

# Introduction to data science and network science

---

Data science course LPI, 2024/2025

Liubov Tupikina [liubov.tupikina@cri-paris.org](mailto:liubov.tupikina@cri-paris.org)  
LPI, Bell labs France

## Contacts

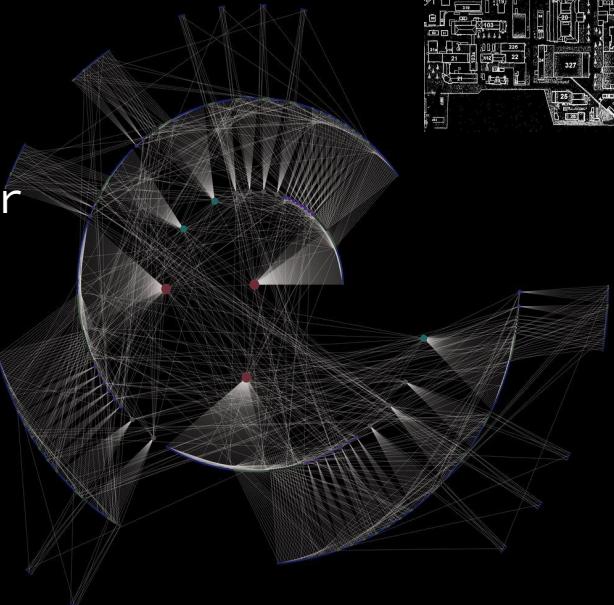
Affiliate researcher at LPI,  
Researcher, Nokia Bell Labs, France

## Main topics:

Explainable AI (geometry of embeddings, analysis of higher-order data)

Geometry of graphs and hypergraphs  
(processes on graphs, transport applications)

Students contact:  
Zahara Farook, Hritika Kathuria



# Outline and connection to other courses

September 2024				October 2024				November 2024				
Day	9h30-12h30	14h-17h	Day	9h30-12h30	14h-17h	Day	9h30-12h30	14h-17h	Day	9h30-12h30	14h-17h	Day
14	Weekend-1		S 6	Weekend-4			W 6	Neurosciences-5	Research methods / CS - 6		F 6	Artificial I
15			M 7	Technologies for learning-1	Python-3	T 7	French Language-8	CCA	S 7			
16			T 8	Exploring Sustainability-4	Statistics-2	F 8	Mental Health with CBT-2	Data Science-6	S 8			
17	Statistics-1	Exploring Sustainability-1	W 9	Neurosciences-2	Research methods / CS - 2	S 9	Weekend-9			M 9	Technologie	
18	Introduction to Data Science-1	Introduction to Data Science-2	T 10	French Language-4	CCA	S 10				T 10	Physics of Fundamental	
19	French Language-1	Co-curricular Activities	F 11	Internship Workshop-1	Exploring Sustainability-5	M 11	National Holiday			W 11	Neuroscie	
20	Introduction to Data Science-3	Introduction to Data Science-4	S 12	Weekend-5			T 12	Technologies for learning-6	Exploring Sustainability-10	T 12	French la	
21	Weekend-2		S 13	W 13	Neurosciences-6	Statistics-7	F 13	Artificial I				
22			M 14	Technologies for learning-2	Data Science-1	T 14	French Language-9	CCA	S 14			
23	Discovery Days		T 15	Statistics-3	Exploring Sustainability-6	F 15	Masterclass	Data Science-7	S 15			
24	Exploring Sustainability-2		W 16	Neurosciences-3	Research methods / CS - 3	S 16	Weekend-10			M 16	Technologie	
25			T 17	French Language-5	CCA	S 17				T 17	Physics of Fundamental	
26	French Language-2	CCA	F 18	Data Science-2	Inclusion & Diversity workshop	M 18	Technologies for learning-7	Data Science-8	W 18	Mental Hea		
27		Python-1	S 19	Weekend-6			T 19	Physics of the cells-1/ Fundamentals of Learning-1	Exploring Sustainability-11	T 19	French la	
28	Weekend-3		S 20	W 20	Neurosciences-7	Research methods / CS - 7	F 20	Artificial I				
29			M 21	Technologies for learning-3	Data Science-3	T 21	French language-10	CCA	S 21			
30	Research methods / CS - 1	Workshop In&Di: Le Paris Noir	T 22	Statistics-4	Exploring Sustainability-7	F 22	Mental Health with CBT-3	Statistics-8	S 22			
			W 23	Research methods / CS - 4	Internship Workshop-2	S 23	Weekend-11			M 23		
			T 24	French Language-6	CCA	S 24				T 24		

# Orientation every class



# Outline of Introduction to data science

## Syllabus and Agenda:

### 18th September:

- **Morning:** Elements of statistics for data analysis: building intuition with a dataset)
- **Afternoon:** Introduction to data science, network science

### 20th September

- **Morning:** Foundations on AI for data science: from theory to practice on data fitting, embedding, modeling
- **Afternoon:** Spatial data analysis, Data and Network visualization

# Outline of Introduction to data science

## Syllabus and Agenda:

**18th September:**

- **Morning:** Elements of statistics for data analysis: building intuition with a dataset)
- **Afternoon:** Introduction to data science, network science

**20th September**

- **Morning:** Foundations on AI for data science: from theory to practice on data fitting, embedding, modeling
- **Afternoon:** Spatial data analysis, Data and Network visualization

# Outline of Introduction to network science

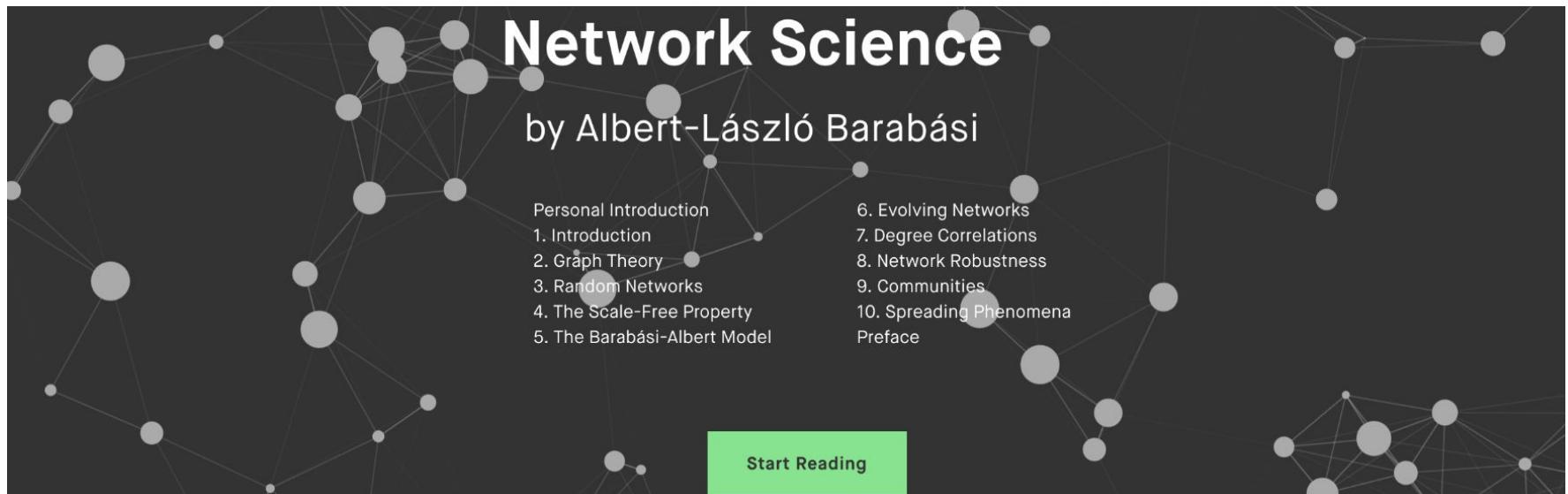
1. Introduction to networks
2. Practical part: notebooks

[Assignment](#) for the 20th September

# Resources and libraries for the course

Standard libraries (Python): numpy, matplotlib, scikit learn, seaborn

Network libraries: networkx, osmnx (open street data)



# Resources and libraries for the course

Standard libraries (Python): numpy, matplotlib, scikit learn, seaborn

Network libraries: networkx, osmnx (open street data)

## Support materials

- Big data course Marc and Liubov <https://github.com/Big-data-course-CRI/>
- Correlaid, Complex system conference CSS 2023 and TidyTuesday  
<https://github.com/rfordatascience/tidytuesday>
- Network science book <http://networksciencebook.com/>
- Network repository <https://networks.skewed.de/>
- Visualisation tools <https://gephi.org/users/download/>
- Network datasets <https://www.complex-networks.net/datasets.html#chap8>

# Resources and libraries for the course

The screenshot shows a GitHub repository page for 'materials\_big\_data\_cri\_2024\_2025'. The repository is public and has 5 commits. It contains files like 'README.md', 'day 1 networks and hypergraphs', 'day 2 foundations AI', and 'README'. The 'About' section describes it as a repository to share example codes and materials of lectures. It has 0 stars, 2 watchers, and 0 forks. There are sections for 'Releases', 'Packages', and 'Languages'.

**Code** Issues Pull requests Actions Projects Wiki Security Insights Settings

**materials\_big\_data\_cri\_2024\_2025** Public

main 1 Branch 0 Tags

Go to file Add file Code

Liyubov day 2 foundations of AI 406bbae · 2 days ago 5 Commits

day 1 networks and hypergraphs Add files via upload 2 days ago

day 2 foundations AI day 2 foundations of AI 2 days ago

README.md Update README.md 2 days ago

README

**materials big data cri 2024-2025**

The repository of the course to share example codes and materials of lectures.

About

The repository of the course to share example codes and materials of lectures.

Readme Activity Custom properties

0 stars 2 watching 0 forks

Report repository

Releases

No releases published Create a new release

Packages

No packages published Publish your first package

Languages

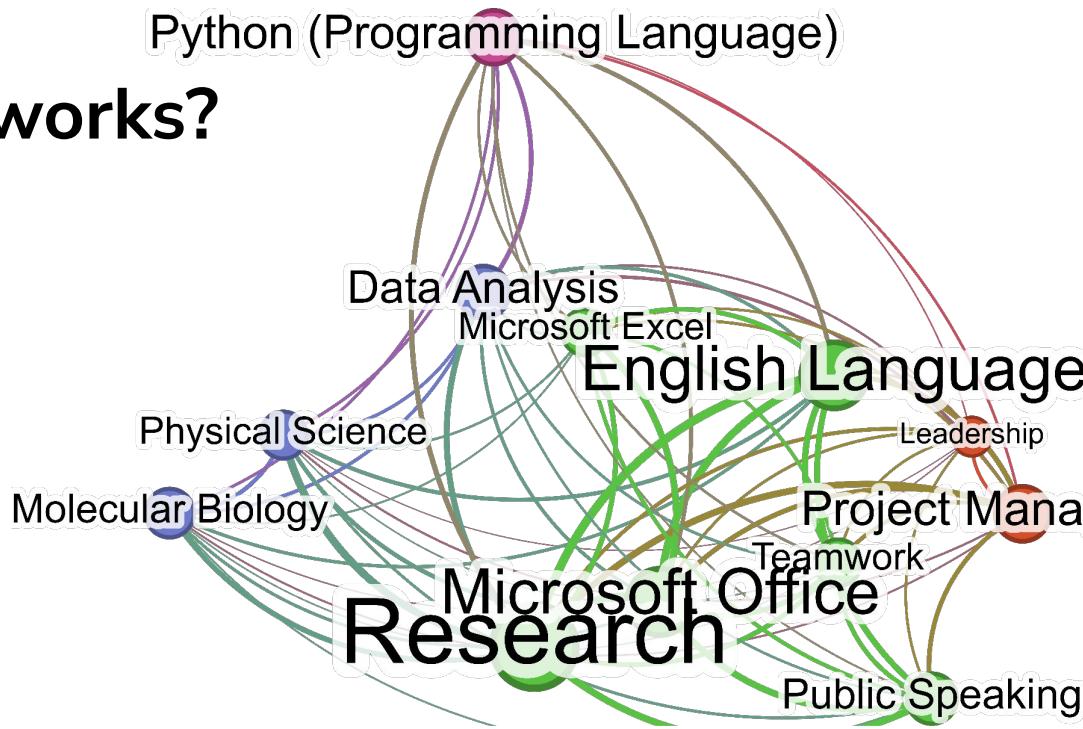
Jupyter Notebook 100.0%

[https://github.com/Big-data-course-CRI/materials\\_big\\_data\\_cri\\_2024\\_2025](https://github.com/Big-data-course-CRI/materials_big_data_cri_2024_2025)

# Why networks?

Figure 7.11

# What are graphs / networks?



# Python (Programming Language)

## Where graphs / networks can be used?

Networks are good to represent data.

What are structures which we can process using networks?

See other examples of data in [Github of the course](#)

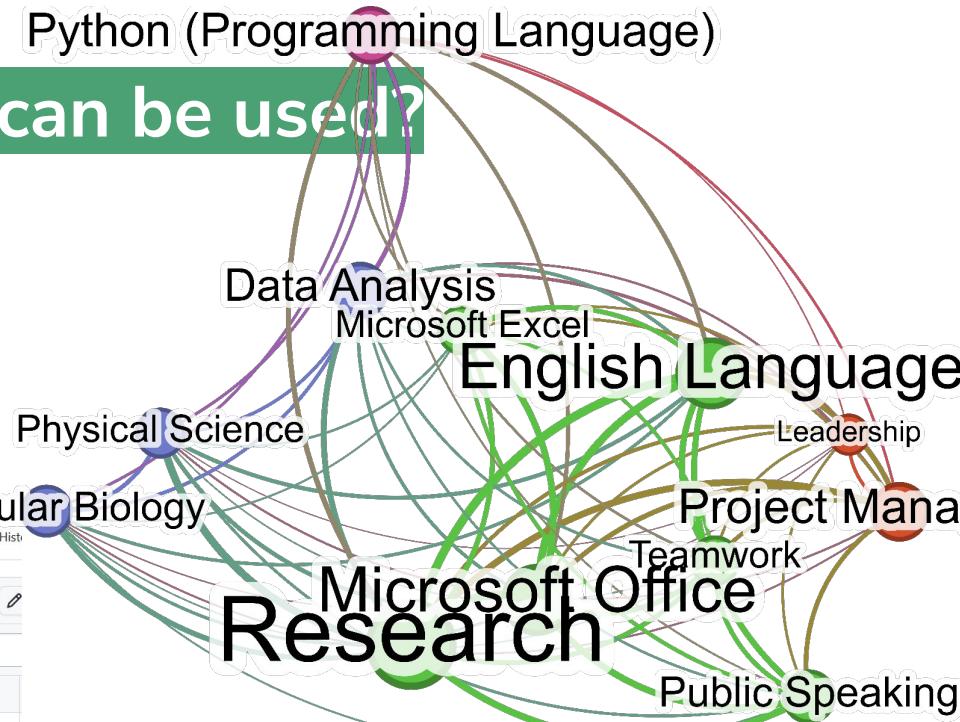
ziqingchery organized\_version

Preview Code Blame 8 lines (8 loc) · 56.3 KB

Raw

Search this file

#	Name	Occupation	City	State	Country
1	7f97716741aea4d227491b5a7d87d4e	Stagiaire at INSTITUT GUSTAVE ROUSSY			France
2	c10674167315089247ea5fa8c98256f6	Initiatrice de projet at Tous Tes Possibles	Paris	Île-de-France	France
3	bf83d3cbf7ebfeeeccdedd9519df56af8	PhD Student at Medical University of Vienna	Österreich		Austria
4	94369f5f833008a9b3d6e9ae7a6a533	Chargée de communication grand public et jeunes at ADEME	Paris	Île-de-France	France
5	67af1831de65104374b77b9a597f4671	Stagiaire UX/ Product Owner at Tylt	Talence	Nouvelle-Aquitaine	France
6	95dc67fb5d78d0c099249d553343451	Research Associate at King's College London			United Kingdom
7	8009f906321941401...800	Project Manager at Microsoft Research	Redmond	Washington	United States

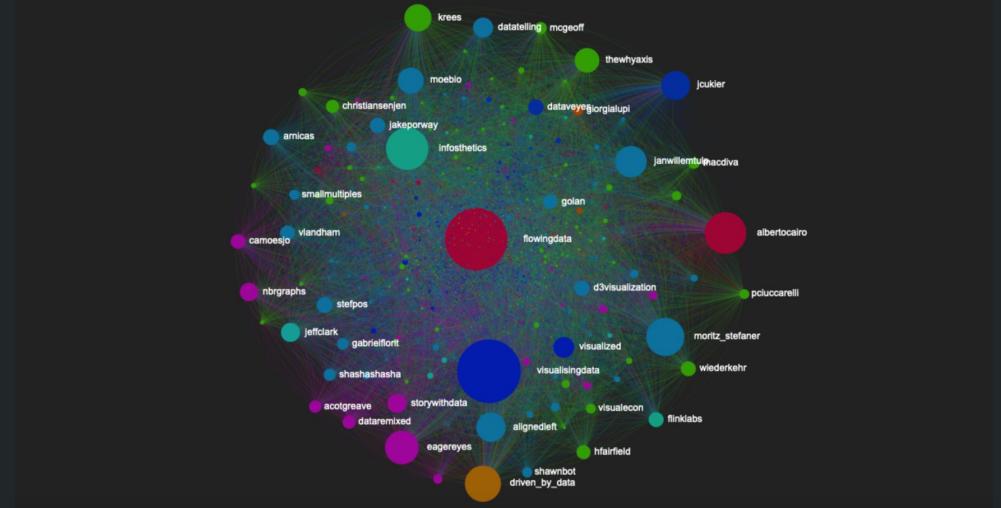


A photograph of a man in a server room, wearing a dark t-shirt and pants, using a broom to sweep up a massive tangle of red, white, and purple network cables that have spilled onto the floor between server racks. The cables are extremely dense and chaotic.

Which data can be represented  
using networks representation?

