

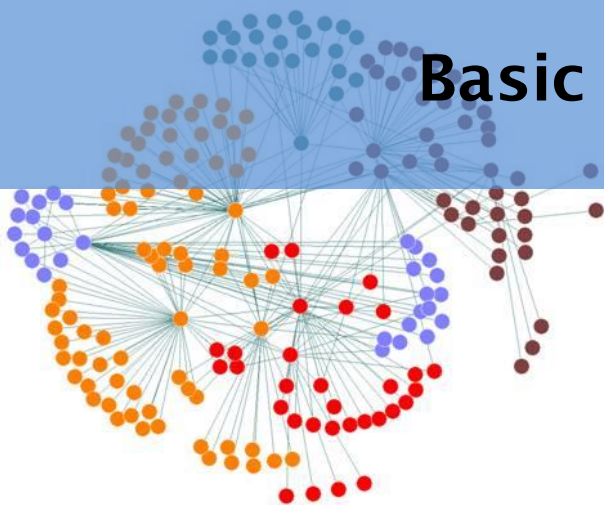


BABEȘ-BOLYAI UNIVERSITY
Faculty of Mathematics and Computer Science



Social Network Analysis

Basic definitions and properties



Camelia Chira

camelia.chira@ubbcluj.ro

Network science

- Properties
 - Interdisciplinary
 - Empirical
 - Quantitative and Mathematical
 - Computational
- Networks really matter

Some basic definitions

- Network:
 - Collection of individuals or entities (nodes)
 - List of pairs of nodes that are linked (edges)
- Networks can represent any binary relationship over individuals
- N = the number of nodes
- The number of possible edges: $N(N-1)/2$
- Degree of a node = number of neighbors
- Degree distribution $P(k)$ – fraction of nodes with degree k
- Distance (or shortest path) between two nodes = length of the shortest path connecting them (what if two nodes are in different components?)
- Diameter of the network = the longest shortest path
- Directed vs undirected networks

Network analysis

- Network structure
 - Size, diameter, connectedness, degree distribution
 - Communities, centrality, important nodes
 - Robustness, vulnerability
- Network dynamics
 - Evolving networks
 - Spreading phenomena
 - Linked with network structure
- Network formation
 - models

Complex network properties

- **Scale free**: power law degree distribution
- **Small world**: small diameter and average path length
- Dense local structure (high clustering coefficient)
- Giant connected component
- Community structure

Example: scientific collaboration

- **Nodes**: researchers (mathematics and computer science)
- **Edges**: link two researchers if they co-author a paper
- Erdős number: distance to **Paul Erdős**
 - Erdős (1913-1996) had 504 co-authors (hub)

MR Erdos Number = 4

Camelia Chira	coauthored with	Gheorghe Ștefănescu	MR2682175
Gheorghe Ștefănescu	coauthored with	Cristian S. Calude	MR2390233
Cristian S. Calude	coauthored with	Solomon Marcus	MR2014322
Solomon Marcus	coauthored with	Paul Erdős ¹	MR0095456

Erdős Number

Erdős number 0 --- 1 person
Erdős number 1 --- 504 people
Erdős number 2 --- 6593 people
Erdős number 3 --- 33605 people
Erdős number 4 --- 83642 people
Erdős number 5 --- 87760 people
Erdős number 6 --- 40014 people
Erdős number 7 --- 11591 people
Erdős number 8 --- 3146 people
Erdős number 9 --- 819 people
Erdős number 10 --- 244 people
Erdős number 11 --- 68 people
Erdős number 12 --- 23 people
Erdős number 13 --- 5 people
The median 5; the mean 4.65.

<https://sites.google.com/a/oakland.edu/jerry-grossman-home-page/home/the-erdoes-number-project/facts-about-erdoes-numbers-and-the-collaboration-graph>

Carmen Chiriac

Data on the entire collaboration graph

There are about **1.9 million authored items** in the Math Reviews database, by a total of about **401,000 different authors**. (This includes all books and papers in MR except those items, such as some conference proceedings, that do not have authors.) **Approximately 62.4% of these items are by a single author, 27.4% by two authors, 8.0% by three authors, 1.7% by four authors, 0.4% by five authors, and 0.1% by six or more authors.** The largest number of authors shown for a single item is in the 20s, but sometimes the author list includes “et al.”, whom we did not count as a real person. **The fraction of items authored by just one person has steadily decreased over time, starting out above 90% in the 1940s and currently standing at under 50%.**

Let B be the bipartite graph whose vertices are papers and authors, with an edge joining a paper with each author of that paper. Then B has about **2.9 million edges**. The **average number of authors per paper is 1.51**, and the **average number of papers per author is 7.21**. Here is the distribution of the number of papers per author. The median is 2, the mean is 7.21, and the standard deviation is 18.02. It is interesting (for tenure review committees?) to note that the 60th percentile is 3 papers, the 70th percentile is 4, the 80th percentile is 8, the 90th percentile is 18, and the 95th percentile is 32. Indeed, **over 42% of all authors in the database have just one paper.**

There are **four authors with more than 700 papers**: Paul Erdős with 1416 (he actually wrote more papers than that, but these are just the ones covered by Math Reviews), Drumi Bainov with 823, SAHARON SHELAH with 760, and Leonard Carlitz with 730. Bainov's Erdős number is 4, SHELAH's is 1, and Carlitz's is 2. The other mathematicians with more than 500 papers listed in MathSciNet (and their Erdős numbers) are Hari M. Srivastava (2), Lucien Godeaux (infinite — actually he wrote only one joint paper), Ravi Agarwal (3), Edoardo Ballico (3), FRANK HARARY (1), Josip E. Pecaric (2), Shigeyoshi Owa (3), and Richard Bellman (2). The most prolific authors listed in the DBLP (dealing with computer science publications) can be found on a [list](#) at their website (DBLP), which is definitely worth exploring.

The **collaboration graph** C has the roughly 401,000 authors as its vertices, with an edge between every pair of people who have a joint publication (with or without other coauthors — but see below for a discussion of “Erdős number of the second kind”, where we restrict links to just two-author papers). Here is a [picture of a small portion of this graph](#). The entire graph has about **676,000 edges**, so the **average number of collaborators per person is 3.36**. (If we were to view C as a multigraph, with one edge between two vertices for each paper in which they collaborated, then there would be about 1,300,000 edges, for an average of 6.55 collaborations per person.) In C there is one large component consisting of about 268,000 vertices. Of the remaining 133,000 authors, **84,000 of them have written no joint papers** (these are **isolated vertices** in C). The average number of collaborators for people who have collaborated is 4.25; the average number of collaborators for people in the large component is 4.73; and the average number of collaborators for people who have collaborated but are not in the large component is 1.65.

Degree Distribution

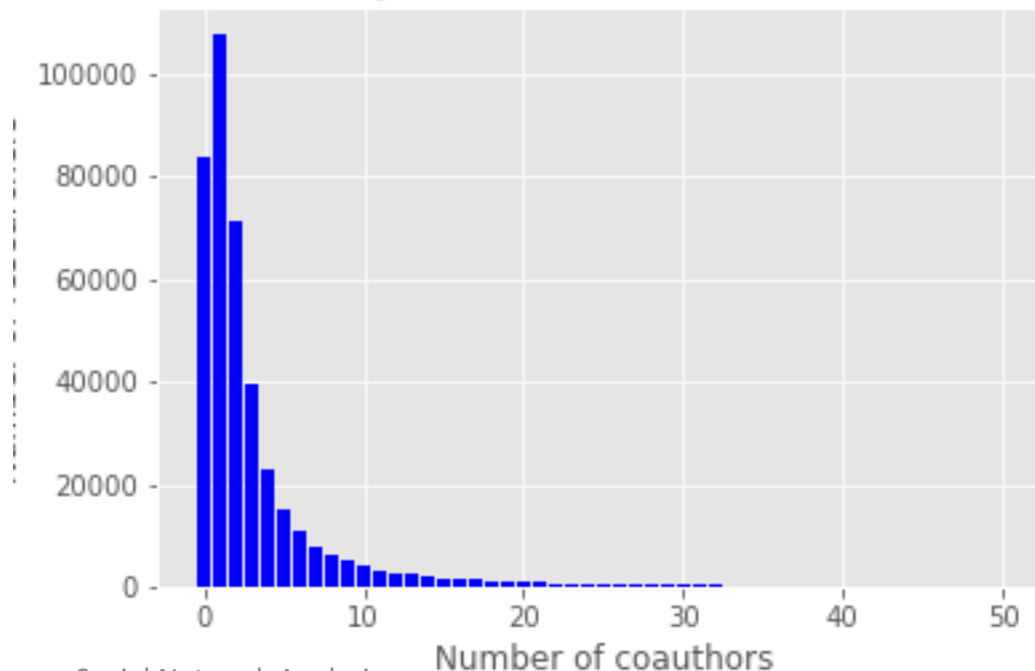
This table shows the distribution of the number of coauthors for the 401,445 authors in the Mathematical Reviews database. For example, 39,574 different mathematicians have exactly 3 coauthors. In other words, this table gives the distribution of the degrees of the vertices in the collaboration graph C . The mean number of coauthors is 4.24819 and the standard deviation is 6.60515.

number of coauthors	count
0	83621
1	107647
2	71452
3	39574
4	22815
5	15205
6	10679
7	7917
8	6255
9	4959
10	4141
11	3283
12	2808
13	2368
14	1993
15	1756
16	1575
17	1311
18	1147

...

Camelia Chira

How many researchers have x co-authors?

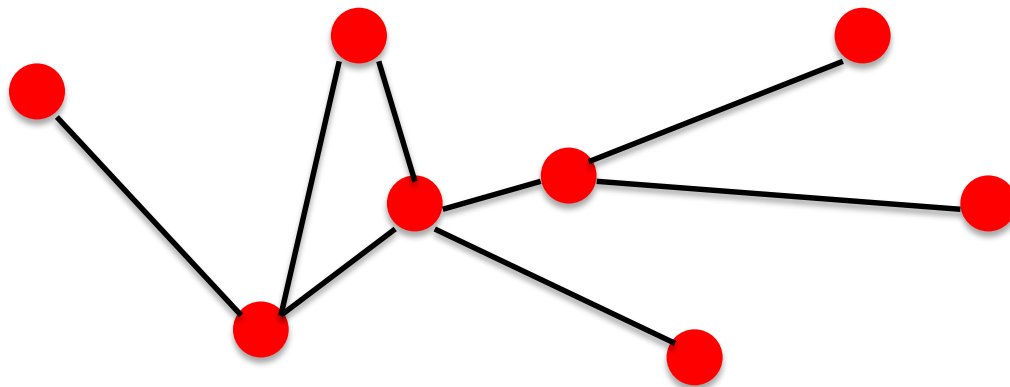


Social Network Analysis

TODAY: Graphs and Networks

- Graph $G(V, E)$, V – nodes, E – edges
- Can be: directed, undirected, weighted
- Connected (if every pair of nodes is connected)
- Connected component
- Path
- Diameter
- Node degree
- Degree distribution

Graphs and Networks



- **components:** nodes, vertices N
- **interactions:** links, edges L
- **system:** network, graph (N,L)

Networks or graphs?

