nodes arrive one at the time
    - Each new node connects to  $m$  existing nodes selected with probability proportional to their degree
    - Let  $[d_1, d_2, \dots, d_t]$  be the degree sequence at time  $t$ . Then the node at  $t+1$  will be connected to node  $i$  with probability

$$p_i = \frac{d_i}{\sum_i d_i}$$

## Barabasi-Albert Model (2/2)

---

- This phenomenon is also known as the **rich get richer** effect
  - E.g., a web page that already has many incoming hyperlinks is likely to get more in the future
- The BA model produces graphs with **power-law** degree distribution  $C_k = k^{-\gamma}$ , where  $\gamma = 3$



- Barabasi-Albert graph
- $n = 100,000$  nodes
- $m = 4$

The BA model holds for several real-world networks (flickr, Delicious, LinkedIn) [Leskovec et al., 2008]

# Network Models and Temporal Evolution

---

- Most of the existing models (e.g., BA) consider that
  - The **number of edges** grows **linearly** with respect to the number of nodes
  - The **diameter increases** based on a factor of **log n** or **log log n**
- In real networks we have observed
  - **Densification power law**
  - **Shrinking diameter**

How to model the temporal evolution of real-world networks?

---

# Kronecker Model of Graphs (1/4)

- Reminder: **Kronecker product** of matrices
  - $A = [a_{ij}]$  an  $n \times m$  matrix
  - $B = [b_{ij}]$  an  $p \times q$  matrix
  - Then  $C = A \otimes B$  is defined as the  $np \times mq$  matrix

$$C = A \otimes B = \begin{pmatrix} a_{1,1}B & a_{1,2}B & \cdots & a_{1,m}B \\ a_{2,1}B & a_{2,2}B & \cdots & a_{2,m}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{n,1}B & a_{n,2}B & \cdots & a_{n,m}B \end{pmatrix}$$

- Intuition:** repeat the Kronecker product between the adjacency matrix of an initial graph to get the final graph

# Kronecker Model of Graphs (2/4)

---

- **Kronecker** model:

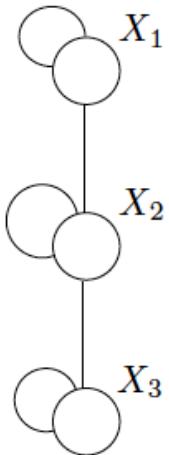
- Start by an initiator adjacency matrix  $\mathbf{A}_1$  of size  $\mathbf{p} \times \mathbf{p}$
- The Kronecker product of two graphs is defined as the Kronecker product of their adjacency matrices
- The Kronecker graph after  $\mathbf{k}$  iterations is defined as the graph with the following adjacency matrix

$$\mathbf{A}_k = \underbrace{\mathbf{A}_1 \otimes \mathbf{A}_1 \otimes \cdots \otimes \mathbf{A}_1}_{k \text{ iterations}} = \mathbf{A}_{k-1} \otimes \mathbf{A}_1$$

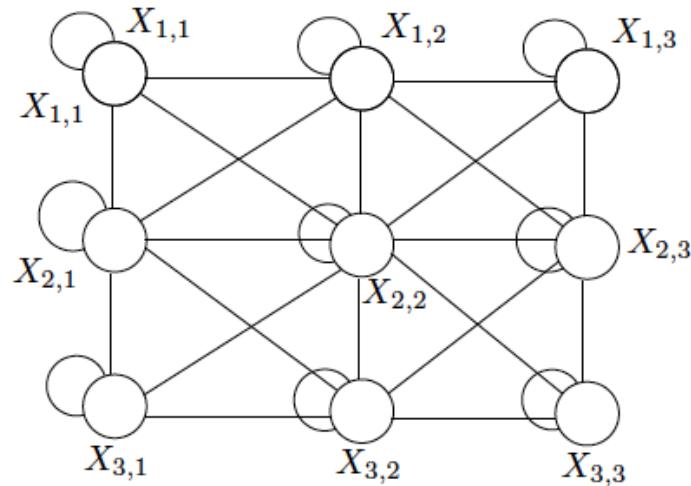
- Each Kronecker multiplication exponentially increases the size of the graph

# Kronecker Model of Graphs (3/4)

---



Graph  $\mathbf{G}_1$

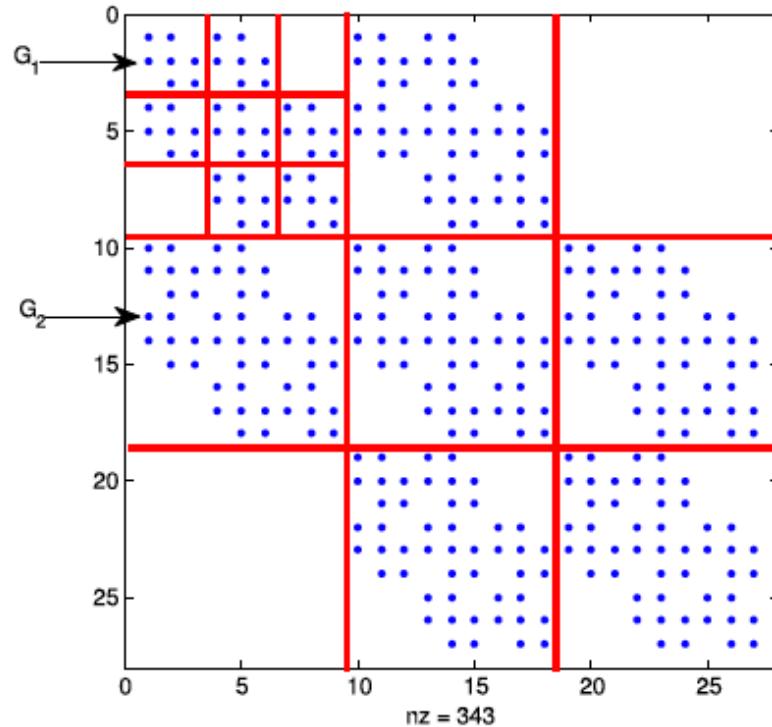


Graph  $\mathbf{G}_2 = \mathbf{G}_1 \boxtimes \mathbf{G}_1$

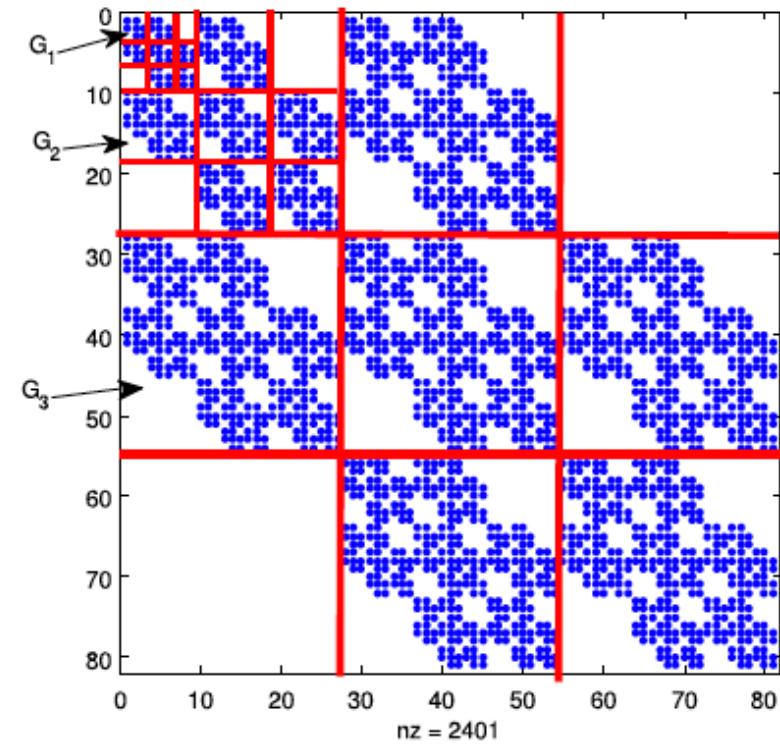
1	1	0
1	1	1
0	1	1

$\mathbf{G}_1$	$\mathbf{G}_1$	0
$\mathbf{G}_1$	$\mathbf{G}_1$	$\mathbf{G}_1$
0	$\mathbf{G}_1$	$\mathbf{G}_1$

# Kronecker Model of Graphs (4/4)



$$(a) A(G_3) = A(G_2) \otimes A(G_1)$$



$$(\beta) A(G_4) = A(G_3) \otimes A(G_1)$$

**Intuition:** Recursion and self-similarity

# Stochastic Kronecker Model

---

- In practice, the **stochastic Kronecker graph** is used
  - Start by an initiator matrix  $\theta$

a	b
c	d

- We obtain a graph with  $n = 2^k$  nodes by repeating  $k$  times the Kronecker product:  $A_{k,\theta} = \theta \otimes \dots \otimes \theta$
- Consider the value  $(i, j)$  of the matrix  $A_{k,\theta}$  as the probability of existence of the edge  $(i, j)$  (applying randomized rounding)
- Typically,  $2 \times 2$  initiator matrices produce good results

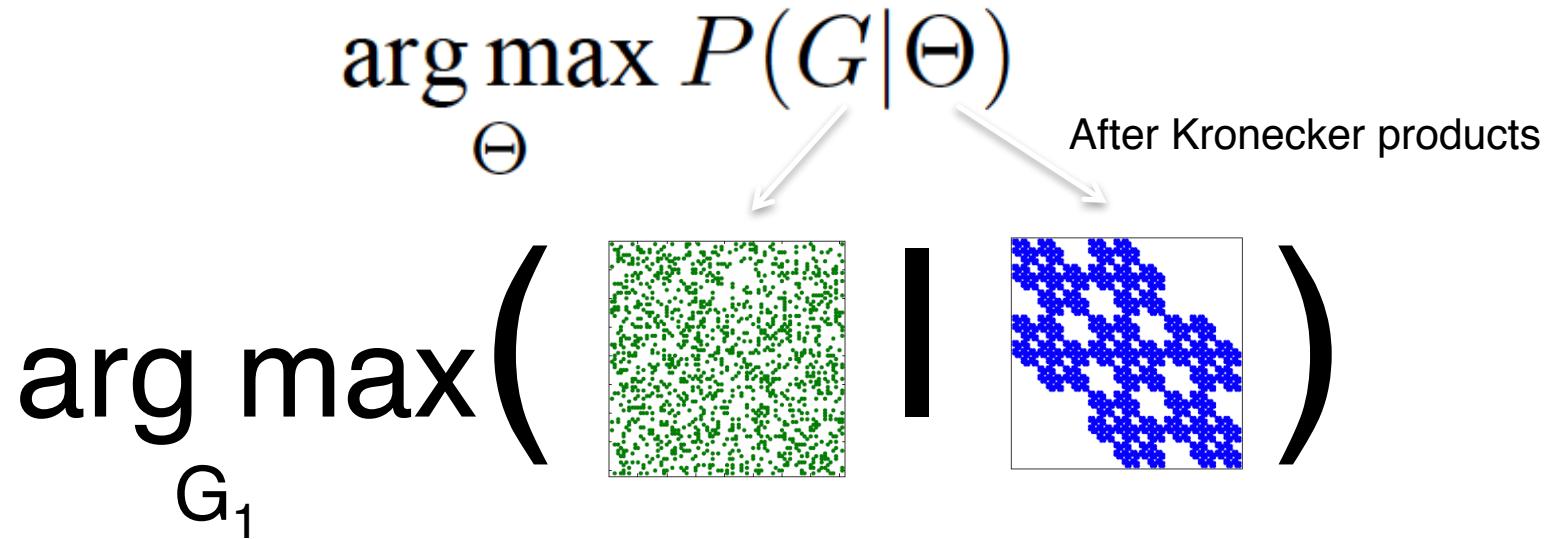
# Generate Realistic Kronecker Graphs

---

- Given a network  $\mathbf{G}$ , how can we find a “good” initiator matrix  $\Theta$ , such that  $\mathbf{A}_G \approx \Theta \boxtimes \dots \boxtimes \Theta$ ?
  - Fit the parameters  $\Theta$  of the model
  - Idea: use **maximum-likelihood estimation**

$$\arg \max_{\Theta} P(G|\Theta)$$

After Kronecker products

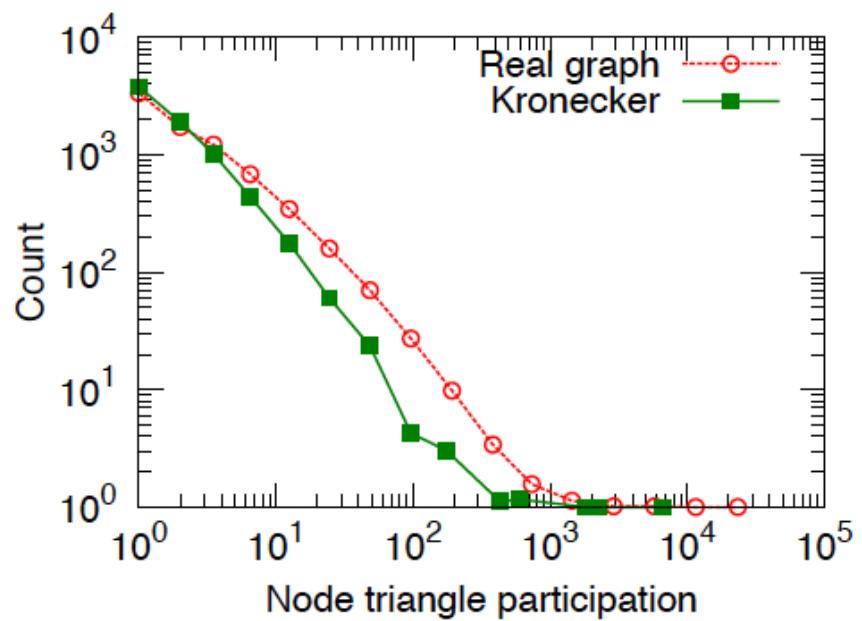
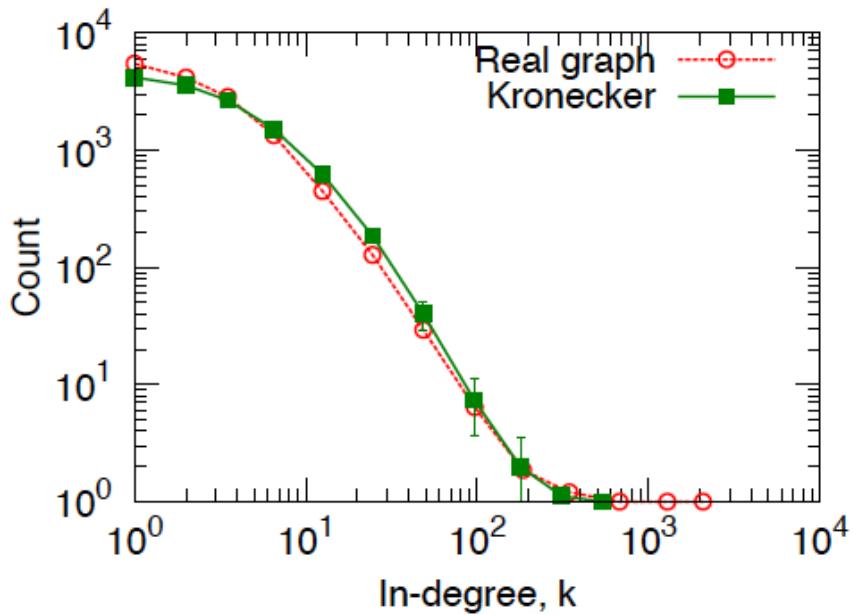
$$\arg \max_{G_1} \left( \begin{array}{c|c} \text{[green noise pattern]} & \text{[blue noise pattern]} \end{array} \right)$$


# Properties of Kronecker Model

---

- The Kronecker (stochastic) graph model is able to reproduce a plethora of properties
  - Power-law degree distribution
  - Small diameter
  - Shrinking diameter
  - Densification power-law
  - Triangle participation
  - ...