<https://www.wired.com/2014/09/coupland-bell-labs/>

# Which data can be represented using networks representation?



Social media: Twitter, Instagram, Facebook data analysis

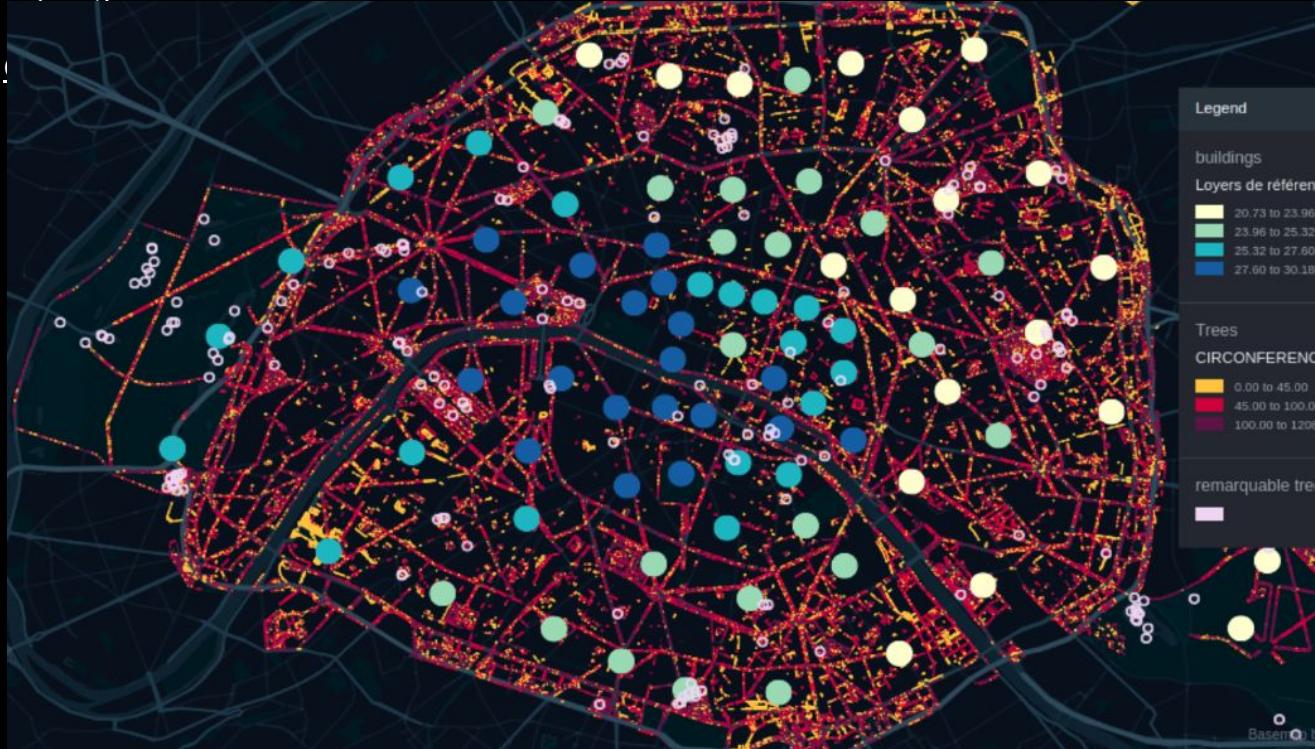
Ask Marc Santolini on his project or check other projects

<https://exploring-data.com/vis/visualisingdata-census-twitter-network/#infosthetics>

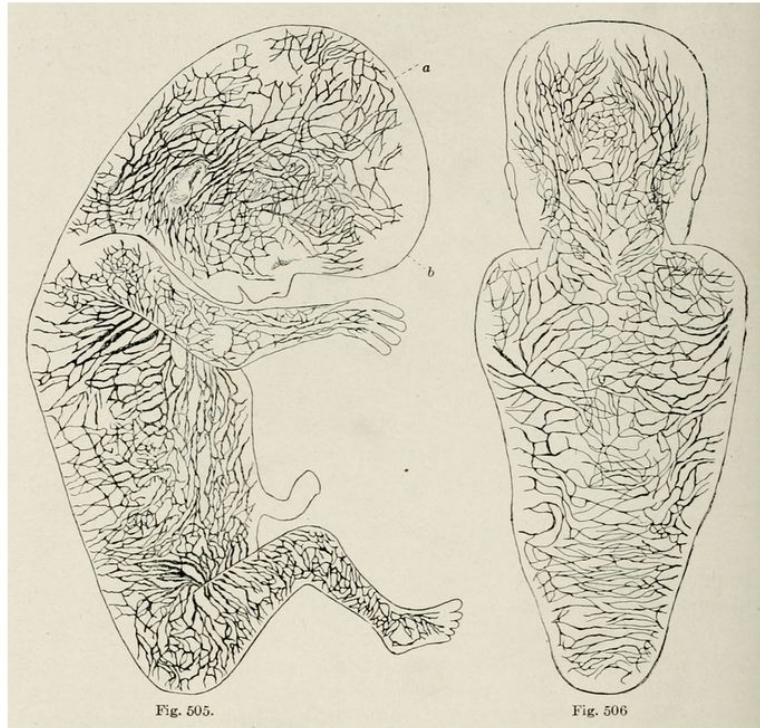
# Networks and hypergraphs in data

Network data: Olympics in Paris, statistics, the Guardian data stories, Humanitarian data HDX

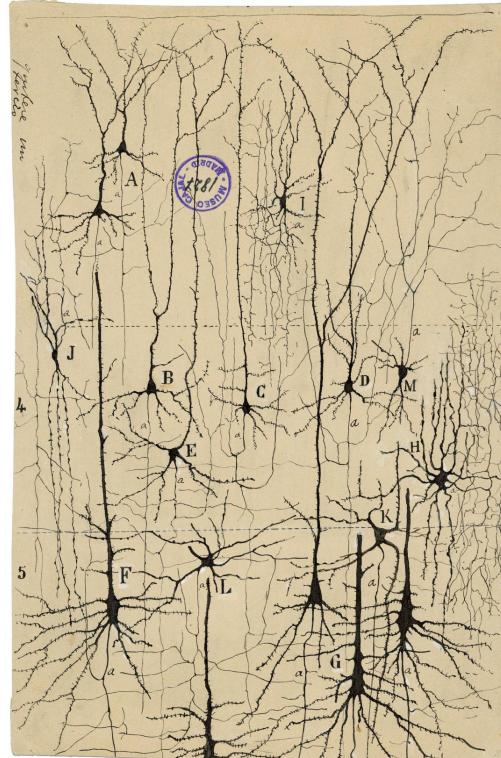
Kepler.gl visualisation



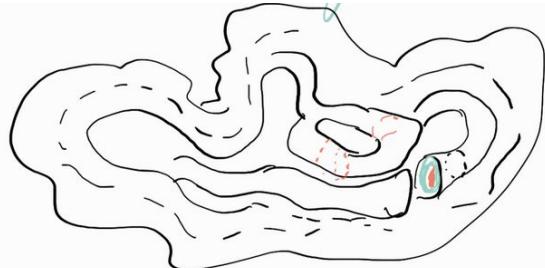
# NETWORKS DESCRIBE HOW THINGS CONNECT AND INTERACT



Distension of the lymphatic vessels in the human foetus, from Franz Kreibel, *Manual of human embryology*, 1910



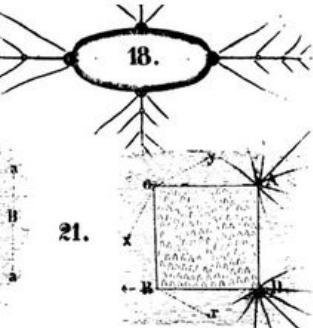
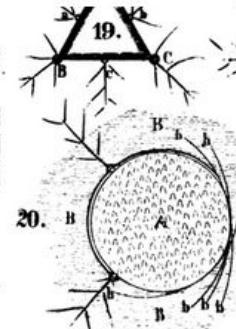
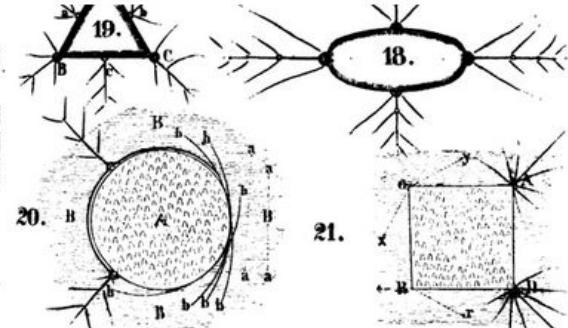
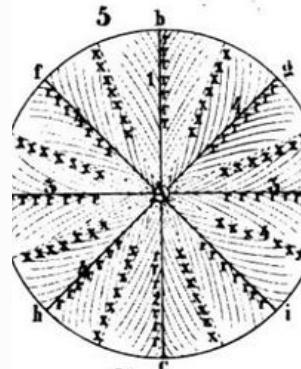
# NETWORKS DESCRIBE HOW THINGS CONNECT AND INTERACT



I think it can  
be approximated  
by manifold  
in  $\mathbb{R}^N$



If we cut it into pieces  
and see what we can  
record from its parts...

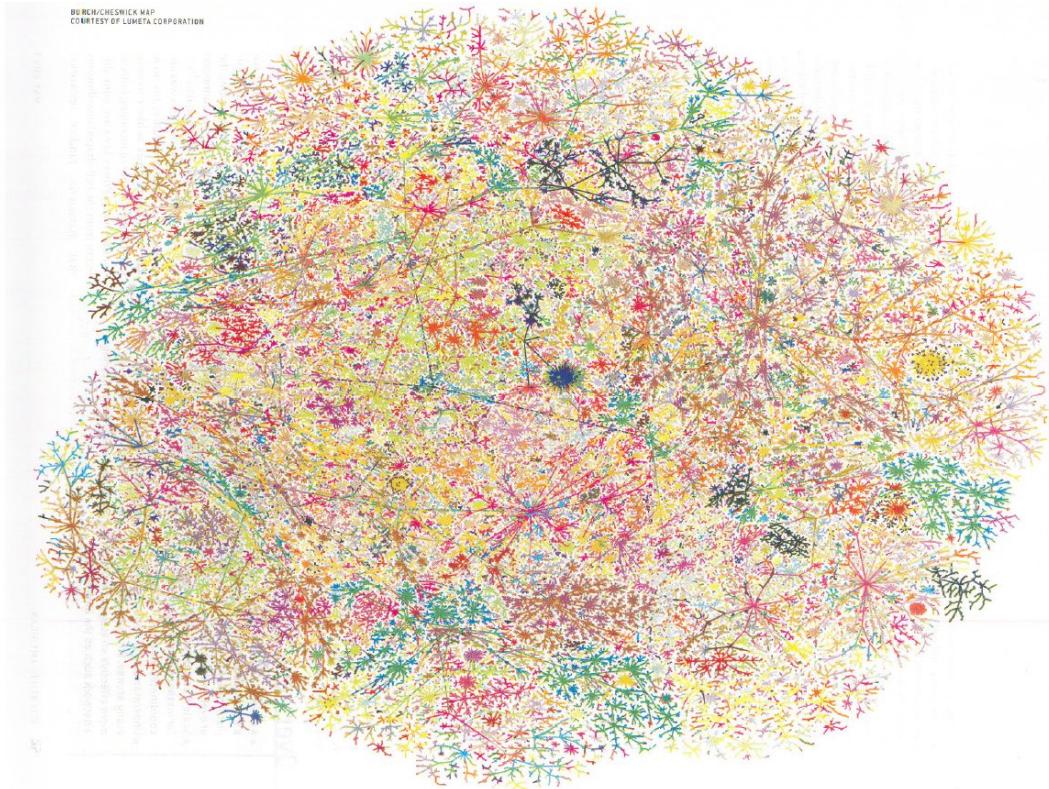




# What is network science?

One idea in network science is that any node can influence other nodes, not only their direct connections.

Such indirect influence happens through some external phenomenon—travel in a transportation network, information transfer in the Internet, vibrations in a spiderweb, etc.—and depends on how the network is connected.  
(P. Holme) <https://petterhol.me/>



# Network science

Network **measures** and network **types**

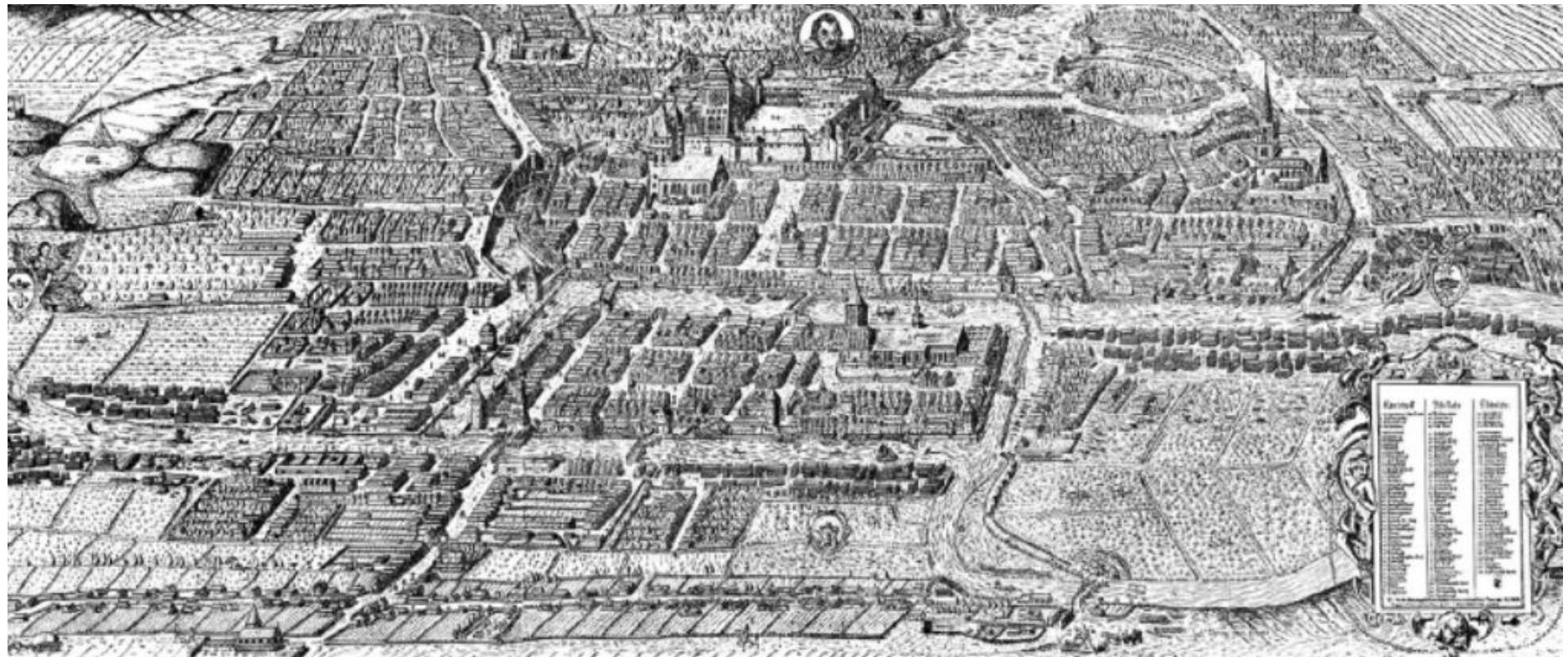
Networks in **time and space**, **dynamics on** networks

Networks from **data**

Figure 7.11

Aaron Koblin's Flight Patterns (2005): visualization of the flight paths of aircraft crossing North America

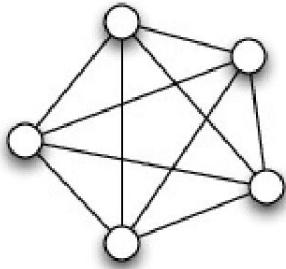
# How did the network science start?



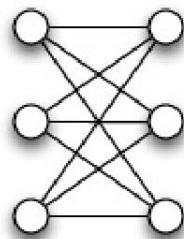
<http://networksciencebook.com/chapter/2#bridges>

# Network science, graph and topology theory

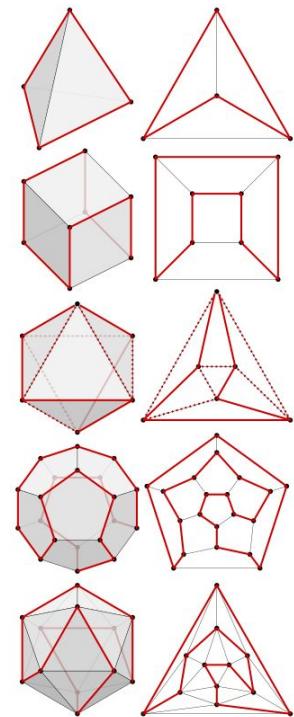
- Graph theory:  
Koenigsberg problem 1736  
Eulerian paths algorithms 1873
- Soft matter physics



$K_5$

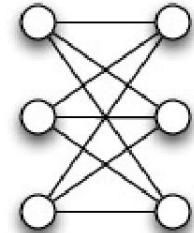
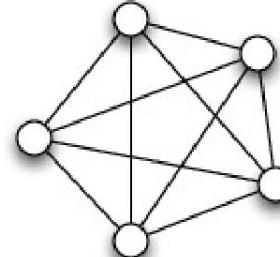


$K_{3,3}$



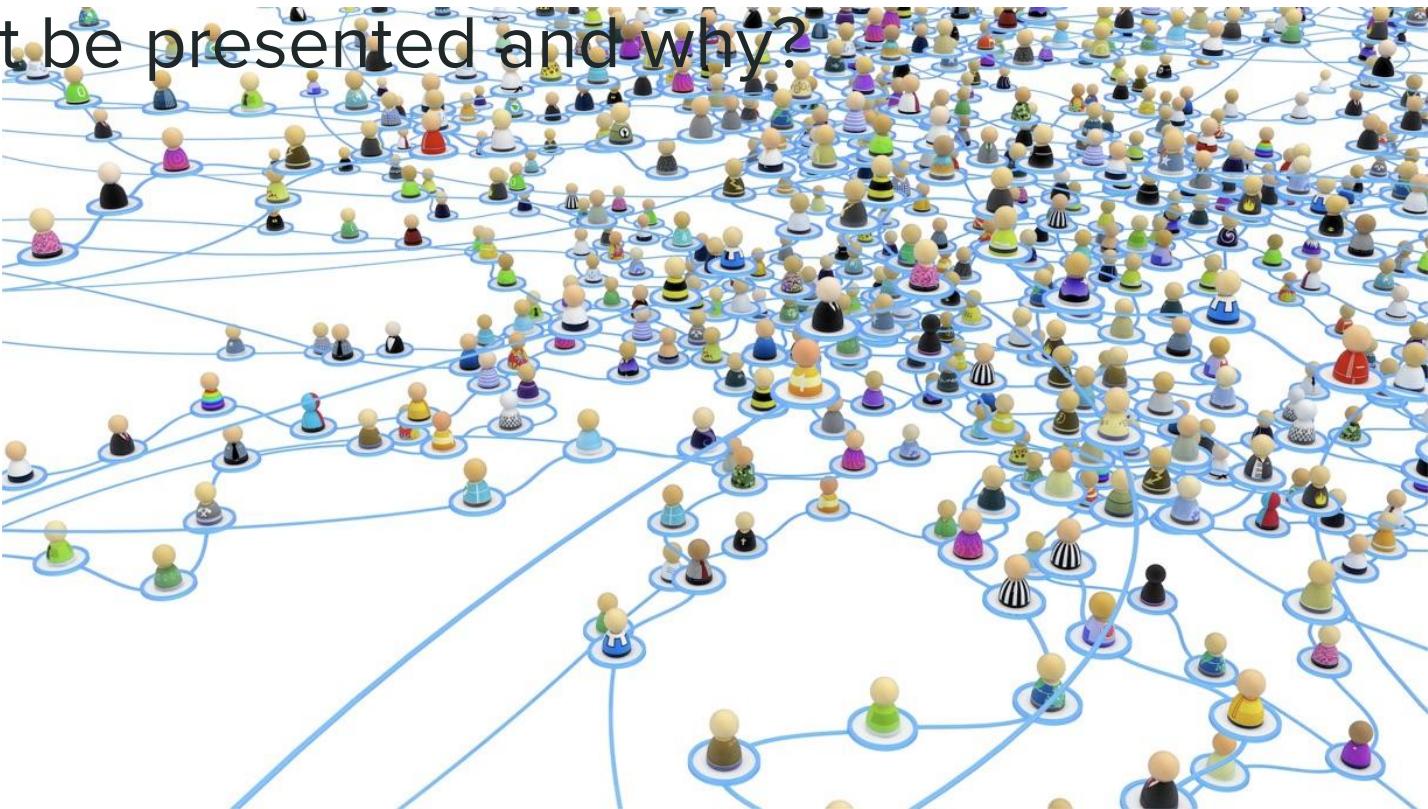
# Network science and graph theory

*Graph* (discrete mathematics),  
a structure made of vertices  $V$   
and edges  $E$  (subset of two vertices).