Network Science	Graph Theory
Network	Graph
Node	Vertex
Link	Edge

Real systems

Mathematical
representation of
networks

WWW - a network of
documents linked by
URLs

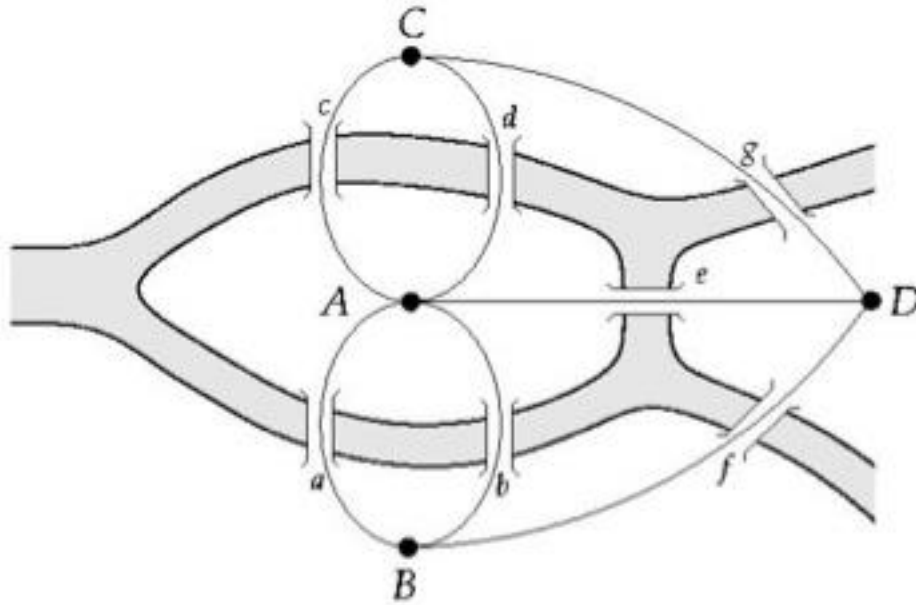
Society - a network of individuals
linked by family / friendship /
professional ties

THE BRIDGES OF KONIGSBERG



Can one walk across the seven bridges and never cross the same bridge twice?

THE BRIDGES OF KONIGSBERG

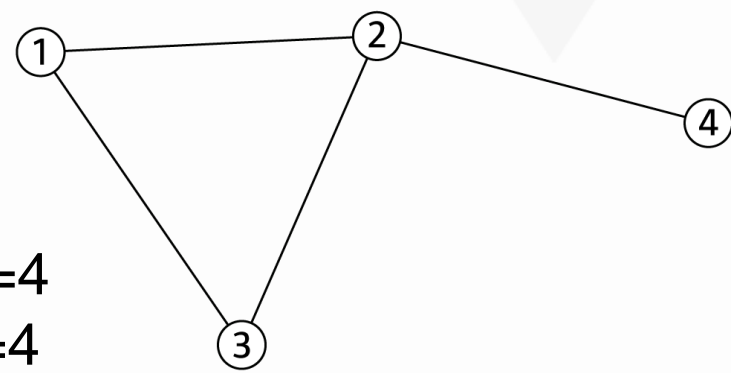
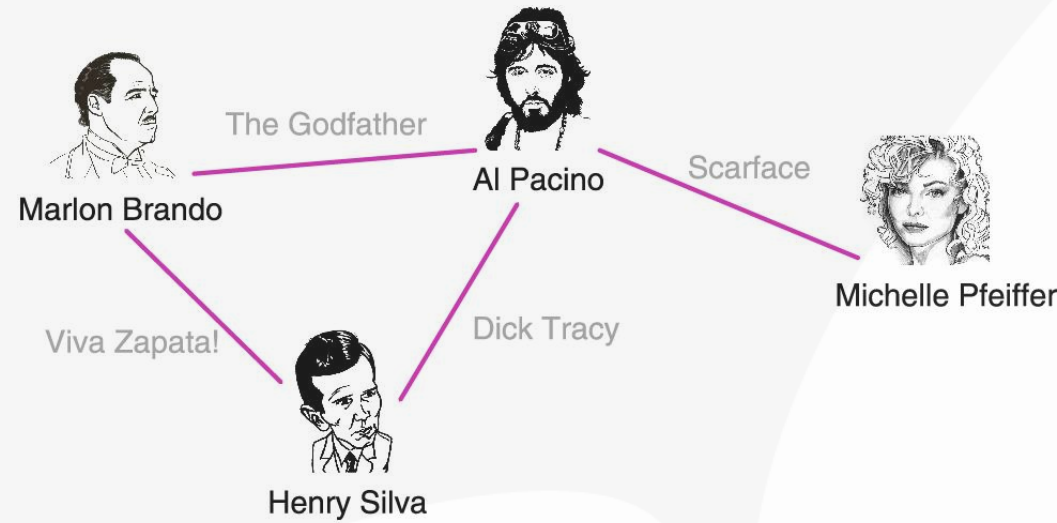
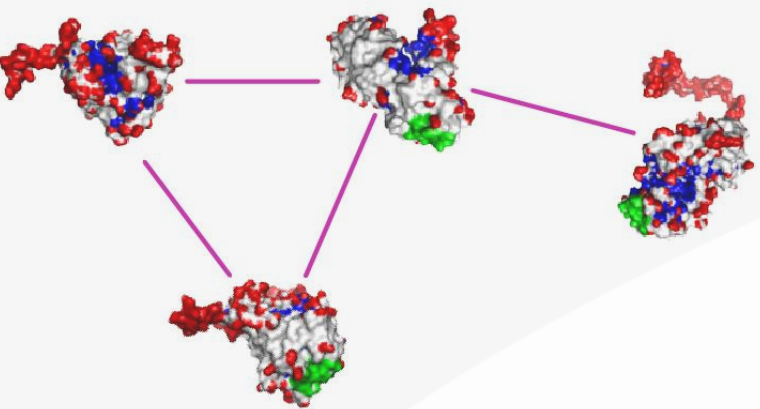
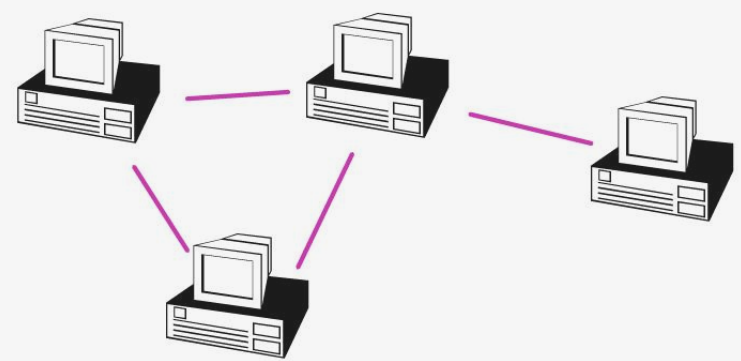


Can one walk across the seven bridges and never cross the same bridge twice?

1735: Euler's theorem:

- (a) If a graph has more than two nodes of odd degree, there is no path.
- (b) If a graph is connected and has no odd degree nodes, it has at least one path.

A COMMON LANGUAGE



$N=4$
 $L=4$

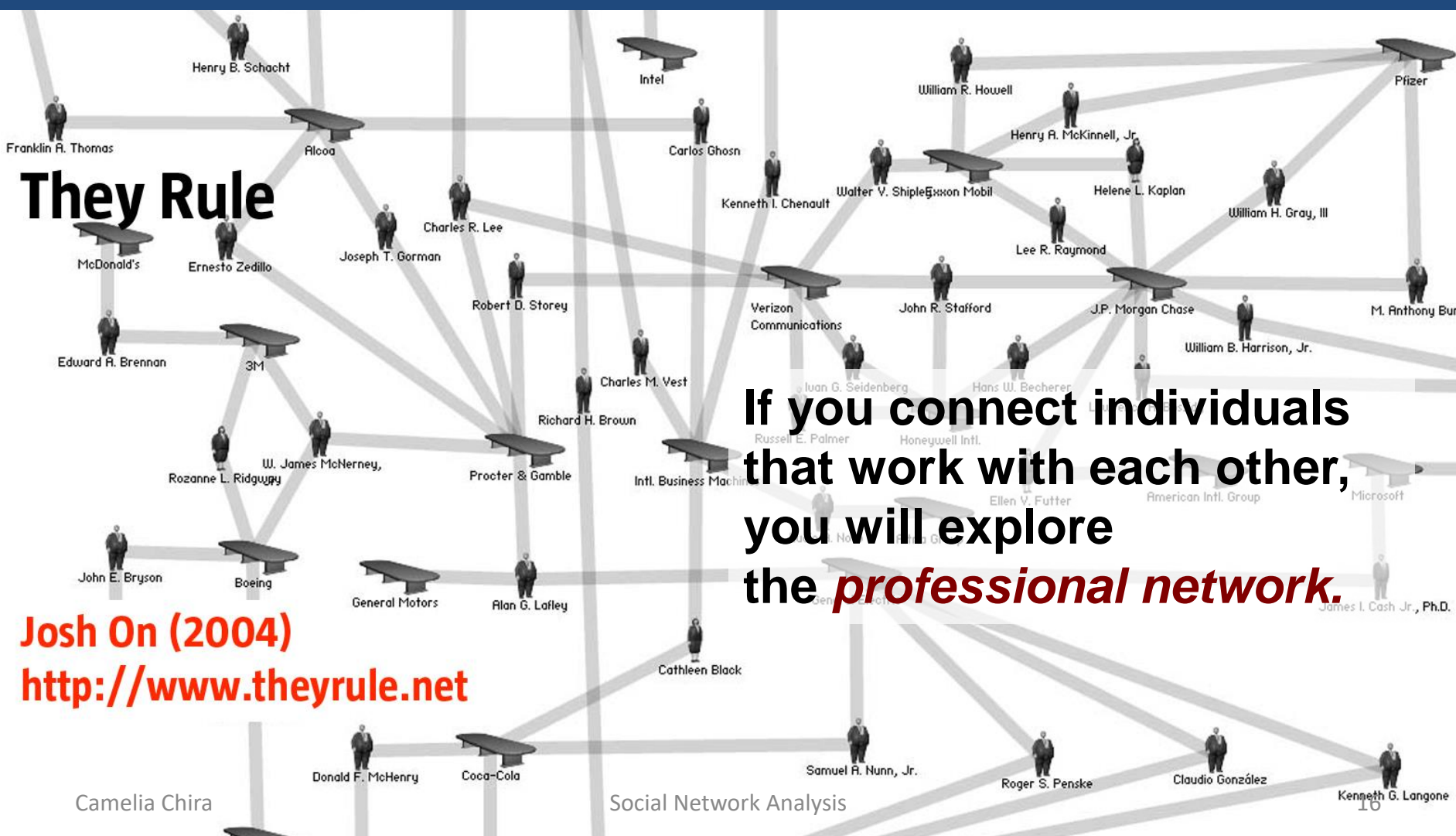
CHOOSING A PROPER REPRESENTATION

They Rule

Josh On (2004)

<http://www.theyrule.net>

**If you connect individuals
that work with each other,
you will explore
the *professional network*.**



The structure of adolescent romantic and sexual networks

If you connect those that have a romantic and sexual relationship, you will be exploring the *sexual networks*.