# Example: Fitting Kronecker Model to a Graph



Blog-to-Blog network

# References

---

- J. Leskovec. Modeling Large Social and Information Networks. Tutorial at ICML, 2009.
- J. McAuley. Data Mining and Predictive Analytics, UCSD, 2015.
- D. Easley and J. Kleinberg. Networks, Crowds, and Markets: Reasoning About a Highly Connected World. Cambridge University Press, 2010.
- J. Leskovec, D. Chakrabarti, J. Kleinberg, C. Faloutsos, Z. Ghahramani. Kronecker Graphs: An approach to modeling networks. JMLR, 2010.

## **2. Community evaluation measures**

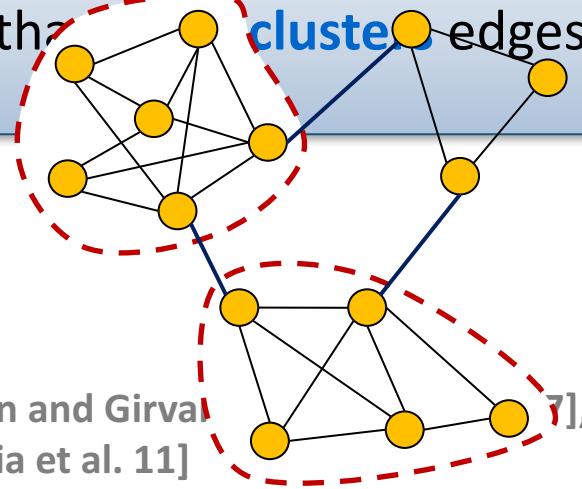
---

3.

# Basics

- The notion of **community structure** captures the tendency of nodes to be organized into modules (communities, clusters, groups)
  - Members within a community are **more similar** among each other
- Typically, the communities in graphs (networks) correspond to **densely connected** entities (nodes)

A community corresponds to a group of nodes with more **intra-cluster** edges than **cluster** edges

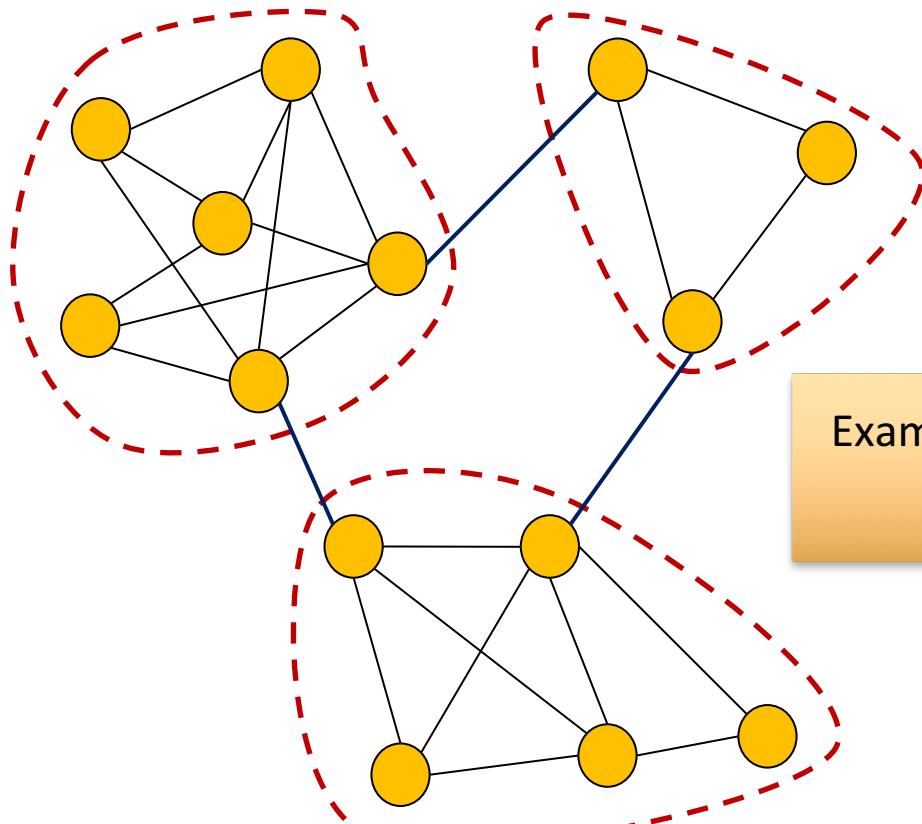


Example graph  
with three  
communities

[Newman '03], [Newman and Girvan '04], [Girvan and Newman '02], [Lancichinetti et al. '08], [Lancichinetti et al. '09], [Fortunato '10],  
[Danon et al. '05], [Coscia et al. 11]

# Schematic representation of communities

---



Example graph with three  
communities

# Community detection in graphs

---

- How can we extract the inherent communities of graphs?
- Typically, a two-step approach
  1. Specify a **quality measure** (evaluation measure, objective function) that quantifies the desired properties of communities
  2. Apply **algorithmic techniques** to assign the nodes of graph into communities, optimizing the objective function
- Several measures for quantifying the quality of communities have been proposed
- They mostly consider that communities are set of nodes with many edges between them and few connections with nodes of different communities
  - Many possible ways to formalize it

# Community evaluation measures

---

## ■ Focus on

- Intra-cluster edge density (# of edges within community),
- Inter-cluster edge density (# of edges across communities)
- Both two criteria

## ■ We group the community evaluation measures according to

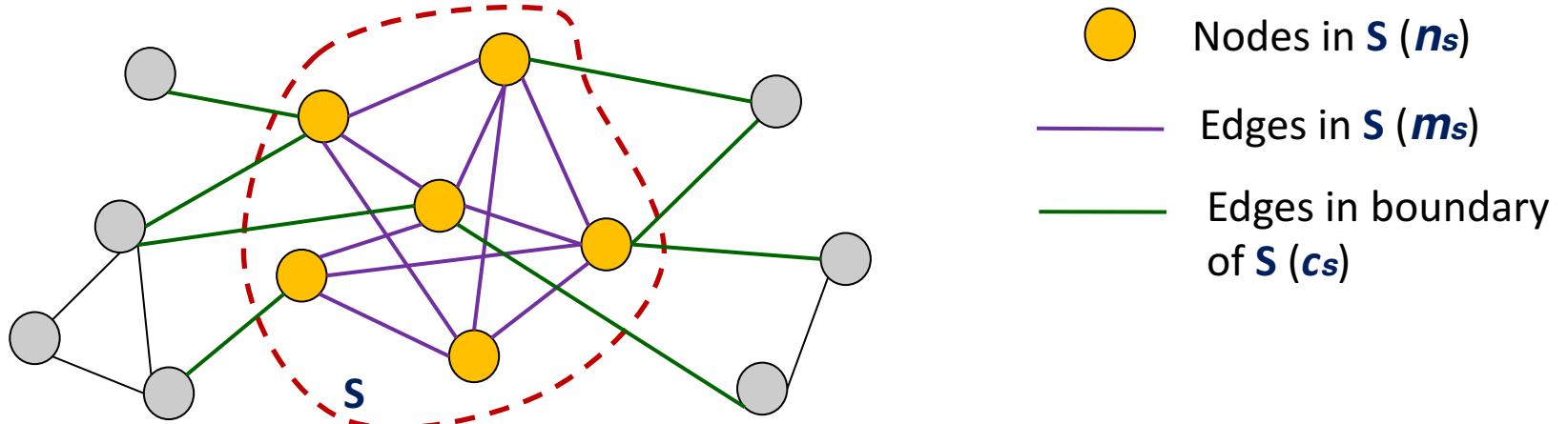
- Evaluation based on **internal** connectivity
- Evaluation based on **external** connectivity
- Evaluation based on **internal and external** connectivity
- Evaluation based on **network model**

[Leskovec et al. '10], [Yang and Leskovec '12], [Fortunato '10]

---

# Notation

- $\mathbf{G} = (\mathbf{V}, \mathbf{E})$  is an undirected graph,  $|\mathbf{V}| = n$ ,  $|\mathbf{E}| = m$
- $\mathbf{S}$  is the set of nodes in the cluster
- $n_s = |\mathbf{S}|$  is the number of nodes in  $\mathbf{S}$
- $m_s$  is the number of edges in  $\mathbf{S}$ ,  $m_s = |\{(u,v) : u \in S, v \in S\}|$
- $c_s$  is the number of edges on the boundary of  $\mathbf{S}$ ,  $c_s = |\{(u,v) : u \in S, v \notin S\}|$
- $d_u$  is the degree of node  $u$
- $f(\mathbf{S})$  represent the clustering quality of set  $\mathbf{S}$

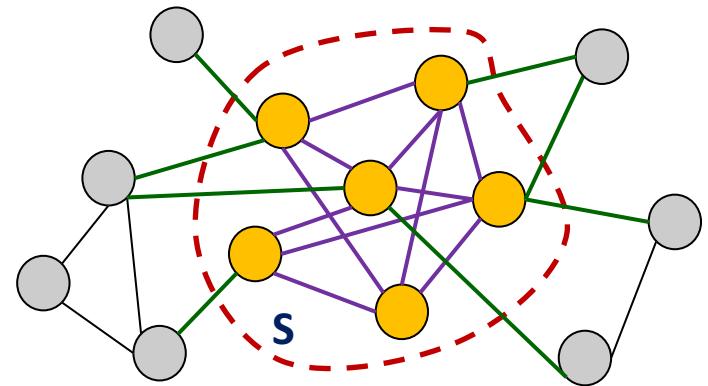


# Evaluation based on internal connectivity (1)

## ■ Internal density [Radicchi et al. '04]

$$f(S) = \frac{m_s}{n_s(n_s - 1)/2}$$

Captures the internal edge density of community  $S$



## ■ Edges inside [Radicchi et al. '04]

$$f(S) = m_s$$

Number of edges between the nodes of  $S$

# Evaluation based on external connectivity

## ■ Expansion [Radicchi et al. '04]

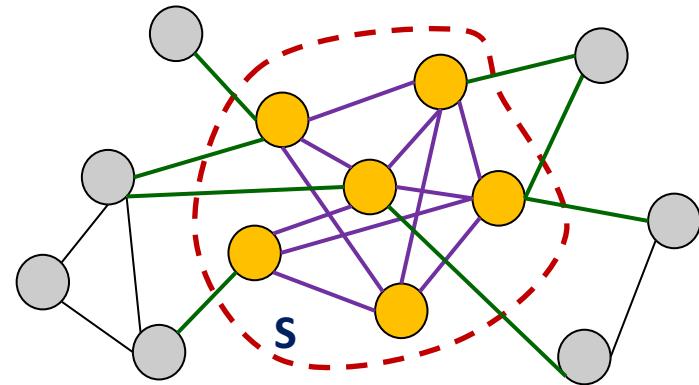
$$f(S) = \frac{c_s}{n_s}$$

Measures the number of edges per node that point outside  $S$

## ■ Cut ratio [Fortunato '10]

$$f(S) = \frac{c_s}{n_s(n - n_s)}$$

Fraction of existing edges – out of all possible edges – that leaving  $S$

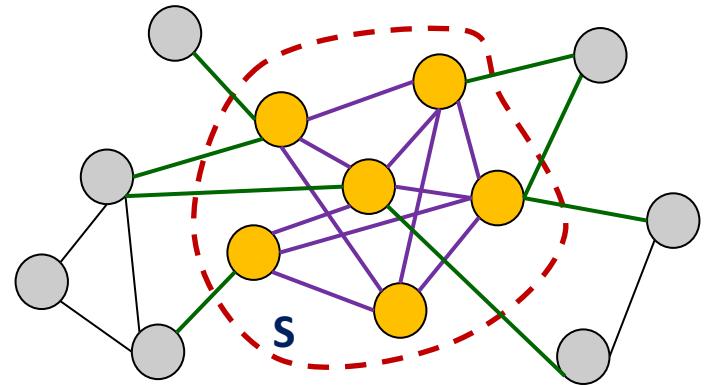


# Evaluation based on internal and external connectivity (1)

## ■ Conductance [Chung '97]

$$f(S) = \frac{c_s}{2m_s + c_s}$$

Measures the fraction of total edge volume that points outside  $S$



## ■ Normalized cut [Shi and Malic '00]

$$f(S) = \frac{c_s}{2m_s + c_s} + \frac{c_s}{2(m - m_s) + c_s}$$

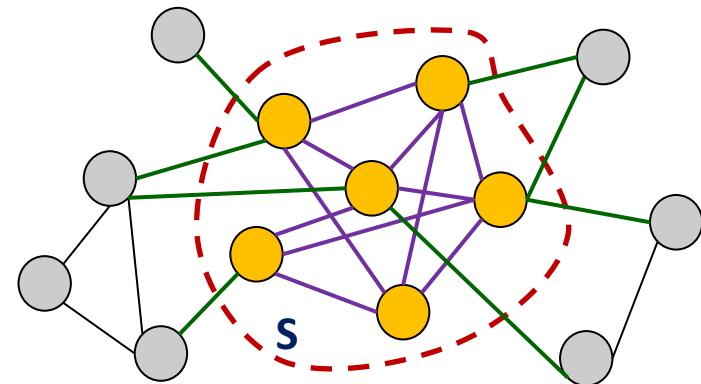
Measures the fraction of total edge volume that points outside  $S$  normalized by the size of  $S$

# Evaluation based on internal connectivity (3)

## ■ Triangle participation ratio (TPR) [Yang and Leskovec '12]

$$f(S) = \frac{|\{u : u \in S, \{(v, w) : v, w \in S, (u, v) \in E, (u, w) \in E, (v, w) \in E\} \neq \emptyset\}|}{n_s}$$

Fraction of nodes in  $S$  that belong to a triangle



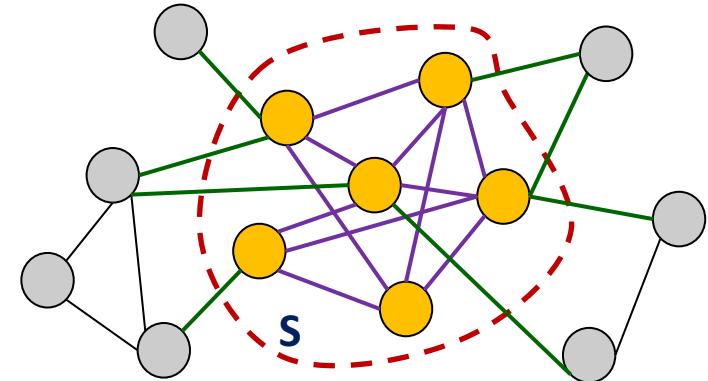
# Evaluation based on network model

## ■ Modularity [Newman and Girvan '04], [Newman '06]

$$f(S) = \frac{1}{4} (m_s - E(m_s))$$

Measures the difference between the number of edges in **S** and the expected number of edges **E(m<sub>s</sub>)** in case of a configuration model