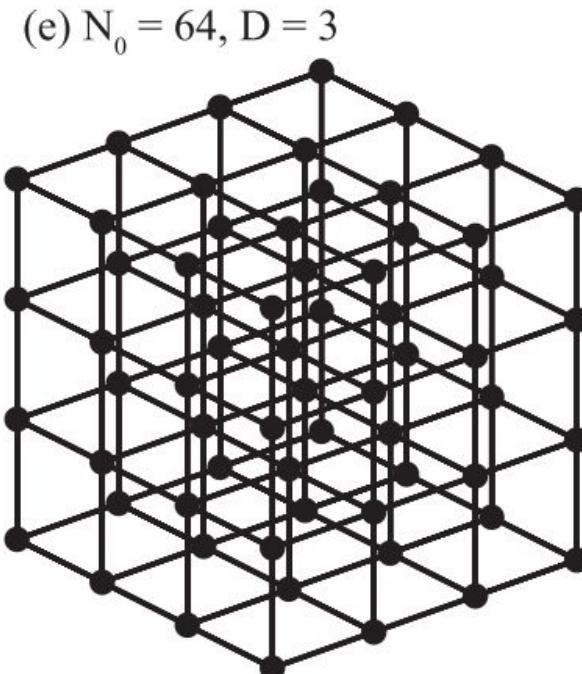
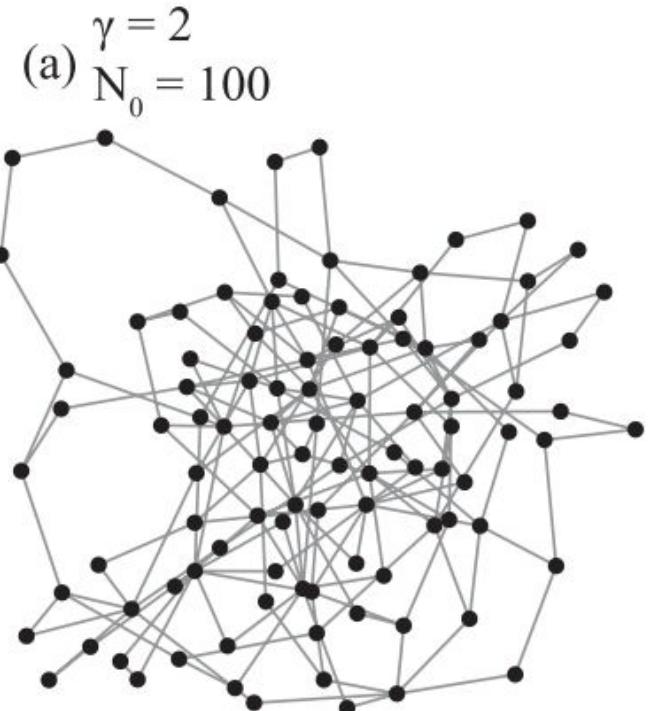
*Networks (from real world)*  
can be represented/encoded as graphs.

Can anything be presented as a network?  
What cannot be presented and why?



# Examples of network representations

M.Stella et al.



# What we will look at in network science?

1. **Network definition and measures**
2. Networks in time and space
3. Networks from data

# Defining a network to the computer

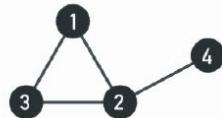
How to represent a network: **edgelist** and **adjacency matrix**.

Adjacency matrix encodes the same information about the network as edgelists.

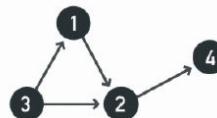
## a. Adjacency matrix

$$A_{ij} = \begin{matrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{matrix}$$

## b. Undirected network



## c. Directed network



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

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

# Network definitions

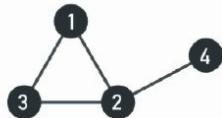
How to represent a network: **edgelist** and **adjacency matrix**.

Adjacency matrix encodes the same information about the network as edgelists.

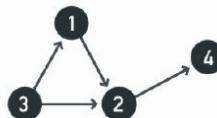
## a. Adjacency matrix

$$A_{ij} = \begin{matrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{matrix}$$

## b. Undirected network



## c. Directed network



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

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

## Questions to check:

How to encode the network data to computer?

What is the most efficient representation of a network for computer?

[Notebook](#) to encode the network

# Network definitions

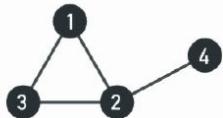
How to represent a network: **edgelist** and **adjacency matrix**.

Adjacency matrix encodes the same information about the network as edgelists.

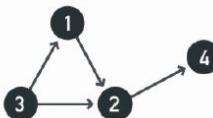
## a. Adjacency matrix

$$A_{ij} = \begin{matrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{matrix}$$

## b. Undirected network



## c. Directed network



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

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

[Notebook](#) to encode the network

```
In [ ]: from collections import Counter
from pprint import pprint
import numpy as np
import matplotlib.pyplot as plt
```

## Edge List

Let us start by defining a list of edges. This will give us our first "dataset" to work with

```
In [ ]: edge_list = [
    ('A', 'B'),
    ('A', 'C'),
    ('A', 'E'),
    ('B', 'C'),
    ('C', 'D'),
    ('C', 'E'),
    ('D', 'E')
]
```

This is a particularly useful representation as many datasets are distributed in this (or a closely related) format. From this list, we can easily measure the number of edges that constitute our network. It's main limitations are that it has no way to explicitly take into account disconnected nodes (it only accounts for nodes that are part of edges) and no indication on whether it is directed or not.

```
In [ ]: number_edges = len(edge_list)
print(number_edges)
```

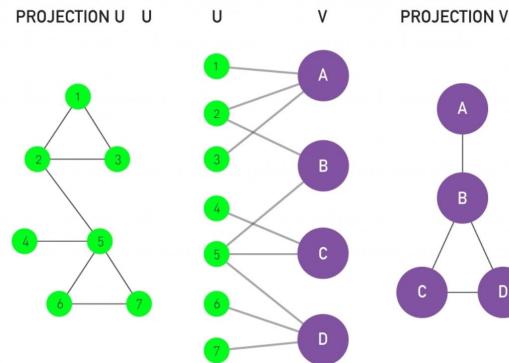
To get the number of node is a bit trickier. We must go edge by edge and keep track of all new nodes. For efficiency, we use a set to automatically remove duplicates

# Network types based on links, nodes properties

How to encode some more information which we want to include into our network?

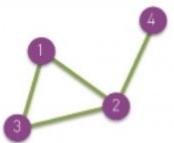
E.g. bipartite networks, or layed networks?

Try it yourself in the [Notebook](#)  
Try it with networkX precoded library [Notebook here](#)



# Network types based on links, nodes properties

a. Undirected

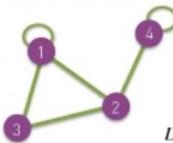


$$A_{ij} = \begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

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

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

b. Self-loops

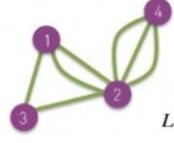


$$A_{ij} = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}$$

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

$$L = \frac{1}{2} \sum_{i,j=1, i \neq j}^N A_{ij} + \sum_{i=1}^N A_{ii} \quad ?$$

c. Multigraph  
(undirected)

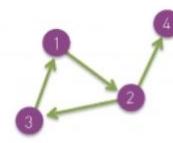


$$A_{ij} = \begin{pmatrix} 0 & 2 & 1 & 0 \\ 2 & 0 & 1 & 3 \\ 1 & 1 & 0 & 0 \\ 0 & 3 & 0 & 0 \end{pmatrix}$$

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

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

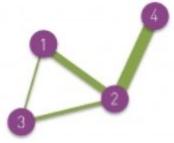
d. Directed



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

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

e. Weighted  
(undirected)

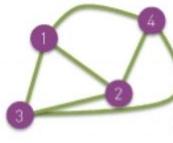


$$A_{ij} = \begin{pmatrix} 0 & 2 & 0.5 & 0 \\ 2 & 0 & 1 & 4 \\ 0.5 & 1 & 0 & 0 \\ 0 & 4 & 0 & 0 \end{pmatrix}$$

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

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

f. Complete Graph  
(undirected)

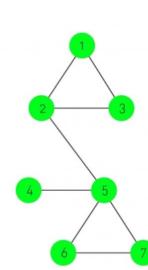


$$A_{ij} = \begin{pmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

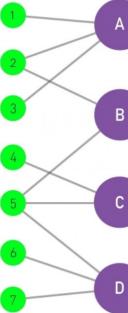
$$A_{ii} = 0 \quad A_{i \neq j} = 1$$

$$L = L_{\max} = \frac{N(N-1)}{2} \quad \langle k \rangle = N - 1$$

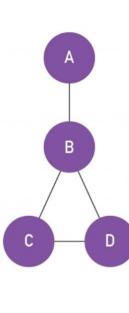
PROJECTION U



PROJECTION V



PROJECTION V



# Network types based on links, nodes properties

## Contact

Mailing list  
Issue tracker  
Source



NetworkX is a Python package for the creation, manipulation, and study of the structure, dynamics, and functions of complex networks.

## Releases

Stable (notes)

3.3 — April 2024  
[Documentation](#)

Latest (notes)

3.4 development  
[Documentation](#)

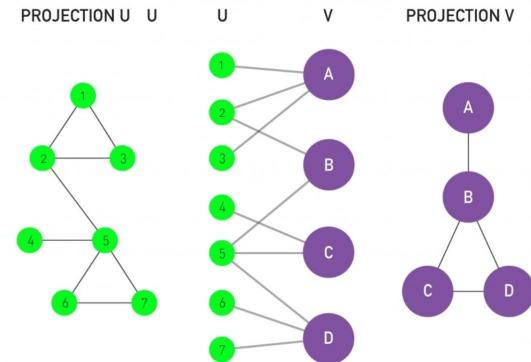
Archive

## Software for complex networks

- Data structures for graphs, digraphs, and multigraphs
- Many standard graph algorithms
- Network structure and analysis measures
- Generators for classic graphs, random graphs, and synthetic networks
- Nodes can be "anything" (e.g., text, images, XML records)
- Edges can hold arbitrary data (e.g., weights, time-series)
- Open source [3-clause BSD license](#)
- Well tested with over 90% code coverage
- Additional benefits from Python include fast prototyping, easy to teach, and multi-platform

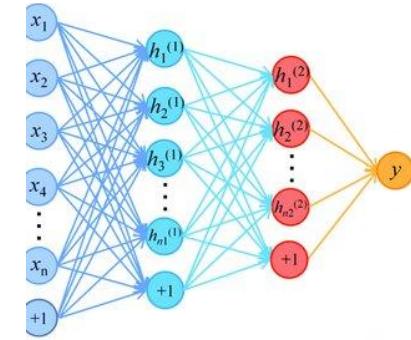
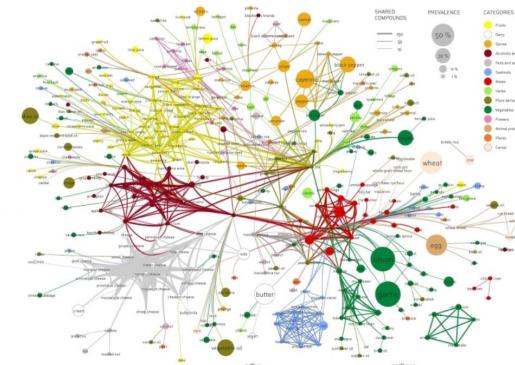
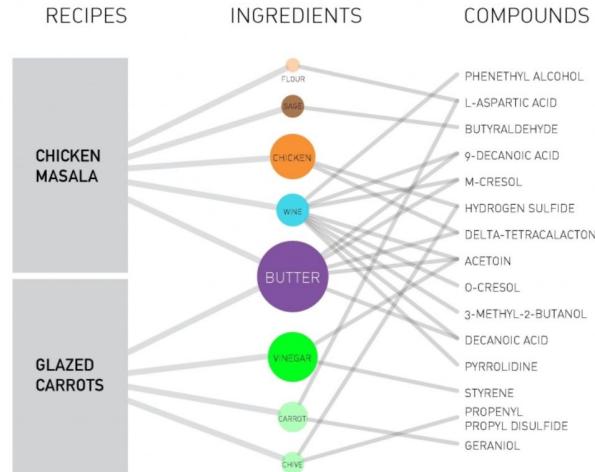
©2014-2024, NetworkX developers.

<https://networkx.org/documentation/stable/tutorial.html>



# Network types based on links, nodes properties

## bipartite networks

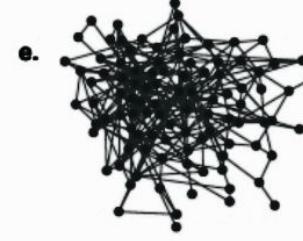
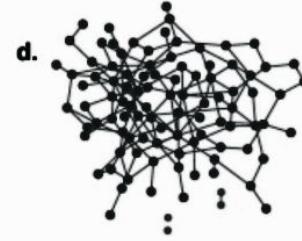


# Network measures

Main idea is to characterise their properties.

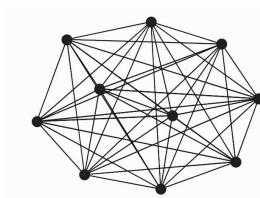
**Local measures** for each node (degree)

**Global measures** for the whole network (density - number of links normalised by number of links in a complete graph)



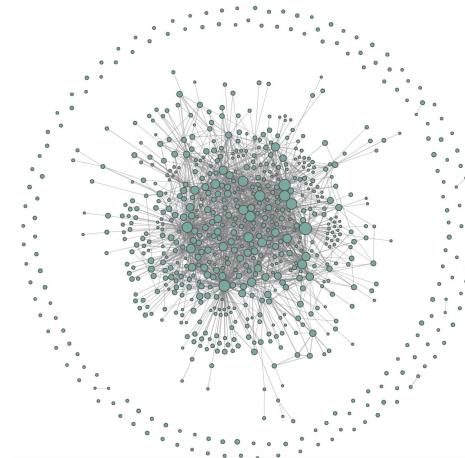
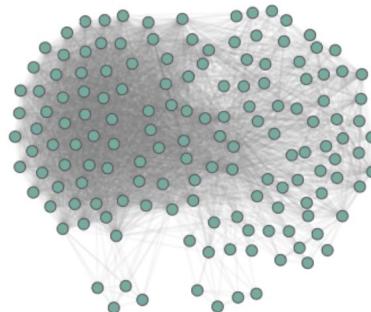
low

high



# Network measures and layout

**Local measures** for each node. **Global measures** for the whole network.  
Layouts of the same network (left, right).



# Network statistics



# Network statistics

What are nodes with highest centrality?

Size       $n = 34$

Volume      $m = 78$

Loop count     $l = 0$

Triangle count     $t = 45$

Square count     $q = 154$

Maximum degree     $d_{\max} = 17$

Average degree     $d = 4.588$

Size of Large Connected Component  $N = 34$

Diameter     $\delta = 5$

Median distance     $\delta_M = 2$

Mean distance     $\delta_m = 2.443$

Gini coefficient     $G = 0.385$

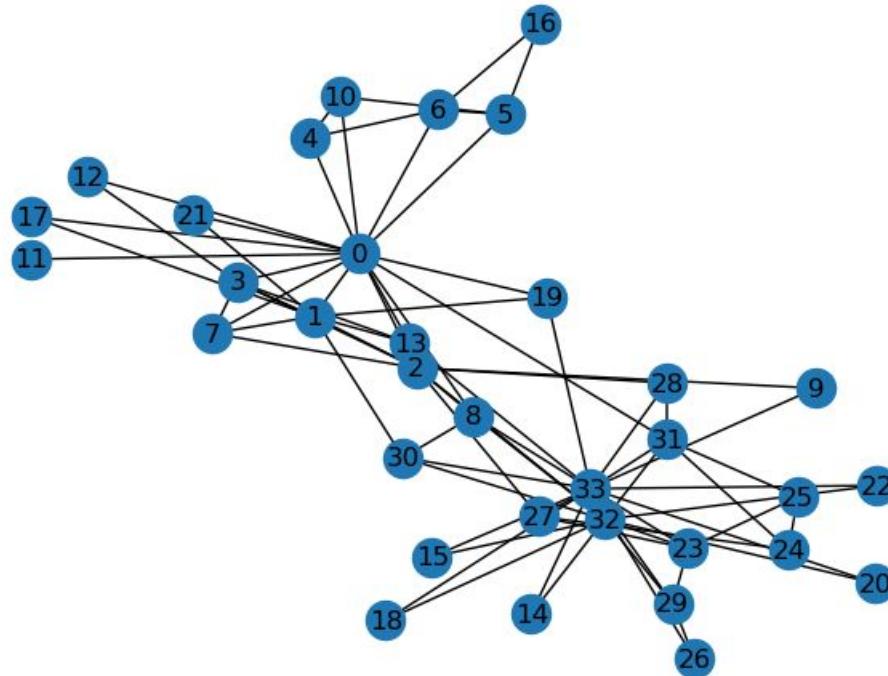
Power law exponent     $\gamma = 1.780$



# Quick check-in

What are network measures for this network?

What node would have the highest betweenness centrality?  
What would be the best spreader?

