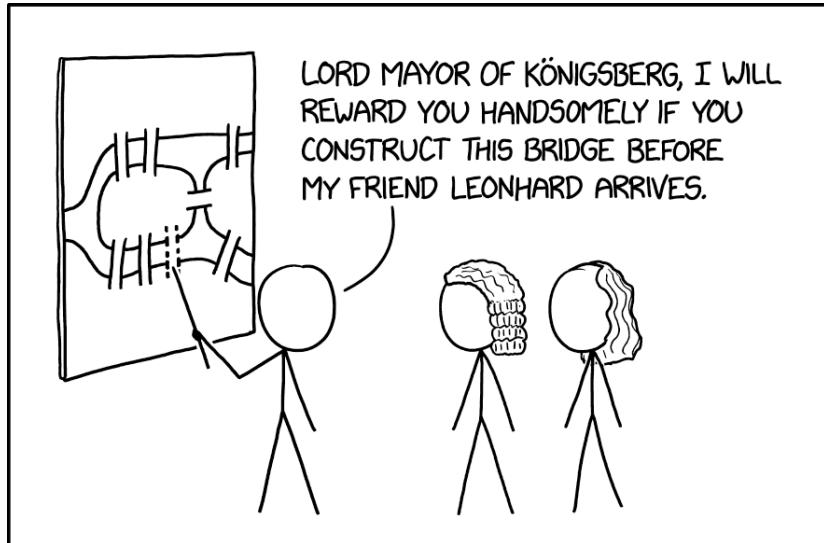


DS8104: Network Science

Class 2: First Connections

a.k.a. graph theory 101

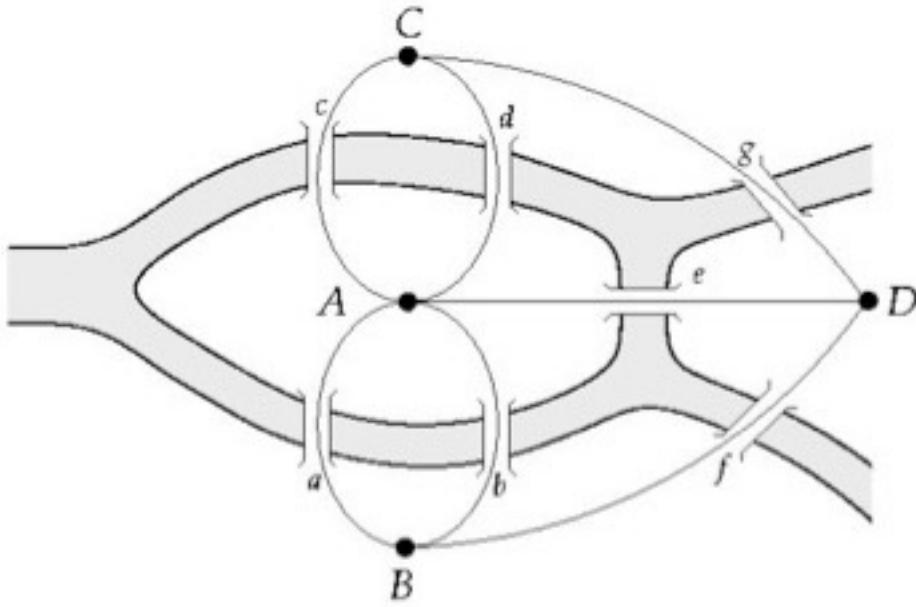


Bridges of Konigsberg



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

Bridges of Konigsberg

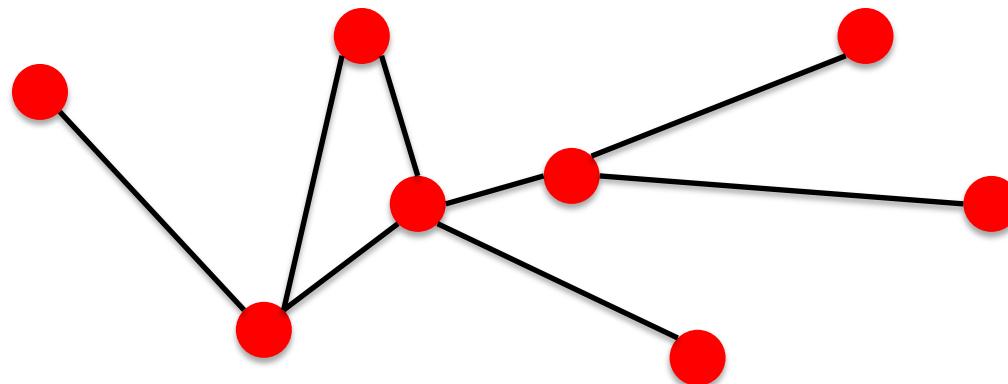


Can one walk across the seven bridges and never cross the same bridge twice?

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.

Basic building blocks



components: nodes, vertices, elements

N

relationships: links, edges, interactions

L

system: network, graph

G:(N,L)

network often refers to real systems

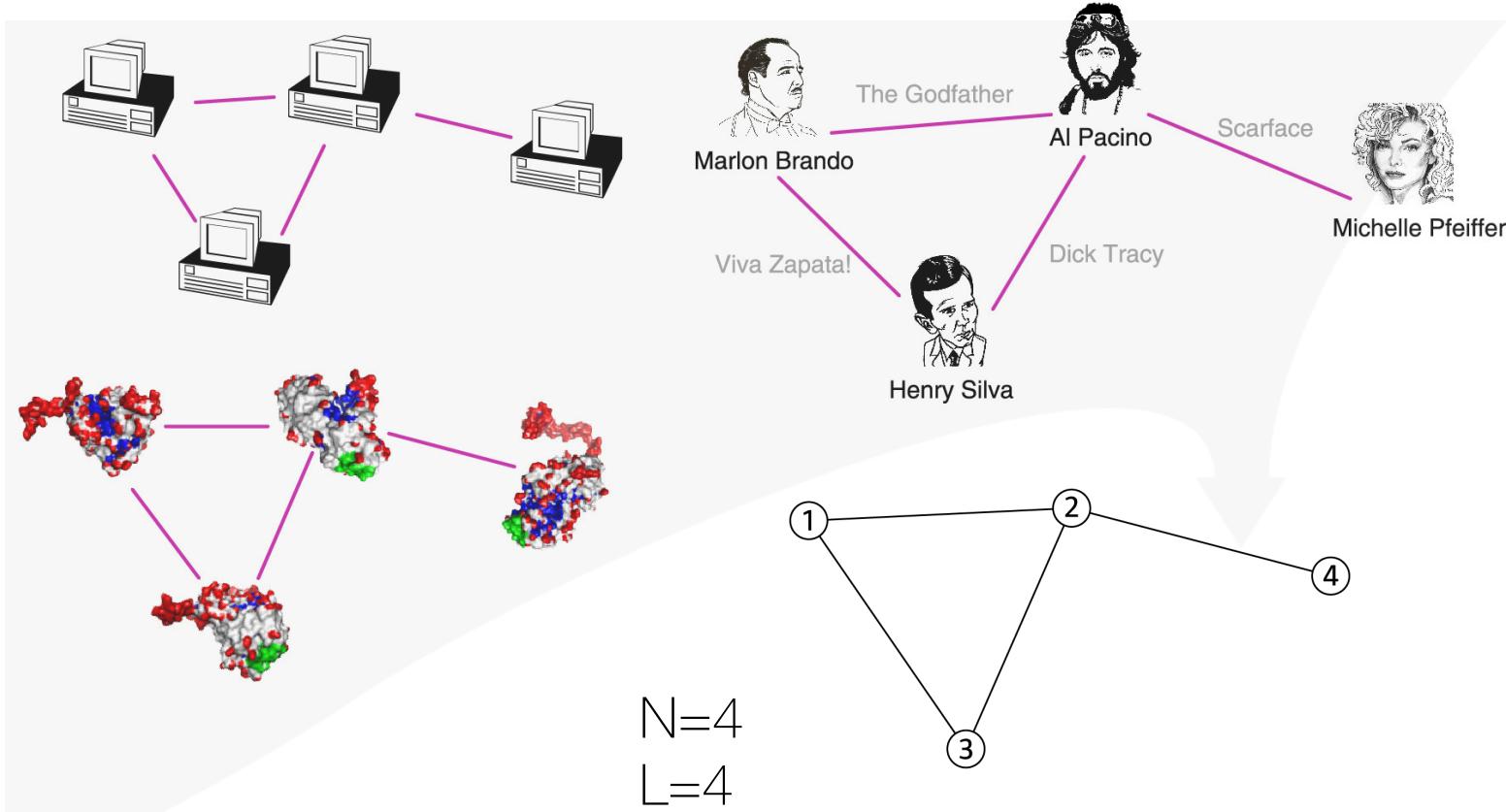
Language: (Network, node, link)

graph refers to the mathematical representation of a network

Language: (Graph, vertex, edge)

We will try to make this distinction whenever it is appropriate, but in most cases we will use the two terms interchangeably.

Networks provide a general framework to capture the pattern underlying a system's relationships



Whats my network?

One aspect of network science is to choose the appropriate set of nodes and relationships.

In some cases there is a unique, unambiguous representation.

In other cases, the representation is by no means unique.

For example, the way we assign the links between a group of individuals will determine the nature of the question we can study.

PERSPECTIVE

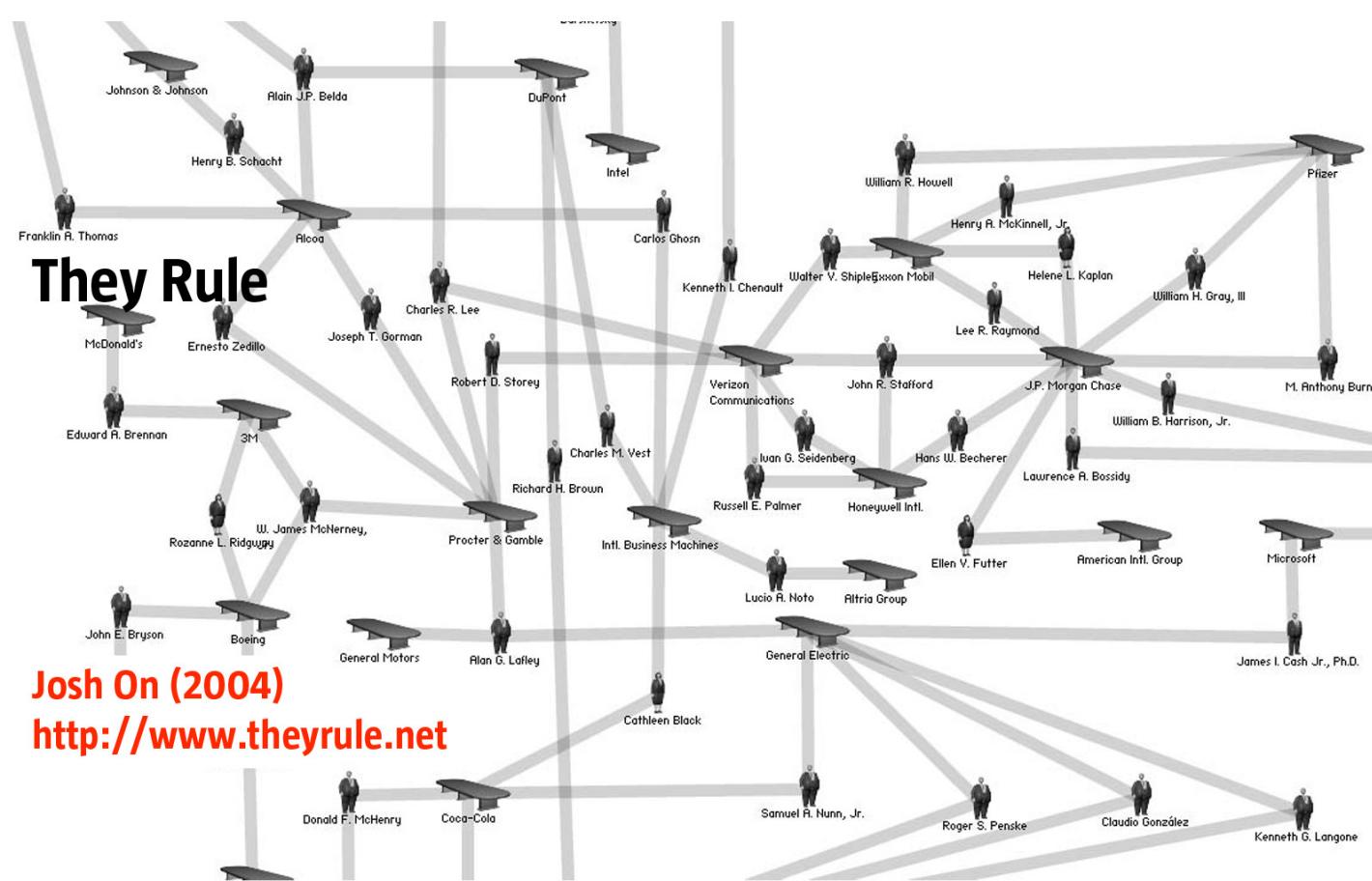
Revisiting the Foundations of Network Analysis

Carter T. Butts

Network analysis has emerged as a powerful way of studying phenomena as diverse as interpersonal interaction, connections among neurons, and the structure of the Internet. Appropriate use of network analysis depends, however, on choosing the right network representation for the problem at hand.