- Typically, a random graph model with the same degree sequence



# Graph clustering

---

# Notations

---

## Given Graph $G=(V,E)$ undirected:

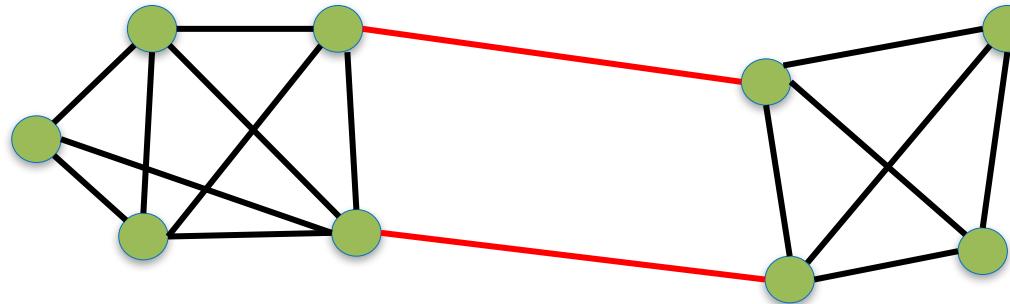
- Vertex Set  $V=\{v_1, \dots, v_n\}$ , Edge  $e_{ij}$  between  $v_i$  and  $v_j$ 
  - we assume weight  $w_{ij} > 0$  for  $e_{ij}$
- $|V|$  : number of vertices
- $d_i$  degree of  $v_i$  :  $d_i = \sum_{v_j \in V} w_{ij}$
- $\nu(V) = \sum_{v_i \in V} d_i$
- for  $A \subset V$   $\overline{A} = V - A$
- Given
  - $A, B \subset V$  &  $A \cap B = \emptyset$ ,  $w(A, B) = \sum_{v_i \in A, v_j \in B} w_{ij}$
- $D$  : Diagonal matrix where  $D(i, i) = d_i$
- $W$  : Adjacency matrix  $W(i, j) = w_{ij}$

# Graph-Cut

---

## ■ For k clusters:

- $cut(A_1, \dots, A_k) = 1/2 \sum_{i=1}^k w(A_i, \overline{A}_i)$ 
  - undirected graph: 1/2 we count twice each edge

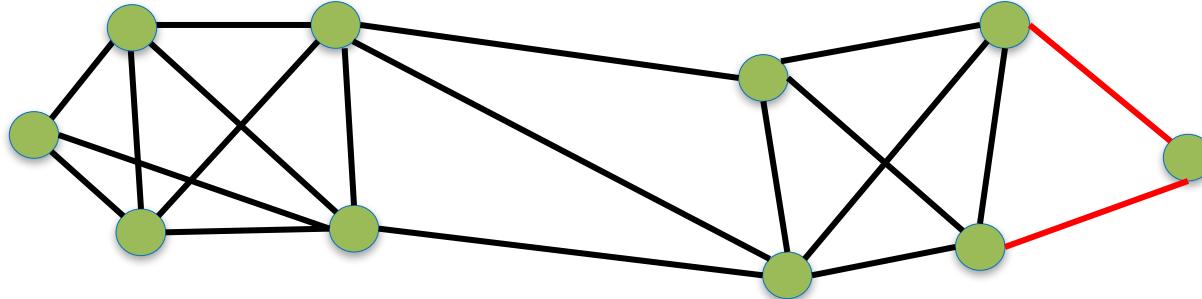


- Min-cut: Minimize the edges' weight a cluster shares with the rest of the graph

# Min-Cut

---

- Easy for  $k=2$  :  $\text{Mincut}(A_1, A_2)$ 
  - Stoer and Wagner: “A Simple Min-Cut Algorithm”
- In practice one vertex is separated from the rest
  - The algorithm is drawn to outliers



# Normalized Graph Cuts

---

- We can normalize by the size of the cluster (size of sub-graph) :

- number of Vertices (Hagen and Kahng, 1992):

$$Ratiocut(A_1, \dots, A_k) = \sum_{i=1}^k \frac{cut(A_i, \overline{A}_i)}{|A_i|}$$

- sum of weights (Shi and Malik, 2000) :

$$Ncut(A_1, \dots, A_k) = \sum_{i=1}^k \frac{cut(A_i, \overline{A}_i)}{v(A_i)}$$

- Optimizing these functions is NP-hard
- Spectral Clustering provides solution to a relaxed version of the above

# From Graph Cuts to Spectral Clustering

---

- For simplicity assume  $k=2$ :

- Define  $f: V \rightarrow \mathbb{R}$  for Graph  $G$  :

$$f_i = \begin{cases} 1 & v_i \in A \\ -1 & v_i \in \bar{A} \end{cases}$$

- Optimizing the original cut is equivalent to an optimization of:

$$\begin{aligned} & \sum_{i,j=1}^n w_{ij} (f_i - f_j)^2 \\ &= \sum_{v_i \in A, v_j \in \bar{A}} w_{ij} (1 + 1)^2 + \sum_{v_i \in \bar{A}, v_j \in A} w_{ij} (-1 - 1)^2 \\ &= 8 * \text{cut}(A, \bar{A}) \end{aligned}$$

# Graph Laplacian

---

- How is the previous useful in Spectral clustering?

$$\begin{aligned} & \sum_{i,j=1}^n w_{ij}(f_i - f_j)^2 \\ &= \sum_{i,j=1}^n w_{ij}f_i^2 - 2 \sum_{i,j=1}^n w_{ij}f_i f_j + \sum_{i,j=1}^n w_{ij}f_j^2 \\ &= \sum_{i,j=1}^n d_i f_i^2 - 2 \sum_{i,j=1}^n w_{ij}f_i f_j + \sum_{i,j=1}^n d_j f_j^2 \\ &= 2 \left( \sum_{i,j=1}^n d_{ii}f_i^2 - \sum_{i,j=1}^n w_{ij}f_i f_j \right) \\ &= 2(\mathbf{f}^T \mathbf{D} \mathbf{f} - \mathbf{f}^T \mathbf{W} \mathbf{f}) = 2\mathbf{f}^T (\mathbf{D} - \mathbf{W}) \mathbf{f} = 2\mathbf{f}^T \mathbf{L} \mathbf{f} \end{aligned}$$

- $\mathbf{f}$ : a single vector with the cluster assignments of the vertices
- $\mathbf{L} = \mathbf{D} - \mathbf{W}$  : the Laplacian of a graph

# Properties of L

---

- L is
  - Symmetric
  - Positive
  - Semi-definite
- The smallest eigenvalue of L is 0
  - The corresponding eigenvector is  $\mathbf{1}$
- L has n non-negative, real valued eigenvalues
  - $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$

## Two Way Cut from the Laplacian

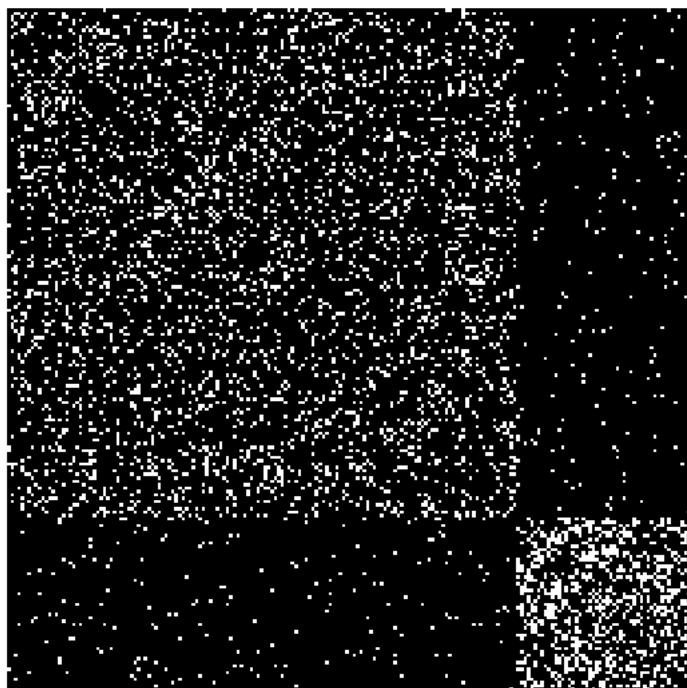
---

- We could solve  $\min_f f^T L f$  where  $f \in \{-1,1\}^n$
- NP-Hard for discrete cluster assignments
  - Relax the constraint to  $f \in R^n$  :  
$$\min_f f^T L f \text{ subject to } f^T f = n$$
- The solution to this problem is given by:
  - (**Rayleigh-Ritz Theorem**) the eigenvector corresponding to smallest eigenvalue: 0 and the corresponding eigenvector (full of 1s) offers no information
- We use the second eigenvector as an approximation
  - $f_i > 0$  the vertex belongs to one cluster ,  $f_i < 0$  to the other

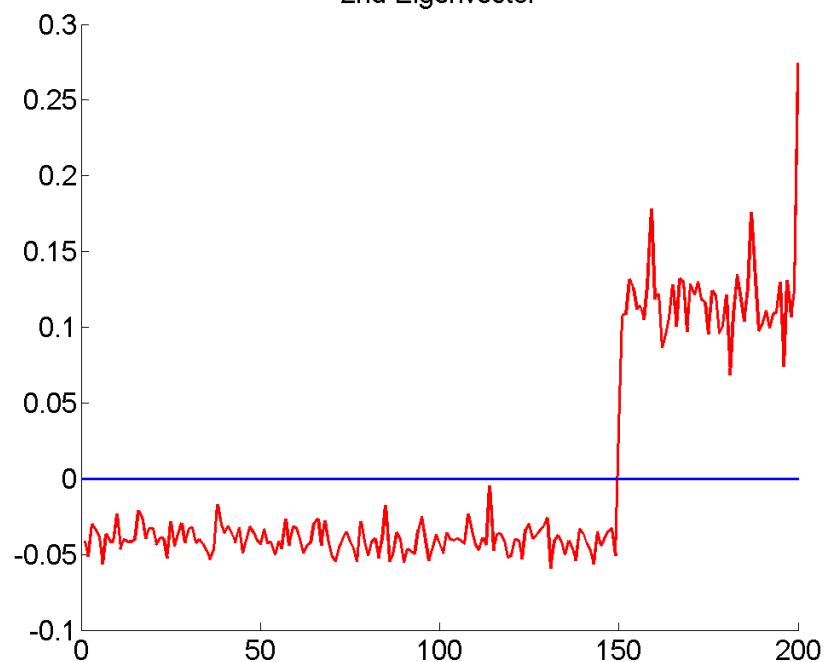
# Example

---

Adjacency Matrix

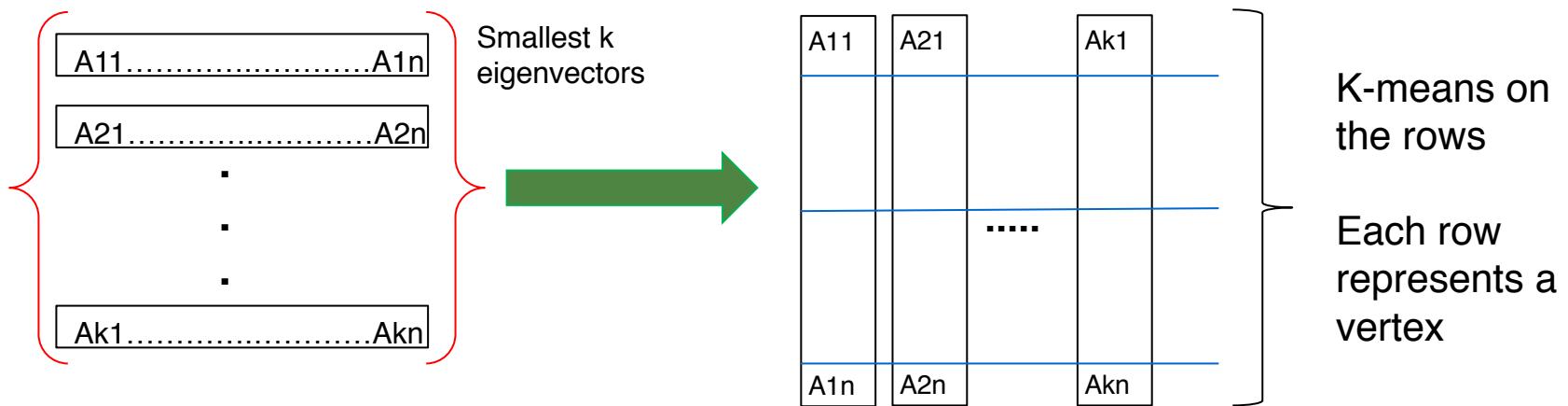


2nd Eigenvector



# Multi-Way Graph Partition

- The cluster assignment is given by the smallest k eigenvectors of  $L$
- The real values need to be converted to cluster assignments
  - We use k-means to cluster the rows
  - We can substitute  $L$  with  $L_{sym}$



# References – Graph clustering

---

- Ulrike von Luxburg, A Tutorial on Spectral Clustering, Statistics and Computing, 2007
- Davis, C., W. M. Kahan (March 1970). The rotation of eigenvectors by a perturbation. III. SIAM J. Numerical Analysis 7
- Shi, Jianbo, and Jitendra Malik. "Normalized cuts and image segmentation, "*Pattern Analysis and Machine Intelligence, IEEE Transactions on* (2000).
- Mechthild Stoer and Frank Wagner. 1997. A simple min-cut algorithm. *J. ACM*
- Ng, Jordan & Weiss, K-means algorithm on the embedded eigen-space, NIPS 2001
- Hagen, L. Kahng, , "New spectral methods for ratio cut partitioning and clustering," *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on* , 1992

# References (modularity)

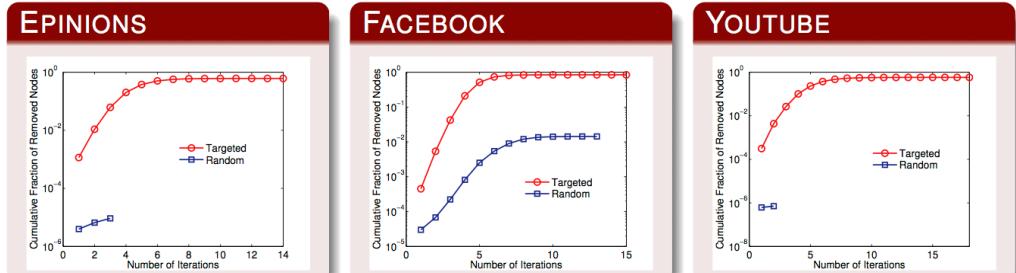
---

- M.E.J. Newman and M. Girvan. Finding and evaluating community structure in networks. *Physical Review E* 69(02), 2004.
  - M.E.J. Newman. Modularity and community structure in networks. *PNAS*, 103(23), 2006.
  - S.E. Schaeffer. Graph clustering. *Computer Science Review* 1(1), 2007.
  - S. Fortunato. Community detection in graphs. *Physics Reports* 486 (3-5), 2010.
  - M. Coscia, F. Giannotti, and D. Pedreschi. A classification for community discovery methods in complex networks. *Statistical Analysis and Data Mining* 4 (5), 2011.
  - A. Arenas, J. Duch, A. Fernandez, and S. Gomez. Size reduction of complex networks preserving modularity. *New J. Phys.*, 9(176), 2007.
  - M. Girvan and M.E.J. Newman. Community structure in social and biological networks. *PNAS* 99(12), 2002.
  - U. Brandes, D. Delling, M. Gaertler, R. Gorke, M. Hoefer, Z. Nikoloski, and D. Wagner. On Modularity Clustering. *IEEE TKDE* 20(2), 2008.
  - M.E.J. Newman. Fast algorithm for detecting community structure in networks. *Phys. Rev. E* 69, 2004.
  - A. Clauset, M.E.J. Newman, and C. Moore. Finding community structure in very large networks. *Phys. Rev. E* 70, 2004.
-