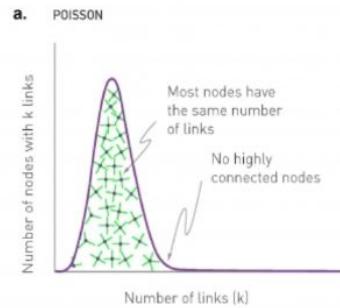
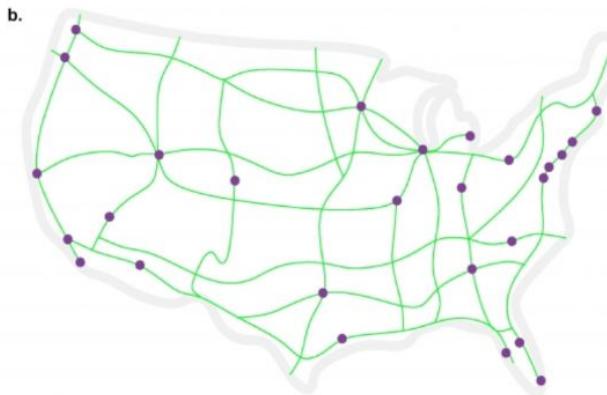


# Network statistics

**Degree measure** - is a local measure to characterise how many nodes each node is connected to.

How to look into degree for N nodes?

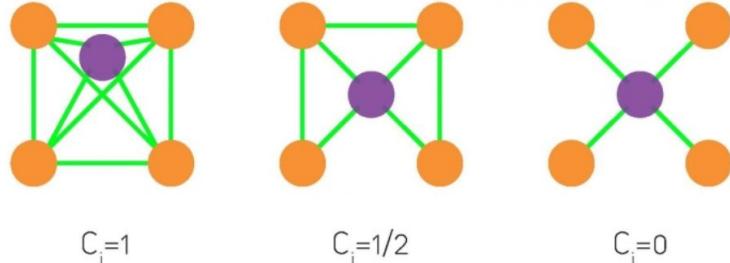
Looking into the degree distribution: plotting how many nodes have degree= $k$ .



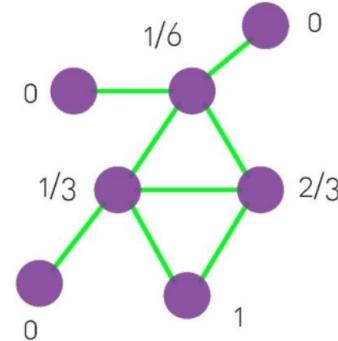
# Network statistics

**Clustering coefficient** is a measure of the degree to which nodes in a graph tend to cluster together.

a.



b.



$$\langle C \rangle = \frac{13}{42} \approx 0.310$$
$$C_{\Delta} = \frac{3}{8} = 0.375$$

# Network statistics

**Clustering coefficient** is a measure of the degree to which nodes in a graph tend to cluster together.

a.

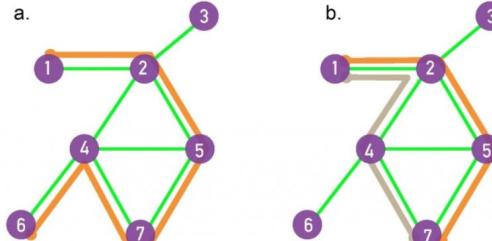
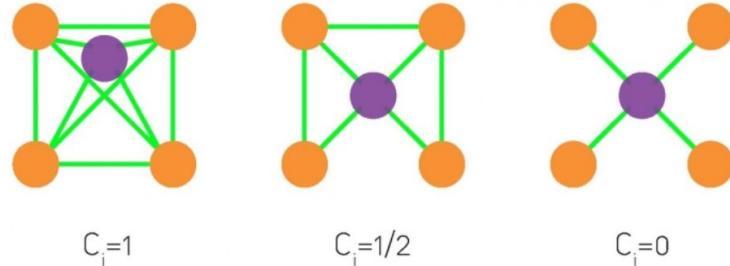


Image 2.12

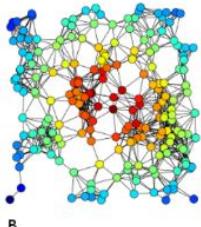
## Paths

- A path between nodes  $i_0$  and  $i_n$  is an ordered list of  $n$  links  $P = \{(i_0, i_1), (i_1, i_2), (i_2, i_3), \dots, (i_{n-1}, i_n)\}$ . The length of this path is  $n$ . The path shown in orange in (a) follows the route 1 → 2 → 5 → 7 → 4 → 6, hence its length is  $n = 5$ .

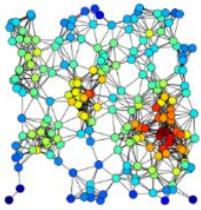
**Shortest path finding between nodes are used in many algorithms for networks.**

Example path between node 1 and node 6 in a graph is then encoded as a sequence of nodes, e.g. (1,2,5,7,4,6). One of the most known shortest path algorithm is Dijkstra's algorithm (1956).

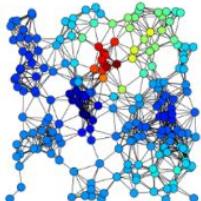
# Network measures



B



D



F

**Closeness centrality**

**Degree centrality**

**Katz centrality**

TABLE 2: Definitions of network science terms and variables.

Term/variable	Definition
$N$	number of nodes, $N$ , in graph
$E$	number of edges, $E$ , in graph
network density	ratio of the number of edges to the maximum number of possible edges $\frac{2E}{N(N-1)}$
$d(n_i, n_j)$	shortest path between node $i$ and node $j$ $d(n_i, n_j)$ where $n_i, n_j \in N$
shortest path length, $L$	average length of shortest path between pairs of nodes $L = \frac{1}{N(N-1)} * \sum_{i,j} d(n_i, n_j)$
$D$	largest shortest path between nodes $D = \max_{n_i \in N, n_j \in N} d(n_i, n_j)$
centrality	inverse of the sum of the length of the shortest paths between node $i$ and all other nodes in the graph $C_i = \frac{1}{\sum_j d(n_i, n_j)}$
degree, $\langle k \rangle$	number of edges attached to node $i$ average number of edges per node in network $\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i$
clustering coefficient, $c_i$	number of edges between the neighbors of node $i$ divided by the maximum number of edges between those neighbors $c_i = \frac{2 e_{j,k} }{k_i(k_i - 1)}$ where $n_j, n_k \in N_D$ , $e_{jk} \in E$
clustering coefficient, $\langle C \rangle$	average clustering coefficient of nodes in the network $\langle C \rangle = \frac{1}{N} \sum_{i=1}^N c_i$
clarity, $Q$	proportion of edges that fall within subgroups of nodes minus the expected proportion if edges were randomly distributed, range $[-1, 1]$
efficiency, $E_G$	measure of how efficiently information is exchanged in the network $E_G = \frac{1}{n(n-1)} \sum_{i,j \in N} \frac{1}{d(n_i, n_j)}$
connected component	largest group of nodes in the network that are connected to each other in a single component
distribution, $P(k)$	probability distribution of node degrees in the network power-law exponent for the degree distribution
erdos structure	network with short average path lengths and relatively high clustering coefficient (relative to a random graph with similar density)
network	network with a degree distribution that is power-law distributed

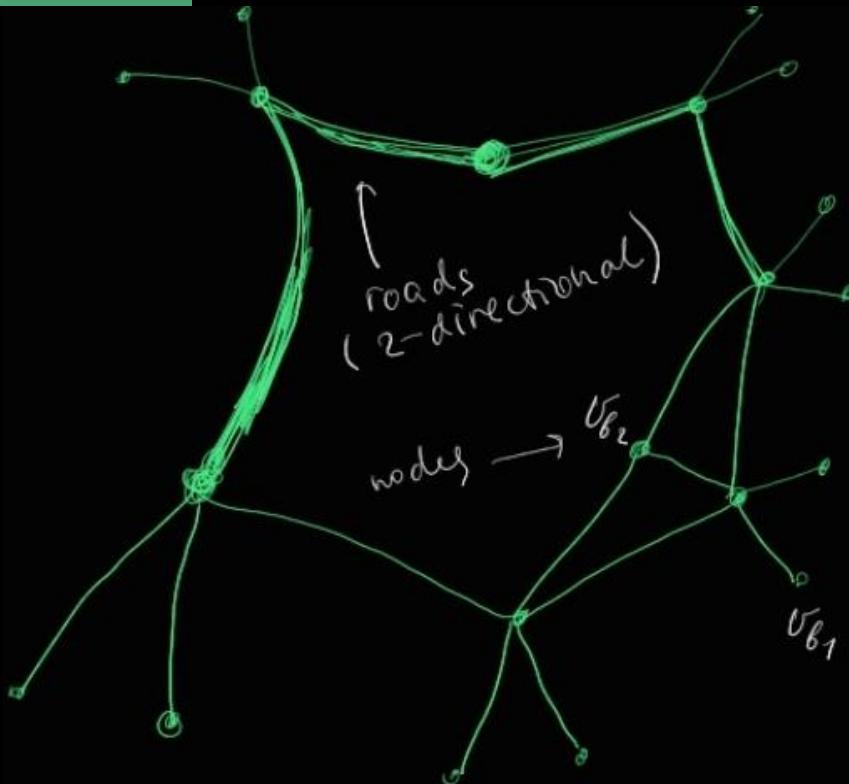
# Network measures

Most of the measures can be estimated directly using networkx python library.

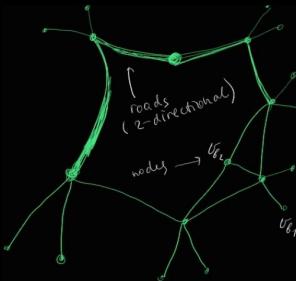
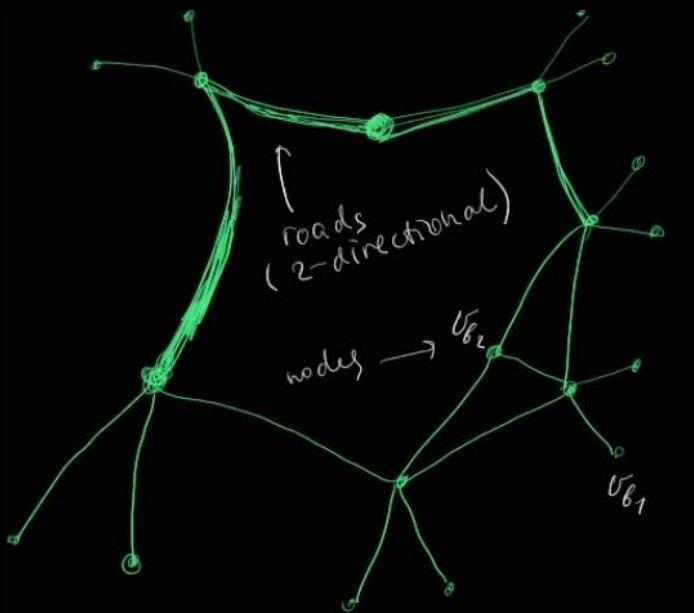
TABLE 2: Definitions of network science terms and variables.

Term/variable	Definition
$N$	number of nodes, $N$ , in graph
$E$	number of edges, $E$ , in graph
network density	ratio of the number of edges to the maximum number of possible edges $\frac{2E}{N(N-1)}$
distance, $d(n_i, n_j)$	shortest path between node $i$ and node $j$ $d(n_i, n_j)$ where $n_i, n_j \in N$
average shortest path length, $L$	average length of shortest path between pairs of nodes $L = \frac{1}{N(N-1)} \cdot \sum_{i,j} d(n_i, n_j)$
diameter, $D$	largest shortest path between nodes $D = \max_{n_i \in N, n_j \in N} d(n_i, n_j)$
closeness centrality	inverse of the sum of the length of the shortest paths between node $i$ and all other nodes in the graph $C_i = \frac{1}{\sum_j d(n_i, n_j)}$
degree, $k_i$	number of edges attached to node $i$
average degree, $\langle k \rangle$	average number of edges per node in network $\langle k \rangle = \frac{1}{N} \sum_{i=1}^N k_i$
local clustering coefficient, $c_i$	number of edges between the neighbors of node $i$ divided by the maximum number of edges between those neighbors $c_i = \frac{2 e_{ji} }{k_i(k_i - 1)}$ where $n_j, n_k \in N$ , $e_{jk} \in E$
average clustering coefficient, $\langle C \rangle$	average clustering coefficient of nodes in the network $\langle C \rangle = \frac{1}{N} \sum_{i=1}^N c_i$
modularity, $Q$	proportion of edges that fall within subgroups of nodes minus the expected proportion if edges were randomly distributed, range $[-1, 1]$
average efficiency, $E_G$	measure of how efficiently information is exchanged in the network $E_G = \frac{1}{n(n-1)} \sum_{i \neq j, i, j \in N} \frac{1}{d(n_i, n_j)}$
largest connected component	largest group of nodes in the network that are connected to each other in a single component
degree distribution, $P(k)$	probability distribution of node degrees in the network
$\gamma$	power-law exponent for the degree distribution
Small world structure	network with short average path lengths and relatively high clustering coefficient (relative to a random graph with similar density)
scale-free network	network with a degree distribution that is power-law distributed

# Getting intuition to work with networks

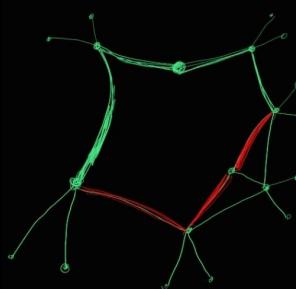


# Getting intuition to work with networks



Given city graph  $G(V, E)$   
we want to characterize  
its bottleneck nodes  $v_b \in V$ .

② Topological characterization  
of bottleneck node  $v_b \in V$ .  
given that for  $\forall v_i, v_j \in V$   
 $b(v_b) = \frac{s_{u_b}(v_i, v_j)}{s(v_i, v_j)}$ , where  $s_{u_b}(v_i, v_j)$  is number of shortest paths from  $v_i$  to  $v_j$  via node  $v_b$ .  
 $b(v_{b_1}) \ll b(v_{b_2})$ ;  $b(v_b)$  is non-local.



③ Given weight on each edge denoted  $\{w_{ij}(t)\}$   
we characterize how bad on average this  
node is in terms of slowing down  $d(v_b)$ .

$$d(v_b) = \sum_{i \neq b} \langle w_{ij}(t) \rangle / \max(w_{ij}(t))$$



$d(v_b)$  is local measure not like betweenness.  
Another measure is deviation from average e.g. if deviation in node  
adjacent to  $v_b$  is high, then, it's a bottleneck.  
 $d(v_b) = \sum \text{div}(w_{ij}(t)) / \max(\text{div}(w_{ij}(t)))$

④ Given time spent on passing on each edge  $t_{ij}(t)$   
we can find contribution of node  $v_b$  to the expected  
path deviation in terms of time.

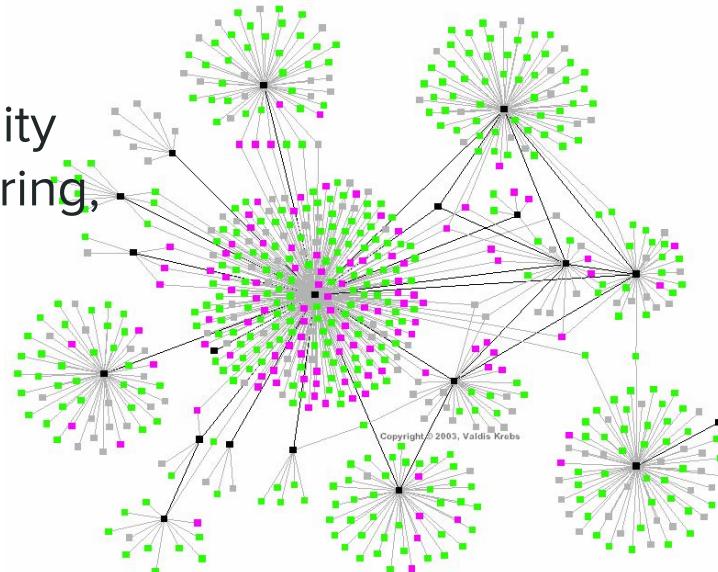
# Getting intuition to work with networks

Bottlenecks - high betweenness centrality

Outliers - distributions of degree, clustering,

Betweenness

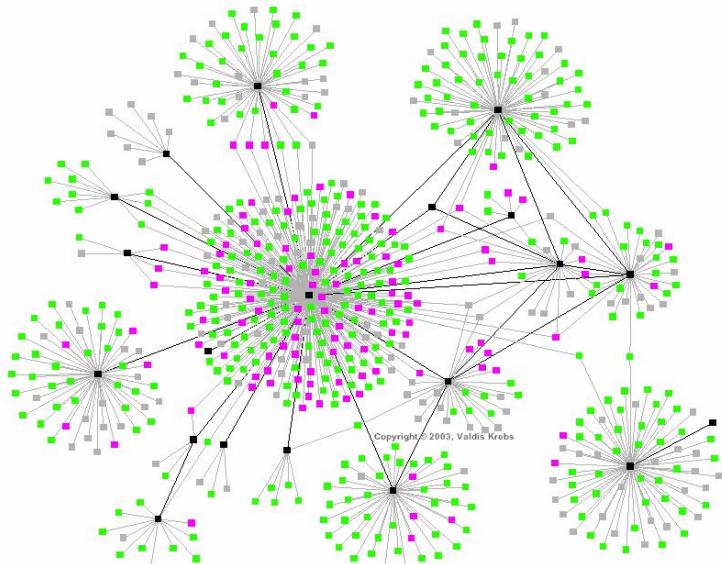
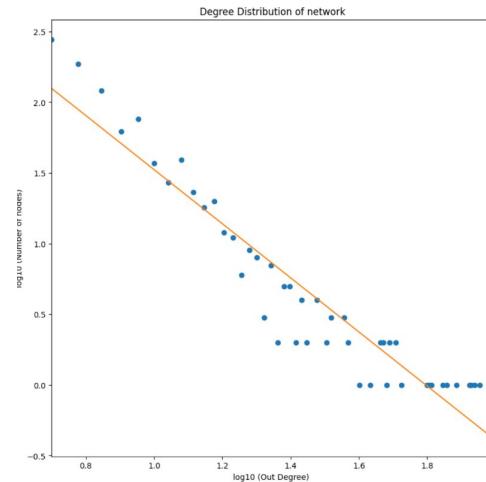
Hubs - nodes with high degree