24 JULY 2009 VOL 325 SCIENCE www.sciencemag.org

One reading this week provides a deeper perspective on this issue.



They Rule

Josh On (2004)

<http://www.theyrule.net>

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

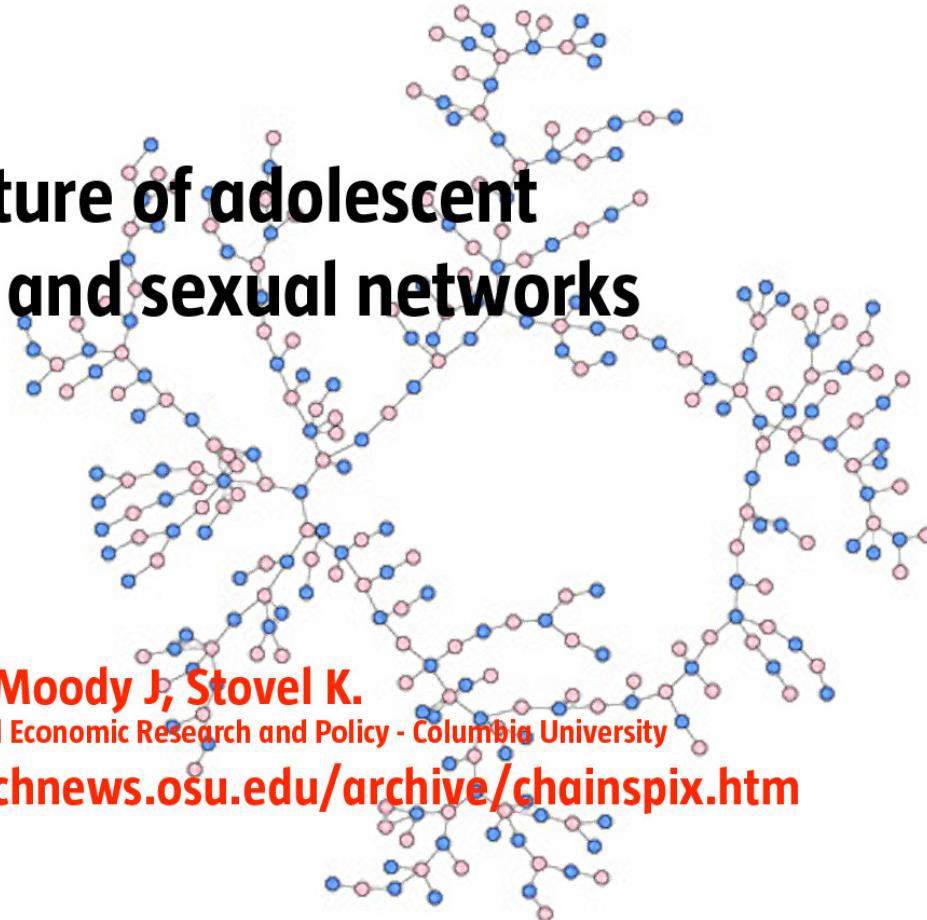
If you connect those that have a romantic and sexual relationship,
you will be exploring their **sexual network**.

The structure of adolescent romantic and 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>



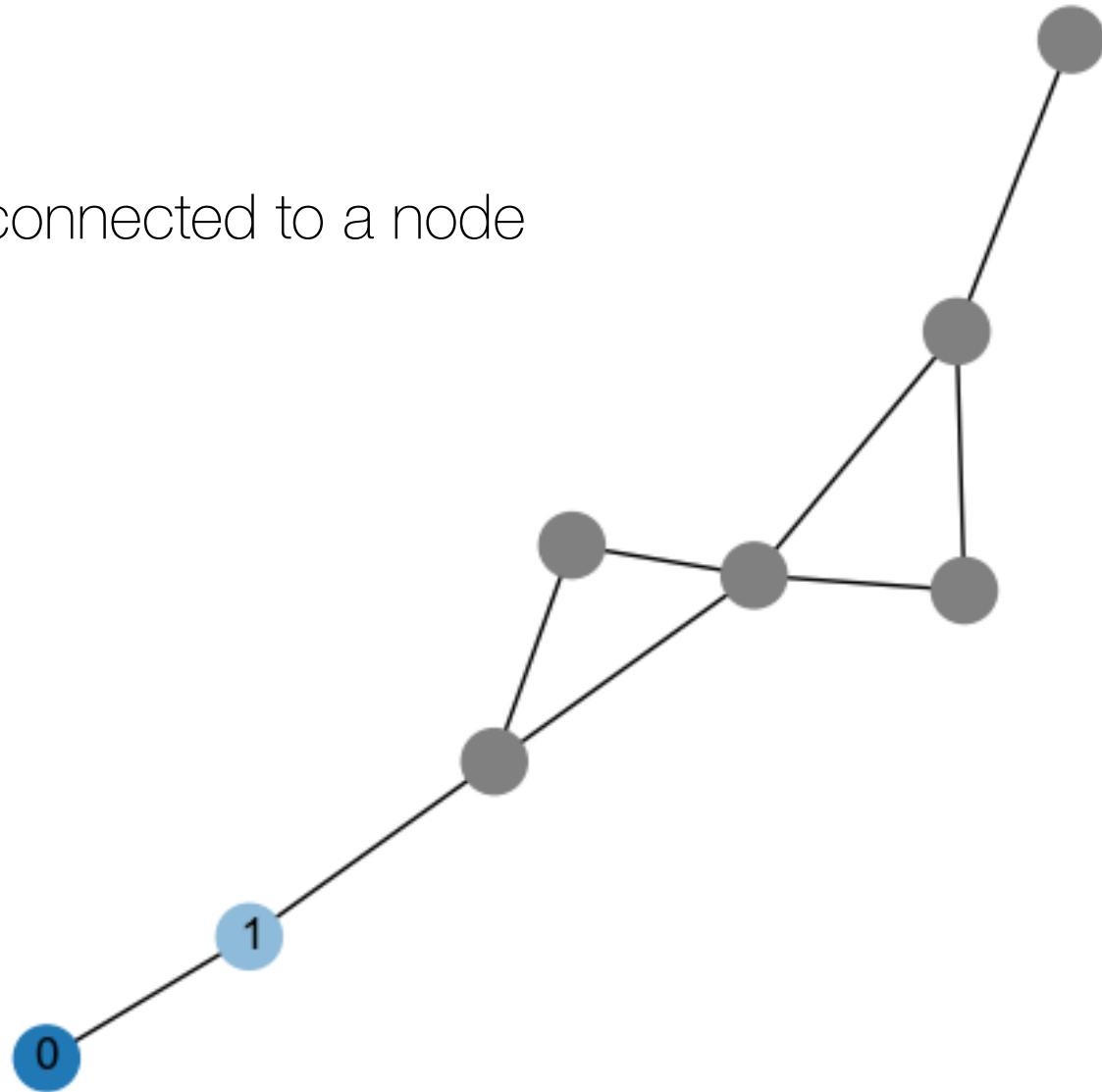
If you connect individuals based on their first name (all Peters connected to each other), you will be exploring what?

But this is a network, nevertheless.

Node degree

the number of links connected to a node

$$k_0 = 1$$

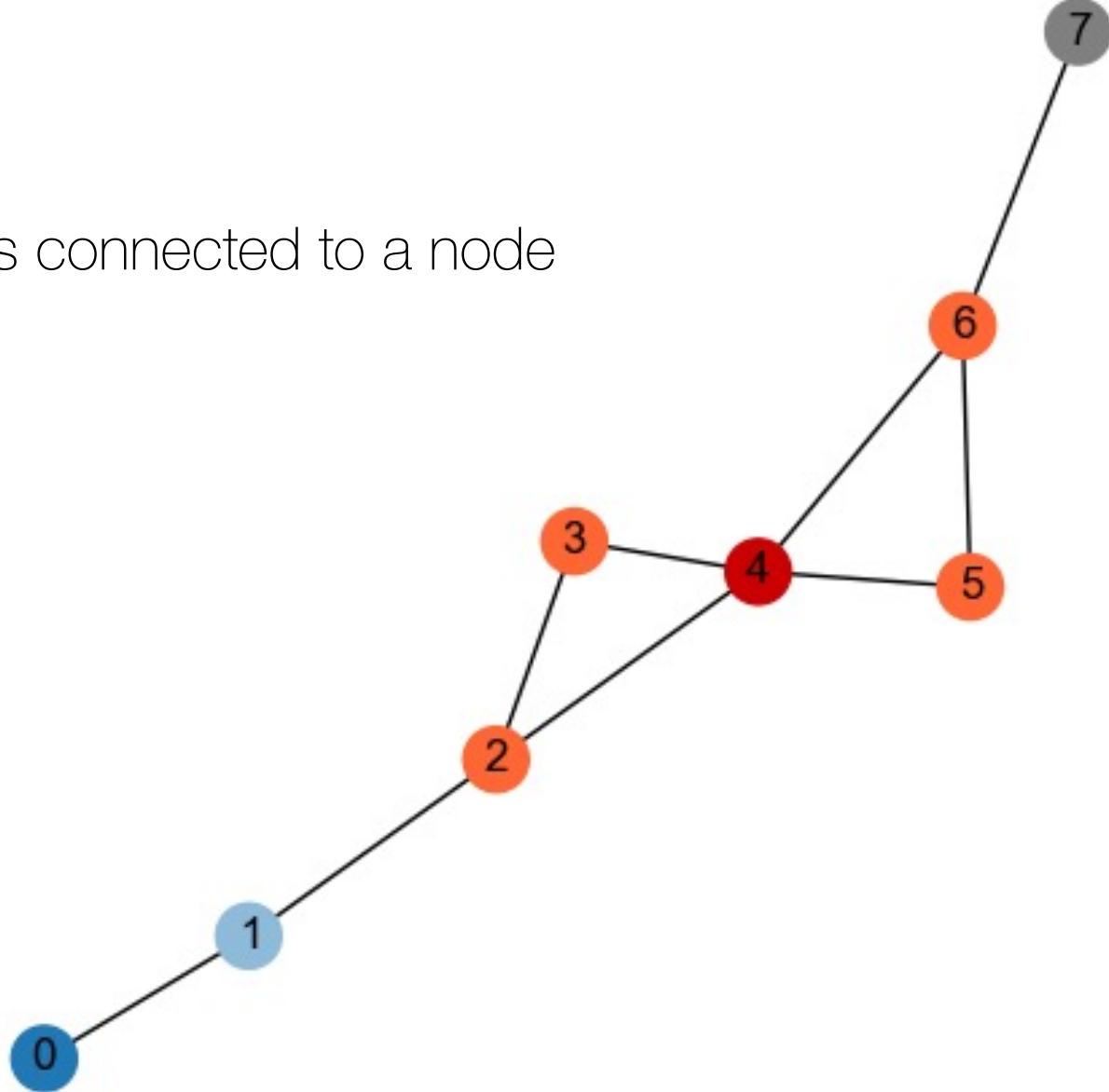


Node degree

the number of links connected to a node

$$k_0 = 1$$

$$k_4 = 4$$

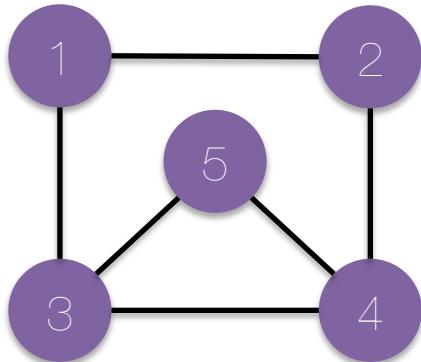


Average degree

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

$$\langle k \rangle \equiv \frac{2L}{N}$$

Graph Representations



Adjacency
Matrix

	0	1	2	3	4
0	0	1	1	0	0
1	1	0	0	1	1
2	1	0	0	0	1
3	0	1	0	0	1
4	0	1	1	1	0

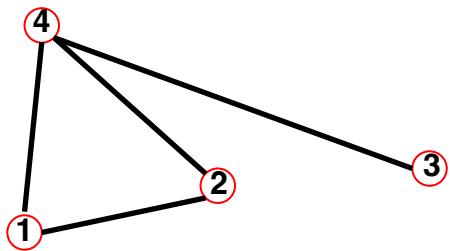
source target

0	1
0	2
1	3
1	4
2	4
3	4

Edge List

Adjacency
List

```
node : list of neighbors
{
    "0": [1, 2],
    "1": [0, 3, 4],
    "2": [0, 4],
    "3": [1, 4],
    "4": [1, 2, 3]
}
```