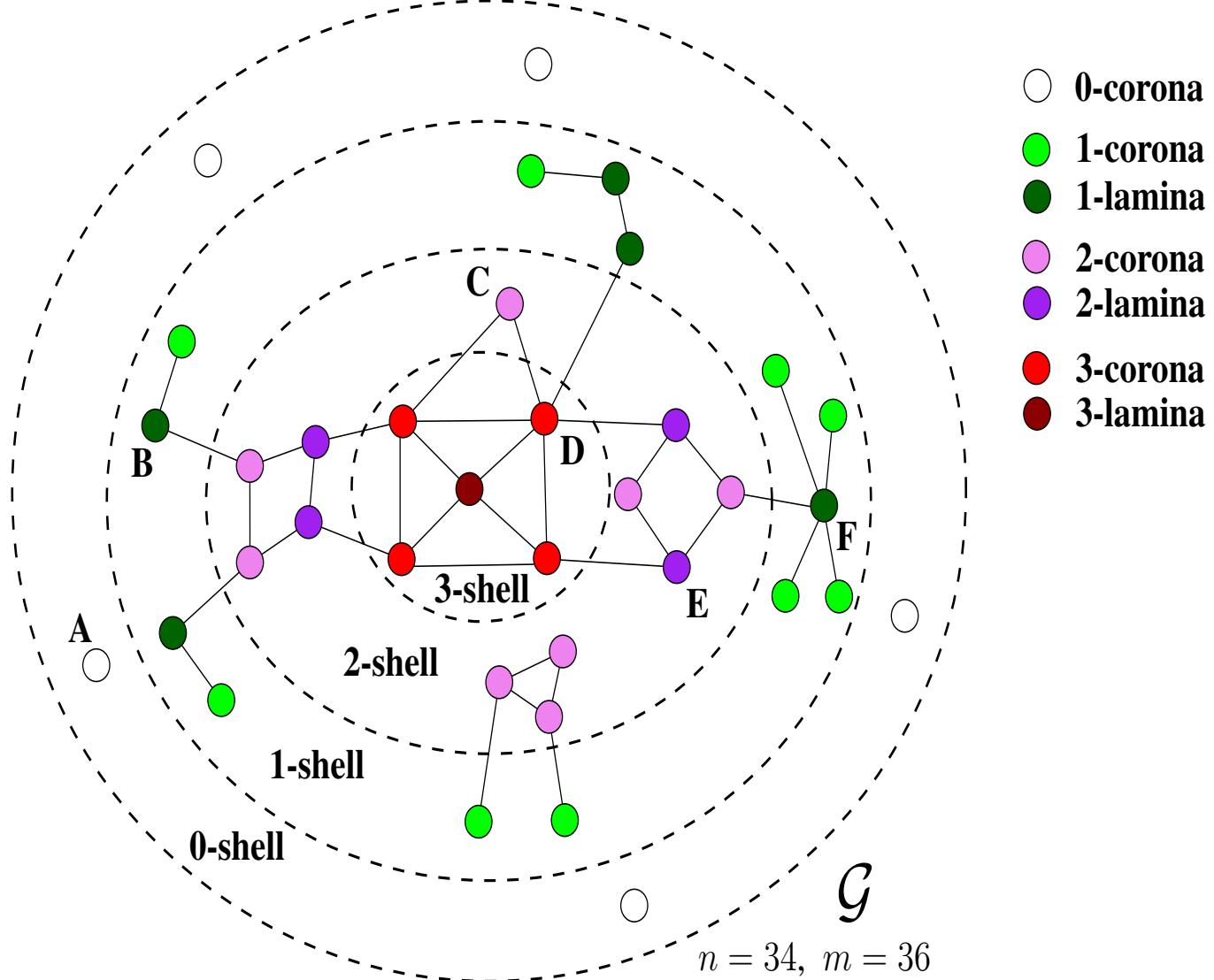
# Graph Mining with degeneracy

## ■ Community detection & evaluation

- Identifying groups of users highly collaborating among them



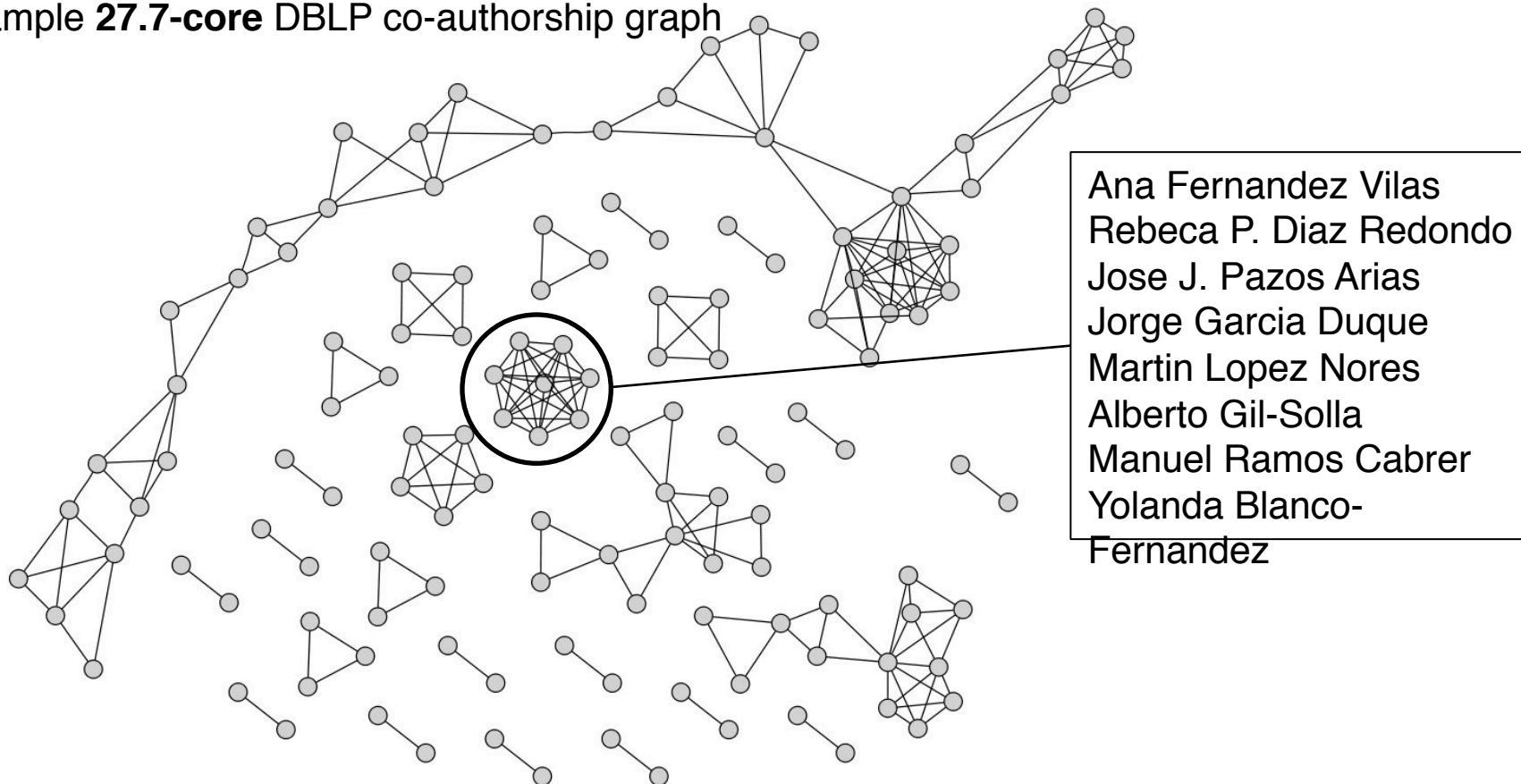
# Graph Mining – k-core concept



# Community detection and evaluation

---

Example **27.7-core** DBLP co-authorship graph



# Community detection and evaluation

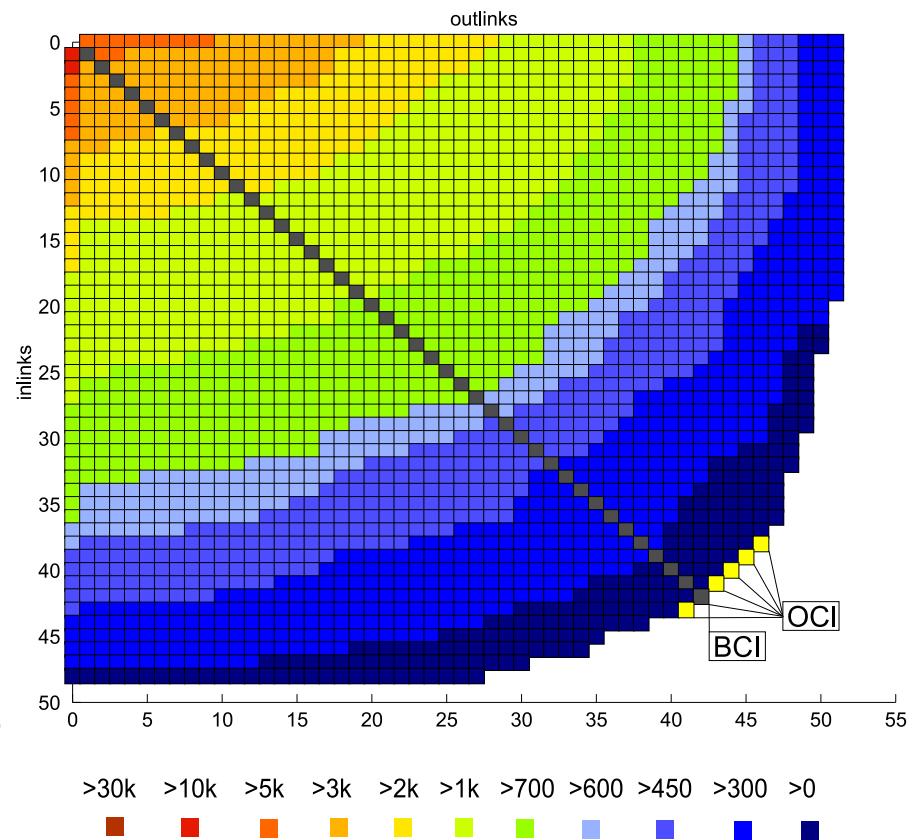
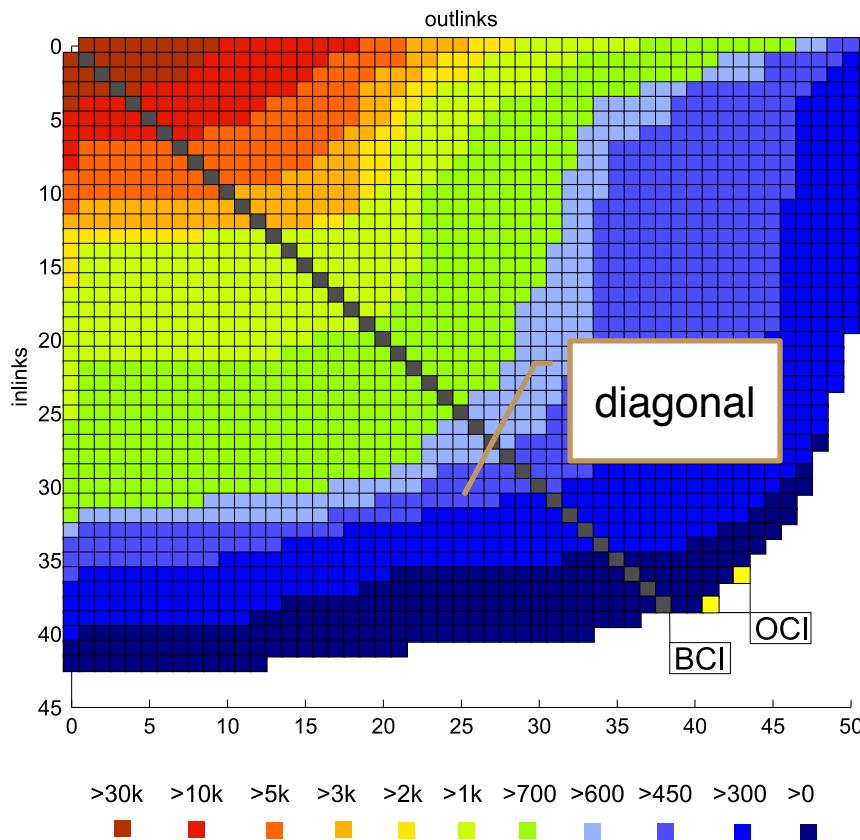
---

## Degeneracy in directed graphs



- Directed graphs:
  - WIKI - graph
  - DBLP & ARXIV – Citation graph
- Is there a degeneracy notion for directed graphs?
- We extend the k-core concept in directed graphs by applying a limit on in/out edges respectively
- Trade off between in/out edges can give us a more specific view of the cohesiveness and the “social” behavior

# D-core matrix Wikipedia & DBLP



**Wikipedia**  
The extreme D-core(38,41) contains 237  
pages

**DBLP**  
One of the extreme D-cores(38,46) contains  
188 authors

# The Extreme DBLP citation graph D-core

---

José A. Blakeley  
Hector Garcia-Molina  
Abraham Silberschatz  
Umeshwar Dayal  
Eric N. Hanson  
Jennifer Widom  
Klaus R. Dittrich  
Nathan Goodman  
Won Kim  
Alfons Kemper  
Guido Moerkotte  
Clement T. Yu  
M. Tamer Å Zsu  
Amit P. Sheth  
Ming-Chien Shan  
Richard T. Snodgrass  
David Maier  
Michael J. Carey  
David J. DeWitt  
Joel E. Richardson  
Eugene J. Shekita  
Waqar Hasan  
Marie-Anne Neimat  
Darrell Woelk  
Roger King  
Stanley B. Zdonik  
Lawrence A. Rowe  
Michael Stonebraker  
Serge Abiteboul  
Richard Hull  
Victor Vianu  
Jeffrey D. Ullman  
Michael Kifer  
Philip A. Bernstein  
Vassos Hadzilacos  
Elisa Bertino  
Stefano Ceri  
Georges Gardarin