# Getting intuition to work with networks

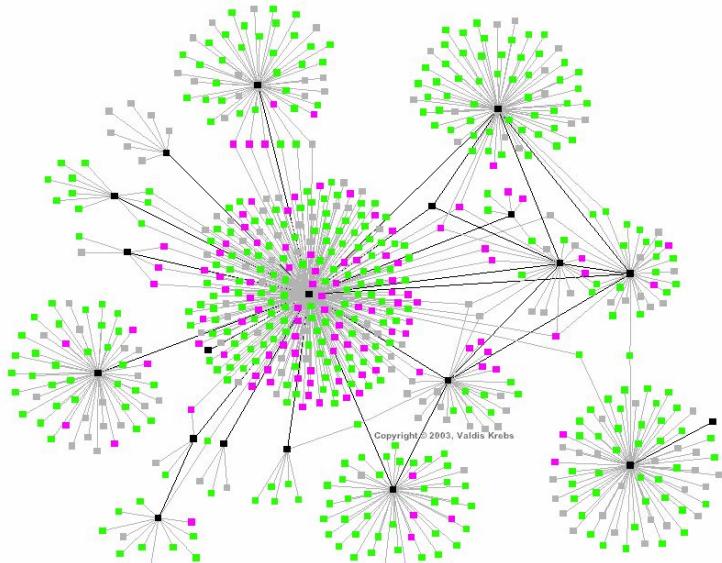
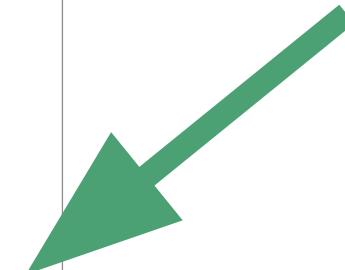
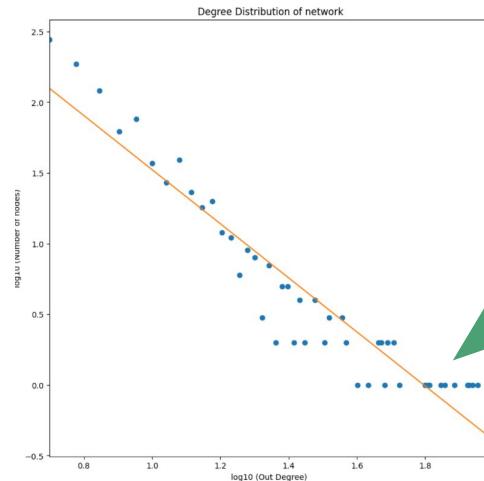
Bottlenecks - high betweenness centrality

Outliers - distributions of degree, clustering,  
betweenness



# Getting intuition to work with networks

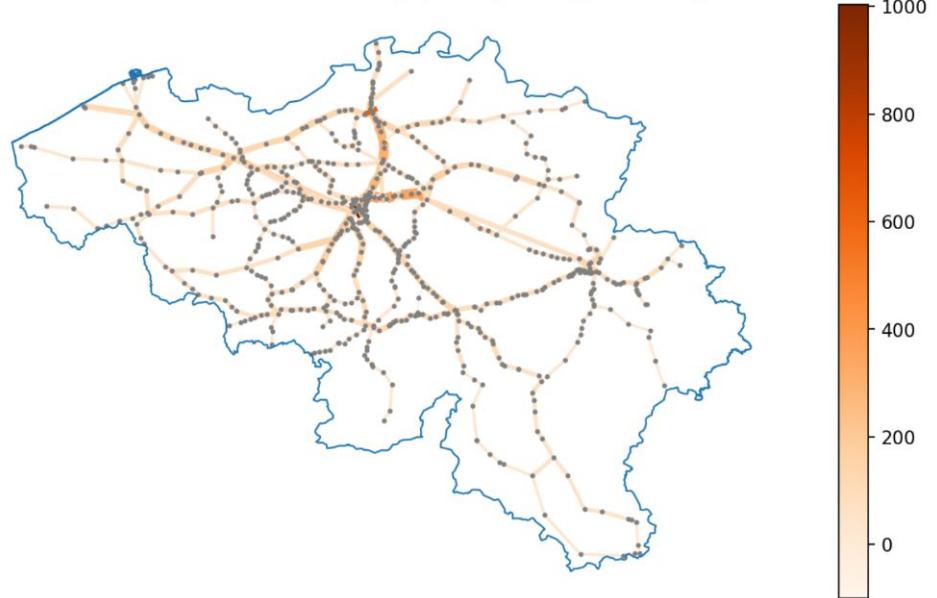
Outliers - from distributions of degree, clustering, betweenness  
Notebooks from [github](#)



# Getting intuition to work with networks

Dataset of train network

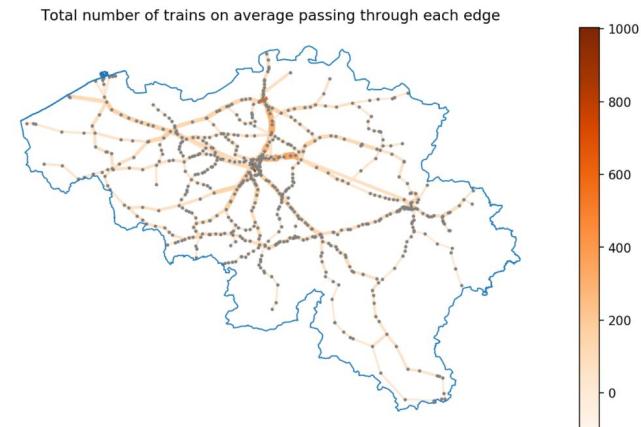
Total number of trains on average passing through each edge



# Getting intuition to work with networks

Global networks measures: given any network centrality we can estimate global (or average network measure).

Dataset of train network

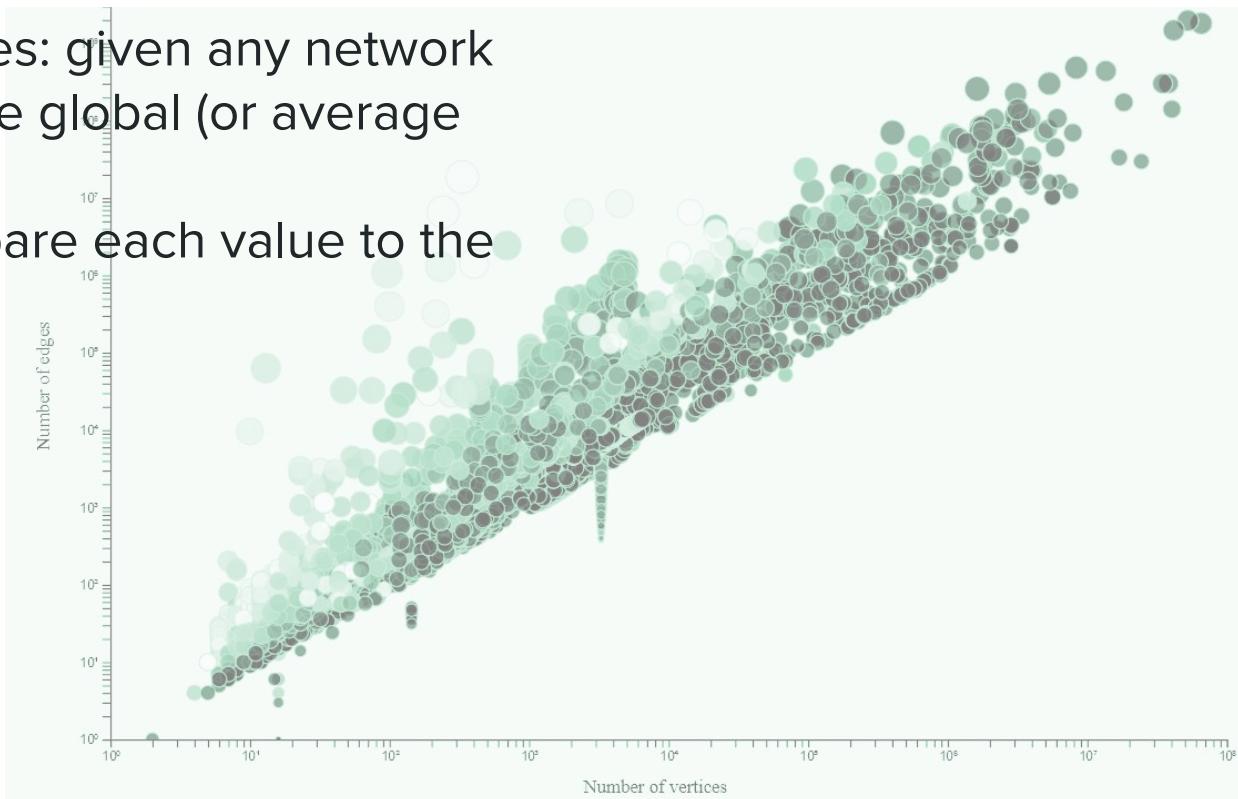


# Choosing a network to consider

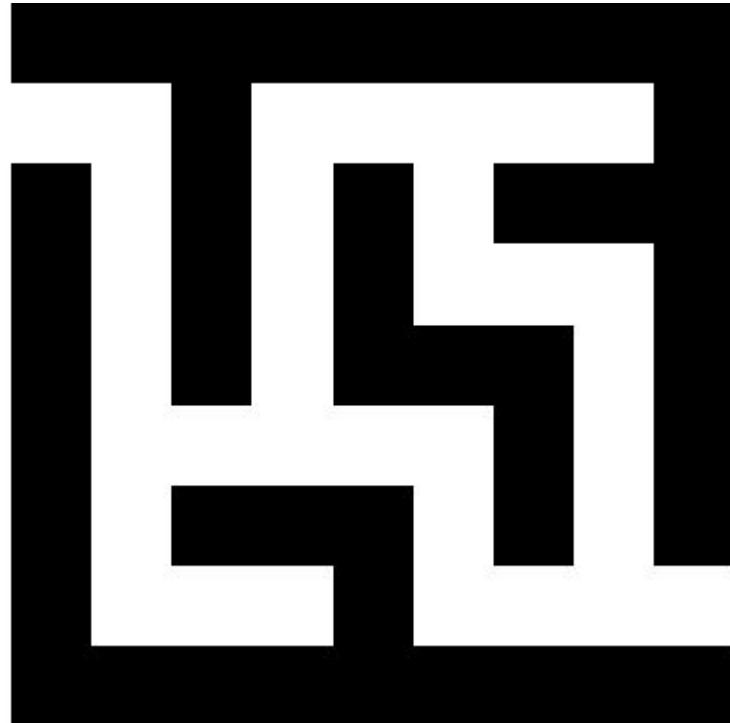
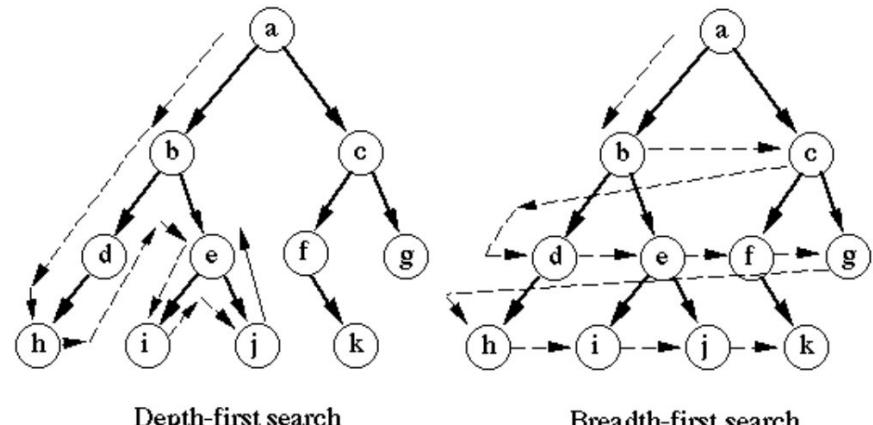
Global networks measures: given any network centrality we can estimate global (or average network measure).

The main idea is to compare each value to the average.

## Datasets



# Algorithms on networks

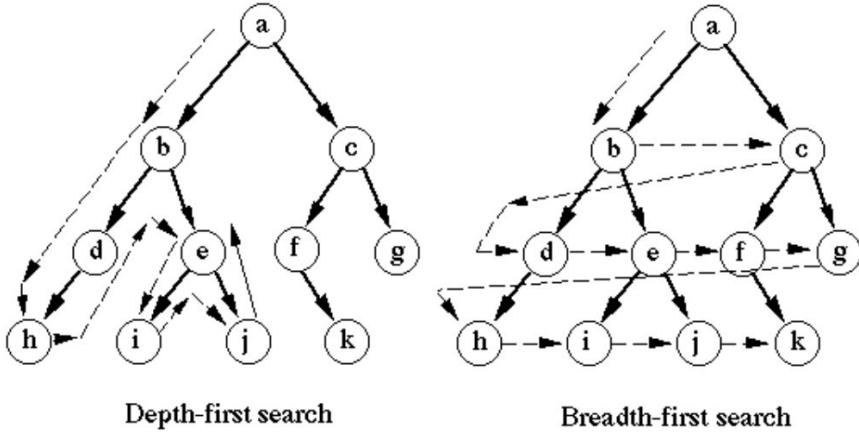


# Algorithms on networks

BFS algorithm

Pseudocode idea

Or why your code on  
large networks is so slow?



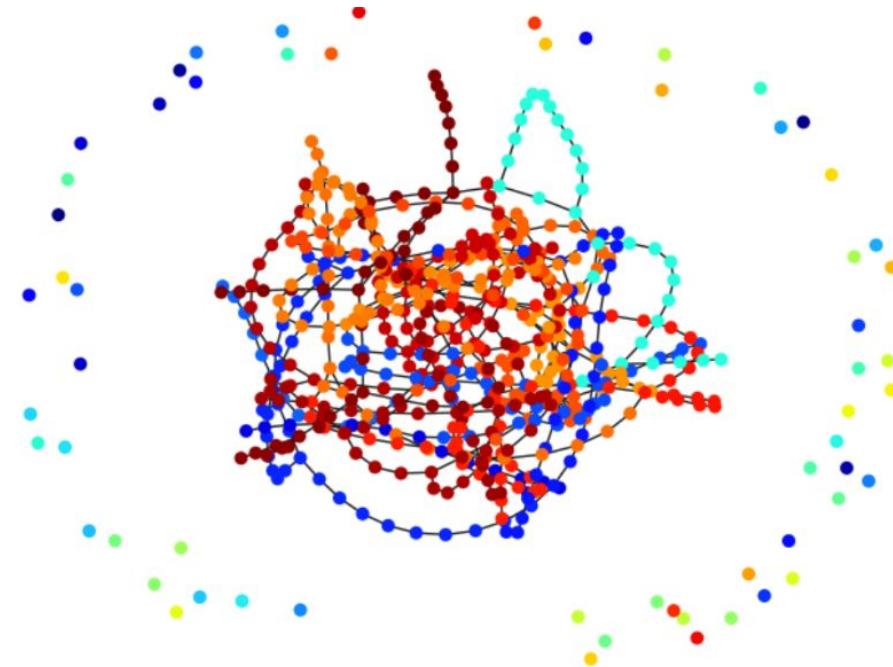
# Algorithms on networks

Louvain algorithm

Pseudocode idea

Communities in networks

See code in [classroom](#)



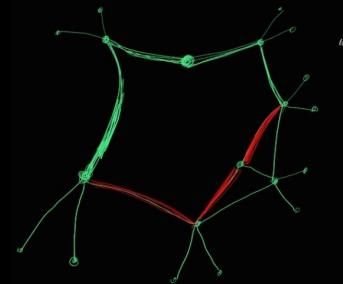
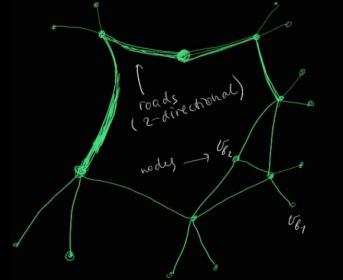
# Hands-on notebooks

Github

[https://github.com/Big-data-course-CRI/materials\\_big\\_data\\_cri\\_2024\\_2025/tree/main/day%201%20networks%20and%20hypergraphs](https://github.com/Big-data-course-CRI/materials_big_data_cri_2024_2025/tree/main/day%201%20networks%20and%20hypergraphs)

Notebook on networks generation and basic network measures

<https://colab.research.google.com/drive/1WmwG30LMmkSoOP5Uc7YIXWIFtM68TIm6?usp=sharing>



# Hands-on notebooks

## Github

[https://github.com/Big-data-course-CRI/materials big data cri 2024 2025/tree/main/day%201%20networks%20and%20hypergraphs](https://github.com/Big-data-course-CRI/materials_big_data_cri_2024_2025/tree/main/day%201%20networks%20and%20hypergraphs)

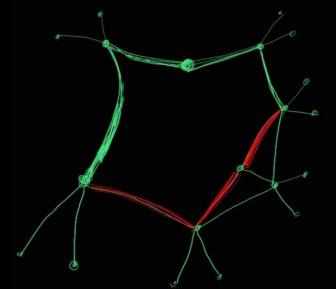
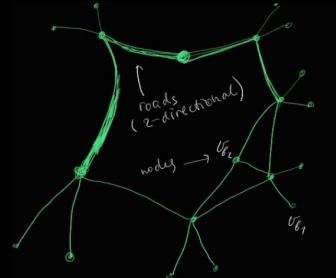
Extra notebook on geopandas and networks

<https://drive.google.com/file/d/1PVFUTS0CuFj4whjUqt7SDROqXb6QQpRK/view?usp=sharing>

Extra notebook on hypergraphs

[https://colab.research.google.com/drive/1bc633d1b5tBtletFJ57nWFYPV\\_JrGahO?usp=sharing](https://colab.research.google.com/drive/1bc633d1b5tBtletFJ57nWFYPV_JrGahO?usp=sharing)

Extra notebook on neural networks



# Random networks



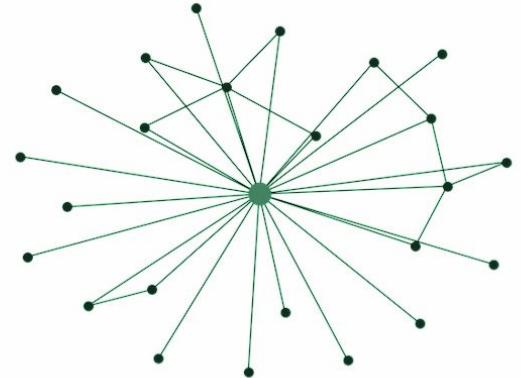
# Random networks: model by Erdős (1913-1996) and Rényi (1921-1970)

1. Create N nodes
2. Connect each pair of N labeled nodes with probability p. You can do it yourself by tossing a coin each time.