Bearman PS, Moody J, Stovel K.

Institute for Social and Economic Research and Policy - Columbia University

<http://researchnews.osu.edu/archive/chainspix.htm>

If you connect individuals based on their first name (*all Peters connected to each other*), you will be exploring what?

It is a network, nevertheless.

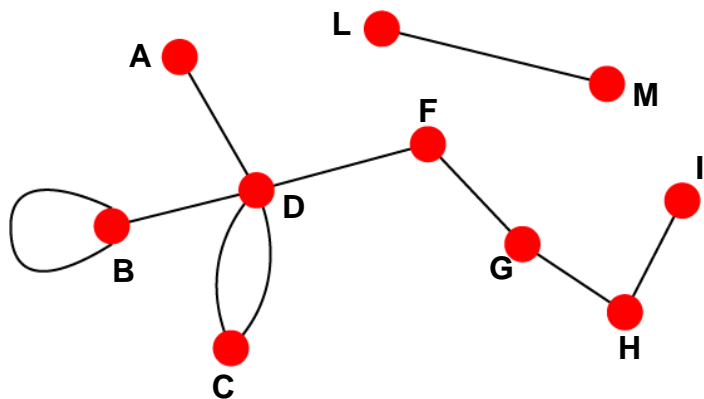
Definitions and properties

- Degrees
 - Node degree
 - Degree distribution
- Paths and distances
 - Distance (shortest path) between two nodes
 - Diameter of the network
 - Average distance (or average path length)
- Connectivity
 - Connected components
 - Clustering coefficient of a node
 - Average clustering coefficient of the network

Undirected

Links: undirected (*symmetrical*)

Graph:



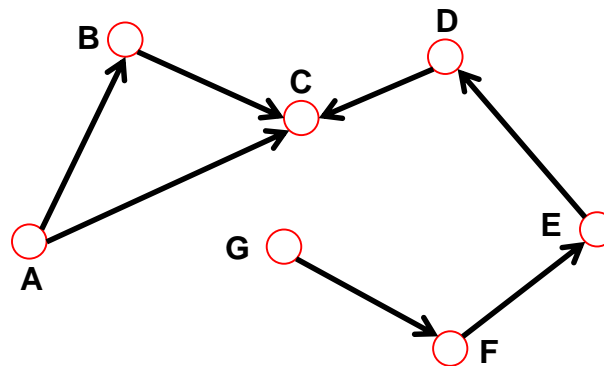
Undirected links :

coauthorship links
Actor network
protein interactions

Directed

Links: directed (*arcs*).

Digraph = directed graph:



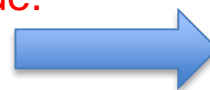
Directed links :

URLs on the www
phone calls
metabolic reactions

Node degree

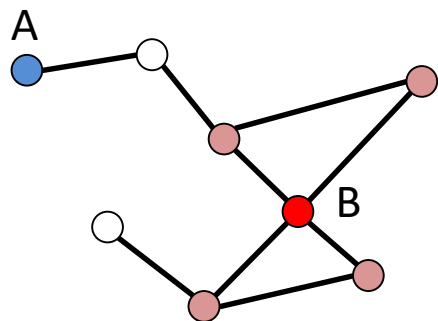
Node degree: the number of links connected to the node.

k_i - the degree of i^{th} node in the network



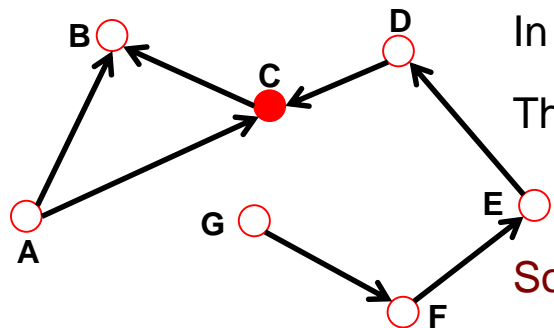
$$L = \frac{1}{2} \sum_{i=1}^N k_i$$

Undirected



$$k_A = 1 \quad k_B = 4$$

Directed



In *directed networks* we can define an **in-degree** and **out-degree**.

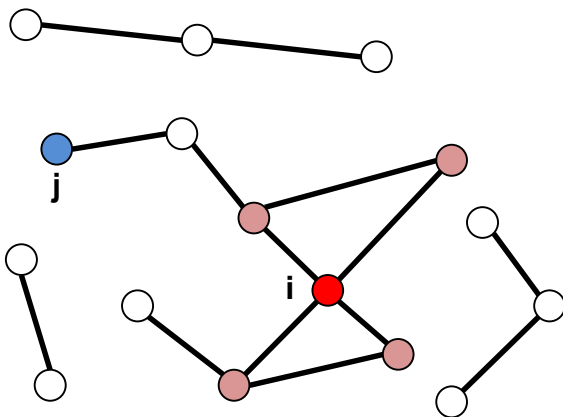
The (total) degree is the sum of in- and out-degree.

$$k_C^{\text{in}} = 2 \quad k_C^{\text{out}} = 1 \quad k_C = 3$$

Source: a node with $k^{\text{in}} = 0$; **Sink**: a node with $k^{\text{out}} = 0$.

Average degree

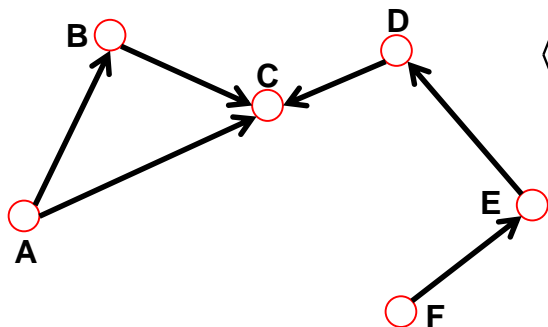
Undirected



$$\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i \quad \langle k \rangle = \frac{2L}{N}$$

N – the number of nodes in the graph

Directed



$$\langle k^{in} \rangle \equiv \frac{1}{N} \sum_{i=1}^N k_i^{in}, \quad \langle k^{out} \rangle \equiv \frac{1}{N} \sum_{i=1}^N k_i^{out}, \quad \langle k^{in} \rangle = \langle k^{out} \rangle$$

$$\langle k \rangle = \frac{L}{N}$$

Reference NETWORKS

Network	Nodes	Links	Directed / Undirected	N	L	⟨K⟩
Internet	Routers	Internet connections	Undirected	192,244	609,066	6.34
WWW	Webpages	Links	Directed	325,729	1,497,134	4.60
Power Grid	Power plants, transformers	Cables	Undirected	4,941	6,594	2.67
Mobile-Phone Calls	Subscribers	Calls	Directed	36,595	91,826	2.51
Email	Email addresses	Emails	Directed	57,194	103,731	1.81
Science Collaboration	Scientists	Co-authorships	Undirected	23,133	93,437	8.08
Actor Network	Actors	Co-acting	Undirected	702,388	29,397,908	83.71
Citation Network	Papers	Citations	Directed	449,673	4,689,479	10.43
E. Coli Metabolism	Metabolites	Chemical reactions	Directed	1,039	5,802	5.58
Protein Interactions	Proteins	Binding interactions	Undirected	2,018	2,930	2.90

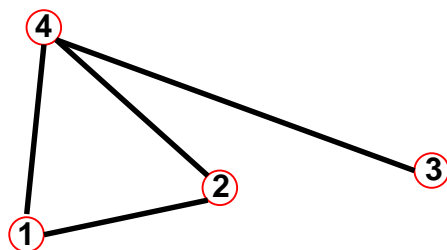
Network representation

ADJACENCY MATRIX

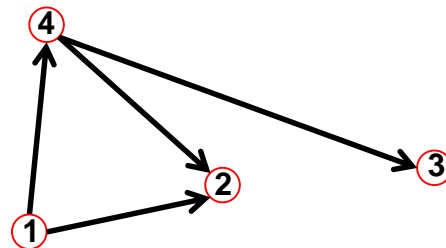
$A_{ij}=1$ if there is a link between node i and j

$A_{ij}=0$ if nodes i and j are not connected to each other.

$$k_i = \sum_{j=1}^N A_{ji} = \sum_{i=1}^N A_{ji}$$



$$A_{ij} = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$



$$A_{ij} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

Note that for a directed graph (right) the matrix is not symmetric.

$A_{ij} = 1$ if there is a link pointing from node j and i

$A_{ij} = 0$ if there is no link pointing from j to i .

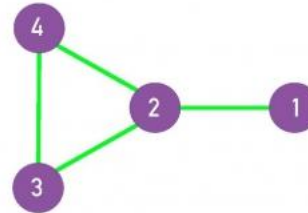
Degree distribution

$P(k)$:
probability that a
randomly chosen node
has degree k

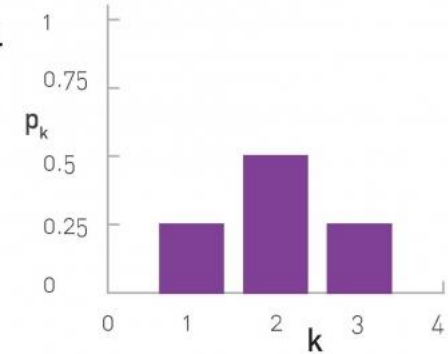
$N_k = \# \text{ nodes with degree } k$

$P(k) = N_k / N$

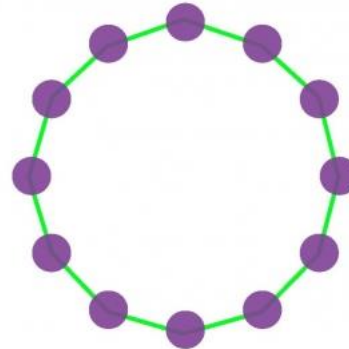
a.



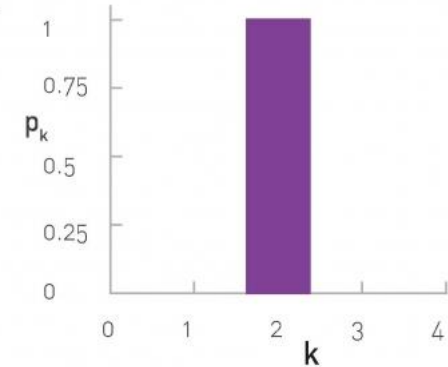
b.