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

$$\begin{aligned} A_{ij} &= A_{ji} \\ A_{ii} &= 0 \end{aligned}$$

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

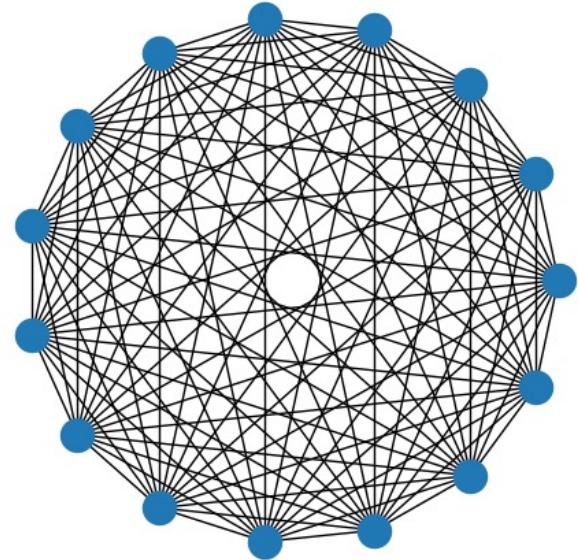
$$k_j = \sum_{i=1}^N A_{ij}$$

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

Dense graphs

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 links $L=L_{\max}$ is called a **complete graph**

$$\langle k \rangle = k_i = N - 1$$

Most observed networks have far fewer links (they are **sparse**)

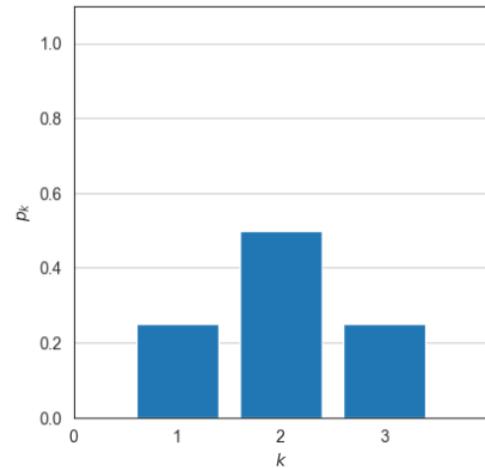
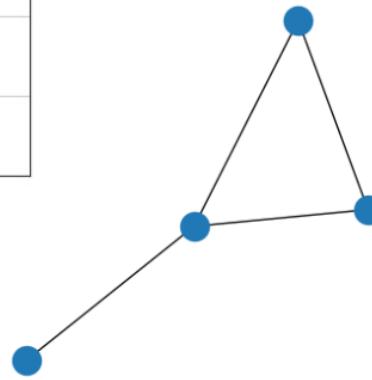
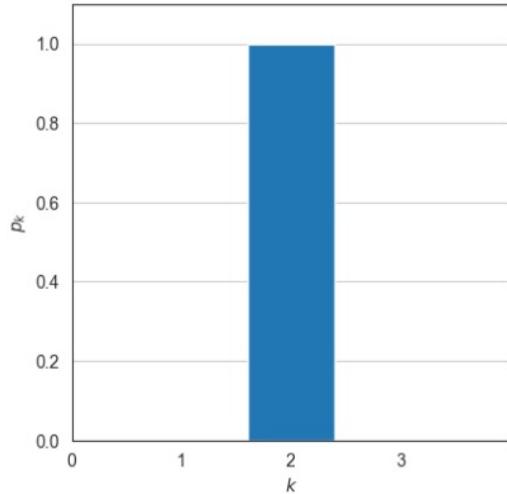
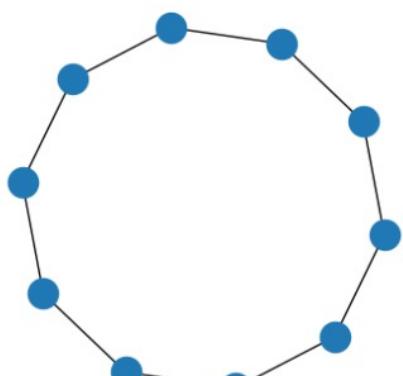
$$\langle k \rangle \ll N - 1$$

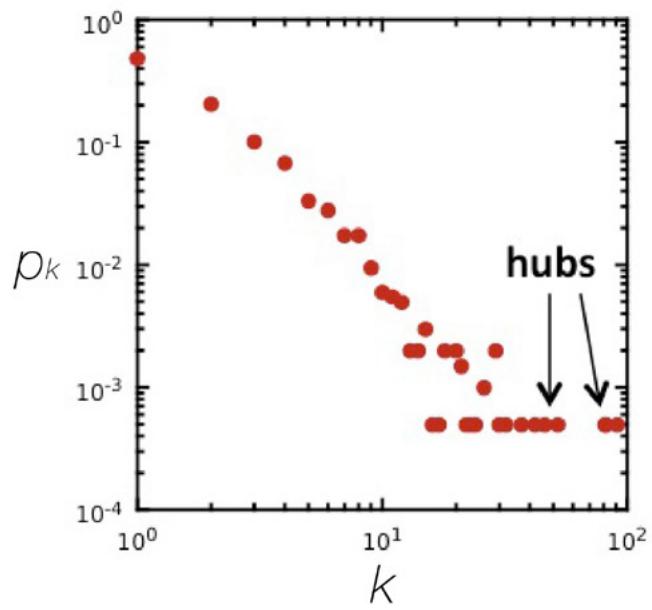
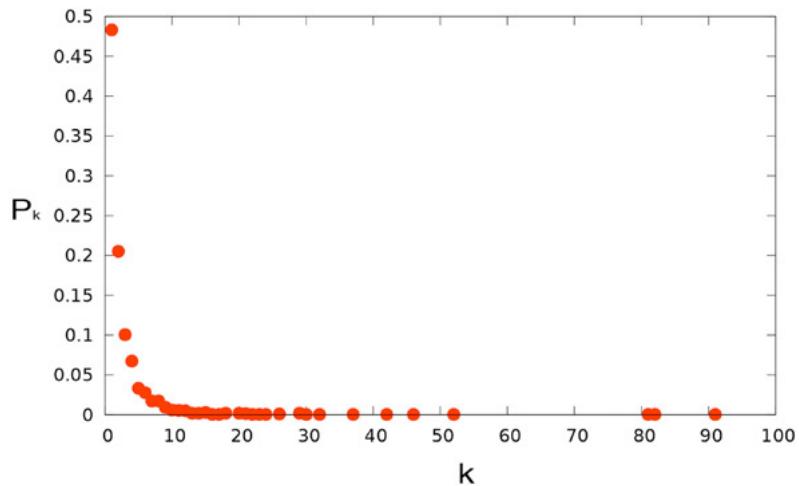
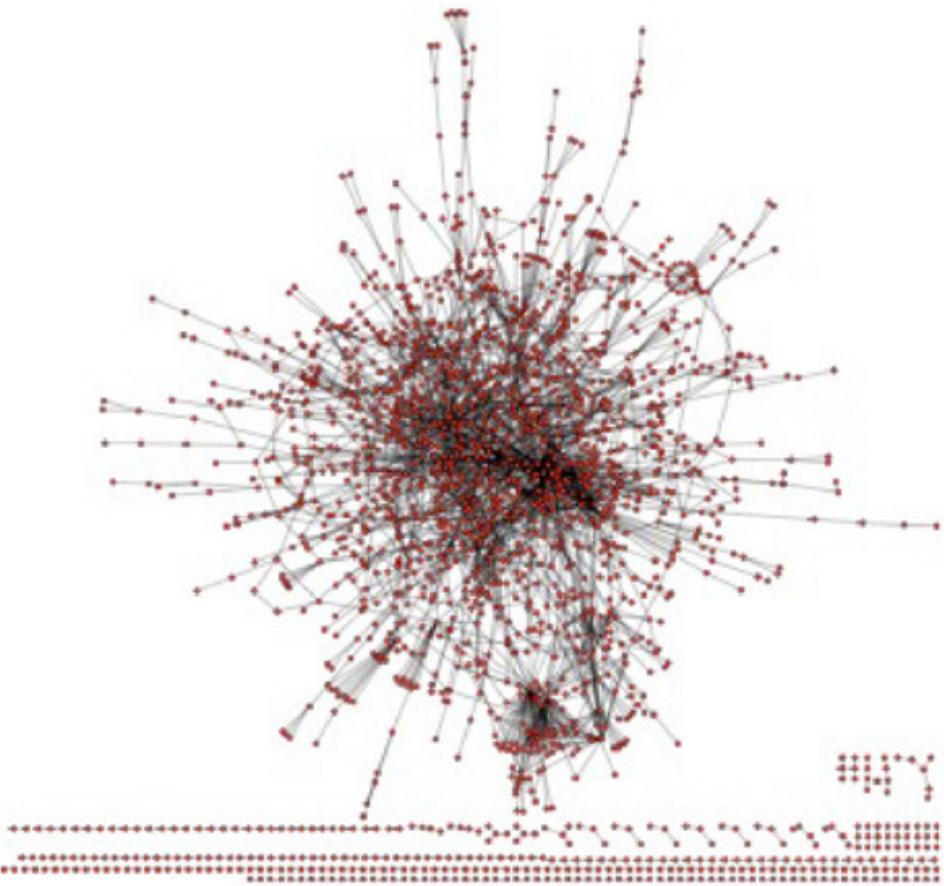
Network	N	L	$\langle k \rangle$
UVA Facebook	17196	789321	91.8
Protien Interaction	2018	2930	2.9
IMDB Co-acting	702388	29397908	83.71
Powergrid	4941	6594	2.67
E coli. metabolism	1039	4741	9.13

Real networks are sparse

Degree distribution

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





Graph Types

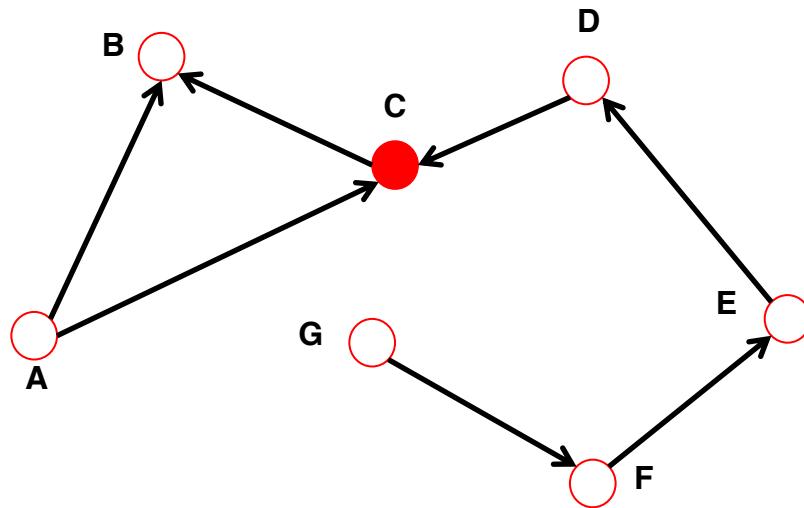
Directed graphs

So far the relationships we have considered are bi-directional (a.k.a. symmetric)

But many relationships are uni-directional

Alex follows Beyonce
(but Beyonce doesn't follow Alex....yet)

In directed networks we can define an **in-degree** and **out-degree**. The (total) degree is the sum of in- and out-degree.



$$k_C^{in} = 2$$

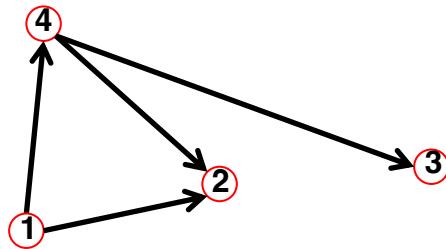
$$k_C^{out} = 1 \quad \langle k^{in} \rangle = \langle k^{out} \rangle$$

$$k_C = 3$$

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

Source: a node with $k^{in} = 0$;

Sink: a node with $k^{out} = 0$.