Corresponding class in networkx:

```
G_er = nx.erdos_renyi_graph(n, p2)
```

# Random networks: model by Erdős (1913-1996) and Rényi (1921-1970)

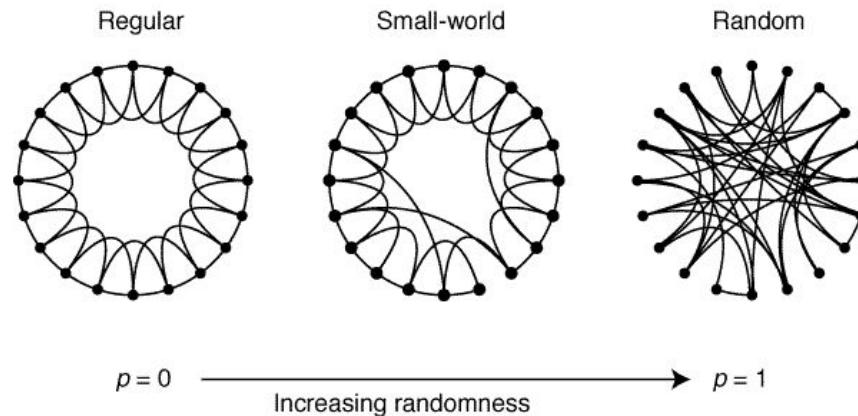
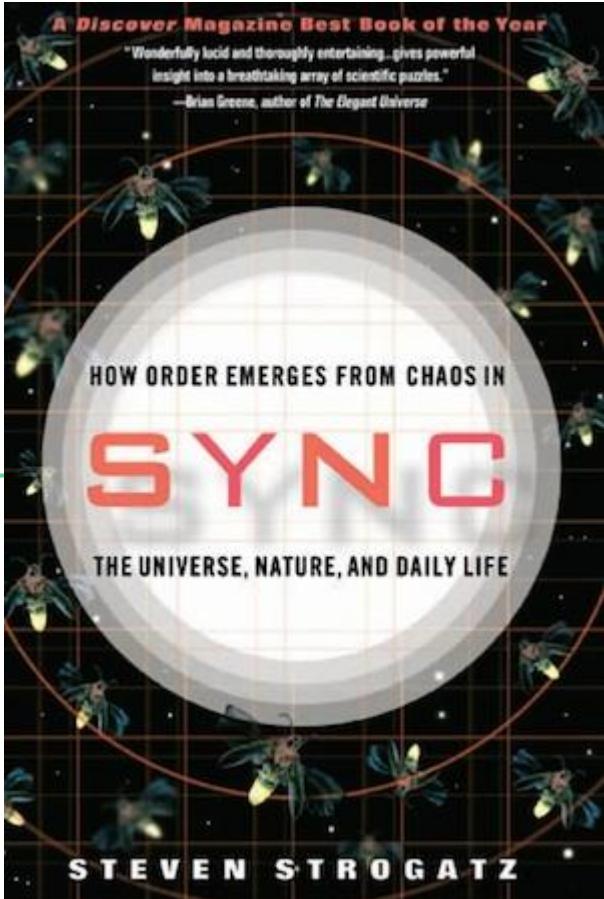


1. Create N nodes
2. Connect each pair of N labeled nodes with probability p. You can do it yourself by tossing a coin each time.

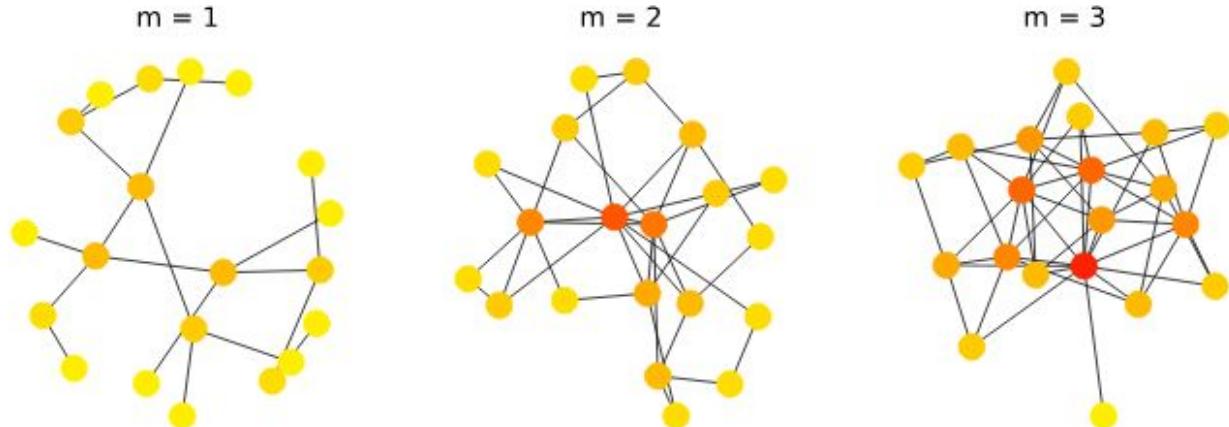
Corresponding class in networkx:

```
G_er = nx.erdos_renyi_graph(n, p2)
```

# Random networks: Watts-Strogatz network



# Random networks: model by Barabasi Albert



A graph of nodes is grown by attaching new nodes each with  $m$  edges that are preferentially attached to existing nodes with high degree.

A. L. Barabási and R. Albert "Emergence of scaling in random networks", Science, 1999.

## Random networks: model Barabasi Albert

Network G on N nodes and m edges preferential attachment.  
Model by Barabasi and Albert creates a random network with  
algorithm:

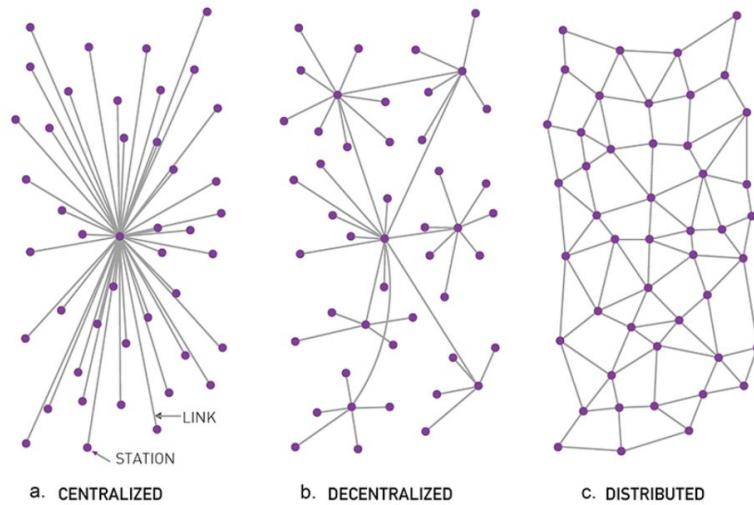
1. Create starting nodes
2. Connect a new node with m edges to existing nodes
3. Repeat (2.) x times for all non existing nodes

# Network and robustness

Network science requires intuition, e.g.  
how to construct a network, such that it  
would have a specific property, e.g.  
robustness, or particular distribution?

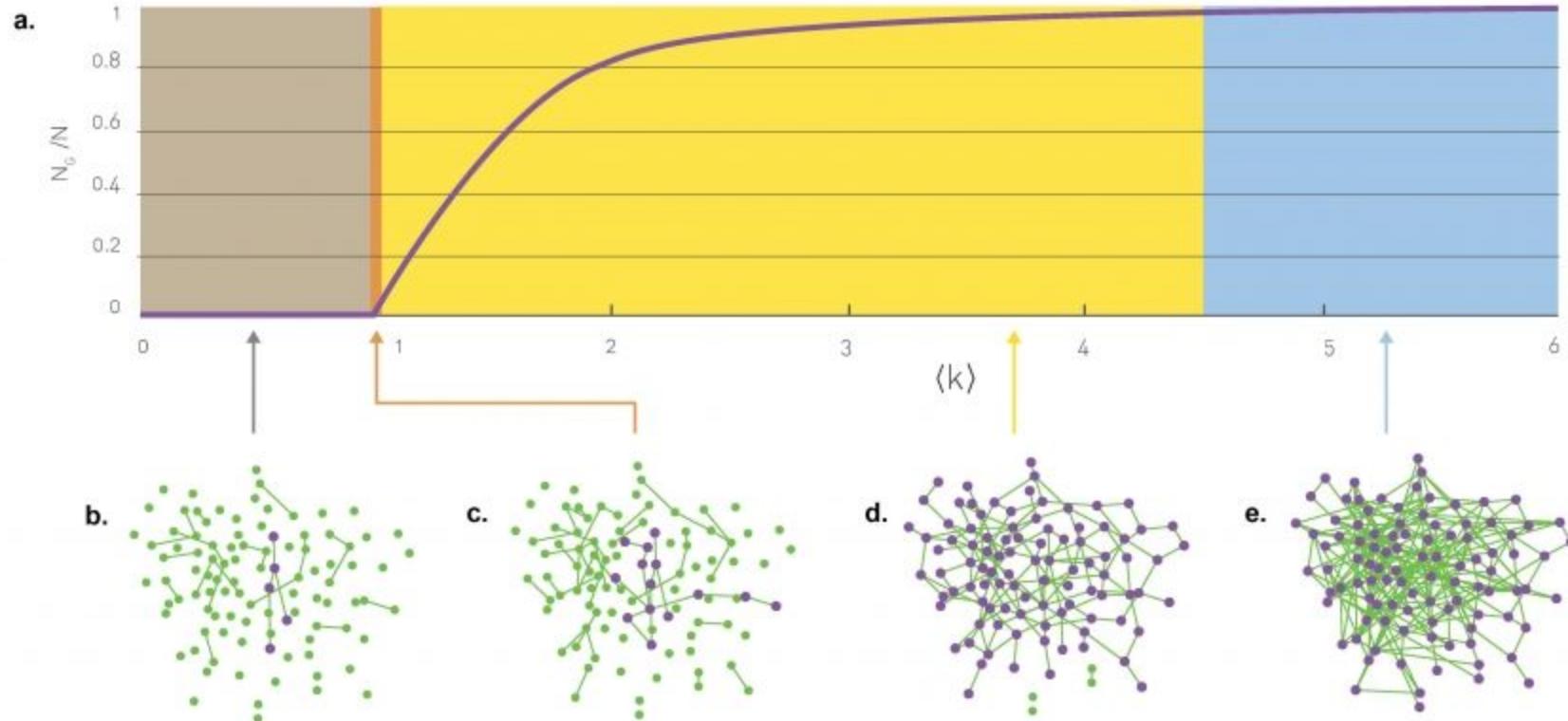
Are random networks robust?

Fig. credits P. Barran



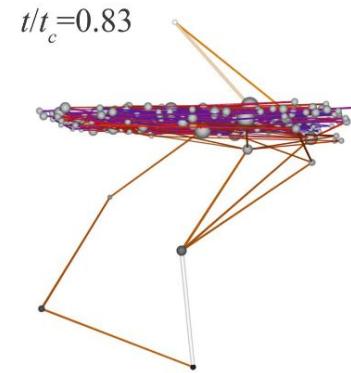
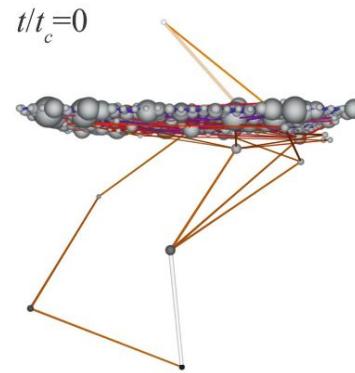
# Example of a universal law at the collective scale

Emergence of a giant component in a network above a threshold of number of links

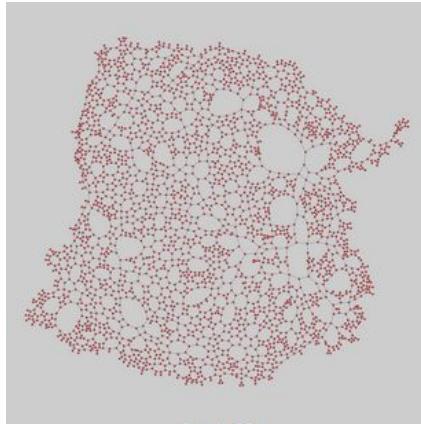


# Networks in time and space

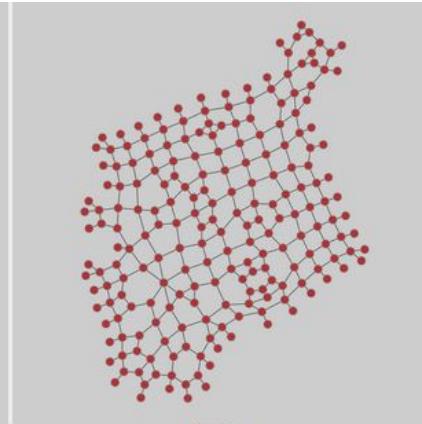
Percolation of networks in time  
Nat.Comm. 2020



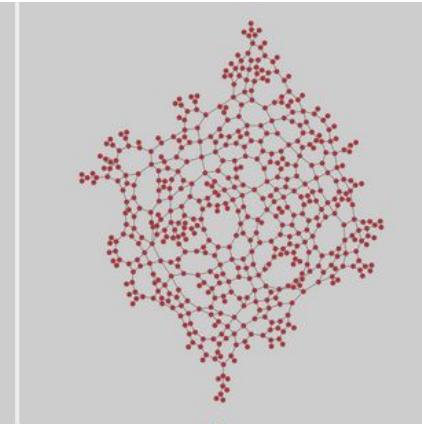
# Networks in time and space



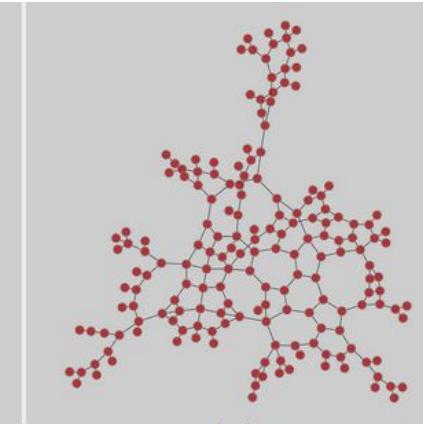
ahmedabad



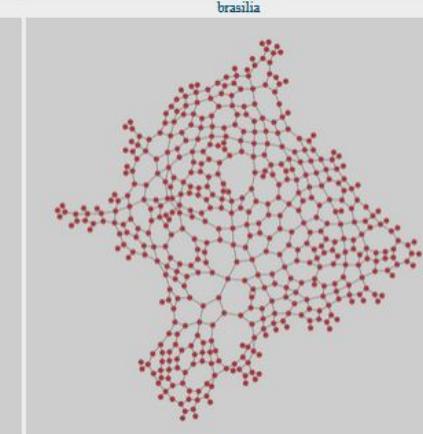
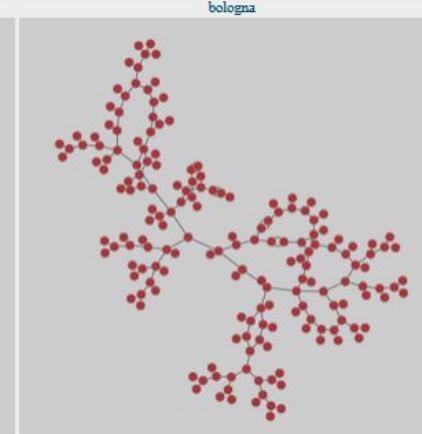
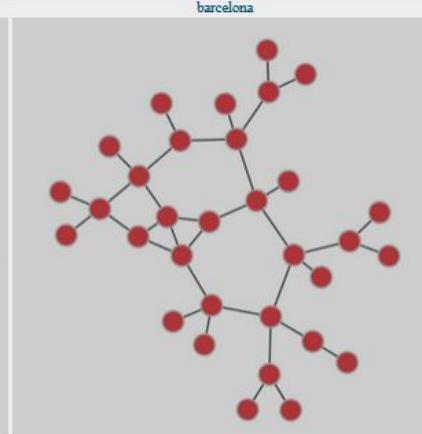
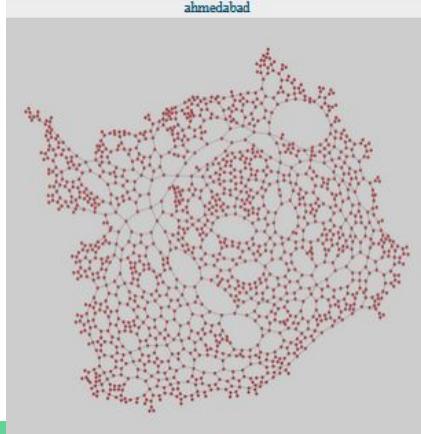
barcelona



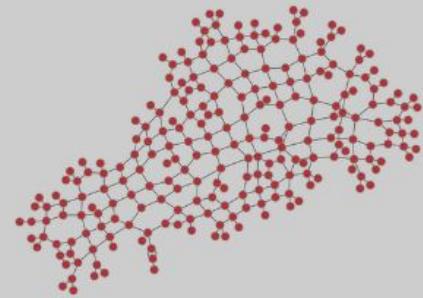
bologna



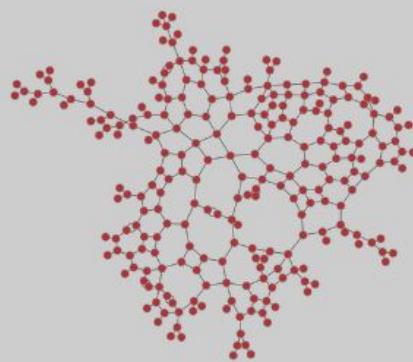
brasilia



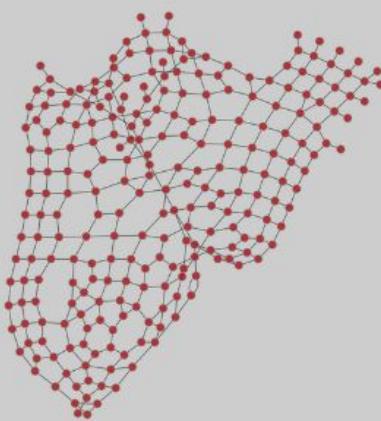
# Networks in time and space



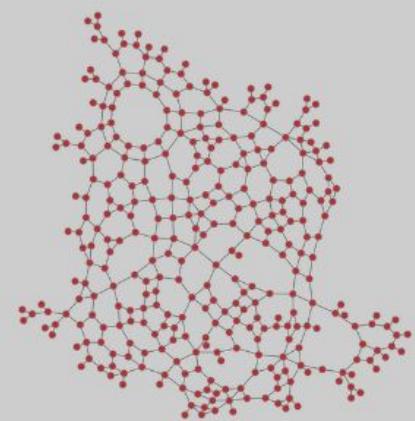
los-angeles



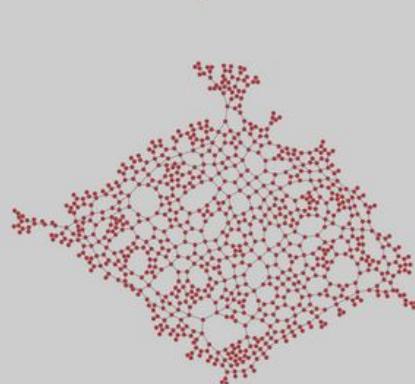
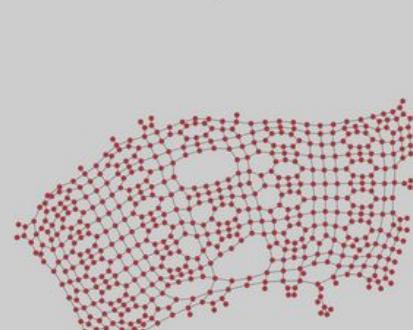
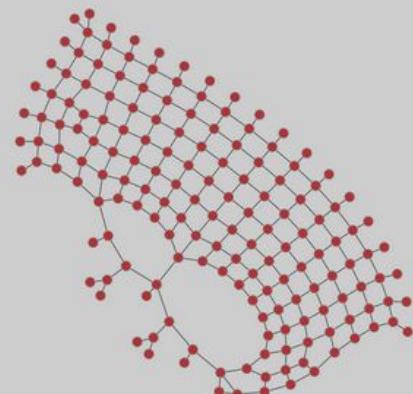
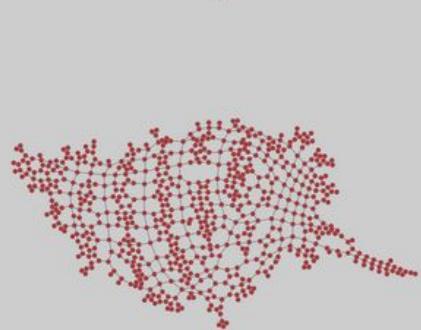
new-delhi



new-york



paris

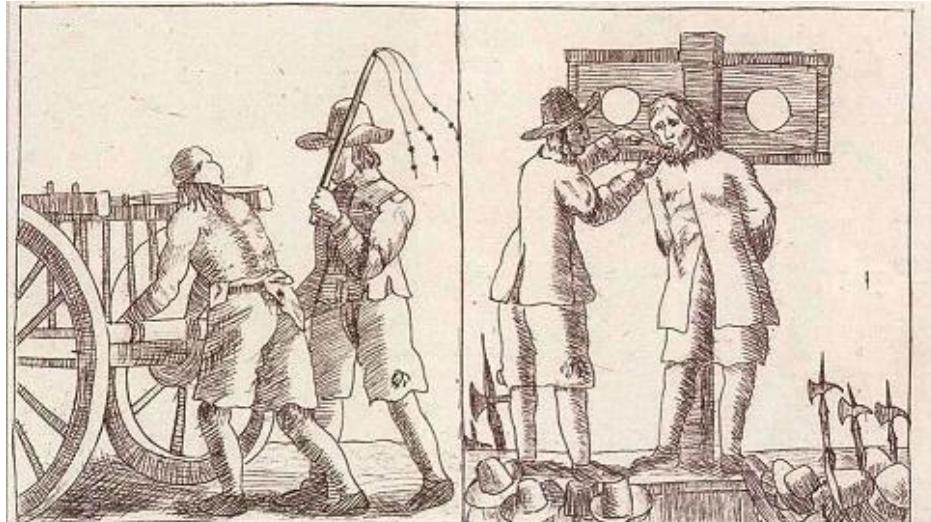


# Quick check-in

Can networks tell a new story about your data?