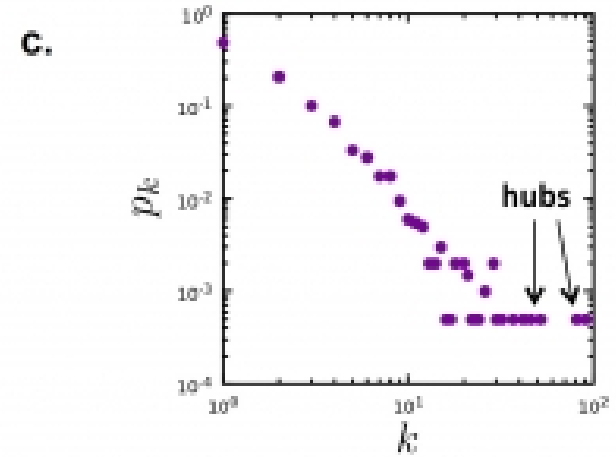
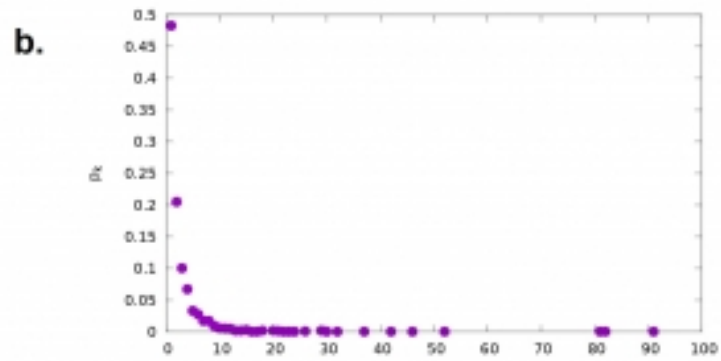
c.



d.

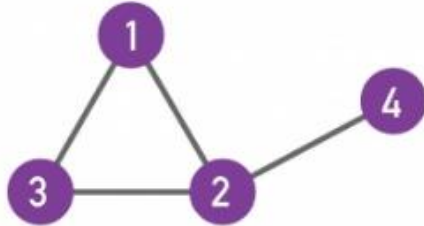


Degree distribution of a real network



ADJACENCY MATRIX AND NODE DEGREES

Undirected



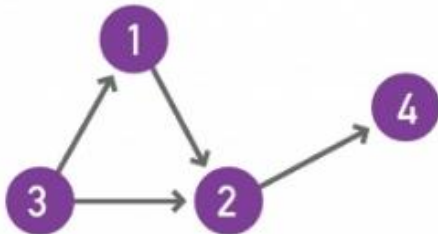
$$A_{ij} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$$A_{ij} = A_{ji} \quad A_{ii} = 0$$

$$k_2 = \sum_{j=1}^4 A_{2j} = \sum_{i=1}^4 A_{i2} = 3$$

$$L = \frac{1}{2} \sum_{i,j=1}^N A_{ij} \quad \langle k \rangle = \frac{2L}{N}$$

Directed



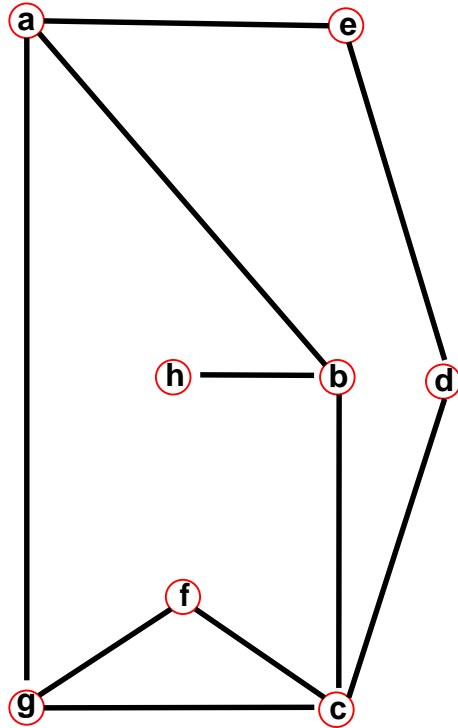
$$A_{ij} = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$$A_{ij} \neq A_{ji} \quad A_{ii} = 0$$

$$k_2^{\text{in}} = \sum_{j=1}^4 A_{2j} = 2, \quad k_2^{\text{out}} = \sum_{i=1}^4 A_{i2} = 1$$

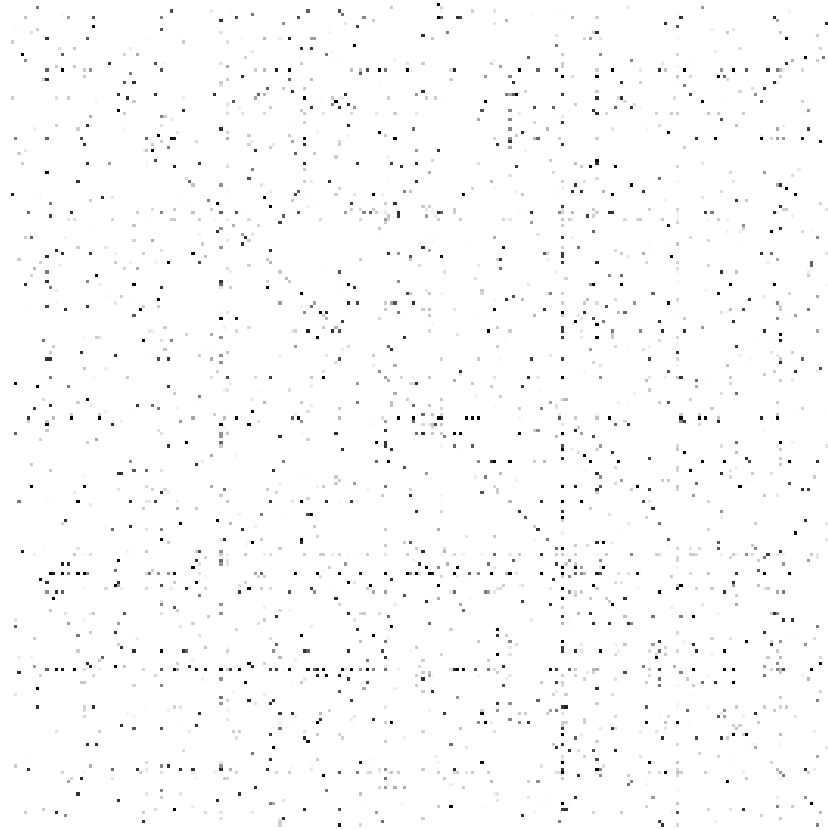
$$L = \sum_{i,j=1}^N A_{ij} \quad \langle k^{\text{in}} \rangle = \langle k^{\text{out}} \rangle = \frac{L}{N}$$

Adjacency matrix



	a	b	c	d	e	f	g	h
a	0	1	0	0	1	0	1	0
b	1	0	1	0	0	0	0	1
c	0	1	0	1	0	1	1	0
d	0	0	1	0	1	0	0	0
e	1	0	0	1	0	0	0	0
f	0	0	1	0	0	0	1	0
g	1	0	1	0	0	0	0	0
h	0	1	0	0	0	0	0	0

Adjacency matrices are sparse



Real networks are sparse

Other network representations

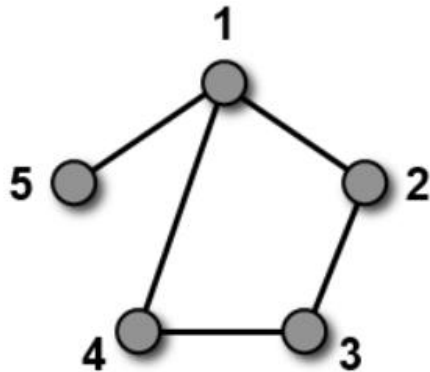
- Edge lists

- (ij) form: only if $\mathbf{A}_{ij}=1$

- Order of rows not important

- Undirected networks: only record links for $i < j$ (L rows)

$$I = \begin{pmatrix} \dots \\ i \\ \dots \end{pmatrix} \quad J = \begin{pmatrix} \dots \\ j \\ \dots \end{pmatrix}$$



$$I = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 3 \end{pmatrix} \quad J = \begin{pmatrix} 2 \\ 4 \\ 5 \\ 3 \\ 4 \end{pmatrix}$$

Other network representations

- Adjacency lists

N+1 elements

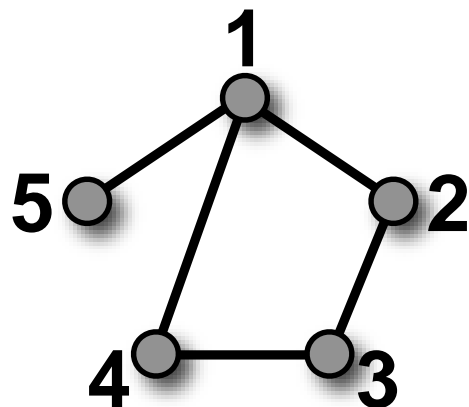
$$r = \begin{pmatrix} 1 \\ k_1 + 1 \\ k_2 + k_1 + 1 \\ \dots \\ 1 + \sum_{i=1}^j k_i \\ \dots \\ 2L + 1 \end{pmatrix}$$

2L elements

$$j = \begin{pmatrix} \dots \\ \dots \\ \dots \\ \dots \\ \dots \\ \dots \\ \dots \end{pmatrix} \begin{matrix} r(1) \\ r(2) - 1 \\ r(2) \\ r(3) - 1 \\ r(j) \\ r(j+1) - 1 \end{matrix} \left. \begin{matrix} k_1 \text{ neighbors of node } 1 \\ k_2 \text{ neighbors of node } 2 \\ k_j \text{ neighbors of node } j \end{matrix} \right\}$$

Other network representations

- Adjacency lists



N+1 elements

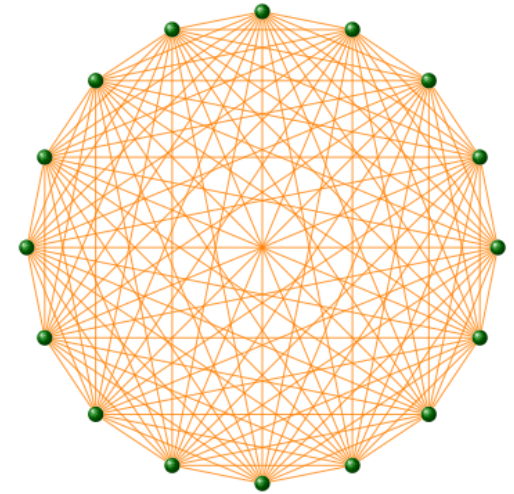
$$r = \begin{pmatrix} 1 \\ 4 \\ 6 \\ 8 \\ 10 \\ 11 \end{pmatrix}$$

2L elements

$$j = \begin{pmatrix} 2 \\ 4 \\ 5 \\ \text{---} \\ 1 \\ 3 \\ \text{---} \\ 2 \\ 4 \\ \text{---} \\ 3 \\ 1 \\ \text{---} \\ 1 \end{pmatrix}$$