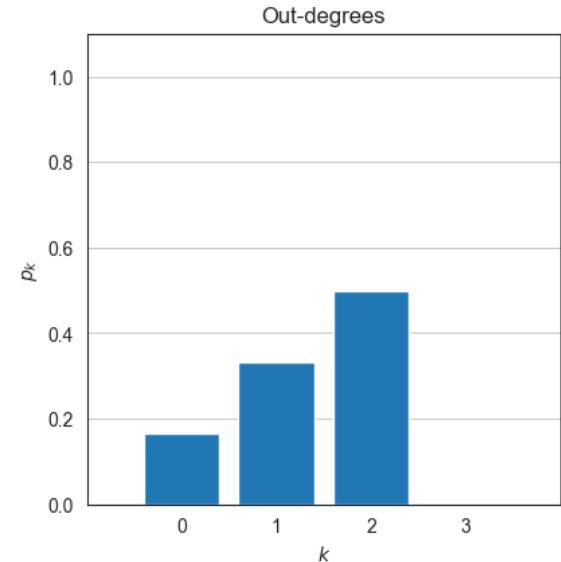
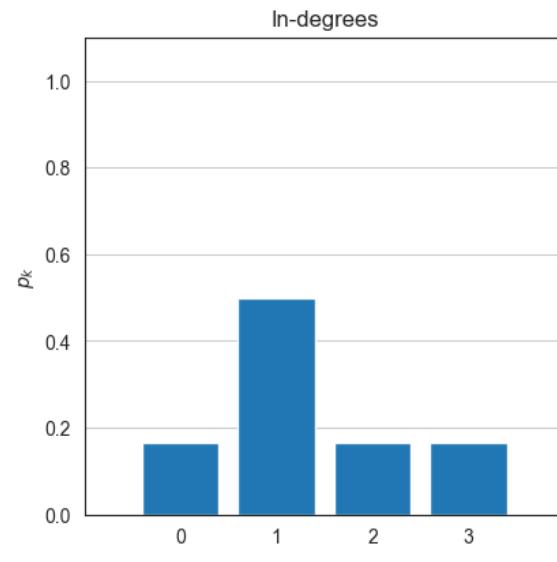
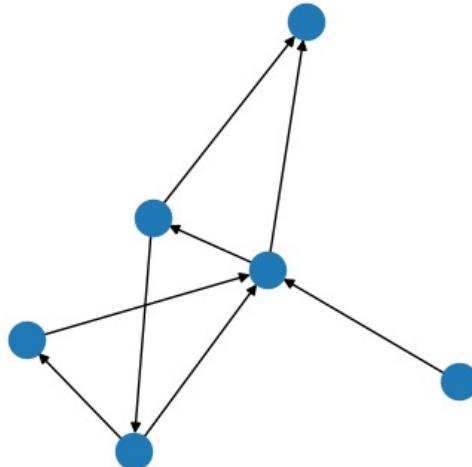
$$A_{ij} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{pmatrix}$$

$$\begin{aligned} A_{ij} &\neq A_{ji} \\ A_{ii} &= 0 \end{aligned}$$

$$k_i^{in} = \sum_{j=1}^N A_{ij}$$

$$k_j^{out} = \sum_{i=1}^N A_{ij}$$

$$L = \sum_{i=1}^N k_i^{in} = \sum_{j=1}^N k_j^{out} = \sum_{i,j} A_{ij}$$

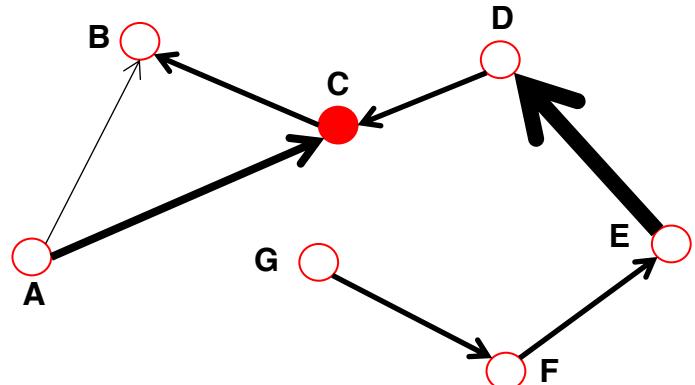


Weighted graphs

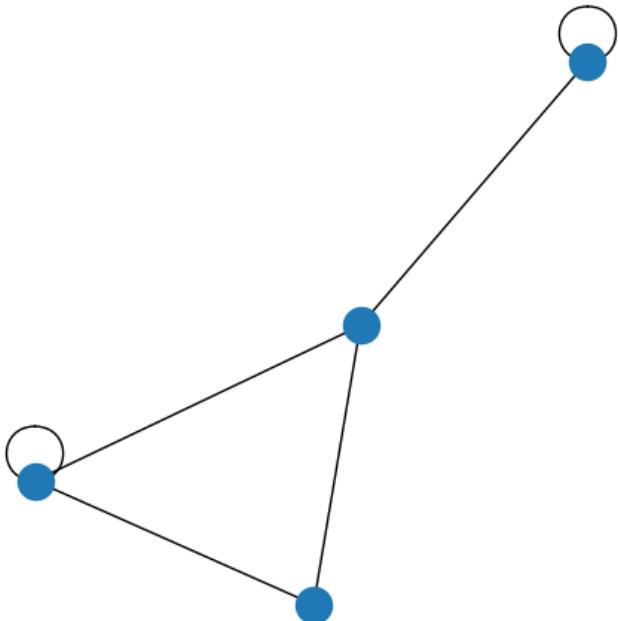
Sometimes, not all edges are equivalent

In many applications, edges carry additional information by an **edge weight**

	Unweighted	Weighted
Adjancy value	$A_{ij} = 1$	$A_{ij} \in \mathbb{R}$



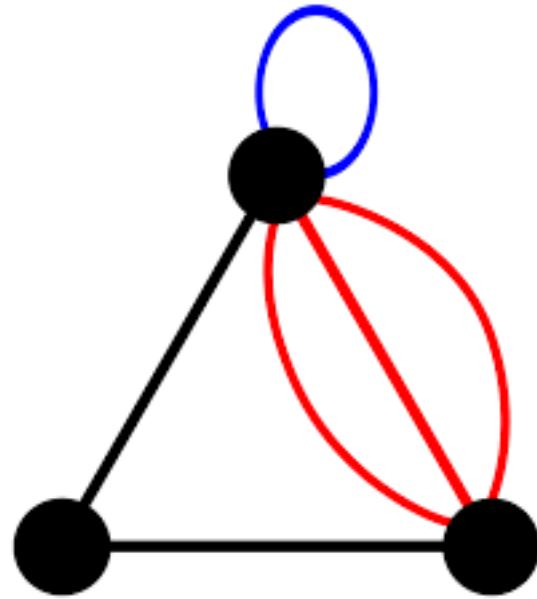
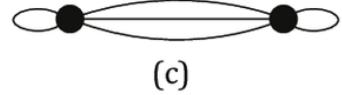
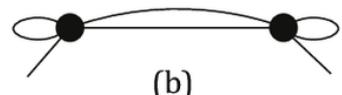
Self interactions



$$A_{ii} \neq 0$$

	0	1	2	3
0	1	1	1	0
1	1	0	1	0
2	1	1	0	1
3	0	0	1	1

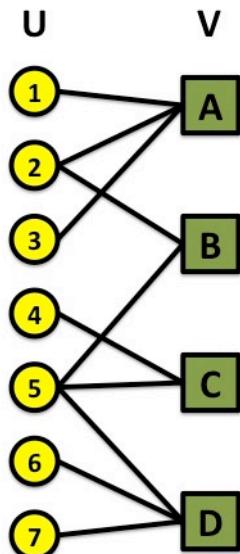
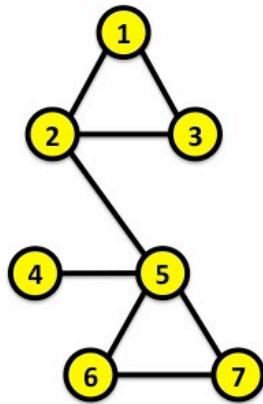
Multigraphs



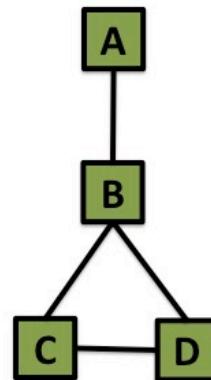
Bipartite Graphs

bipartite graph (or bigraph) is a [graph](#) whose nodes can be divided into two [disjoint sets](#) U and V such that every link connects a node in U to one in V ; that is, U and V are [independent sets](#).

Projection U



Projection V



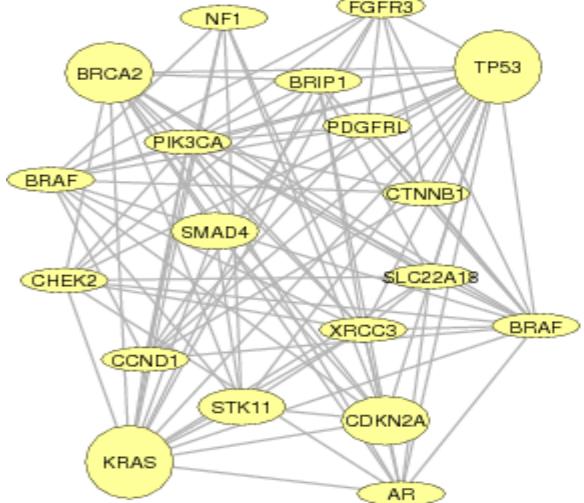
Examples:

Hollywood actor network
Collaboration networks
Disease network (diseasome)