Typical examples:

Quakers, people who belong to a historically Protestant Christian set of denominations known formally as the Religious Society of Friends.



James Nailor Quaker set a howers on the Pillory at Westminster whiped by the Hang man to the old Exchange London. Some dayes after, Stood too howers more on the Pillory at the Exchange and there had his Tongue Bored throug with a hot Iron, & Stigmatalized in the Forehead with the Letter B: Decem: 17 anno Dom: 1656:

# Hypergraphs

---

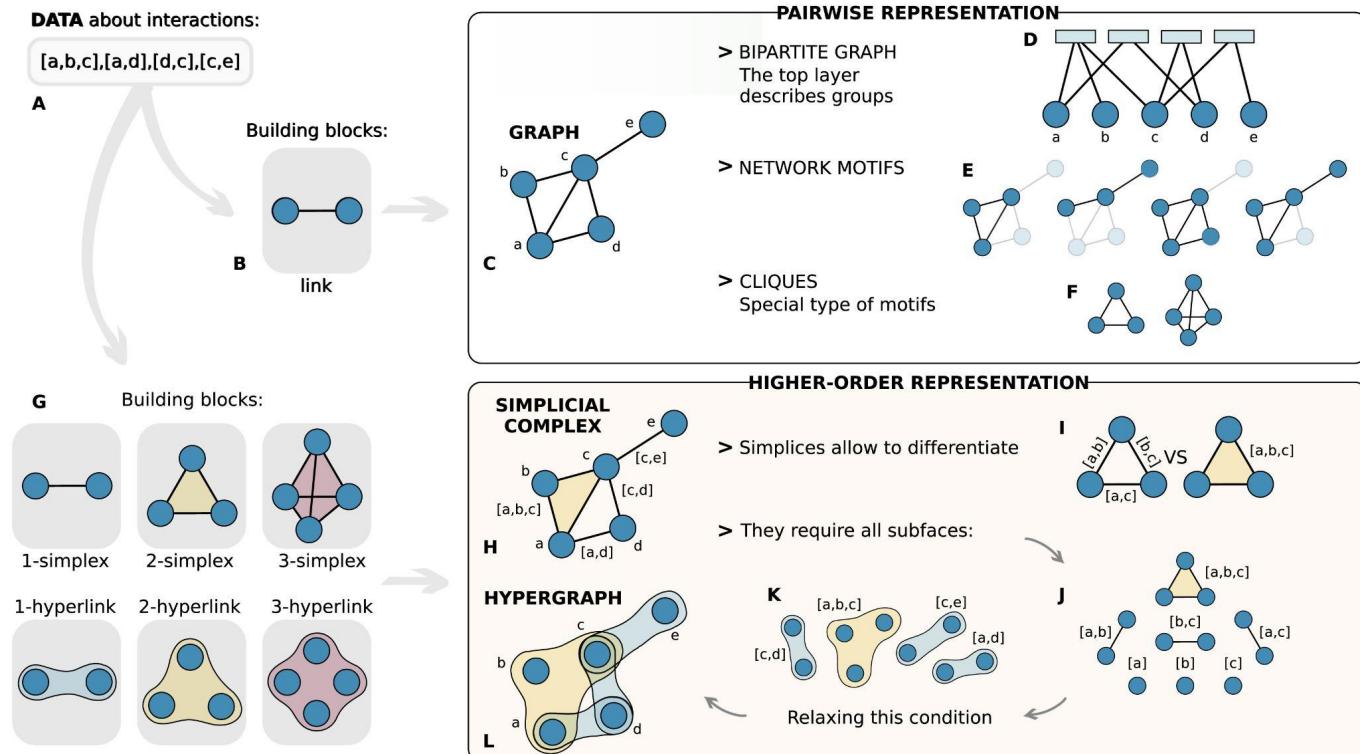
# Non-binary interactions in nature

## Higher-order motif analysis in hypergraphs

Quintino Francesco Lotito, Federico Musciotto, Alberto Montresor & Federico Battiston 

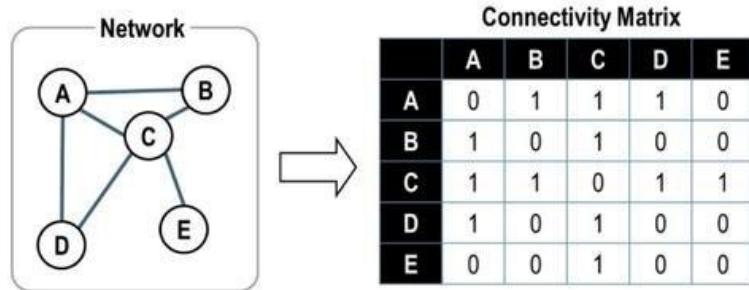
*Communications Physics* 5, Article number: 79 (2022) | [Cite this article](#)

5225 Accesses | 21 Citations | 17 Altmetric | [Metrics](#)



# Main idea and main problematics

Graphs operations, **binary** relations vs. Hyper-graphs operations **higher**-order relations



Simple Connectivity Matrix

$$C_2 = C_1 \cdot C_1^T$$

	A	B	C	D	E
A	3	1	2	1	1
B	1	2	1	2	1
C	2	1	4	1	0
D	1	2	1	2	1
E	1	1	0	1	1

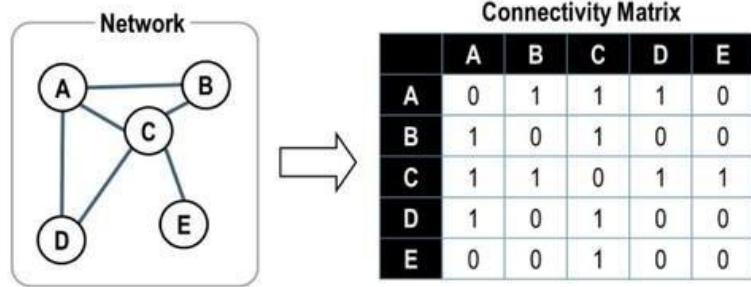
=

	A	B	C	D	E
A	0	1	1	1	0
B	1	0	1	0	0
C	1	1	0	1	1
D	1	0	1	0	0
E	0	0	1	0	0

	A	B	C	D	E
A	0	1	1	1	0
B	1	0	1	0	0
C	1	1	0	1	1
D	1	0	1	0	0
E	0	0	1	0	0

# Main idea and main problematics

Graphs operations, **binary** relations vs. Hyper-graphs operations **higher-order** relations



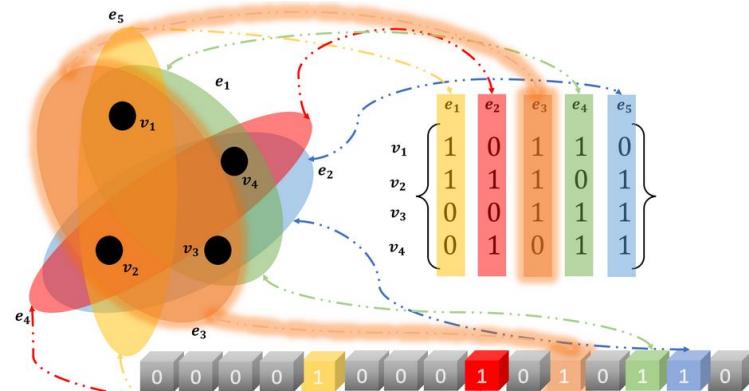
Simple Connectivity Matrix

	C2				
	A	B	C	D	E
A	3	1	2	1	1
B	1	2	1	2	1
C	2	1	4	1	0
D	1	2	1	2	1
E	1	1	0	1	1

=

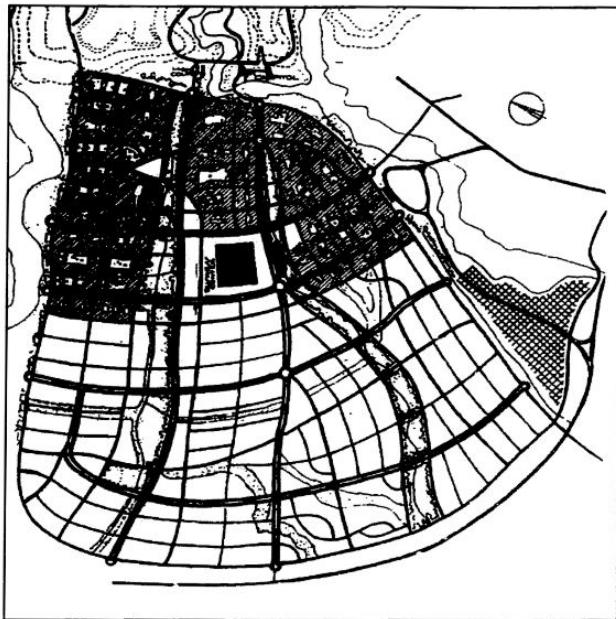
	C1				
	A	B	C	D	E
A	0	1	1	1	0
B	1	0	1	0	0
C	1	1	0	1	X
D	1	0	1	0	0
E	0	0	1	0	0

	CI				
	A	B	C	D	E
A	0	1	1	1	0
B	1	0	1	0	0
C	1	1	0	1	1
D	1	0	1	0	0
E	0	0	1	0	0

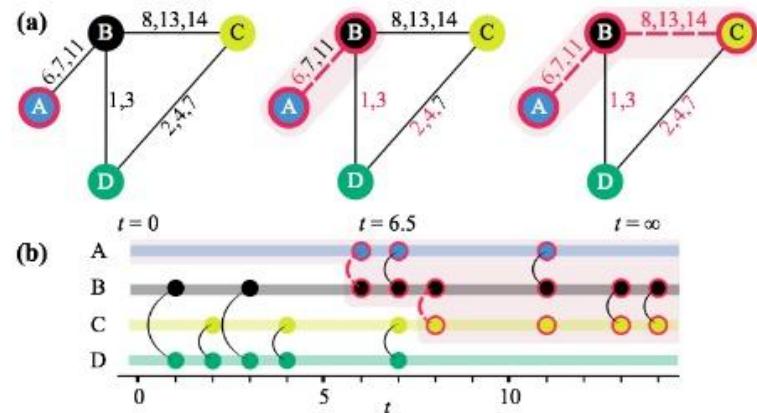


# Networks in time and space

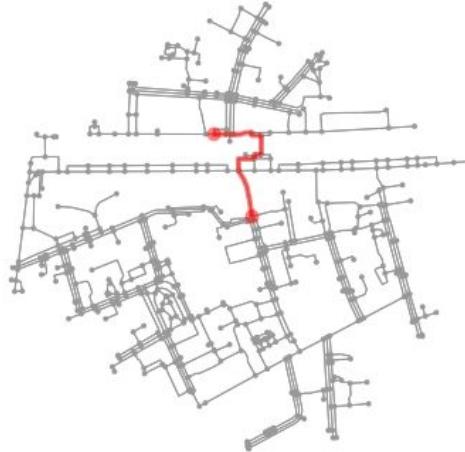
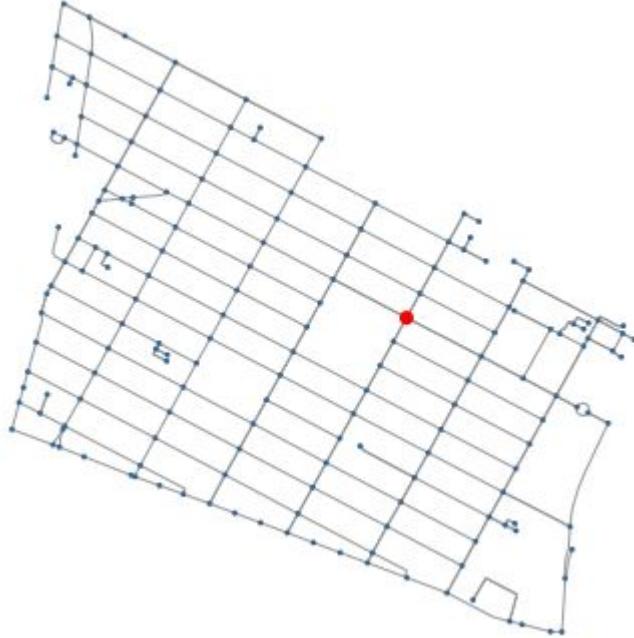
Master Plan for Chandigarh by Albert Mayer RAIC Journal, 1955 (Evenson Norma, Chandigarh, 1966)



Temporality matters:  
reachability issue

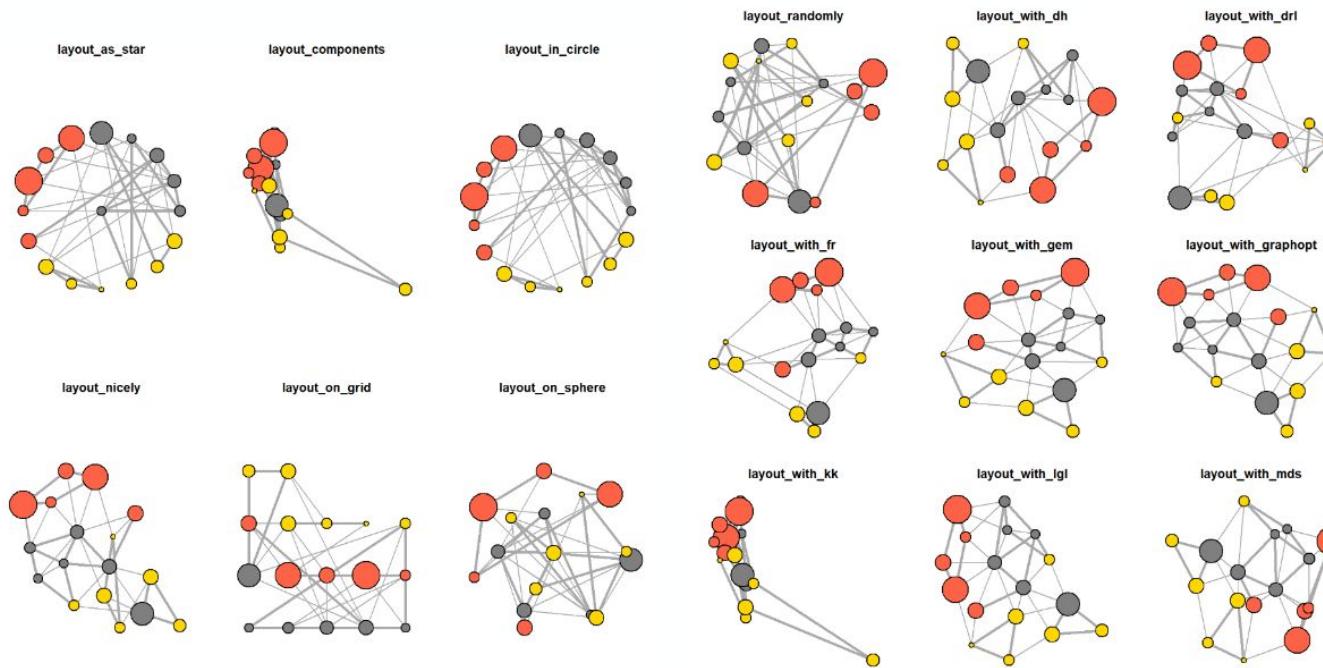


# Networks in time and space



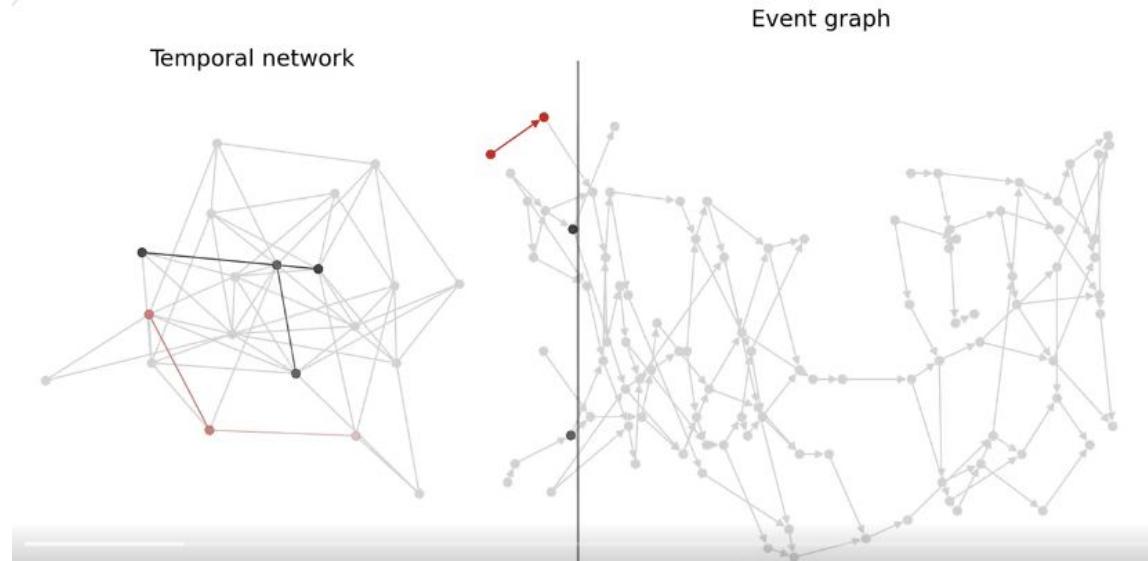
Osmnx for spatial networks  
analysis  
<https://arxiv.org/abs/1010.0302>