The maximum number of links a network of N nodes can have is:

$$L_{\max} = \binom{N}{2} = \frac{N(N-1)}{2}$$



- A graph with degree $L=L_{\max}$ is called a **complete graph**
- Average degree is **$\langle k \rangle = N-1$**

Most networks observed in real systems are sparse:

$$L \ll L_{\max}$$

or

$$\langle k \rangle \ll N-1$$

WWW (ND Sample):	N=325,729;	$L=1.4 \cdot 10^6$	$L_{\max}=10^{12}$	$\langle k \rangle=4.51$
Protein (<i>S. Cerevisiae</i>):	N= 1,870;	$L=4,470$	$L_{\max}=10^7$	$\langle k \rangle=2.39$
Coauthorship (Math):	N= 70,975;	$L=2 \cdot 10^5$	$L_{\max}=3 \cdot 10^{10}$	$\langle k \rangle=3.9$
Movie Actors:	N=212,250;	$L=6 \cdot 10^6$	$L_{\max}=1.8 \cdot 10^{13}$	$\langle k \rangle=28.78$

(Albert, Barabasi, RMP2002)

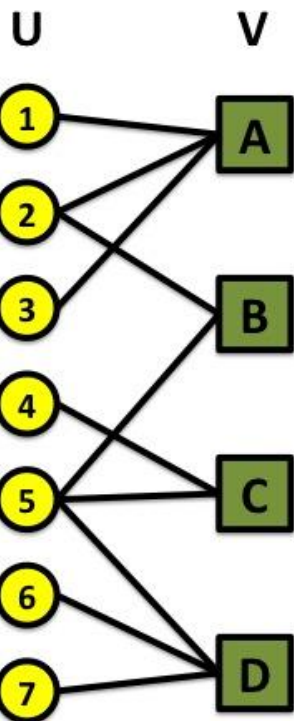
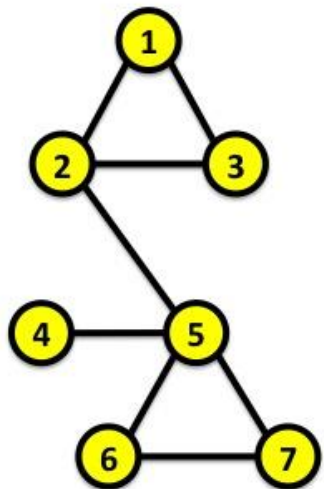
Each link (i, j) has a weight w_{ij}

$$\mathbf{A}_{ij} = \mathbf{w}_{ij}$$

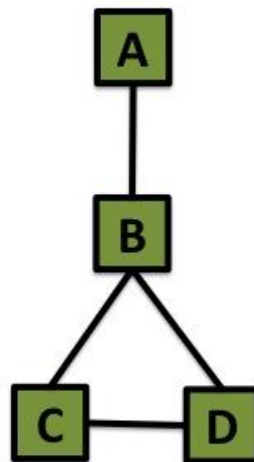
Bipartite graphs (or bigraph)

- Nodes can be divided into two disjoint sets U and V such that every link connects a node in U to one in V

Projection U



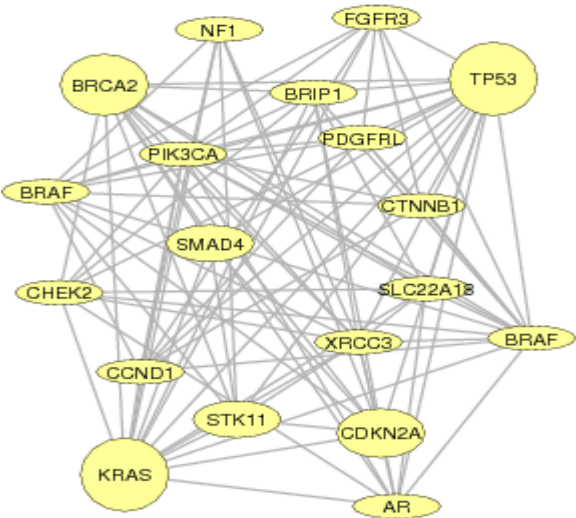
Projection V



Examples:

Hollywood actor network
Collaboration networks
Disease network

Example: GENE NETWORK – DISEASE NETWORK

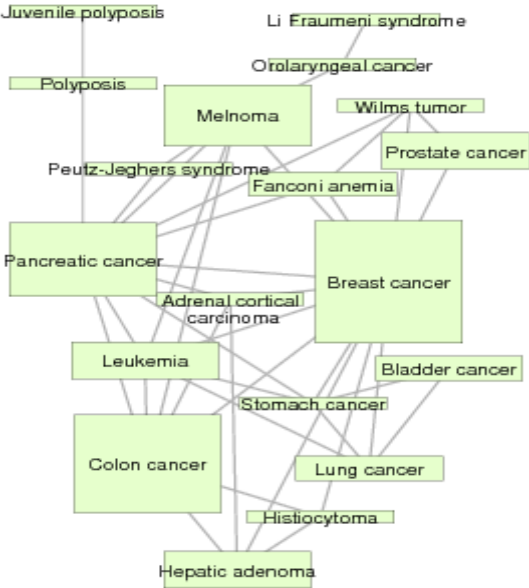
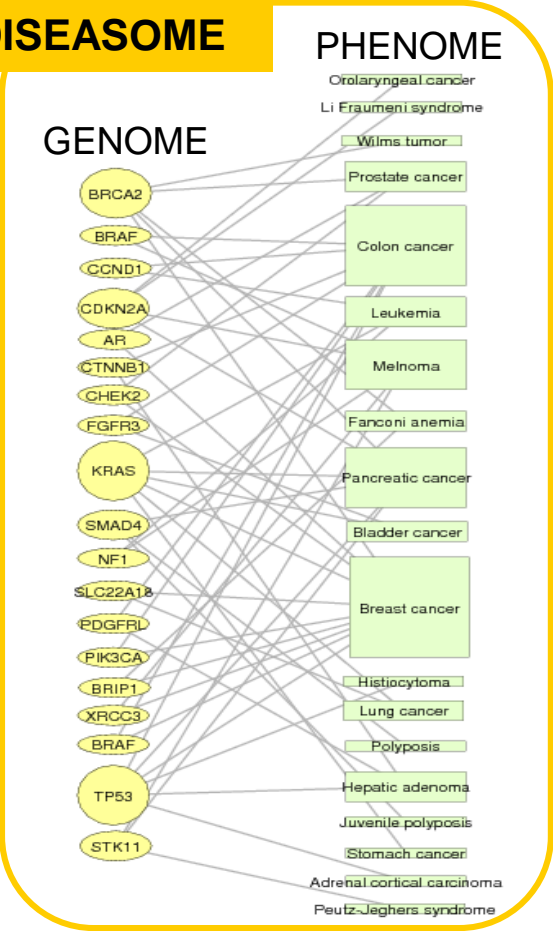