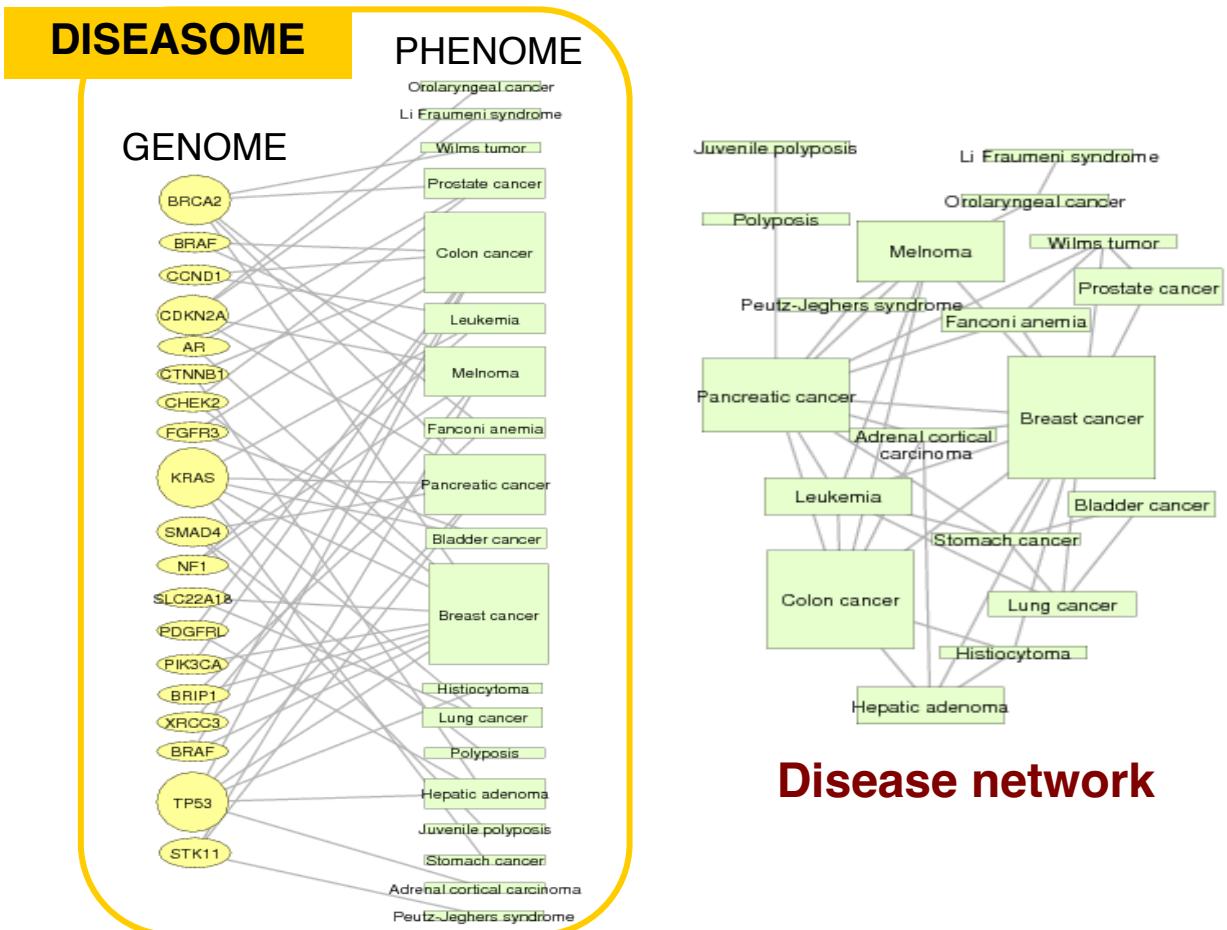
$$BB^T$$

$$B = M_{u,v}$$

$$B^T B$$



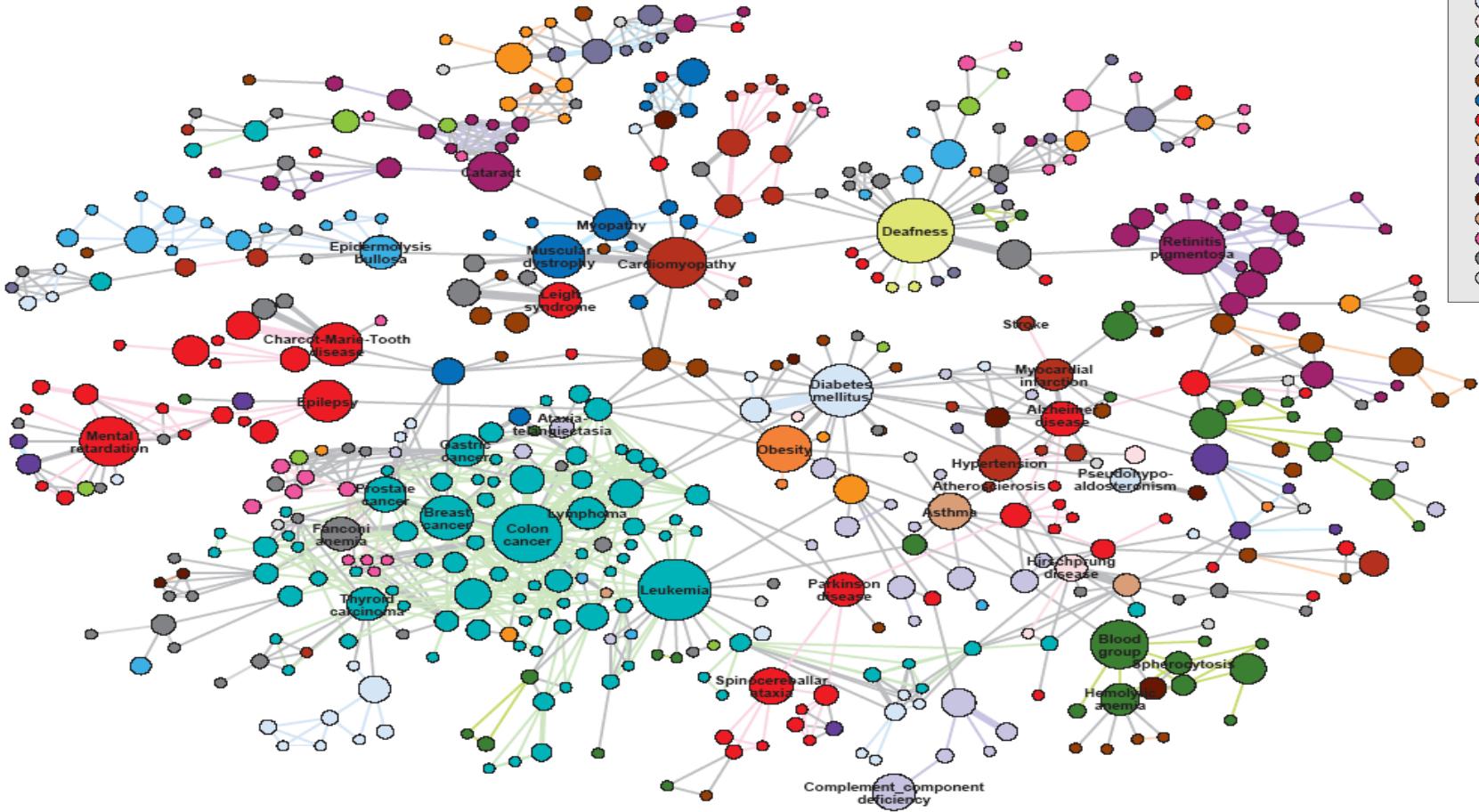
Gene network



Disease network

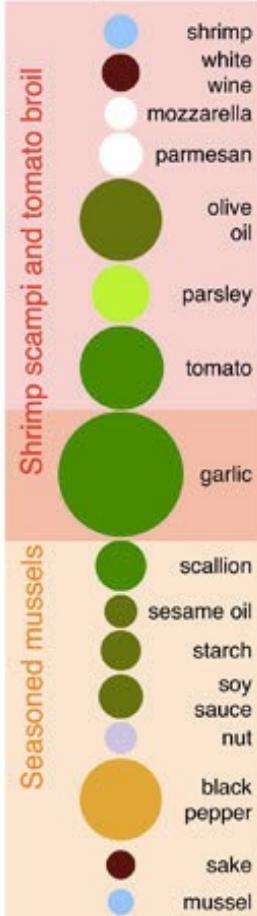
Disorder Class

- Bone
- Cancer
- Cardiovascular
- Connective tissue
- Dermatological
- Developmental
- Ear, Nose, Throat
- Endocrine
- Gastrointestinal
- Hematological
- Immunological
- Metabolic
- Muscular
- Neurological
- Nutritional
- Ophthalmological
- Psychiatric
- Renal
- Respiratory
- Skeletal
- Multiple
- Unclassified



A

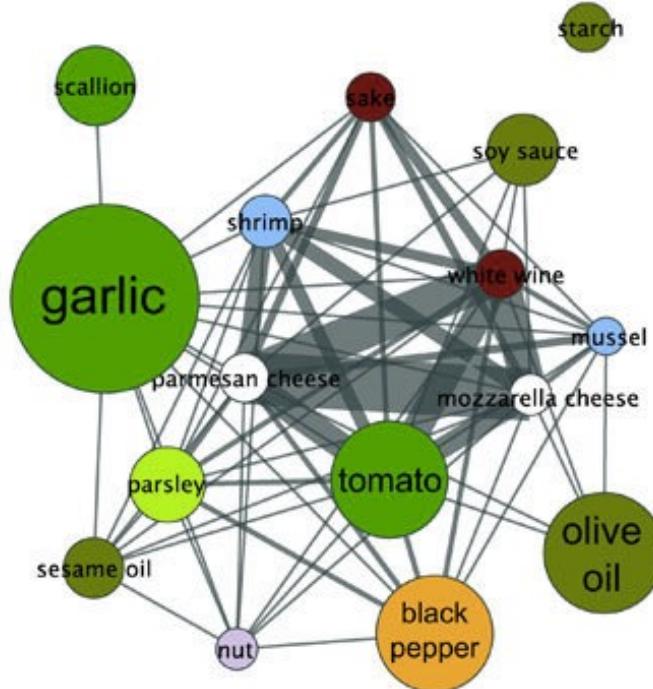
Ingredients



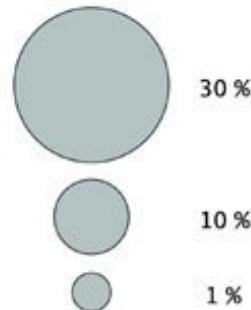
Flavor compounds

1-penten-3-ol
2-hexenal
2-isobutyl thiazole
2,3-diethylpyrazine
2,4-nonadienal
3-hexen-1-ol
4-hydroxy-5-methyl...
4-methylpentanoic acid
acetylpyrazine
allyl 2-furoate
alpha-terpineol
beta-cyclodextrin
cis-3-hexenal
dihydroxyacetone
dimethyl succinate
ethyl propionate
hexyl alcohol
isooamyl alcohol
isobutyl acetate
isobutyl alcohol
lauric acid
limonene (d-,l-, and dl-)
l-malic acid
methyl butyrate
methyl hexanoate
methyl propyl trisulfide
nonanoic acid
phenethyl alcohol
propenyl propyl disulfide
propionaldehyde
propyl disulfide
p-mentha-1,3-diene
p-menth-1-ene-9-al
terpinyl acetate
tetrahydrofurfuryl alcohol
trans, trans-2,4-hexadienal

B Flavor network

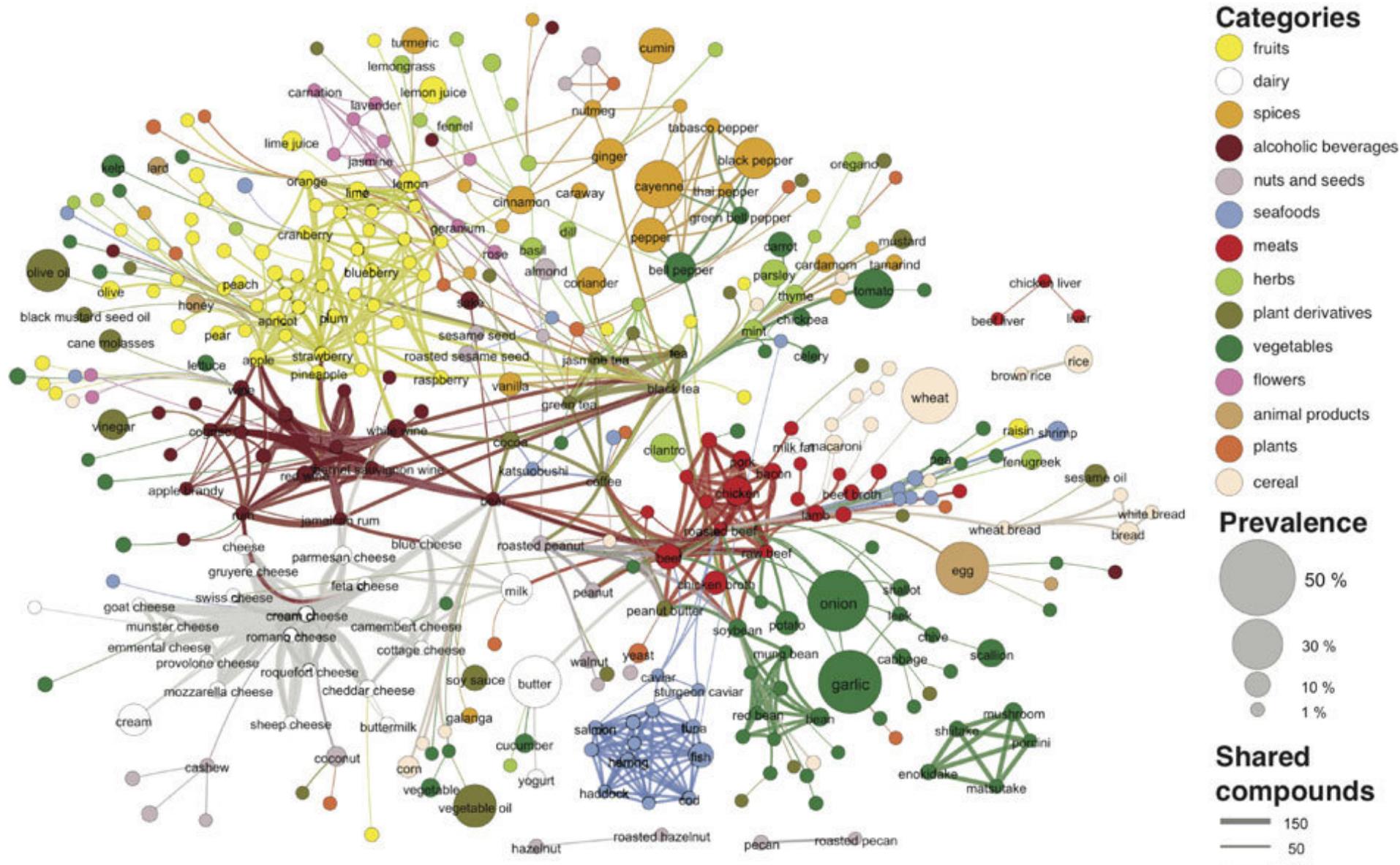


Prevalence



Shared compounds



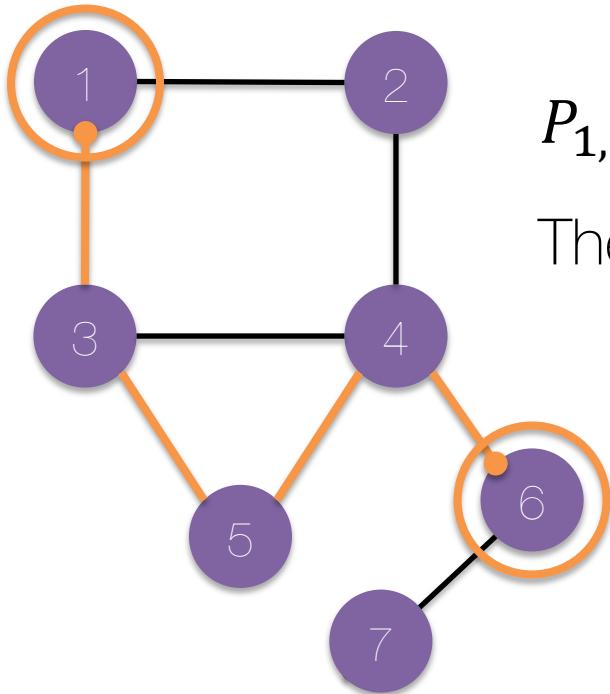


Pathology

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

$$P_{0,n} = (0,1,2, \dots, n)$$

$$P_{0,n} = ((0,1), (1,2), \dots, (n - 1, n))$$

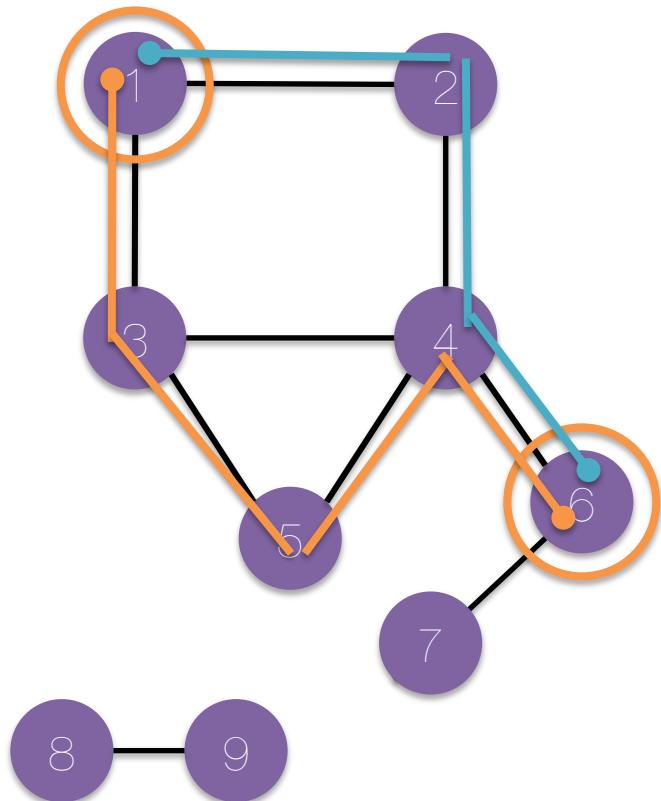


$$P_{1,6} = (1,3,5,4,6)$$

The **path length** is the number of links.