Patrick Valduriez  
Ramez Elmasri  
Richard R. Muntz  
David B. Lomet  
Betty Salzberg  
Shamkant B. Navathe  
Arie Segev  
Gio Wiederhold  
Witold Litwin  
Theo Härdler  
François Bancilhon  
Raghu Ramakrishnan  
Michael J. Franklin  
Yannis E. Ioannidis  
Henry F. Korth  
S. Sudarshan  
Patrick E. O'Neil  
Dennis Shasha  
Shamim A. Naqvi  
Shalom Tsur  
Christos H. Papadimitriou  
Georg Lausen  
Gerhard Weikum  
Kotagiri Ramamohanarao  
Maurizio Lenzerini  
Domenico Saccà  
Giuseppe Pelagatti  
Paris C. Kanellakis  
Jeffrey Scott Vitter  
Letizia Tanca  
Sophie Cluet  
Timos K. Sellis  
Alberto O. Mendelzon  
Dennis McLeod  
Calton Pu  
C. Mohan  
Malcolm P. Atkinson  
Doron Rotem

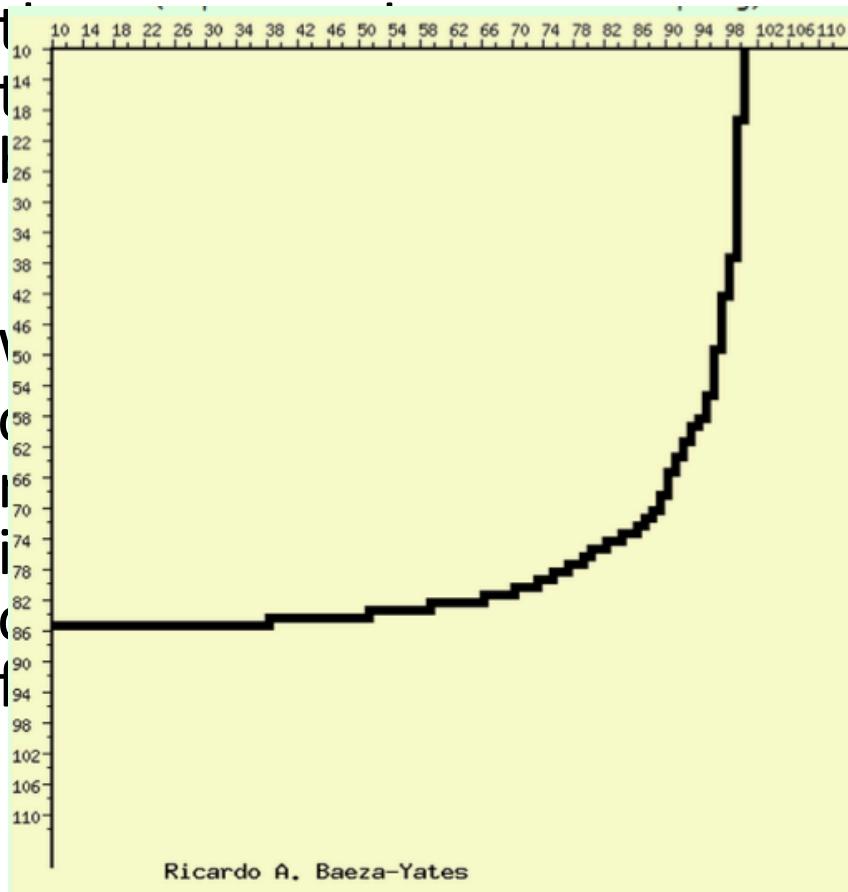
Michel E. Adiba  
Kyuseok Shim  
Goetz Graefe  
Jiawei Han  
Edward Sciore  
Rakesh Agrawal  
Carlo Zaniolo  
V. S. Subrahmanian  
Claude Delobel  
Christophe Lecluse  
Michel Scholl  
Peter C. Lockemann  
Peter M. Schwarz  
Laura M. Haas  
Arnon Rosenthal  
Erich J. Neuhold  
Hans-Jörg Schek  
Dirk Van Gucht  
Hamid Pirahesh  
Marc H. Scholl  
Peter M. G. Apers  
Allen Van Gelder  
Tomasz Imielinski  
Yehoshua Sagiv  
Narain H. Gehani  
H. V. Jagadish  
Eric Simon  
Peter Buneman  
Dan Suciu  
Christos Faloutsos  
Donald D. Chamberlin  
Setrag Khoshafian  
Toby J. Teorey  
Randy H. Katz  
Miron Livny  
Philip S. Yu  
Stanley Y. W. Su  
Henk M. Blanken

Peter Pistor  
Matthias Jarke  
Moshe Y. Vardi  
Daniel Barbară;  
Uwe Deppisch  
H.-Bernhard Paul  
Don S. Batory  
Marco A. Casanova  
Joachim W. Schmidt  
Guy M. Lohman  
Bruce G. Lindsay  
Paul F. Wilms  
Z. Meral Özsoyoglu  
Gultekin Özsoyoglu  
Kyu-Young Whang  
Shahram Ghandeharizadeh  
Tova Milo  
Alon Y. Levy  
Georg Gottlob  
Johann Christoph Freytag  
Klaus Küspert  
Louiqa Raschid  
John Mylopoulos  
Alexander Borgida  
Anand Rajaraman  
Joseph M. Hellerstein  
Masaru Kitsuregawa  
Sumit Ganguly  
Rudolf Bayer  
Raymond T. Ng  
Daniela Florescu  
Per-Åke Larson  
Hongjun Lu  
Ravi Krishnamurthy  
Arthur M. Keller  
Catriel Beeri  
Inderpal Singh Mumick  
Oded Shmueli

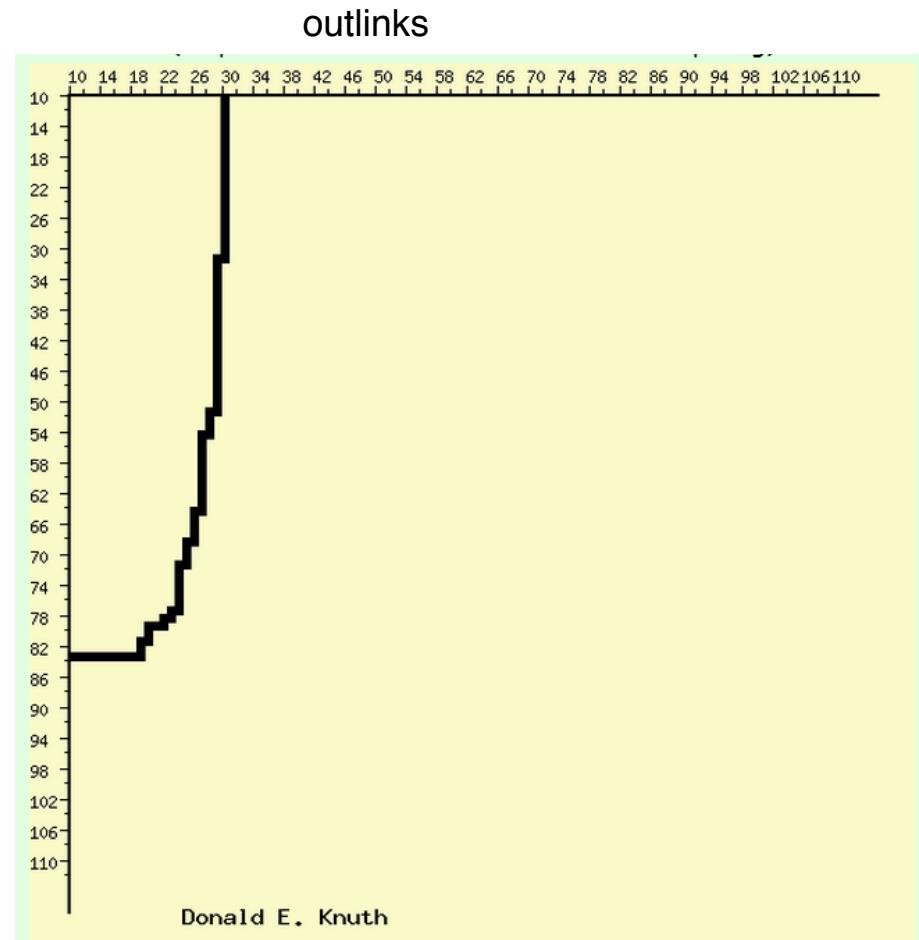
George P. Copeland  
Peter Dadam  
Susan B. Davidson  
Donald Kossmann  
Christophe de Maindreville  
Yannis Papakonstantinou  
Kenneth C. Sevcik  
Gabriel M. Kuper  
Peter J. Haas  
Jeffrey F. Naughton  
Nick Roussopoulos  
Bernhard Seeger  
Georg Walch  
R. Erbe  
Balakrishna R. Iyer  
Ashish Gupta  
Praveen Seshadri  
Walter Chang  
Surajit Chaudhuri  
Divesh Srivastava  
Kenneth A. Ross  
Arun N. Swami  
Donovan A. Schneider  
S. Seshadri  
Edward L. Wimmers  
Kenneth Salem  
Scott L. Vandenberg  
Dallan Quass  
Michael V. Mannino  
John McPherson  
Shaul Dar  
Sheldon J. Finkelstein  
Leonard D. Shapiro  
Anant Jhingran  
George Lapis

# D-Core frontier for individuals

- The frontier of an individual: defined by

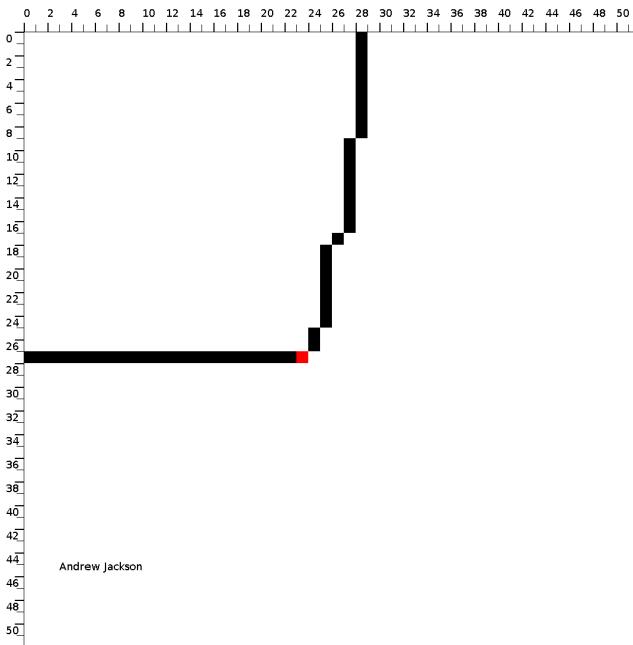


Ricardo A. Baeza-Yates

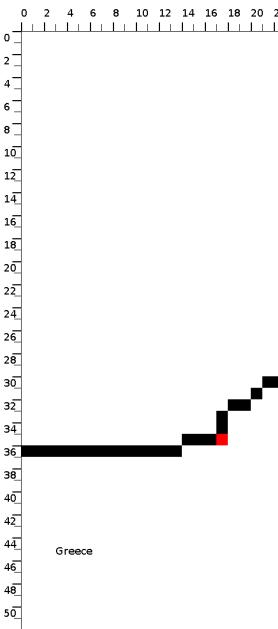


Donald E. Knuth

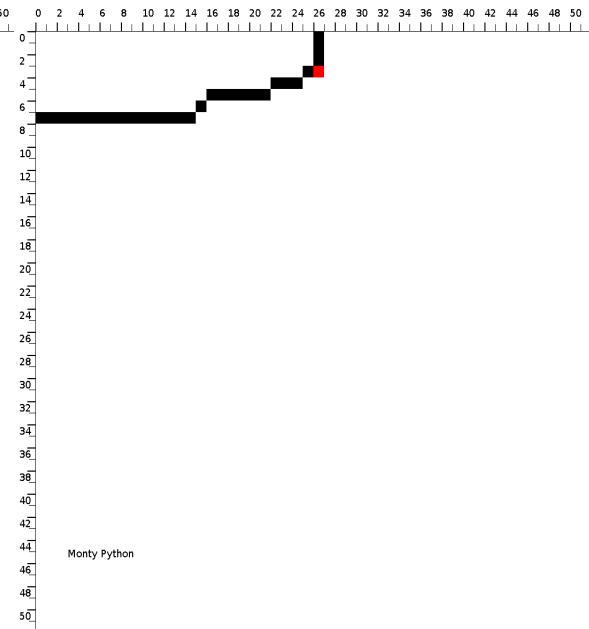
# Thematic D-core frontiers - Wikipedia



“Andrew Jackson”

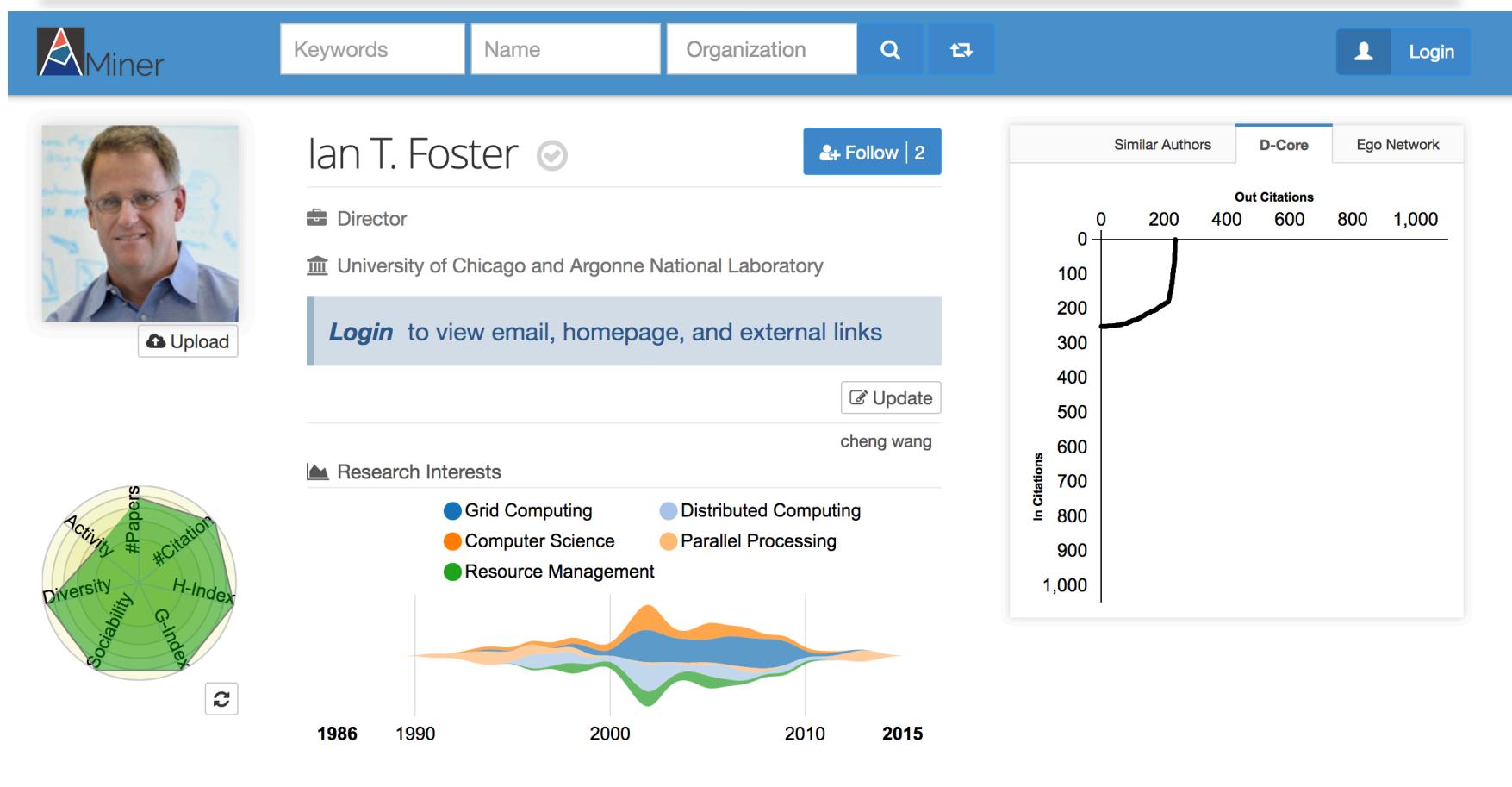


“Greece”