Gene network

DISEASOME

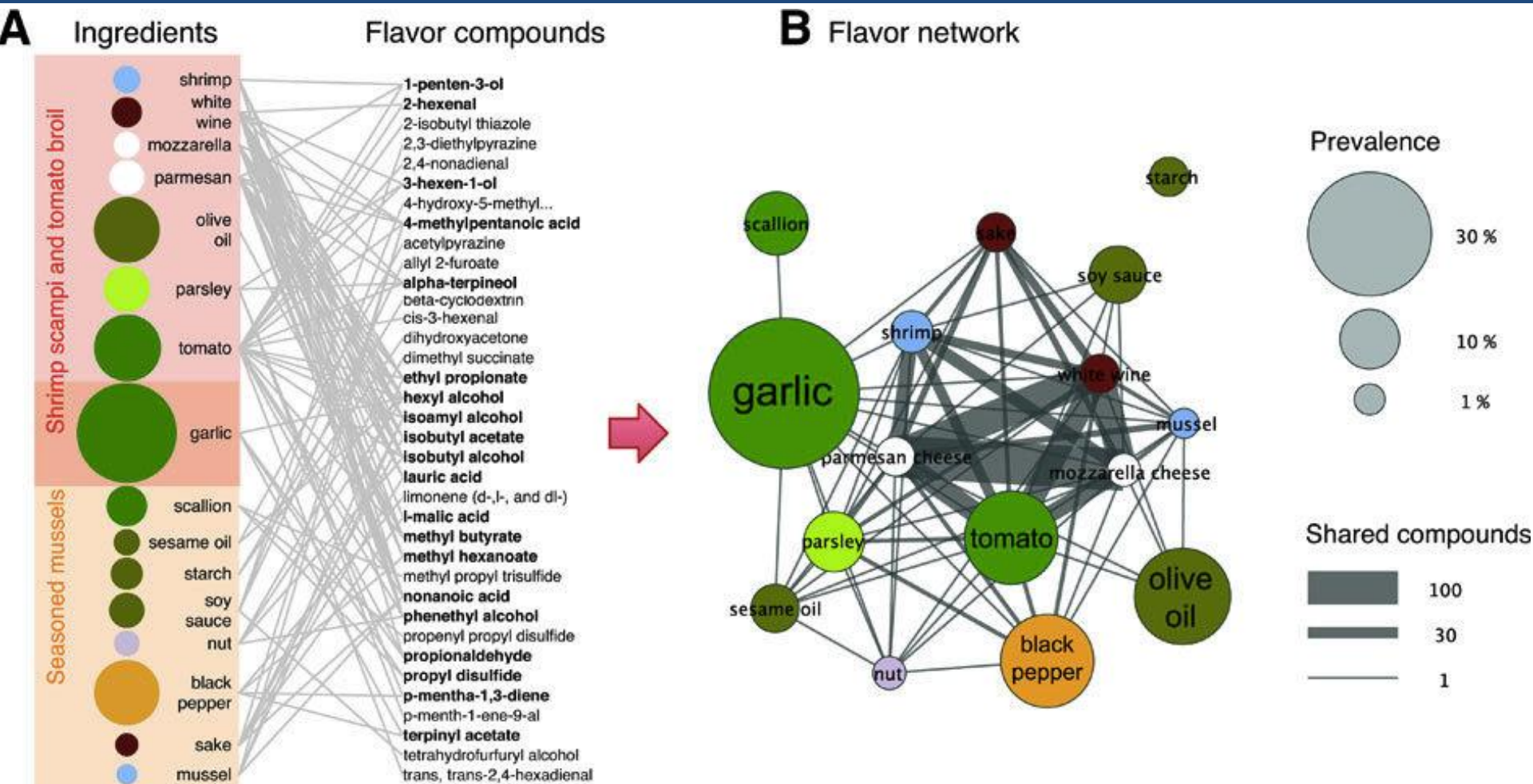
PHENOME

GENOME

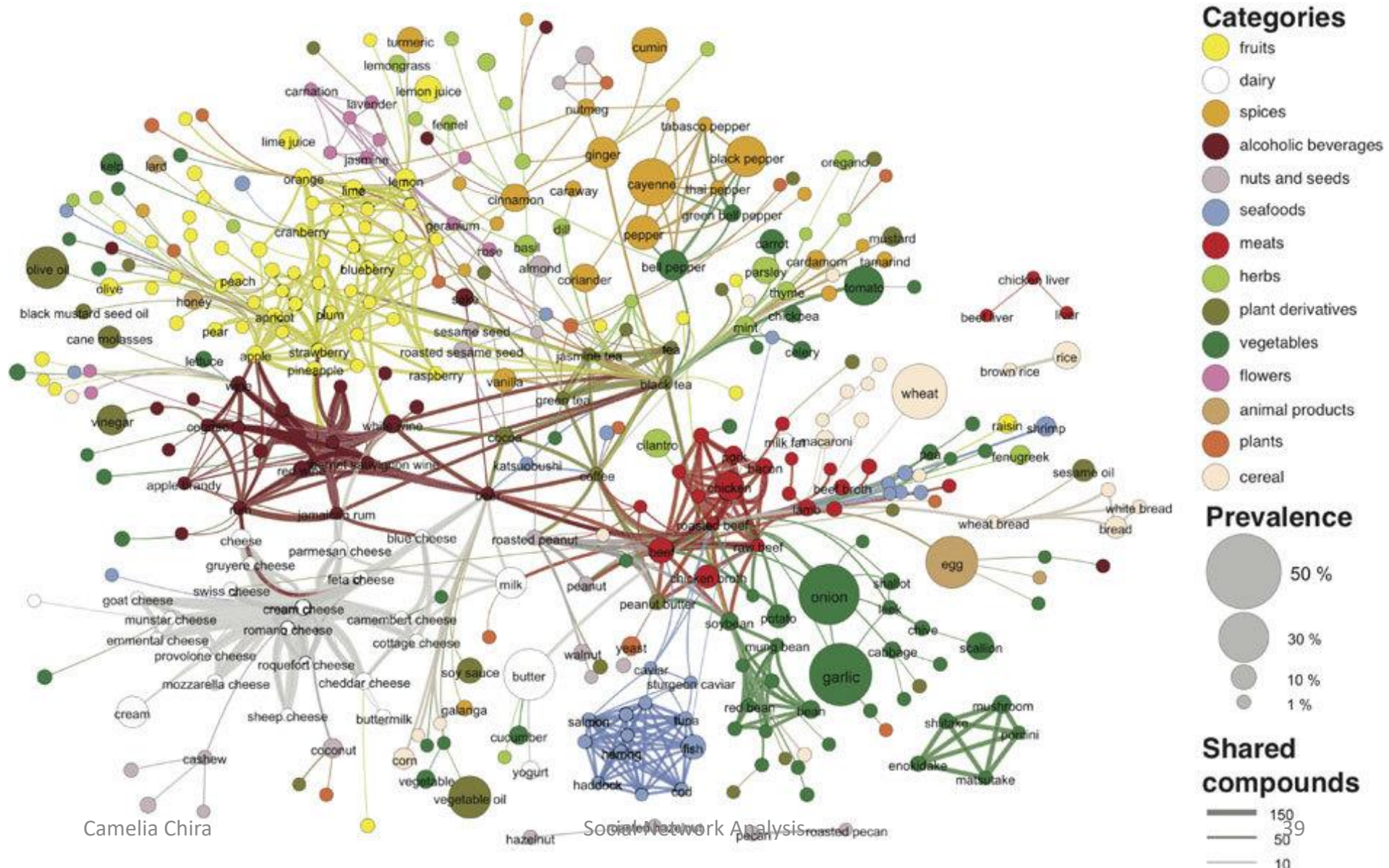


Disease network

Example: Ingredient-Flavor Bipartite Network



Y.-Y. Ahn, S. E. Ahnert, J. P. Bagrow, A.-L. Barabási, Flavor network and the principles of food pairing, Scientific Reports 196, (2011).



Definitions and properties

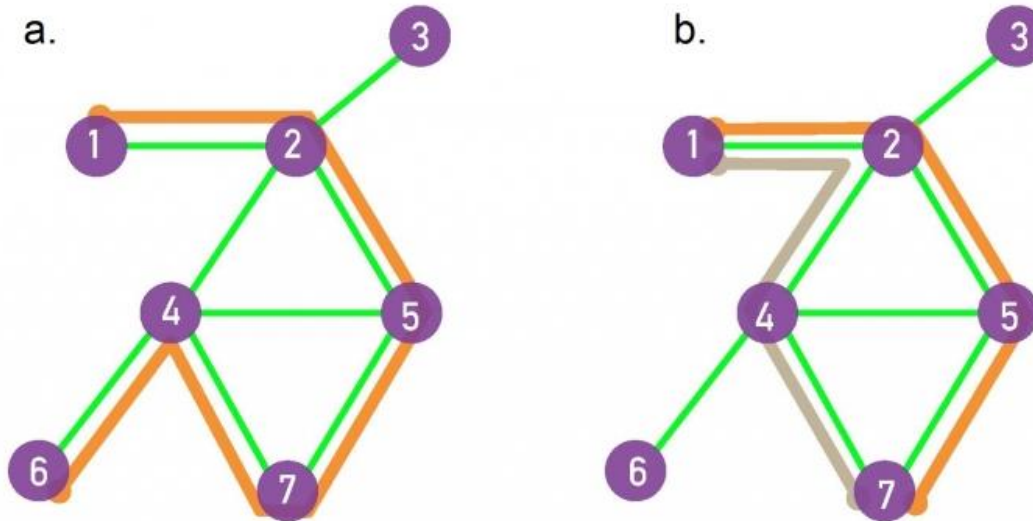
- Degrees
 - Node degree
 - Degree distribution
- Paths and distances
 - Distance (shortest path) between two nodes
 - Diameter of the network
 - Average distance (or average path length)
- Connectivity
 - Connected components
 - Clustering coefficient of a node
 - Average clustering coefficient of the network

Paths in a network

- A **path** is a sequence of nodes in which each node is adjacent to the next one

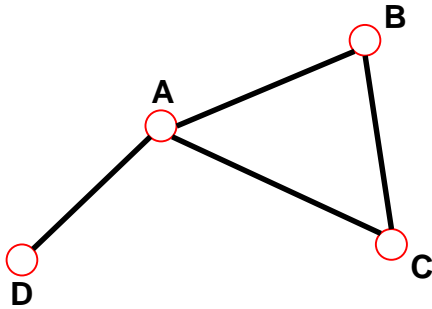
P_{i_0, i_n} of length n between nodes i_0 and i_n is an ordered collection of $n+1$ nodes and n links

$$P_n = \{i_0, i_1, i_2, \dots, i_n\} \quad P_n = \{(i_0, i_1), (i_1, i_2), (i_2, i_3), \dots, (i_{n-1}, i_n)\}$$



- In a directed network, the path can follow only the direction of an arrow.

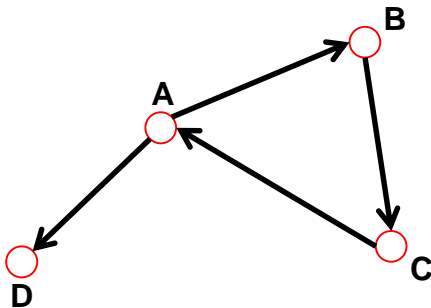
Distance between two nodes



$d(u,v)$ = shortest path between u and v

The *distance (shortest path, geodesic path)* between two nodes is defined as the *number of edges along the shortest path* connecting them.

➤ *If the two nodes are disconnected, the distance is infinity.*



In **directed graphs** each path needs to follow the direction of the arrows.

Number of paths between two nodes

N_{ij} , number of paths between any two nodes i and j :

Length $n=1$: If there is a link between i and j , then $A_{ij}=1$ and $A_{ij}=0$ otherwise.

Length $n=2$: If there is a path of length two between i and j , then $A_{ik}A_{kj}=1$, and $A_{ik}A_{kj}=0$ otherwise.

The number of paths of length 2:

$$N_{ij}^{(2)} = \sum_{k=1}^N A_{ik} A_{kj} = [A^2]_{ij}$$

Length n : In general, if there is a path of length n between i and j , then $A_{ik} \dots A_{lj}=1$ and $A_{ik} \dots A_{lj}=0$ otherwise.

The number of paths of length n between i and j is*

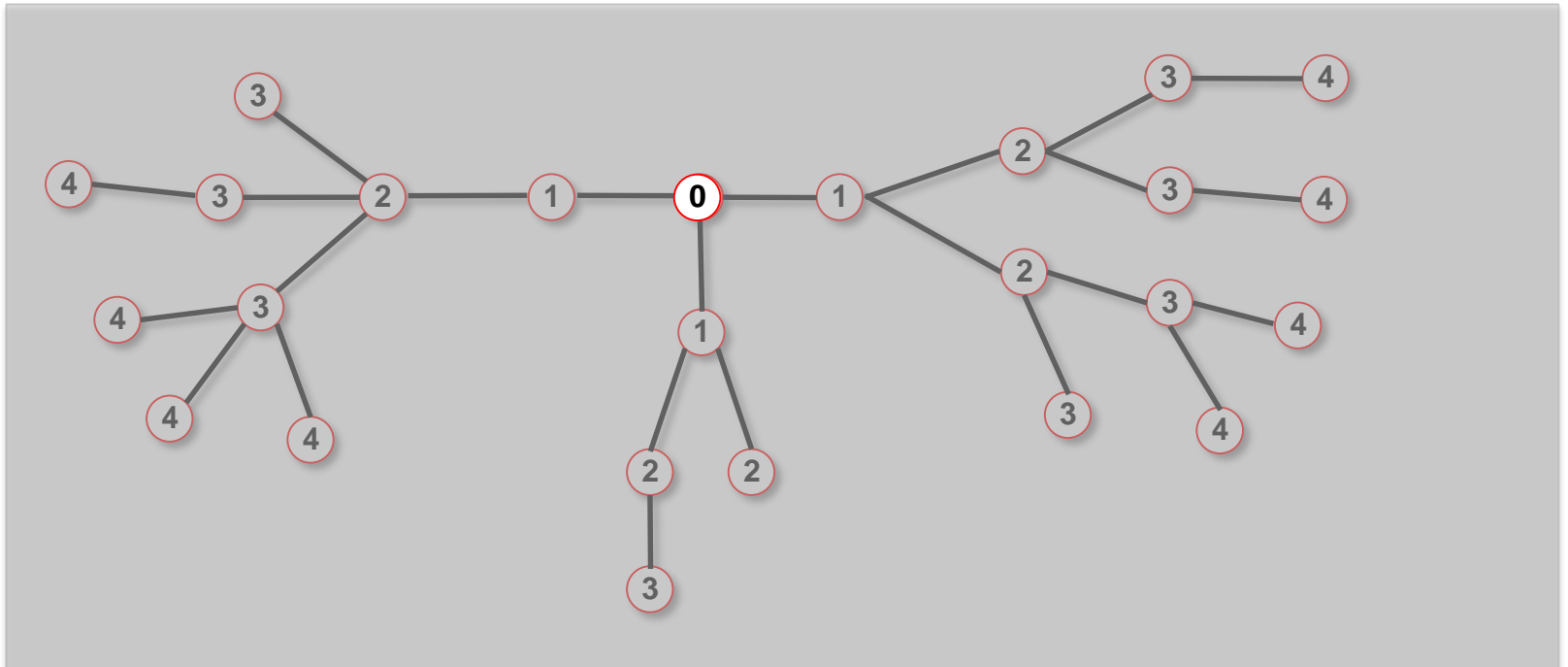
$$N_{ij}^{(n)} = [A^n]_{ij}$$

* holds for both directed and undirected networks.