Pathology

There are typically many paths between any two nodes



$$P_{1,6} = (1,3,5,4,6) \quad P_{1,6} = (1,2,4,6)$$

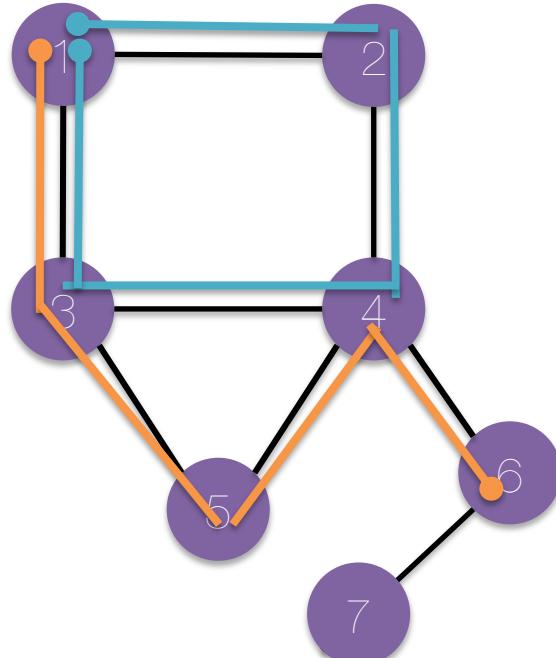
The **shortest path** is the path with the fewest number of links.

The **distance** between two nodes is the length of the shortest path.

If no paths exist between two nodes, the distance is infinite.

Pathology

There are typically many paths between any two nodes

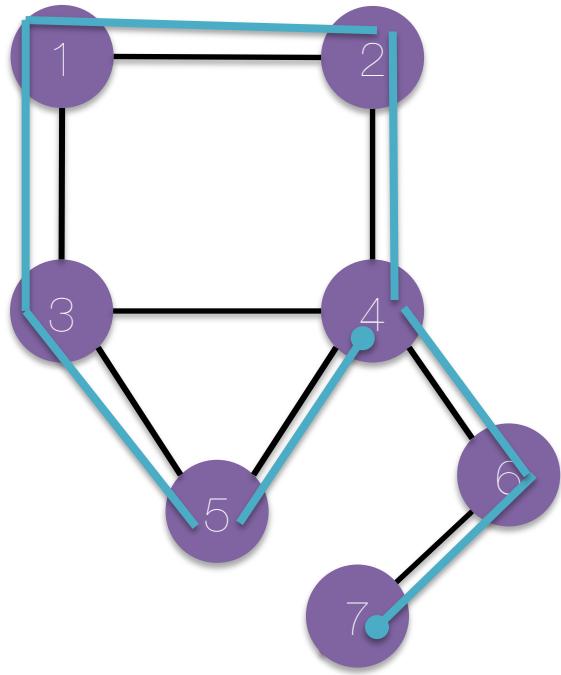


$$P_{1,1} = (1,2,4,3,1)$$

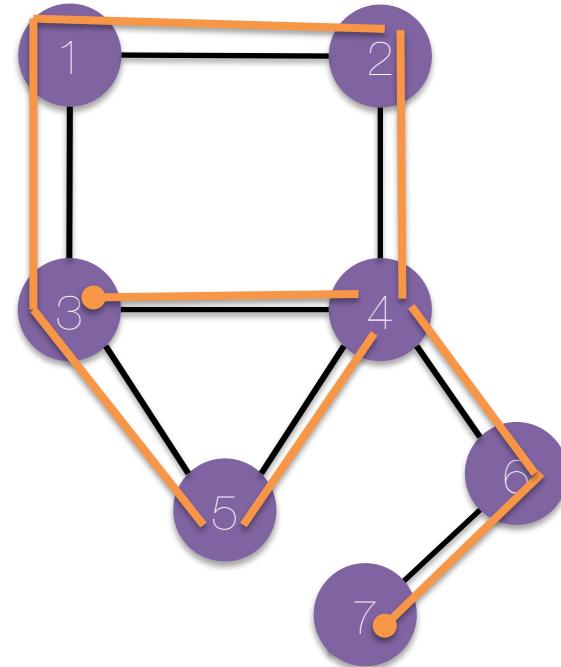
A **cycle** is a path that starts and stops at the same node.

A **self-avoiding path** is a path that does not intersect itself.

Pathology

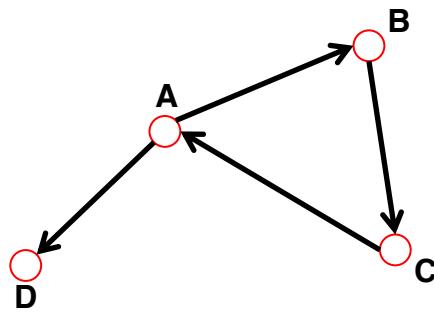


A **Hamilton path** is a path that touches each node exactly once.



A **Euclidean path** is a path that touches each link exactly once.

Pathology



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

Thus in a digraph the distance from node A to B (on an AB path) is generally different from the distance from node B to A (on a BCA path).

Pathology

The graph **diameter** is the maximum distance between any two nodes in the graph (maximum shortest path length).

The **average distance** in a graph:

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

Pathology

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

N_{ij}^n number of paths between any two nodes i and j

N_{ij}^1 Is there a path of length 1 between nodes i and j

$$N_{ij}^1 = A_{ij}$$

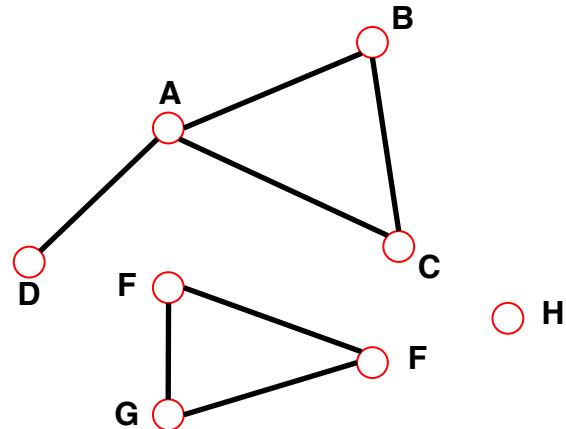
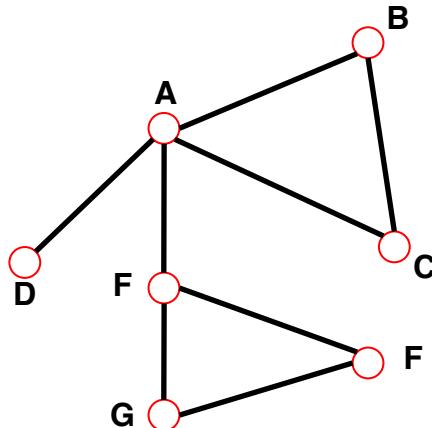
N_{ij}^2 If there is a path of length 2 between nodes i and j,
then there is a node k where $A_{ik}A_{kj} = 1$

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

Connectedness

Connected graph:

(undirected) any two vertices can be joined by a path.

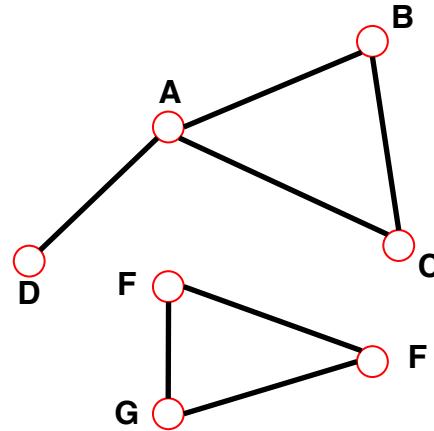


Disconnected graph: is made up of two or more connected graphs (aka connected components).

Connectedness

Largest connected component

(a.k.a. giant component) is the one with the most nodes

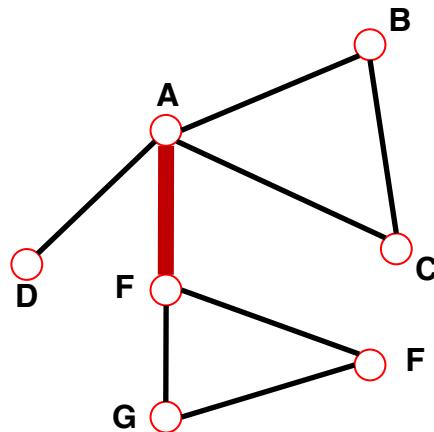


Isolates

are the rest of the components

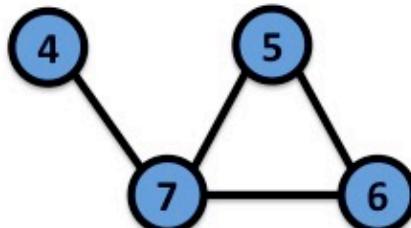
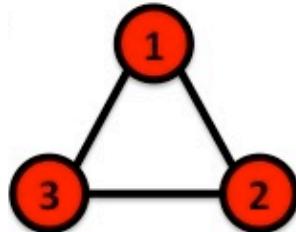
Connectedness

Bridge: an edge whose removal breaks the graph in disconnected components

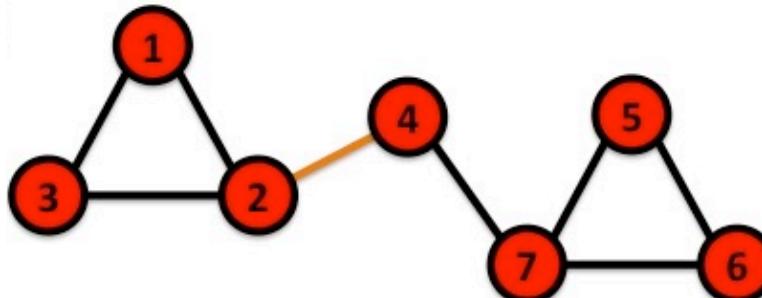


Connectedness

The adjacency matrix of a network with several components can be written in a block-diagonal form, so that nonzero elements are confined to squares, with all other elements being zero:



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

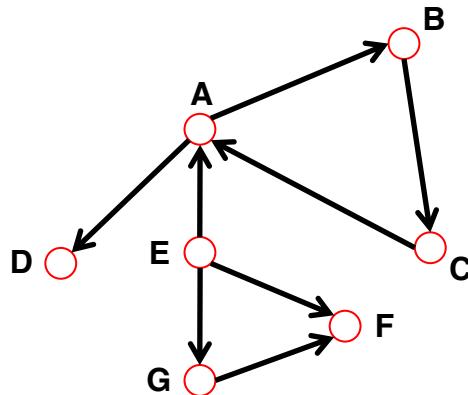


$$\begin{pmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{pmatrix}$$

Connectedness

Strongly connected digraph:

has a directed path between every pair of nodes (and vice versa)



Strongly connected component:

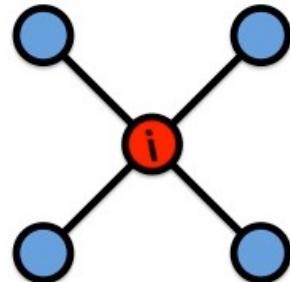
a subgraph with a path between every pair of nodes (and vice versa)

Weakly connected digraph:

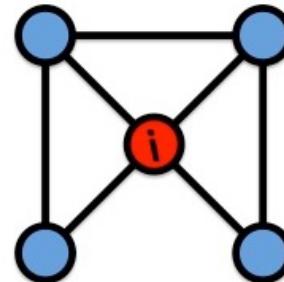
connected if we disregard the edge direction