“Monty Python”

# D-core adopted by aminer.org



<https://cn.aminer.org/profile/ian-t-foster/53f48850dabfaee4dc8b2045>

# Ranking graph nodes – Pagerank

---

# Data are connected!

---

- A.boss = B, B.friends = {C,D,F}
- Social networks, and the web is not just a collection of documents – they form a graph structure
- A link from page *A* to page *B* may indicate:
  - *A* is related to *B*, or
  - *A* is recommending, citing or endorsing *B*
- Links are either
  - referential – *click here and get back home*, or
  - Informational – *click here to get more detail*

# Citation Analysis

---

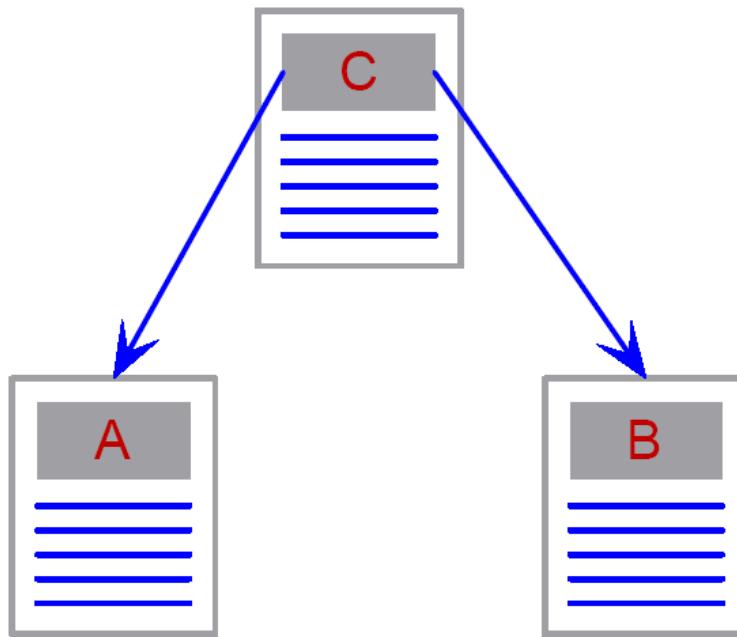
## ■ The **impact factor** of a journal = $A/B$

- $A$  is the number of current year citations to articles appearing in the journal during previous two years.
- $B$  is the number of articles published in the journal during previous two years.

Journal Title	Impact Factor (2002)
J. Mach. Learn. Res.	3.818
IEEE T. Pattern Anal.	2.923
Mach. Learn.	1.944
IEEE Intell. Syst.	1.905
Artif. Intell.	1.703

# Co-Citation

---



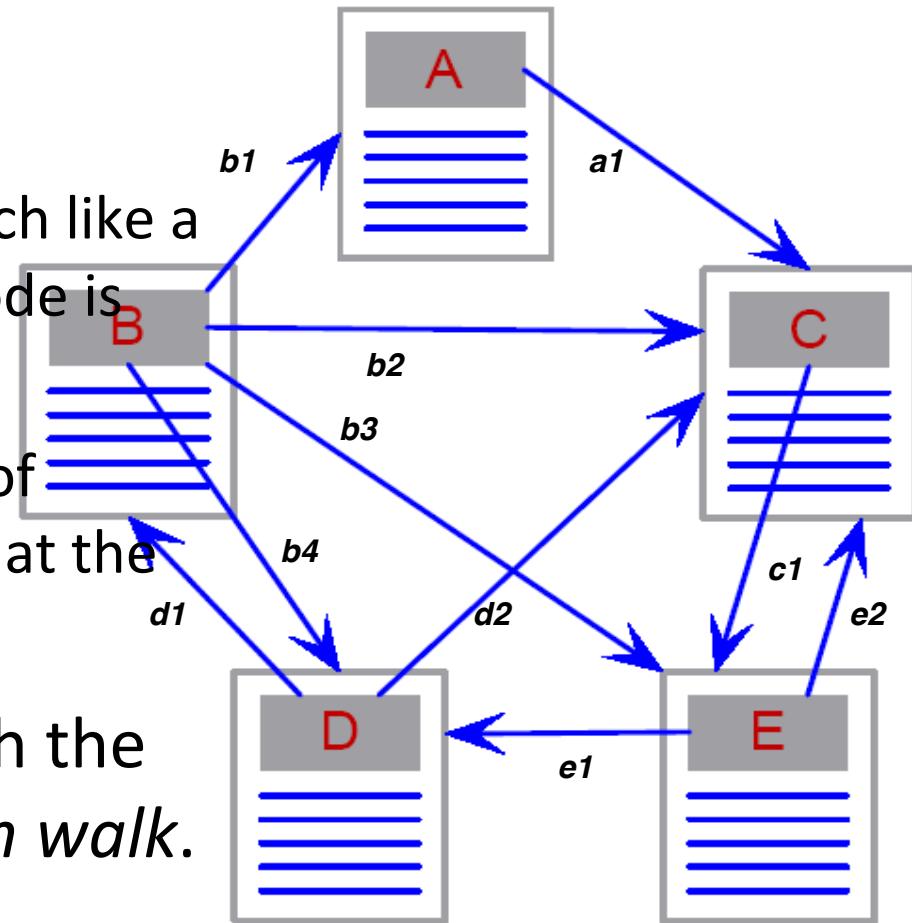
- *A* and *B* are co-cited by *C*, implying that
  - they are related or associated.
- The strength of co-citation between *A* and *B* is the number of times they are co-cited.

# What is a Markov Chain?

A Markov chain has two components:

- A network structure much like a web site, where each node is called a state.
- A transition probability of traversing a link given that the chain is in a state.

A sequence of steps through the chain is called a *random walk*.



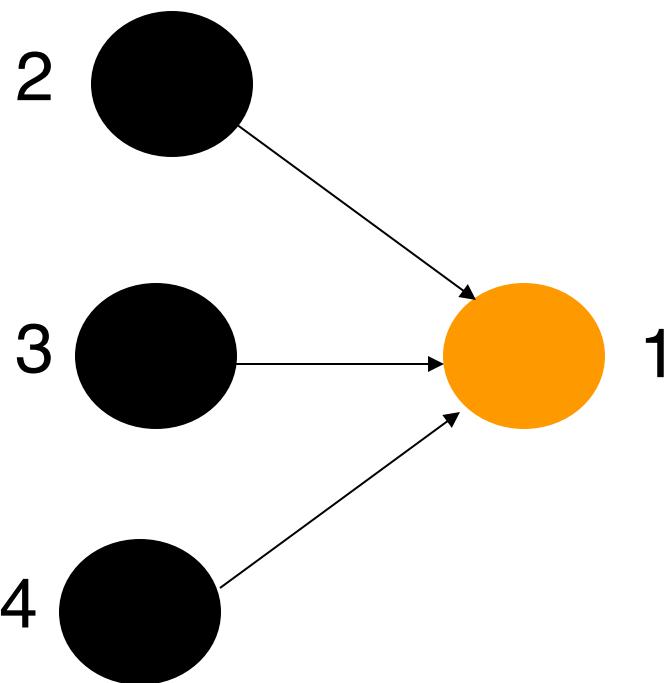
# HITS - Kleinberg's Algorithm

---

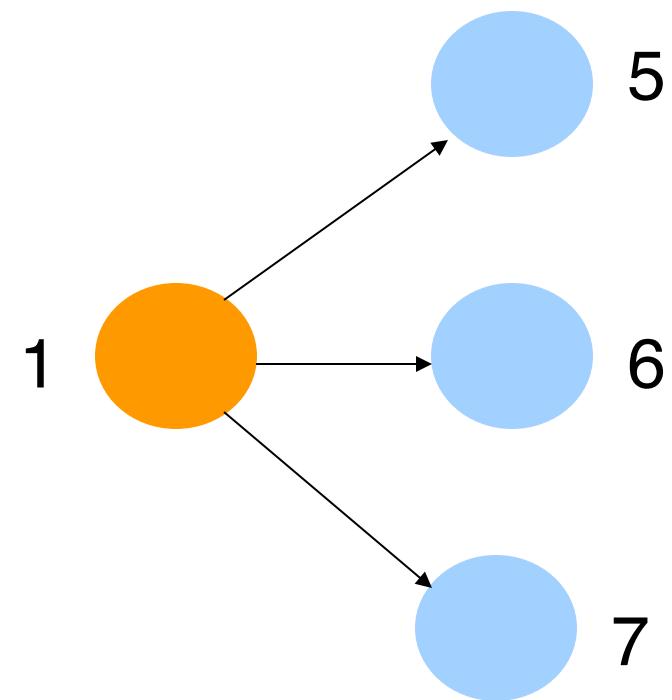
- HITS – Hypertext Induced Topic Selection
- For each vertex  $v \in V$  in a subgraph of interest:
  - $a(v)$  - the authority of  $v$
  - $h(v)$  - the hubness of  $v$
- A site is very authoritative if it receives many citations. Citation from important sites weight more than citations from less-important sites
- Hubness shows the importance of a site. A good hub is a site that links to many authoritative sites

# Authority and Hubness

---



$$a(1) = h(2) + h(3) + h(4)$$



$$h(1) = a(5) + a(6) + a(7)$$

# Authority and Hubness Convergence

---

- Recursive dependency:

$$a(v) \leftarrow \sum_{w \in \text{pa}[v]} h(w)$$

$$h(v) \leftarrow \sum_{w \in \text{ch}[v]} a(w)$$

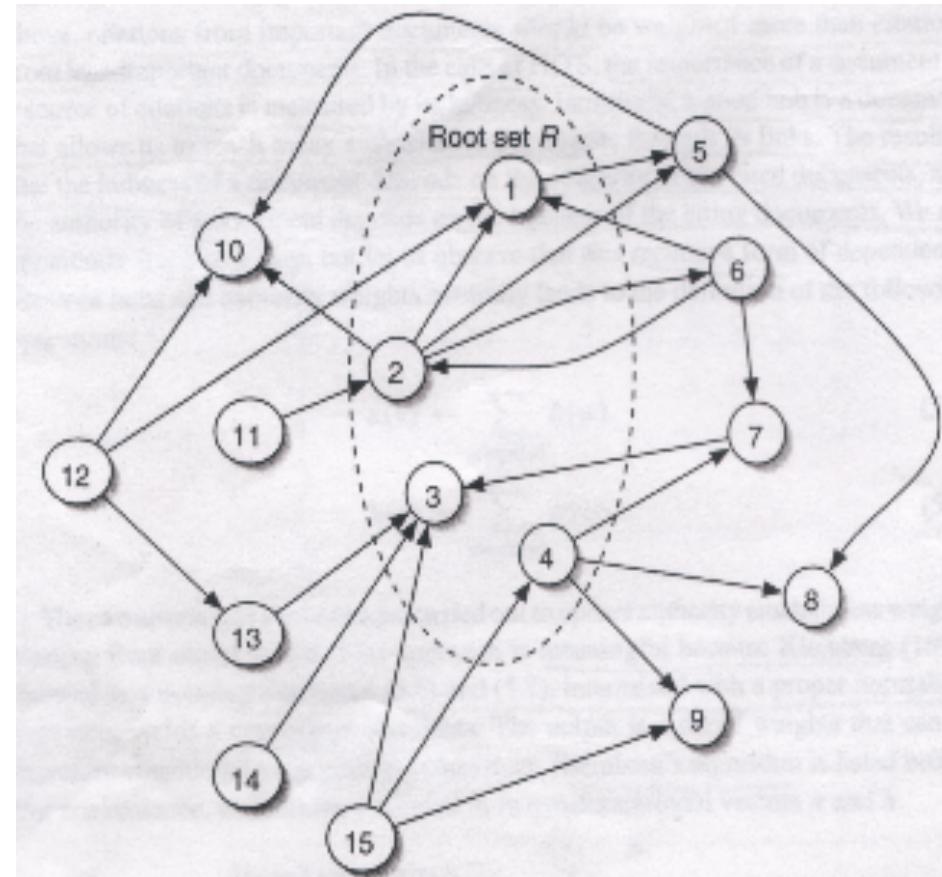
- Using Linear Algebra, we can prove:

$a(v)$  and  $h(v)$  converge

# HITS Example

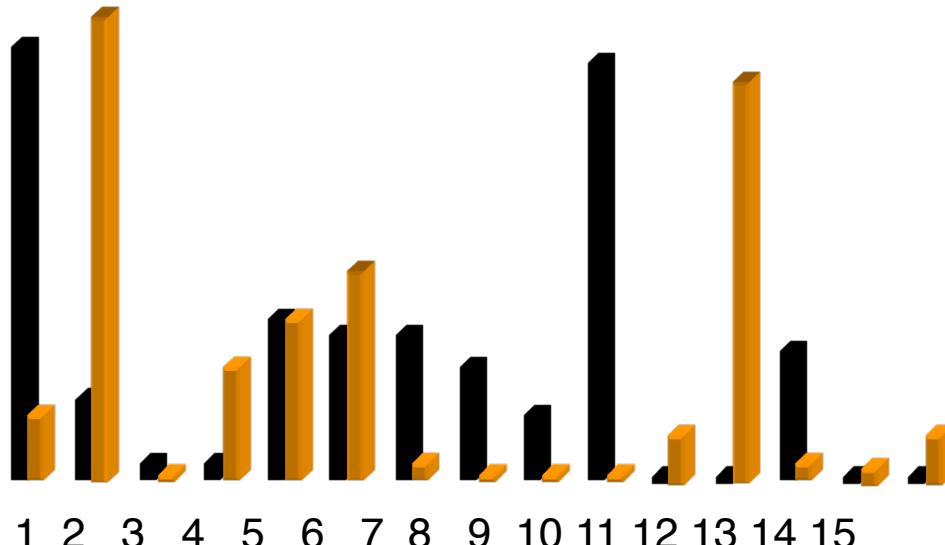
Find a base subgraph:

- Start with a root set  $R \{1, 2, 3, 4\}$
  - $\{1, 2, 3, 4\}$  - nodes relevant to the topic
  - Expand the root set  $R$  to include all the children and a fixed number of parents of nodes in  $R$
- A new set  $S$  (base subgraph) →