Finding distances: Breadth First Search

Distance between node 0 and node 4:

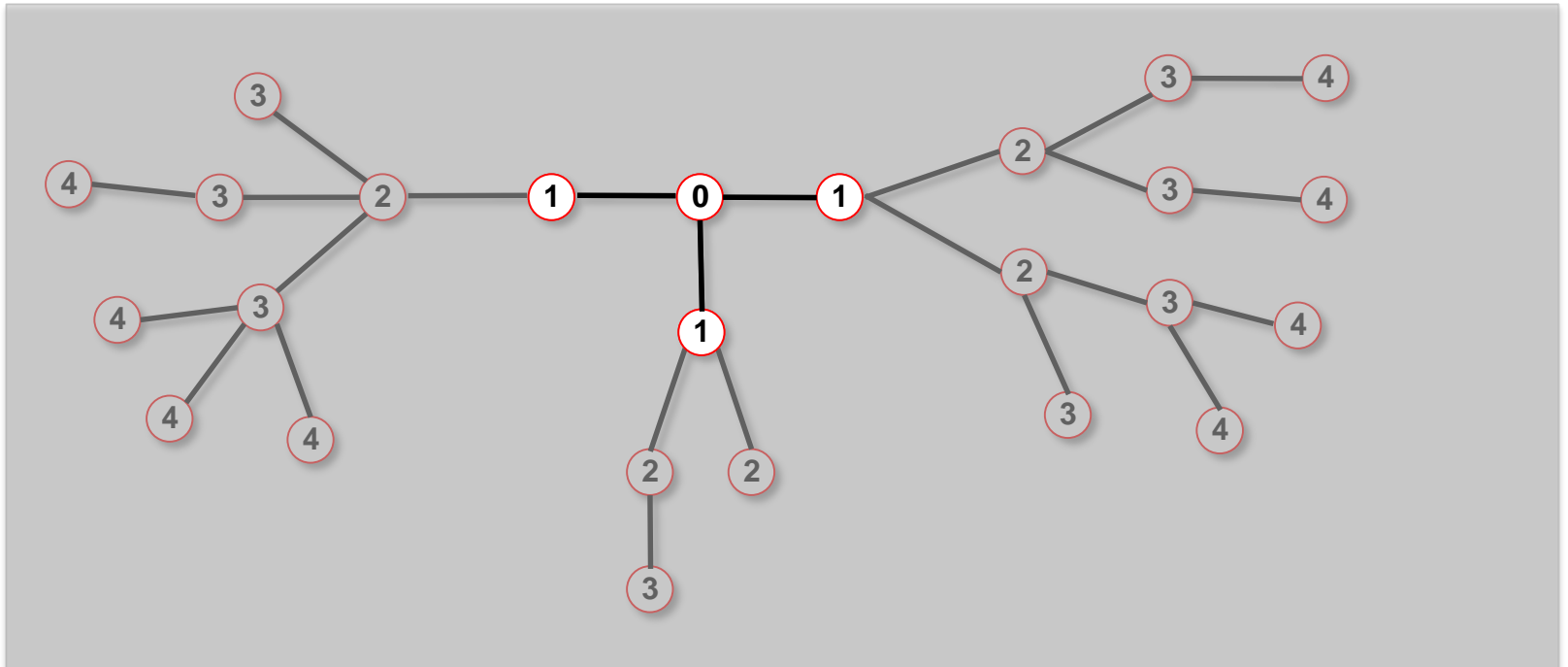
1. Start at 0.



Finding distances: Breadth First Search

Distance between node 0 and node 4:

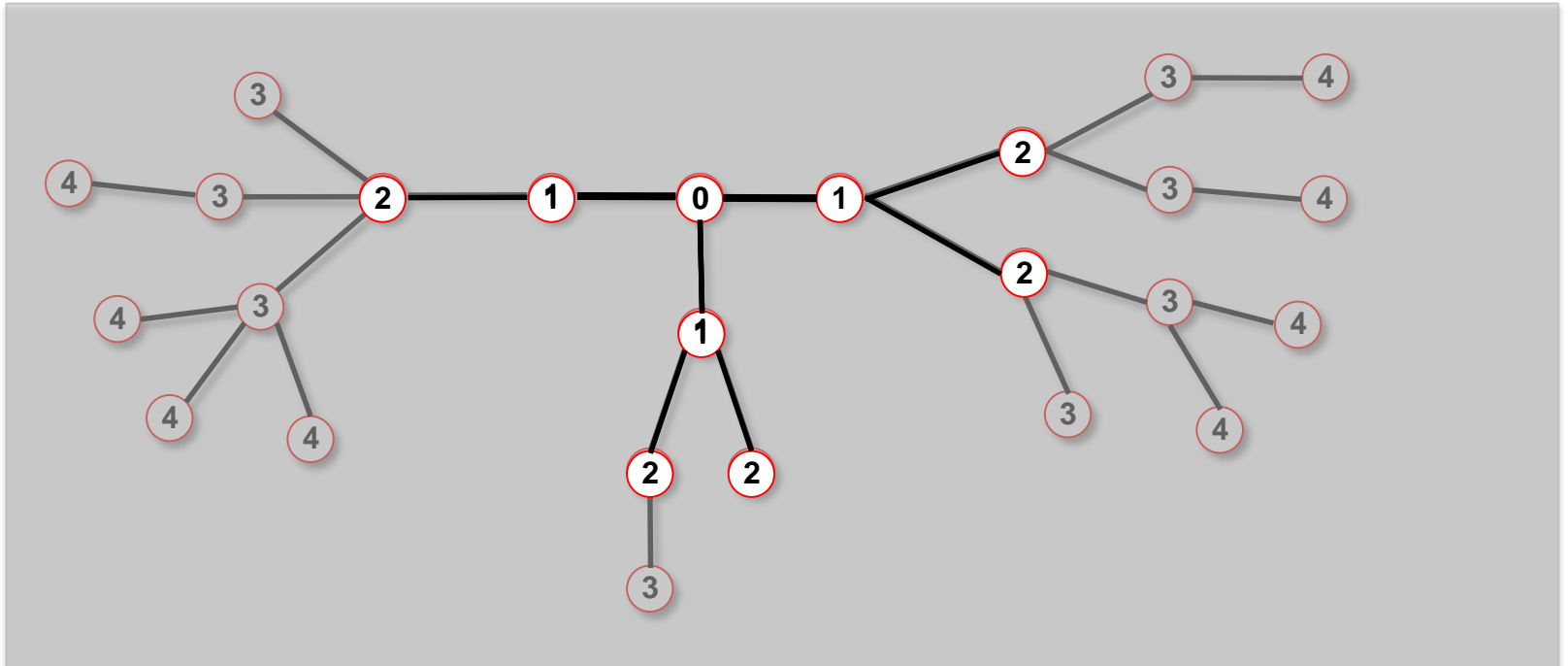
1. Start at 0.
2. Find the nodes adjacent to 1. Mark them as at distance 1. Put them in a queue.



Finding distances: Breadth First Search

Distance between node 0 and node 4:

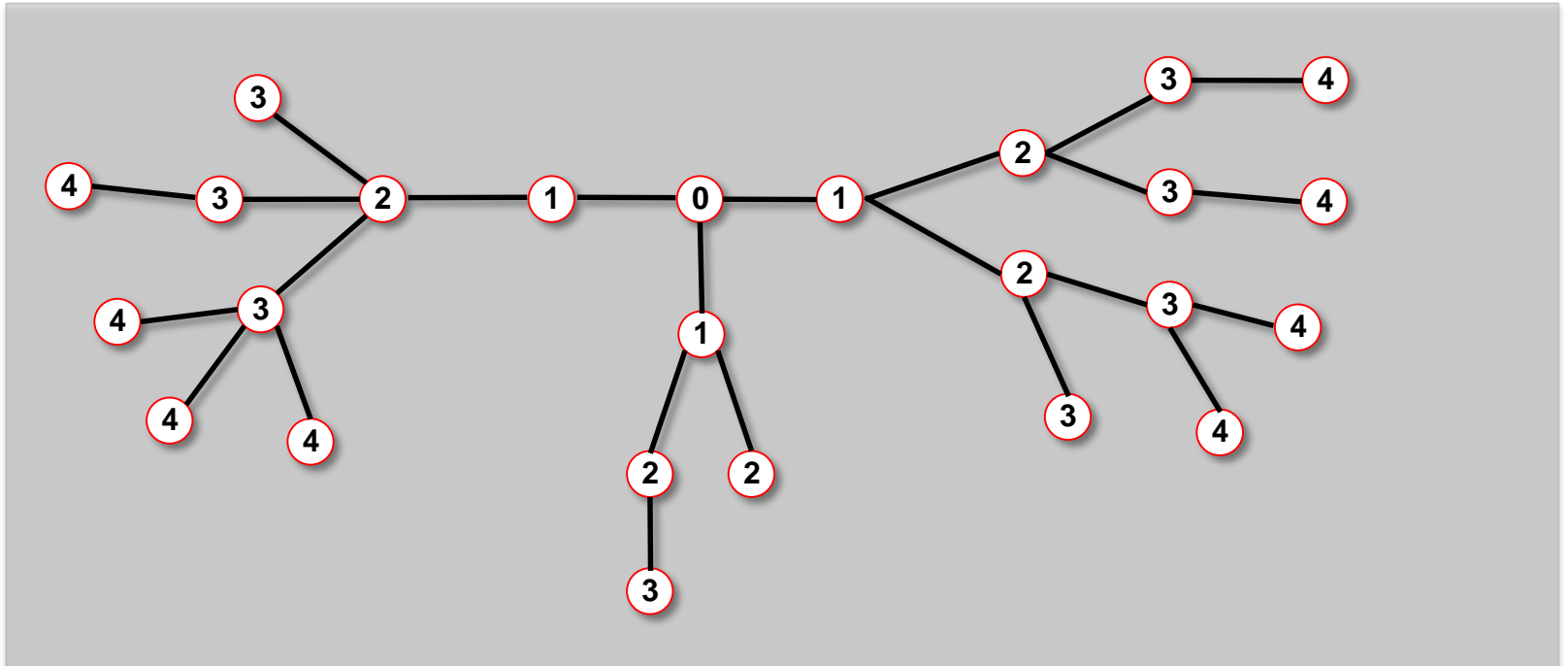
1. Start at 0.
2. Find the nodes adjacent to 0. Mark them as at distance 1. Put them in a queue.
3. Take the first node out of the queue. Find the unmarked nodes adjacent to it in the graph. Mark them with the label of 2. Put them in the queue.