Clustering

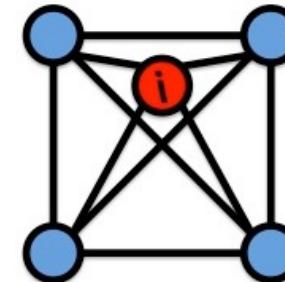
What fraction of your neighbors are also connected to each-other?



$$C_i = 0$$



$$C_i = 1/2$$



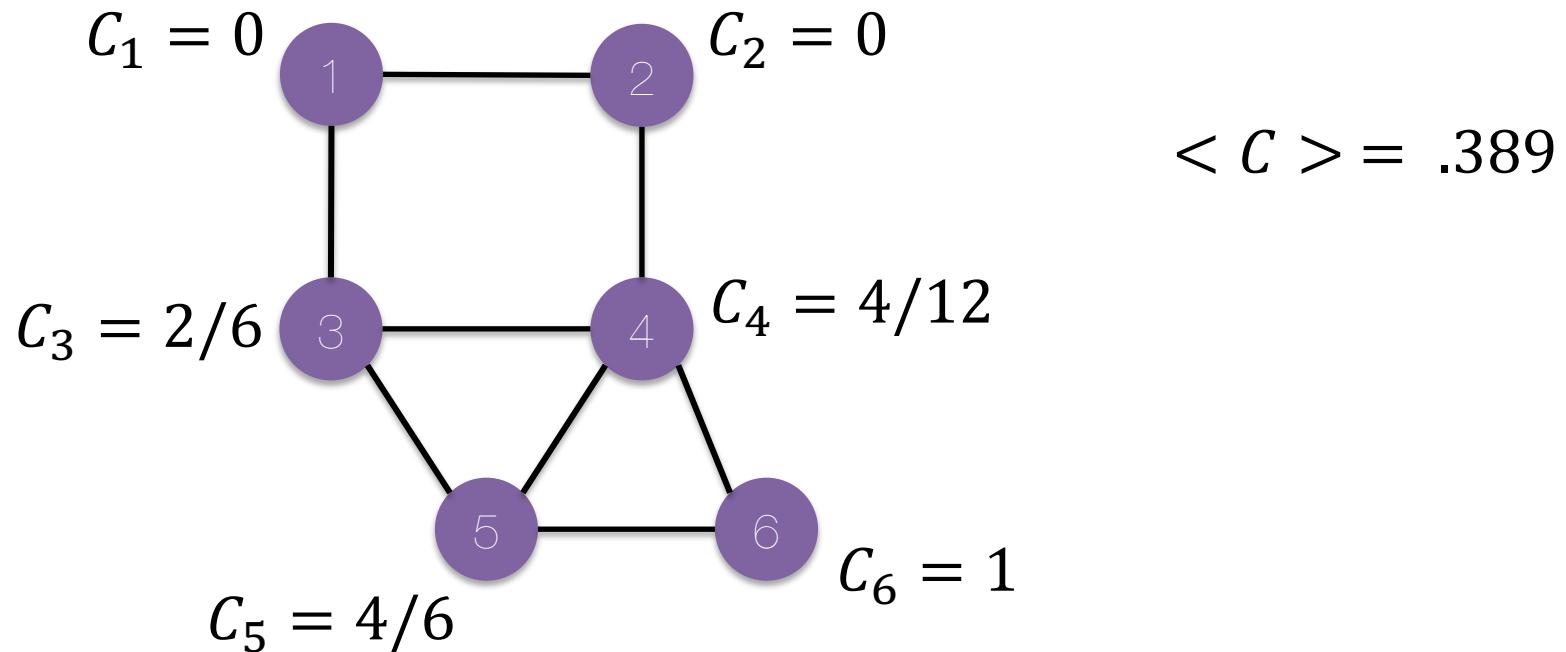
$$C_i = 1$$

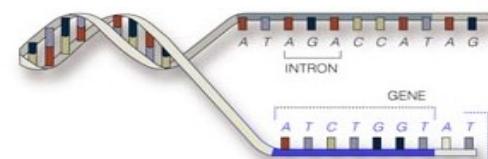
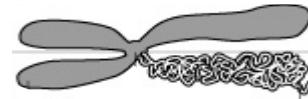
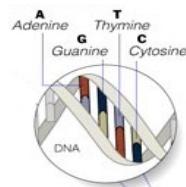
How many **triangles** does node i participate in? $\tau(i)$

Clustering coefficient the fraction of possible triangles in which a node participates:

$$C_i = \frac{\tau(i)}{\tau_{max}(i)} \quad C_i = \frac{2\tau(i)}{k_i(k_i - 1)}$$

Clustering





GENOME

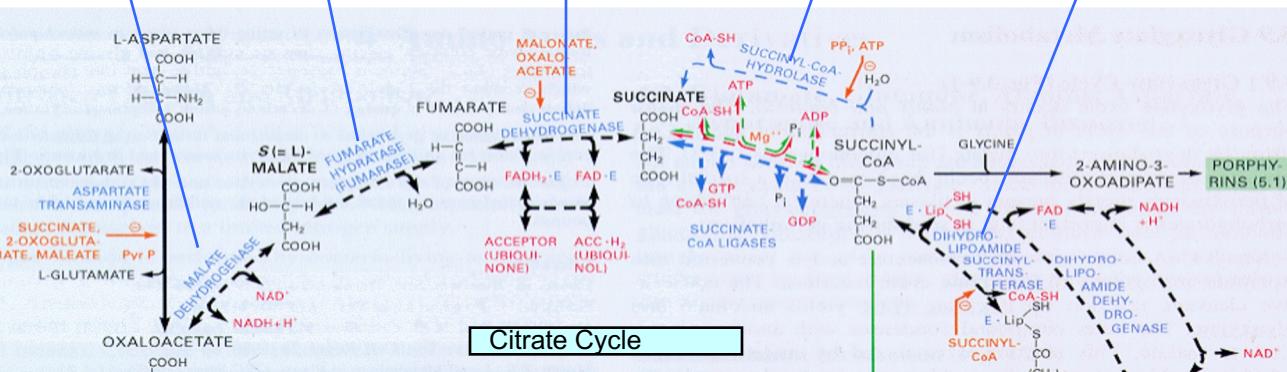
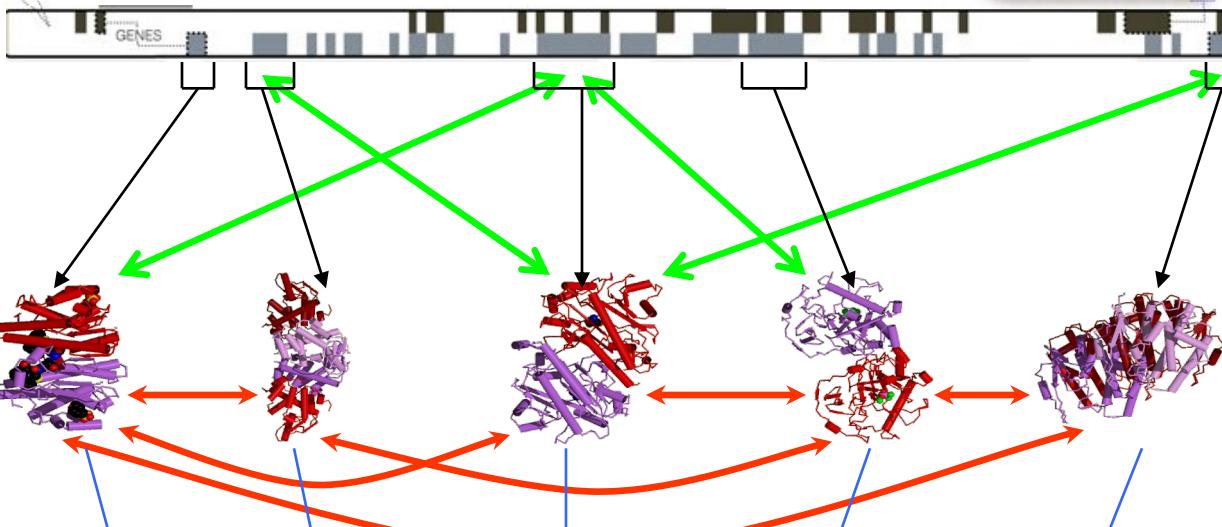
protein-gene interactions