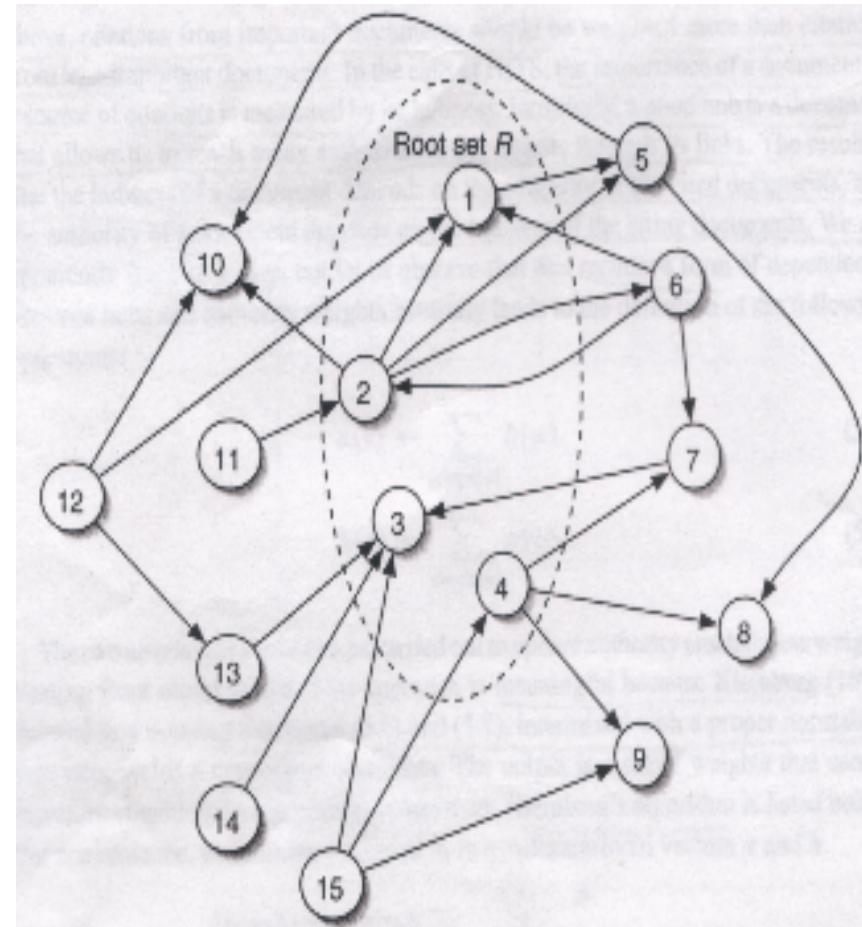


# HITS Example Results

■ Authority  
■ Hubness

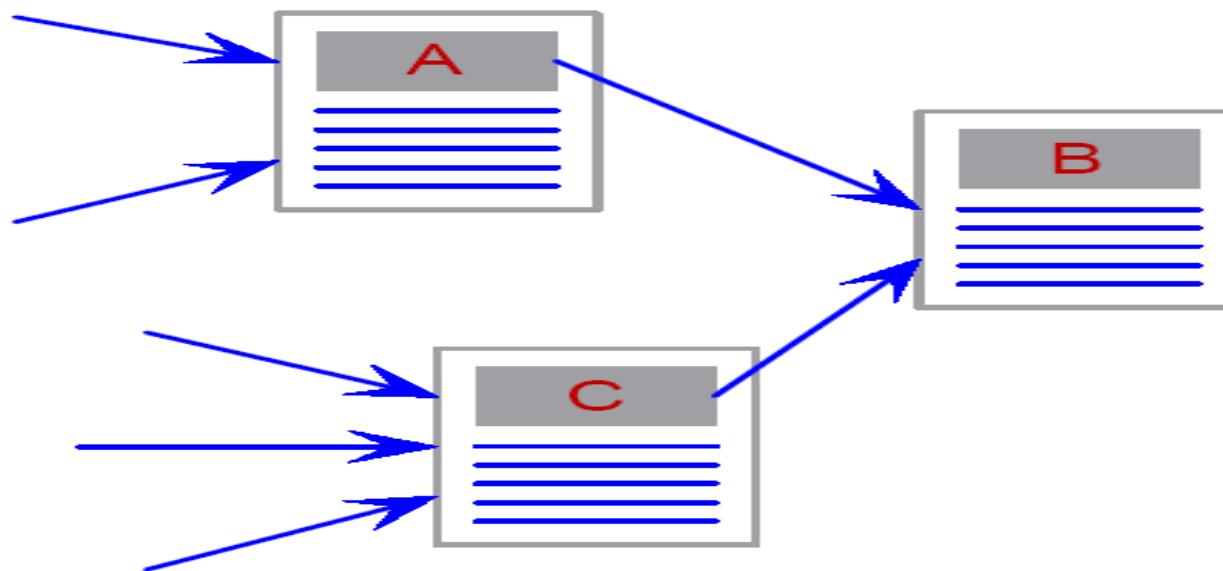


Authority and hubness weights



# PageRank - Motivation

---



- A link from page *A* to page *B* is a **vote** of the author of *A* for *B*, or a **recommendation** of the page.
- The number incoming links to a page is a measure of importance and authority of the page.
- Also take into account the quality of recommendation, so a page is more important if the sources of its incoming links are important.

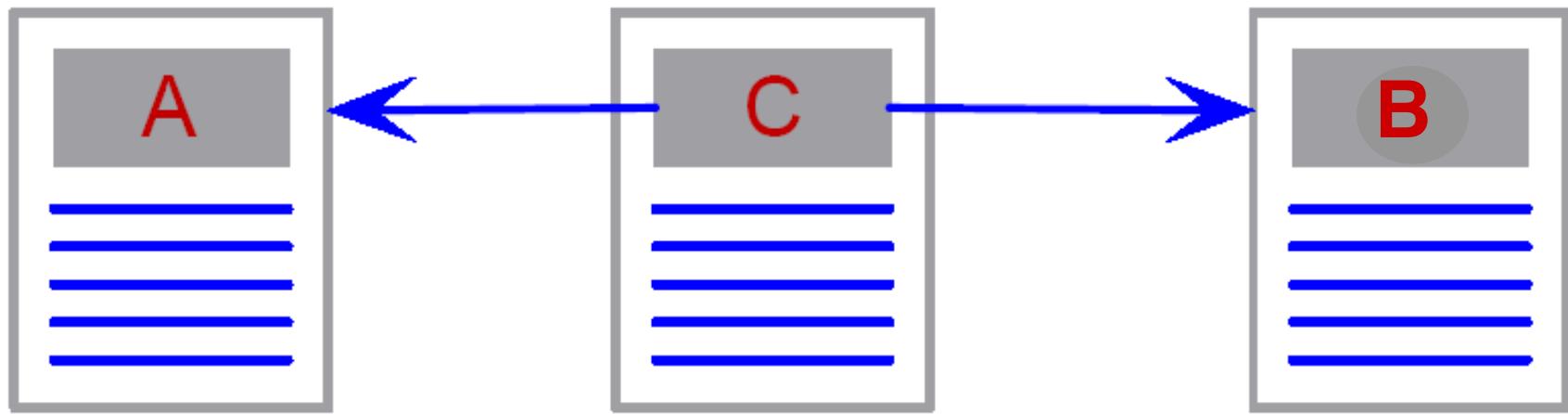
# The Random Surfer

---

- Assume the web is a Markov chain.
- Surfers randomly click on links, where the probability of an outlink from page A is  $1/m$ , where  $m$  is the number of outlinks from A.
- The surfer occasionally gets *bored* and is *teleported* to another web page, say  $B$ , where  $B$  is equally likely to be any page.
- Using the theory of Markov chains it can be shown that if the surfer follows links for long enough, *the PageRank of a web page is the probability that the surfer will visit that page*.

# Dangling Pages

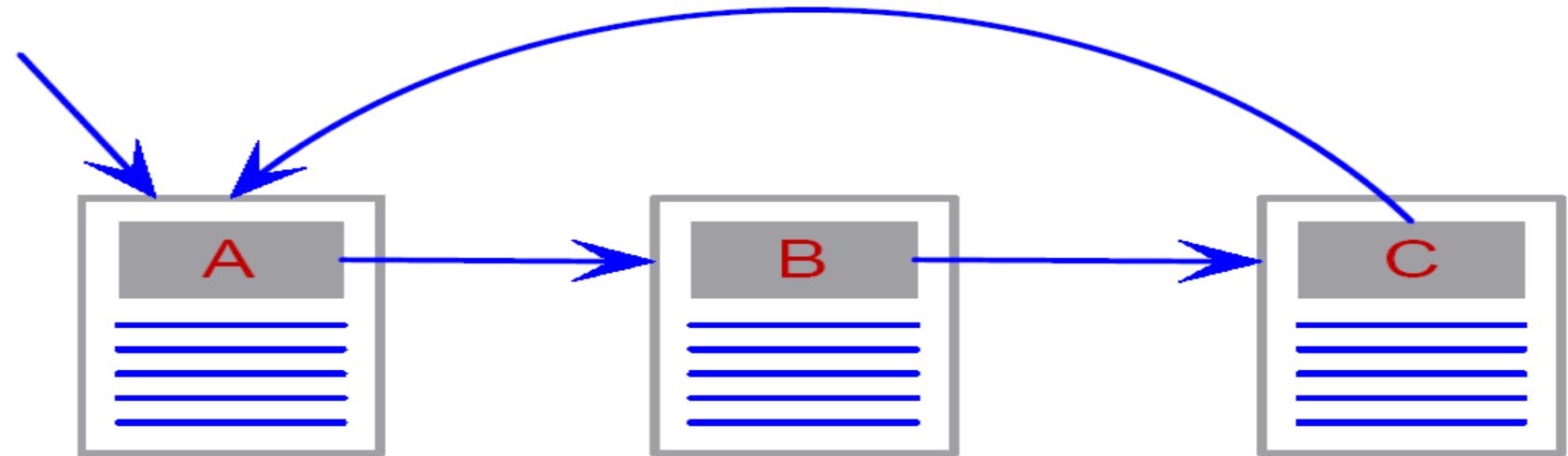
---



■ Problem: *A* and *B* have no outlinks.

Solution: Assume *A* and *B* have links to all web pages with equal probability.

# Rank Sink



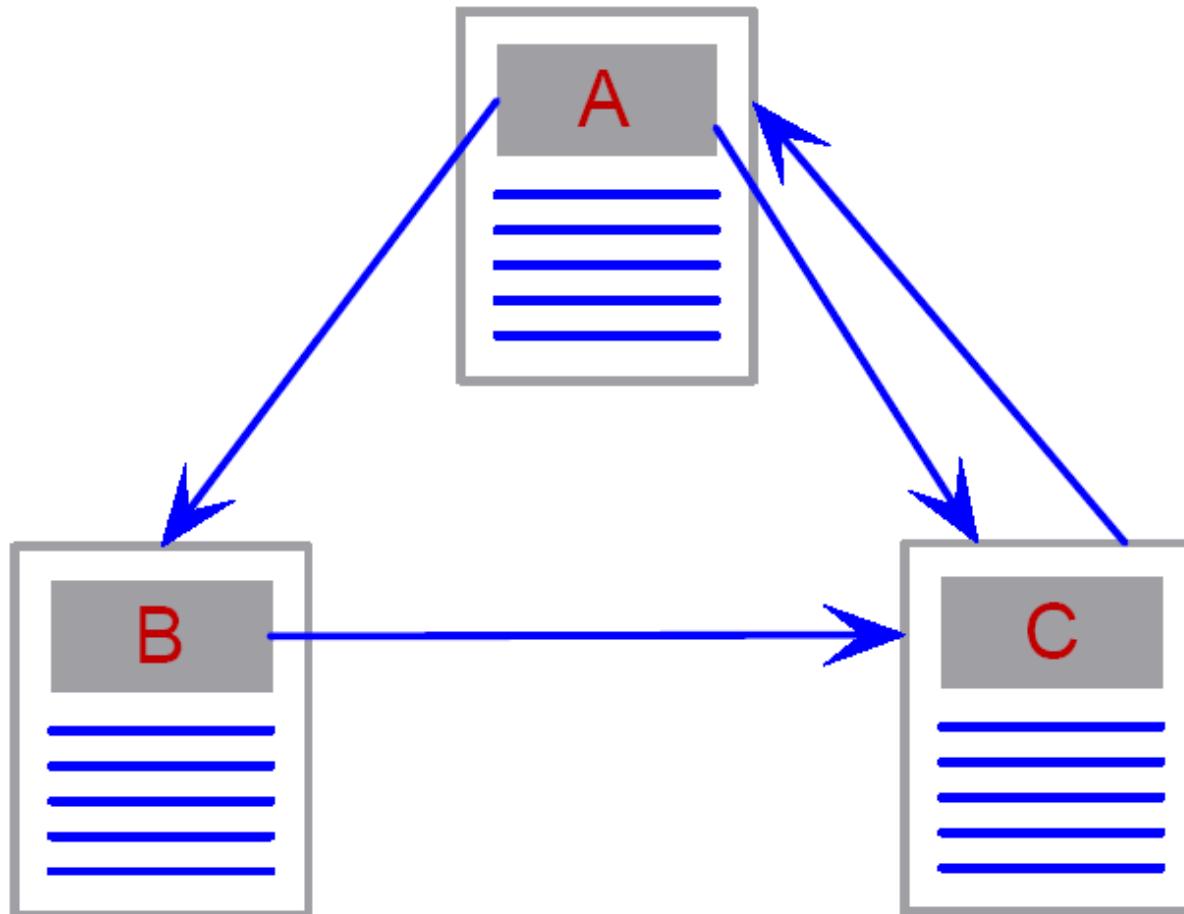
- Problem: Pages in a loop accumulate rank but do not distribute it.
- Solution: Teleportation, i.e. with a certain probability the surfer can jump to any other web page to get out of the loop.

$$PR(P) = \frac{d}{N} + (1-d)\left(\frac{PR(P_1)}{O(P_1)} + \frac{PR(P_2)}{O(P_2)} + \dots + \frac{PR(P_n)}{O(P_n)}\right)$$

- $P$  is a web page
  - $P_i$  are the web pages that have a link to  $P$
  - $O(P_i)$  is the number of outlinks from  $P_i$
  - $d$  is the teleportation probability
  - $N$  is the size of the web
- 
- Difference to HITS
    - HITS takes Hubness & Authority weights
    - The page rank is proportional to its parents' rank, but inversely proportional to its parents' outdegree