Finding distances: Breadth First Search

Distance between node 0 and node 4:

- 1.Repeat until you find node 4 or there are no more nodes in the queue.
- 2.The distance between 0 and 4 is the label of 4 or, if 4 does not have a label, infinity.



Network diameter

Diameter (d_{max}) = the largest distance between any pair of nodes in the network

Average distance

Average path length/distance, $\langle d \rangle$, for a **connected graph**:

$$\langle d \rangle = \frac{1}{2L_{max}} \sum_{i,j \neq i} d_{ij}$$

where d_{ij} is the distance from node i to node j

$$L_{max} = \frac{N(N-1)}{2}$$

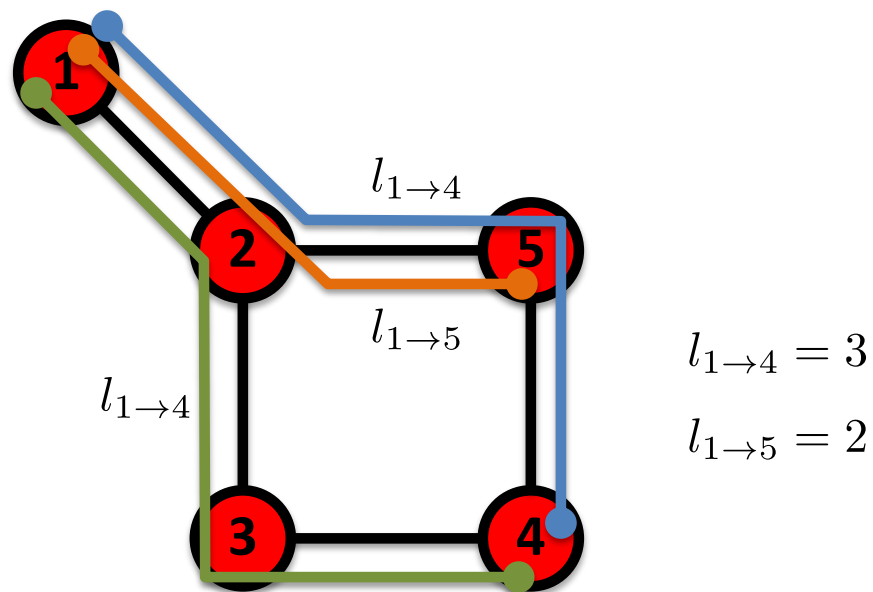
In an *undirected graph* $d_{ij} = d_{ji}$, so we only need to count them once:

$$\langle d \rangle = \frac{1}{L_{max}} \sum_{i,j > i} d_{ij}$$

Average path length for real networks

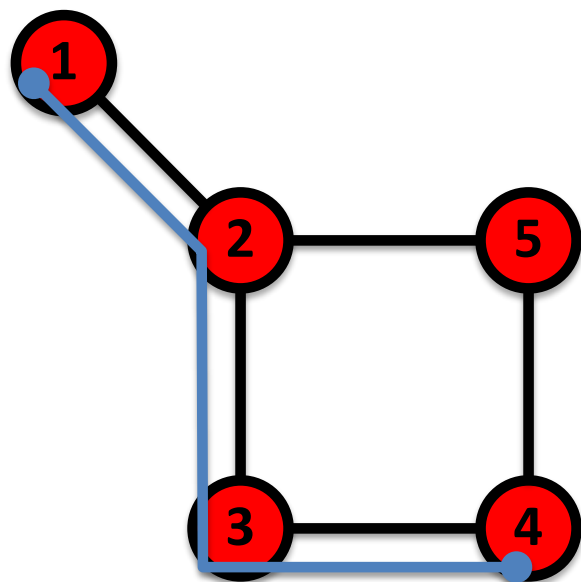
- Assumption: a single connected component
- Compute distance between any two nodes d_{ij}
- **Average path length = avg. over all d_{ij}**
 - Smallest possible average path length? Largest ?
 - ✓ *Should be compared to N (small $\log N$)*
- Real networks: “small” average path length
 - Small is not precisely defined, clearly $\ll N$
 - “Six degrees of separation” (see next lecture)
 - **Examples:**
 - 6, $N \approx 200M$ (Milgram experiment, 1969)
 - 6.5, $N \approx 180M$ (Microsoft messenger network, 2008)
 - 5, $N \approx 721M$ (Facebook social graph, 2012)

Shortest Path (distance)



The path with the shortest length between two nodes (distance).

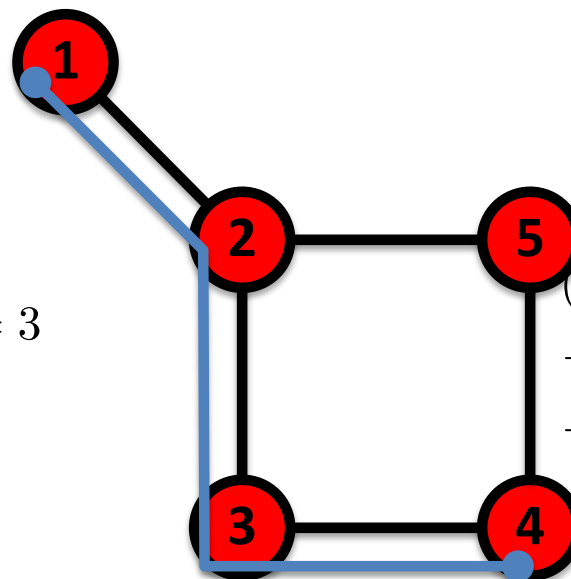
Diameter



The longest shortest path in a graph.

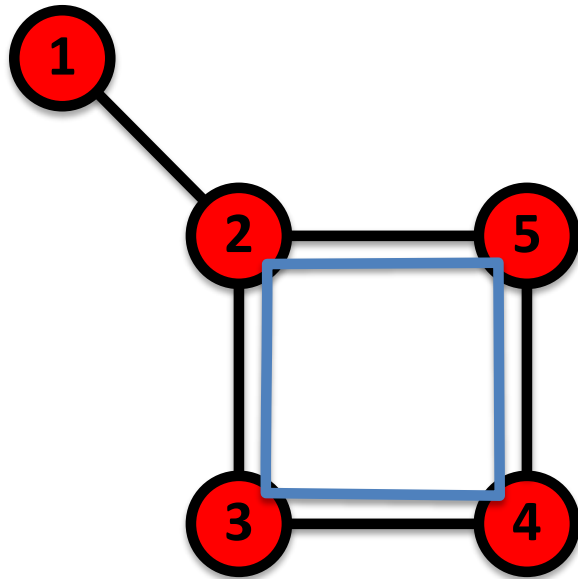
Average Path Length or Average Distance

$$l_{1 \rightarrow 4} = 3$$



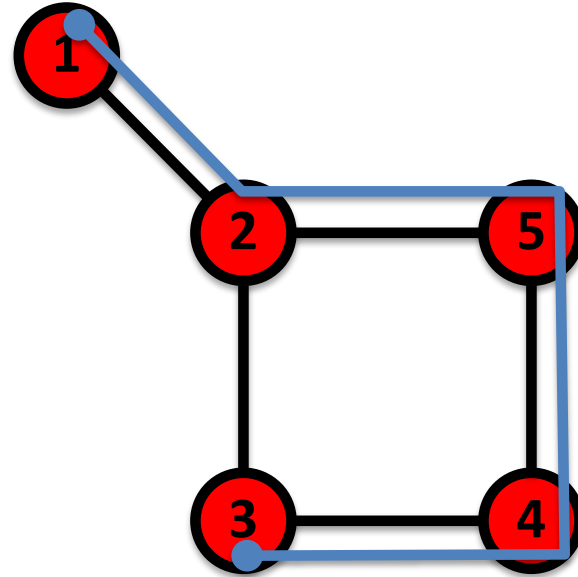
The average of the shortest paths for all pairs of nodes.

Cycle



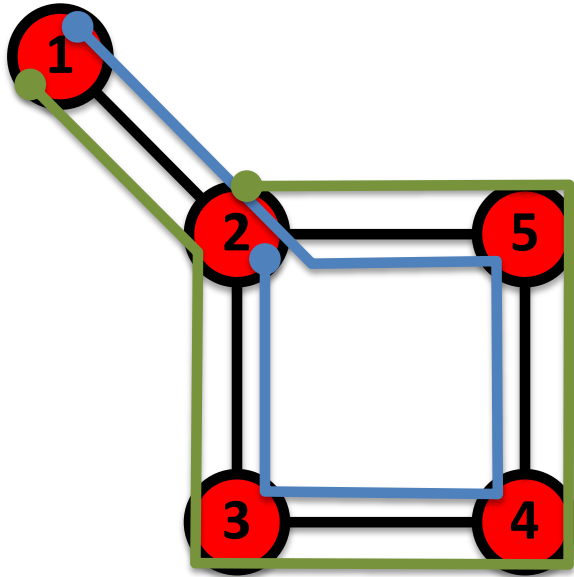
A path with the same start and end node.

Self-avoiding Path



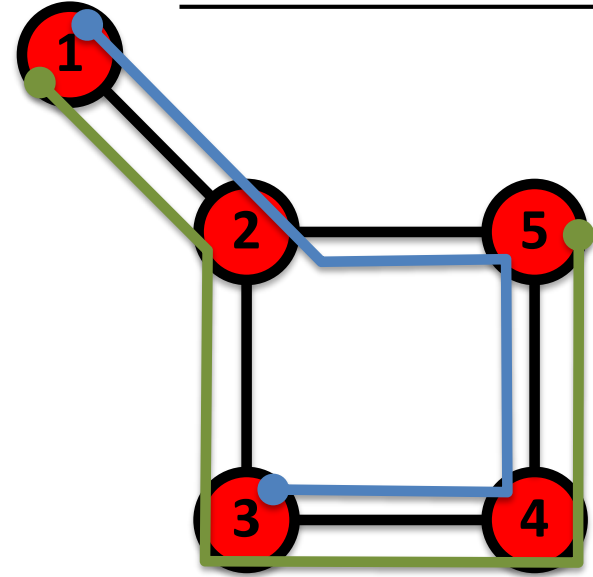
A path that does not intersect itself.

Eulerian Path



A path that traverses each link exactly once.

Hamiltonian Path



A path that visits each node exactly once.

Camelia Chira

<http://cs.ubbcluj.ro/~cchira>

Content based on the textbook:

A.-L. Barabási, Network Science, Cambridge University Press, 2016.

<http://networksciencebook.com/>

References

- A.-L. Barabási, Network Science, Cambridge University Press, 2016. (Chapter 2)
- D. J. Watts and S. H. Strogatz. Collective dynamics of ‘small-world’ networks. Nature, 393:440–442, 1998.
- D. J. Watts, Six Degrees: The Science of a Connected Age, W.W. Norton, 2003.