PROTEOME

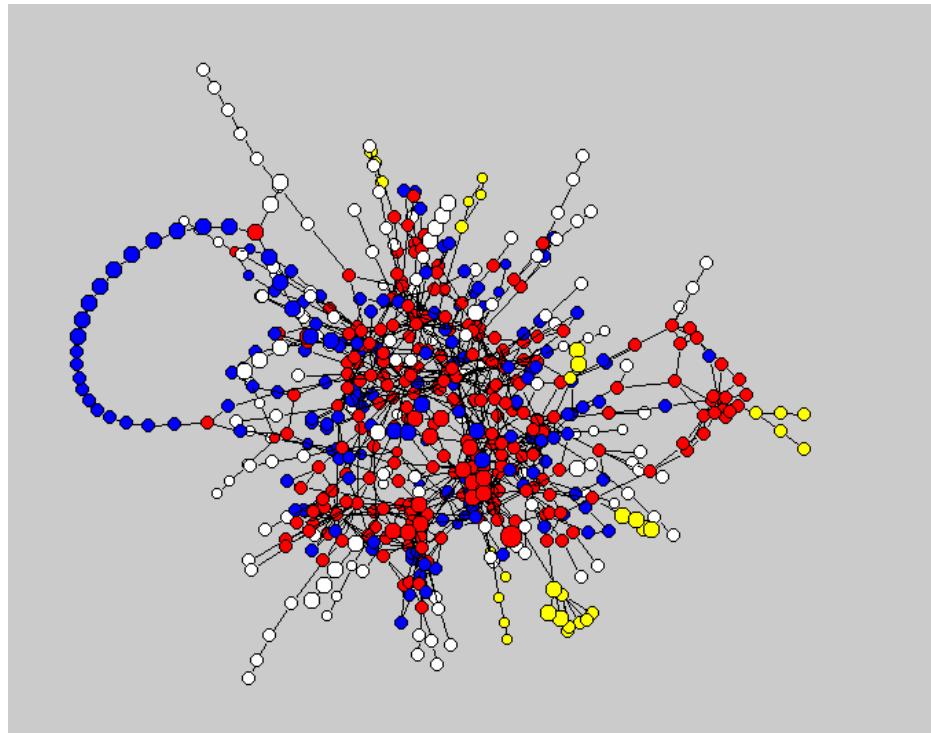
protein-protein interactions

METABOLISM

Bio-chemical reactions

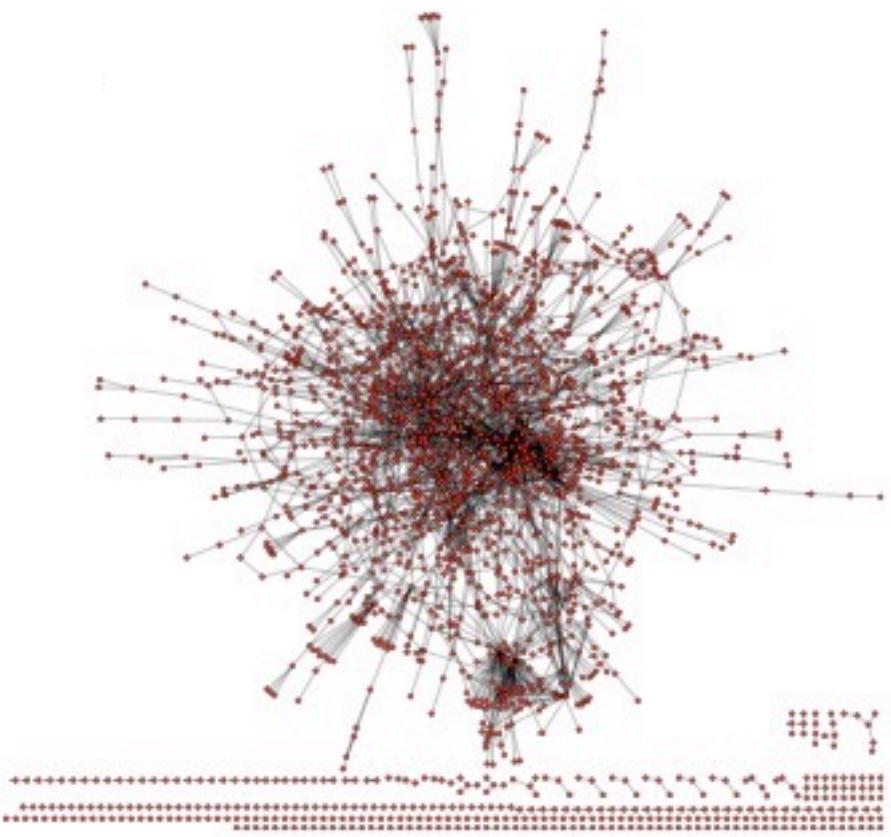


Metabolic Network



Protein Interactions

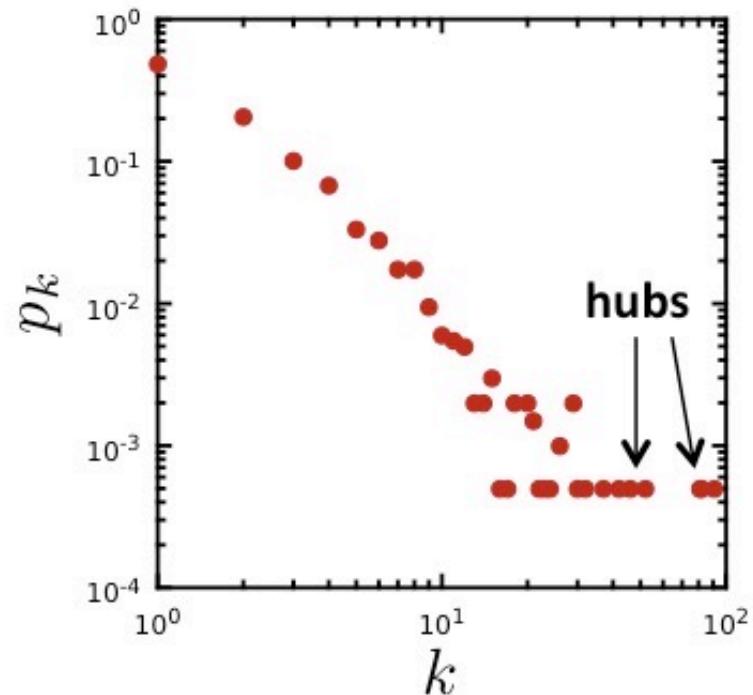




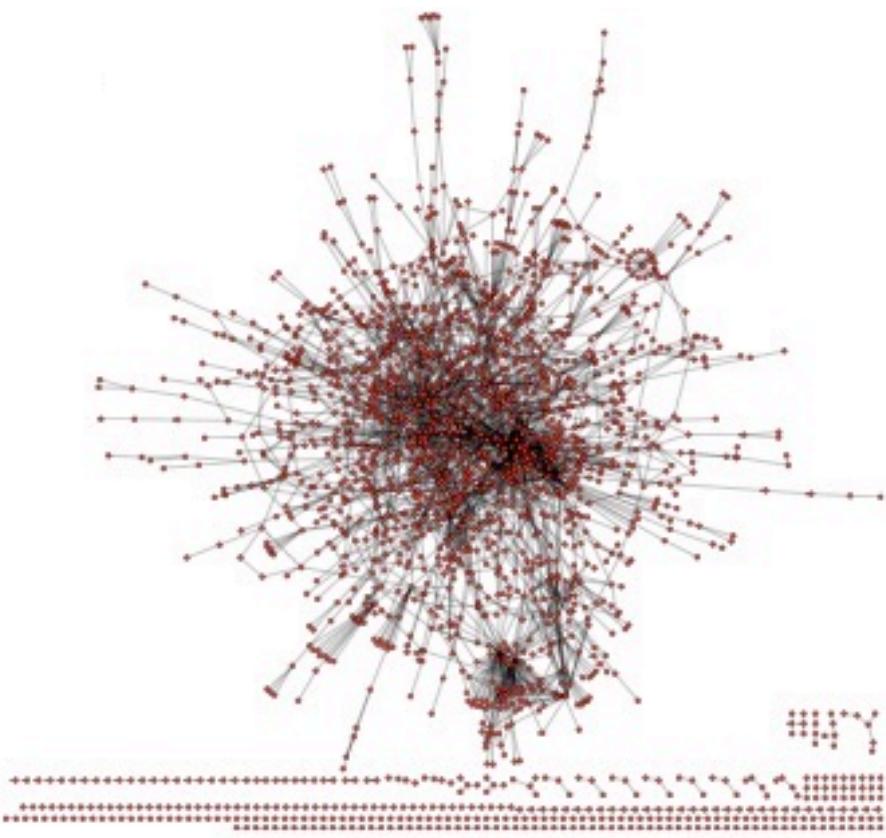
Not connected: 185 components

the largest (giant component)
1,647 nodes

A few hubs connect to ~10%
of the network

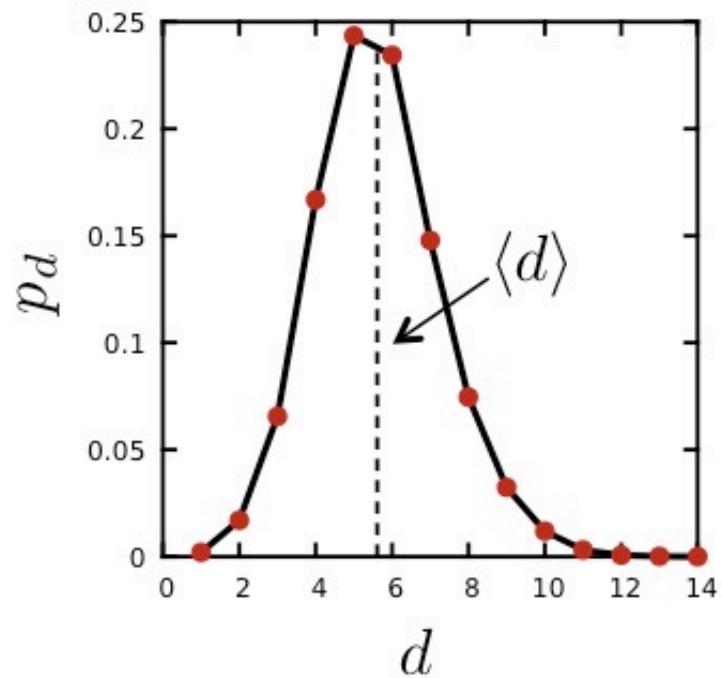


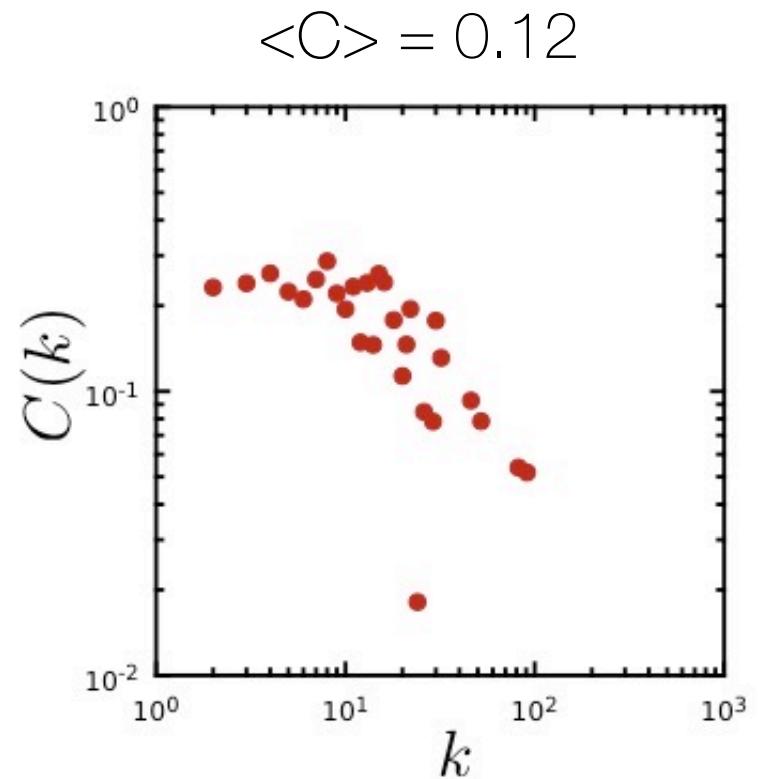
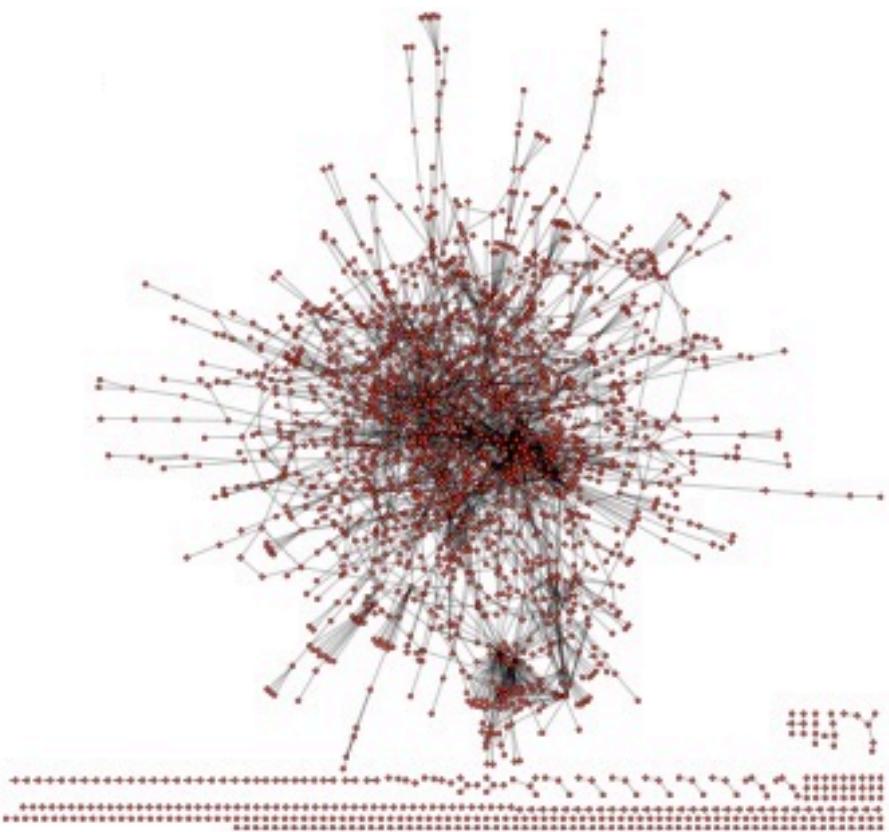
Most nodes have only 1-2
connections $\langle k \rangle = 2.9$



The network's diameter is 14

$$\langle d \rangle = 5.6$$





We can explore how the average clustering coefficient changes with node degree

Weekly Class Readings

Readings

- Chapters 1 & 2 in Easley & Kleinberg, *Networks, crowds, and markets: Reasoning about a highly connected world*. The full book is also here.
- Butts, Carter T. (2009) "Revisiting the foundations of network analysis." *Science* 325, no. 5939: 414-416.

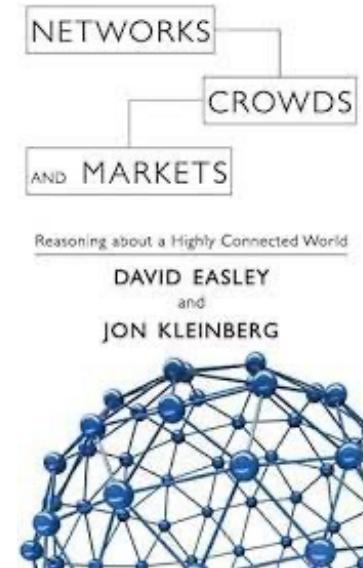
PERSPECTIVE

Revisiting the Foundations of Network Analysis

Carter T. Butts

Network analysis has emerged as a powerful way of studying phenomena as diverse as interpersonal interaction, connections among neurons, and the structure of the Internet. Appropriate use of network analysis depends, however, on choosing the right network representation for the problem at hand.

24 JULY 2009 VOL 325 SCIENCE www.sciencemag.org



Chapters 1 & 2



<https://uvads8104.github.io/content/>

Weekly Class Readings

Next week (1/31 & 2/2):

Readings

-  Milgram, Stanley. (1967) “The small world problem.” *Psychology today* **2**, no. 1: 60-67.
-  Watts, Duncan J., and Steven H. Strogatz. (1998) “Collective dynamics of ‘small world’ networks.” *Nature* **393**, 6684: 440-442.
-  Fosdick, Bailey K., Daniel B. Larremore, Joel Nishimura, and Johan Ugander. (2018) “Configuring random graph models with fixed degree sequences.” *Siam Review* **60**, no. 2: 315-355.
-  Milo, Ron, Shai Shen-Orr, Shalev Itzkovitz, Nadav Kashtan, Dmitri Chklovskii, and Uri Alon. (2002) “Network motifs: simple building blocks of complex networks.” *Science* **298**, no. 5594: 824-827.



<https://uvads8104.github.io/content/>

Discussion moderation

1. Brief summary
 - i. Less than <5 min only!
 - ii. Slides are OK, but not necessary
 - iii. Main Thesis? / Which figure shows this?
2. Contextualize
 - i. What class ideas does the article connect with / draw upon?
3. Questions for discussion
 - i. Do you agree with the findings?
 - ii. Are all results consistent with hypothesis?
 - iii. Anything not clear / vague?
 - iv. Any problems? Were all methods correctly used?
 - v. What would you do differently?
 - vi. Be provocative! (but be respectful!)