# Example Web Graph

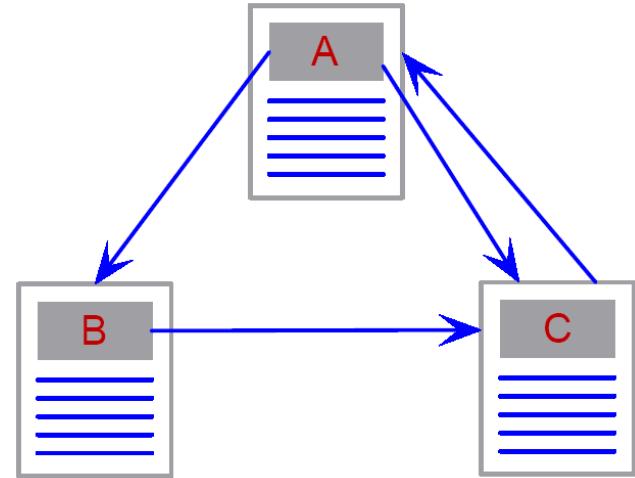
---



# Iteratively Computing PageRank

---

- $d$  is normally set to 0.15
- Set initial  $PR$  values to  $1/3$
- *Solve the following equations iteratively:*



$$PR(A) = 0.15/3 + 0.85PR(C)$$

$$PR(B) = 0.15/3 + 0.85(PR(A)/2)$$

$$PR(C) = 0.15/3 + 0.85(PR(A)/2 + PR(B))$$

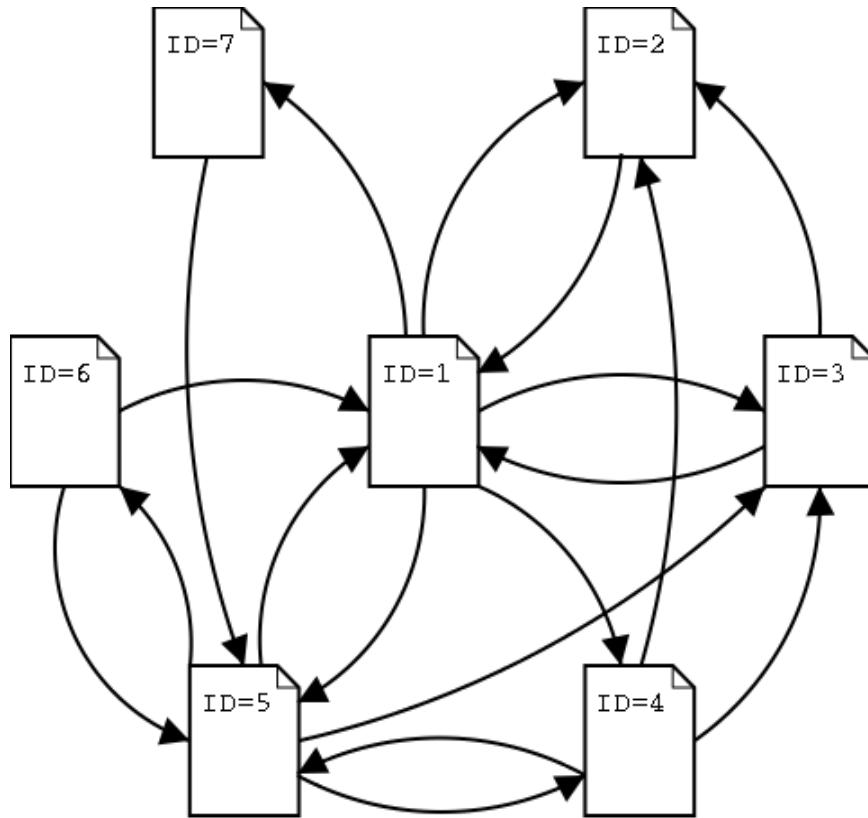
# Example Computation of PR

---

	PR(A)	PR(B)	PR(C)	ERROR
1	0,333333333	0,333333333	0,333333333	
2	0,333333333	0,191666667	0,475	0,04014
3	0,45375	0,191666667	0,354583333	0,029
4	0,351395833	0,24284375	0,405760417	0,01571
5	0,394896354	0,199343229	0,405760417	0,00378
6	0,394896354	0,217830951	0,387272695	0,00068
7	0,379181791	0,217830951	0,402987258	0,00049
8	0,39253917	0,211152261	0,396308569	0,00027
9	0,386862284	0,216829147	0,396308569	6,4E-05
10	0,386862284	0,214416471	0,398721246	1,2E-05
11	0,388913059	0,214416471	0,396670471	8,4E-06
12	0,3871699	0,21528805	0,39754205	4,6E-06
13	0,387910742	0,214547208	0,39754205	1,1E-06
14	0,387910742	0,214862066	0,397227192	2E-07
15	0,387643113	0,214862066	0,397494821	1,4E-07
16	0,387870598	0,214748323	0,397381079	7,8E-08
17	0,387773917	0,214845004	<b>0,397381079</b>	1,9E-08

- Error converges fast, ~10 repetitions
- Page C is the top ranked one

# Matrix Notation



<b>Page ID</b>	<b>OutLinks</b>
1	2,3,4,5,7
2	1
3	1,2
4	2,3,5
5	1,3,4,6
6	1,5
7	5

Adjacency Matrix

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

\* <http://www.kusatro.kyoto-u.com>

# Matrix Notation

---

PageRank: eigenvector of  $\mathbf{A}$  relative to max eigenvalue

$$\mathbf{A} = \mathbf{U} \Lambda \mathbf{U}^T$$

$\mathbf{L}$ : diagonal matrix of eigenvalues  $\{\lambda_1, \dots, \lambda_n\}$

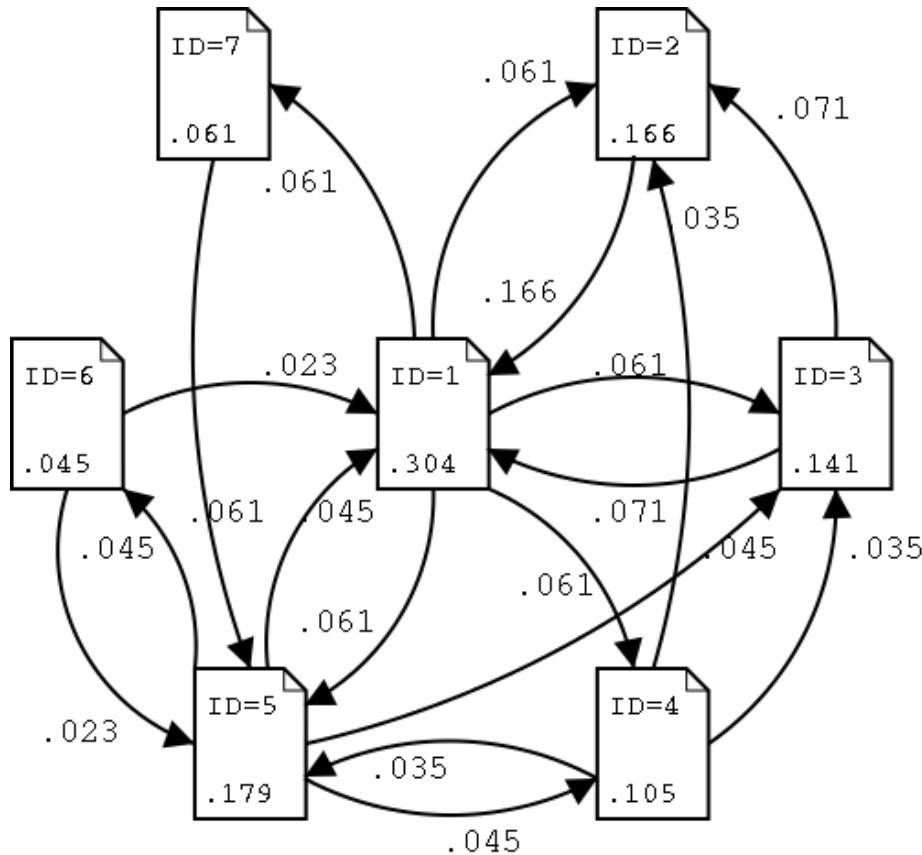
$$\begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{pmatrix} \quad (\mathbf{r}_1 \ \mathbf{r}_2 \ \dots \ \mathbf{r}_n)$$

$\mathbf{U}$ : regular matrix that consists of eigenvectors

Approximation: Power method:  $\mathbf{P}^{i+1} = \mathbf{A} \mathbf{P}^i$

PageRank  $\mathbf{r}_1 = \begin{pmatrix} 0.69946 \\ 0.38286 \\ 0.32396 \\ 0.24297 \\ 0.41231 \\ 0.10308 \\ 0.13989 \end{pmatrix} \xrightarrow{\text{normalized}} \begin{pmatrix} 0.303514 \\ 0.166134 \\ 0.140575 \\ 0.105431 \\ 0.178914 \\ 0.044728 \\ 0.060703 \end{pmatrix}$

# Matrix Notation



PR	ID	OutLink	InLink
<b>0.304</b>	<b>1</b>	<b>2,3,4,5,7</b>	<b>2,3,5,6</b>
<b>0.179</b>	<b>5</b>	<b>1,3,4,6</b>	<b>1,4,6,7</b>
<b>0.166</b>	<b>2</b>	<b>1</b>	<b>1,3,4</b>
<b>0.141</b>	<b>3</b>	<b>1,2</b>	<b>1,4,5</b>
<b>0.105</b>	<b>4</b>	<b>2,3,5</b>	<b>1,5</b>
<b>0.061</b>	<b>7</b>	<b>5</b>	<b>1</b>
<b>0.045</b>	<b>6</b>	<b>1,5</b>	<b>5</b>

- Confirm the result  
# of inlinks from high ranked page  
hard to explain about 5&2, 6&7
- Interesting Topic  
How do you create your homepage  
highly ranked?

# “Rank Sink” Problem

---

- In general, many Web pages have no inlinks/outlinks
- It results in dangling edges in the graph

E.g.

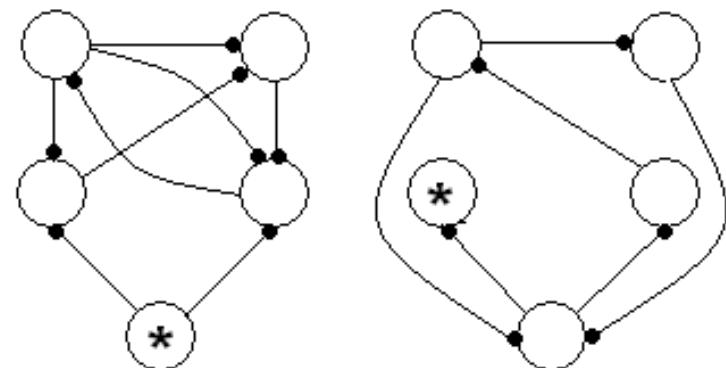
no parent  $\rightarrow$  rank 0

$M^T$  converges to a matrix

whose last column is all zero

no children  $\rightarrow$  no solution

$M^T$  converges to zero matrix



# Modification

---

- Surfer will restart browsing by picking a new Web page at random

$$\mathbf{M} = (\mathbf{B} + \mathbf{E}), \quad e_{vw} = \begin{cases} 0 & \text{if } |ch[v]| > 0 \\ \frac{1}{n} & \text{otherwise} \end{cases}$$

E : escape matrix

M : stochastic matrix

- Still problem?

- It is not guaranteed that **M** is primitive
- If **M** is stochastic and primitive, PageRank converges to corresponding stationary distribution of **M**

# PageRank Algorithm

---

```
PAGERANK( $M$ ,  $n$ ,  $\epsilon$ )
1    $\mathbf{1} \leftarrow [1, \dots, 1] \in \mathbb{R}^n$ 
2    $\mathbf{z} \leftarrow \frac{1}{n}\mathbf{1}$ 
3    $\mathbf{x}_0 \leftarrow \mathbf{z}$ 
4    $t \leftarrow 0$ 
5   repeat
6        $t \leftarrow t + 1$ 
7        $\mathbf{x}_t \leftarrow M^T \mathbf{x}_{t-1}$ 
8        $d_t \leftarrow \|\mathbf{x}_{t-1}\|_1 - \|\mathbf{x}_t\|_1$ 
9        $\mathbf{x}_t \leftarrow \mathbf{x}_1 + d_t \mathbf{z}$ 
10       $\delta \leftarrow \|\mathbf{x}_{t-1} - \mathbf{x}_t\|_1$ 
11      until  $\delta < \epsilon$ 
12  return  $\mathbf{x}_t$ 
```

\* Page et al, 1998

# The Largest Matrix Computation in the World

---

- Computing PageRank can be done via matrix multiplication, where the matrix has several billion rows and columns.
- The matrix is sparse as average number of outlinks is between 7 and 8.
- Setting  $d = 0.15$  or below requires at most 100 iterations to convergence.
- Researchers still trying to speed-up the computation.

# Ranking function web search

---

- Web search engines take into account 100's of features to rank documents assuming a query
- Two important features are
  - The *PageRank* value of the page containing the *query* terms
  - The *relevance* of the term to the specific page
- Given a term  $t$  the score of a document  $d$  is computed as:
- Where  $score_t(d_i) = w_1 \cdot \text{relevance}(t, d_i) + w_2 pr(d_i)$
- In a specific case we used:

$$score_t(d) = (\text{tf/idf}(t, d) \cdot \text{title}(t, d))^{1.5} \cdot pr(d)$$

# References

---

- Amy Nicole Langville, Carl Dean Meyer: Survey: Deeper Inside PageRank. Internet Mathematics 1(3): 335-380 (2003)
- “PageRank Computation and the Structure of the Web: Experiments and Algorithms”, Arvind Arasu, Jasmine Novak, Andrew Tomkins & John Tomlin

# Relevant publications

---

- C. Giatsidis, D. Thilikos, and M. Vazirgiannis, D-cores: Measuring Collaboration of Directed Graphs Based on Degeneracy. *Knowledge and Information Systems Journal*, Springer, 2012.
- C. Giatsidis, K. Berberich, D. M. Thilikos, M. Vazirgiannis, Visual exploration of collaboration networks based on graph degeneracy, *ACM KDD*, 2012.
- C. Giatsidis, D. Thilikos, and M. Vazirgiannis, D-cores: Measuring Collaboration of Directed Graphs Based on Degeneracy. *IEEE ICDM*, 2011,
- C. Giatsidis, D. Thilikos, and M. Vazirgiannis, Evaluating Cooperation in Communities with the k-Core Structure. *ACM/IEEE ASONAM*, 2011.
- F. D. Malliaros and M. Vazirgiannis, To Stay or Not to Stay: Modeling Engagement Dynamics in Social Graphs. *ACM CIKM*, 2013.
- F. D. Malliaros and M. Vazirgiannis, Clustering and Community Detection in Directed Networks: A Survey. *Physics Reports*, 533(4), Elsevier, 2013.
- F. D. Malliaros and M. Vazirgiannis, Vulnerability Assessment in Social Networks under Cascade-based Node Departures, *EPL (Europhysics Letters)*, 11(6), 2015.
- C. Giatsidis, F.D. Malliaros, D. Thilikos, and M. Vazirgiannis, CORECLUSTER: A Degeneracy Based Graph Clustering Framework. *AAAI*, 2014.
- F.D. Malliaros, V. Megalooikonomou and C. Faloutsos. Estimating Robustness in Large Social Graphs. *Knowledge and Information Systems (KAIS)*, Springer, 2015.
- M.-E. G. Rossi, F.D. Malliaros, and M. Vazirgiannis, Spread It Good, Spread It Fast: Identification of Influential Nodes in Social Networks. *WWW*, 2015.

# Relevant publications

---

## Invited Tutorials

- C. Giatsidis, F. D. Malliaros and M. Vazirgiannis, Graph Mining Tools for Community Detection and Evaluation in Social Networks and the Web. *WWW*, Rome, Italy, 2013.
- C. Giatsidis, F. D. Malliaros and M. Vazirgiannis, Advanced graph mining for community evaluation in social networks and the web, *ACM WSDM*, Rio de Janeiro, Brazil, 2013.
- C. Giatsidis, F. D. Malliaros and M. Vazirgiannis, Community Detection and Evaluation in Social and Information Networks. *WISE*, Thessaloniki, Greece, 2014.
- F. D. Malliaros, M. Vazirgiannis and A.N. Papadopoulos, Core Decomposition: Algorithms and Applications, *IEEE/ACM ASONAM*, Paris, France, 2015.
- F. D. Malliaros, A.N. Papadopoulos, Core Decomposition in Graphs: Concepts, Algorithms and Applications. *ICDM*, Barcelona, 2016.

## Demos

<http://graphdegeneracy.org/>