# Networks layout



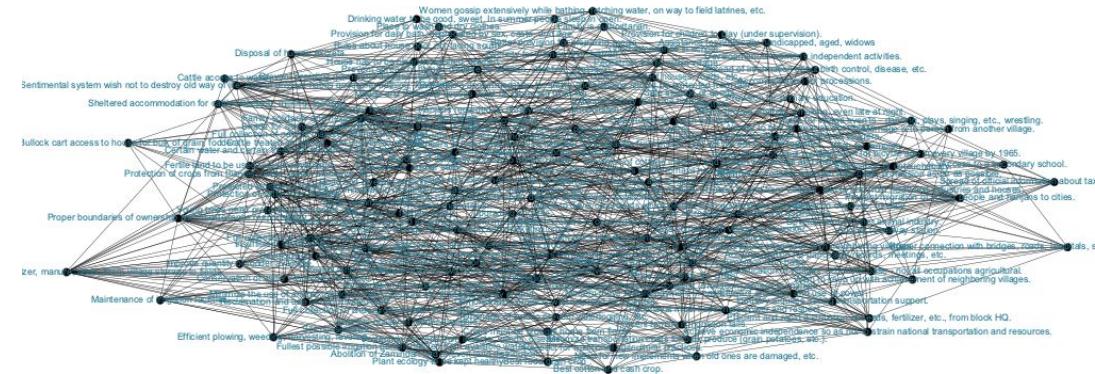
# Networks sonification

Directed percolation in  
temporal networks PRR  
2022



# Networks in time and space

Good resource on spatial networks  
M. Barthelemy “Spatial networks”



Good resource on temporal networks  
P.Holme, J.Saramaki “Temporal networks”  
Holme blog <https://petterhol.me/>

# What we will look at in network science?

1. Network measures and network types
2. Networks in time and space
3. Networks from data

Figure 7.11

Aaron Koblin's *Flight Patterns* (2005): visualization of the flight paths of aircraft crossing North America

# Where can I get network data?

Example:

Highschool: Illinois high school students (1958). A network of friendships among male students in a small high school in Illinois from 1958. 70 nodes, 366 edges.

<https://networks.skewed.de/net/highschool>

Example:

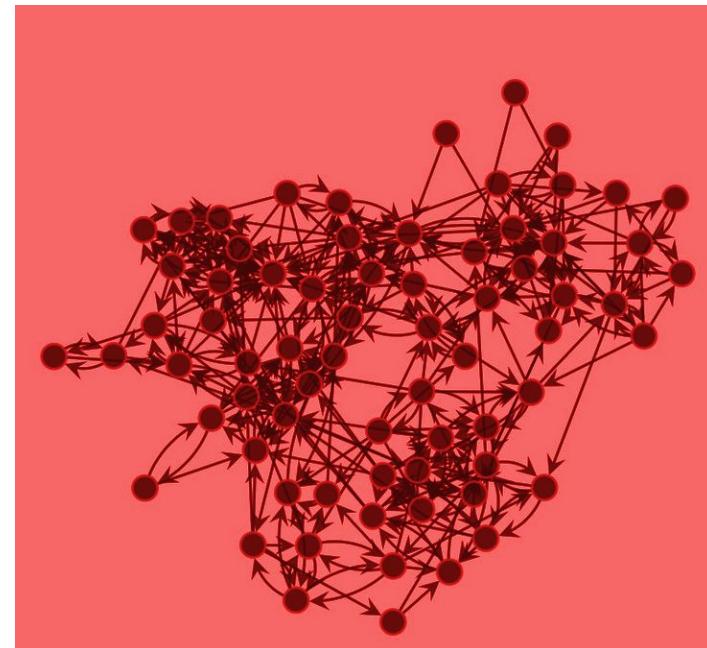
Facebook or wikipedia data

<https://snap.stanford.edu/data/wiki-meta.html>

Syllabus Data Science 2022-2023 ☆ ↗ See new changes

File Edit View Insert Format Tools Extensions Help

Category	Description	Link	Format	Public Datasets	Author
<b>From online repositories (beware, there are a LOT of possibilities in there!)</b>					
ICON network database	Database of 697 network datasets over social, biological, technological, transportation, economic, informational themes. Each dataset contains information on paper, data etc..	<a href="https://icon.colorado.edu/">https://icon.colorado.edu/</a>	No		Liubov Marc
Network repository	Similar to ICON, database of networks	<a href="http://networkrepository.com/">http://networkrepository.com/</a>	No		Liubov Marc



# Social networks analysis

## The Strength of Weak Ties<sup>1</sup>

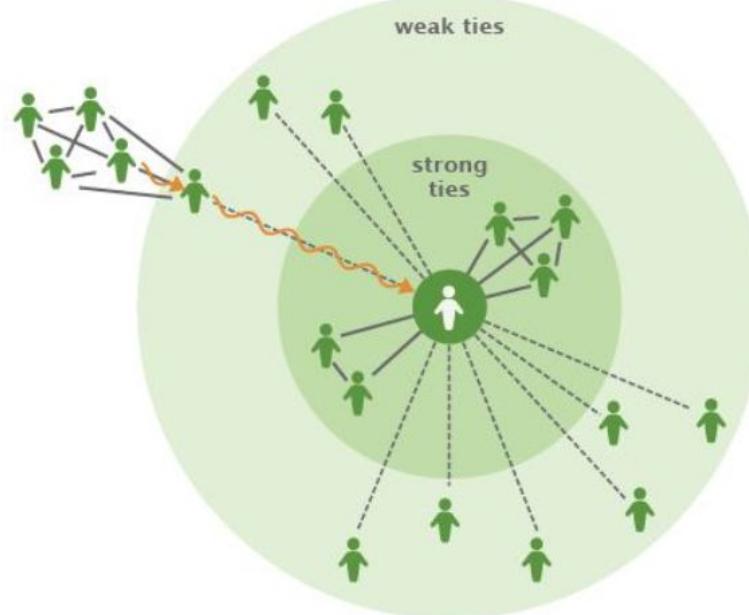
Mark S. Granovetter

*Johns Hopkins University*

Analysis of social networks is suggested as a tool for linking micro and macro levels of sociological theory. The procedure is illustrated by elaboration of the macro implications of one aspect of small-scale interaction: the strength of dyadic ties. It is argued that the degree of overlap of two individuals' friendship networks varies directly with the strength of their tie to one another. The impact of this principle on diffusion of influence and information, mobility opportunity, and community organization is explored. Stress is laid on the cohesive power of weak ties. Most network models deal, implicitly, with strong ties, thus confining their applicability to small, well-defined groups. Emphasis on weak ties lends itself to discussion of relations *between* groups and to analysis of segments of social structure not easily defined in terms of primary groups.

A fundamental weakness of current sociological theory is that it does not relate micro-level interactions to macro-level patterns in any convincing way. Large-scale statistical, as well as qualitative, studies offer a good deal of insight into such macro phenomena as social mobility, community organization, and political structure. At the micro level, a large and increasing body of data and theory offers useful and illuminating ideas about what transpires within the confines of the small group. But how interaction in small groups aggregates to form large-scale patterns eludes us in most cases.

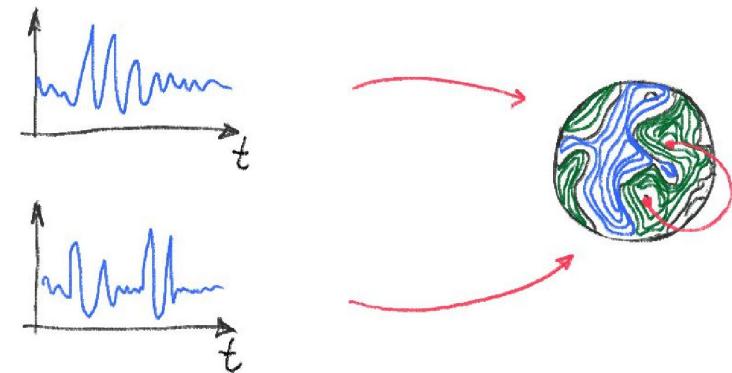
I will argue, in this paper, that the analysis of processes in interpersonal networks provides the most fruitful micro-macro bridge. In one way or another, it is through these networks that small-scale interaction becomes



# How to construct networks from data?

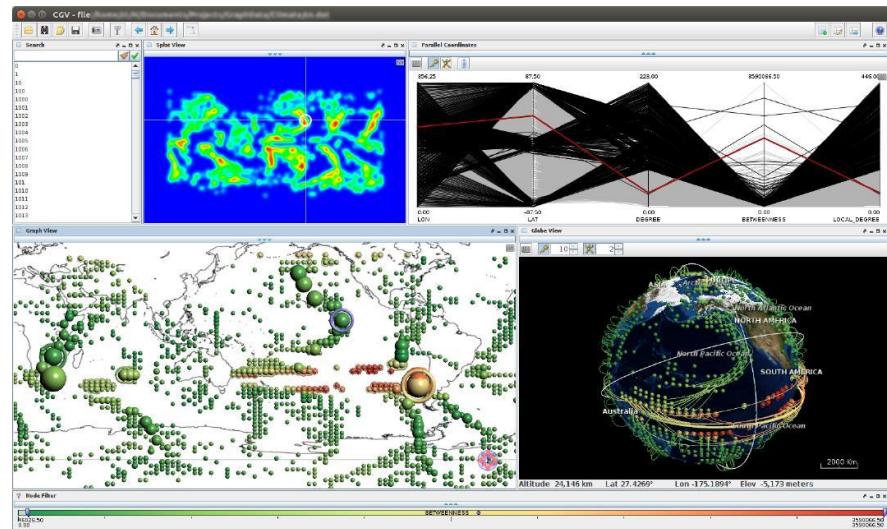
1. Directly build correspondence between links and edges with data (social networks, flights data)
2. Preprocess data (first build correlation from data)

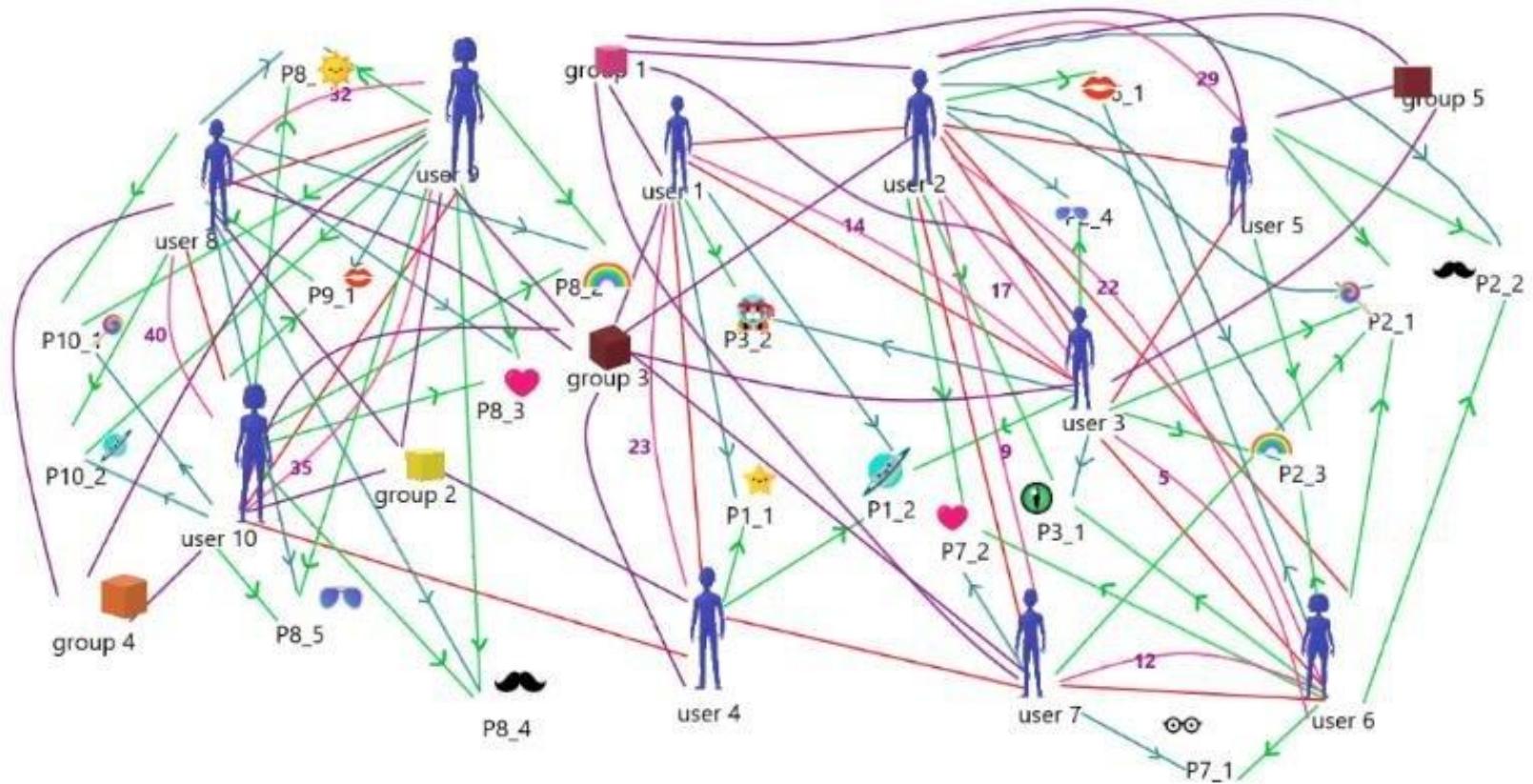
(example of project from [www.pik-potsdam.de](http://www.pik-potsdam.de)  
Climate networks)

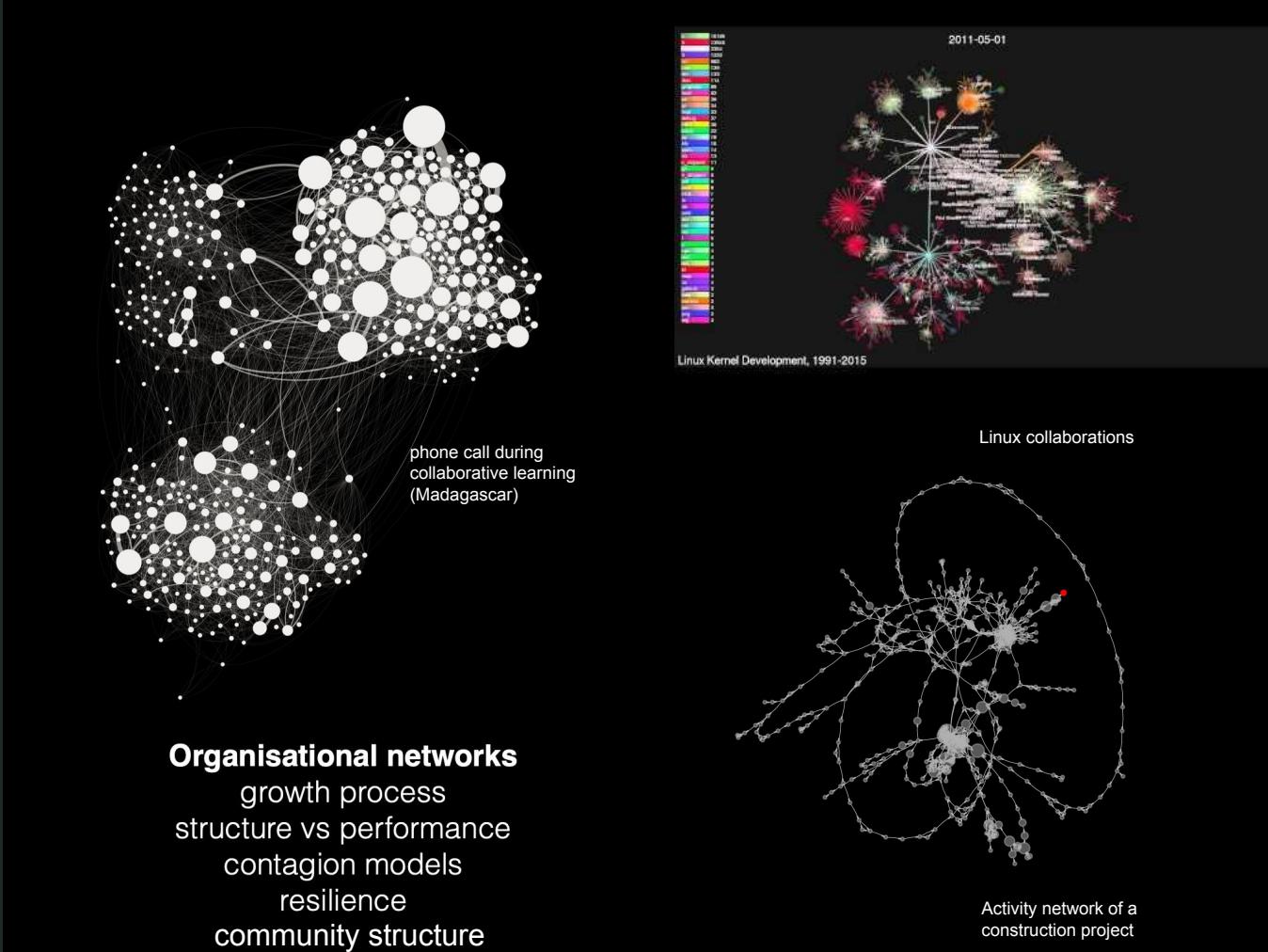


# How to construct networks from data?

1. Directly build correspondence between links and edges with data (social networks, flights data)
2. **Preprocess data (first build correlation from data)**  
Working with data from big systems





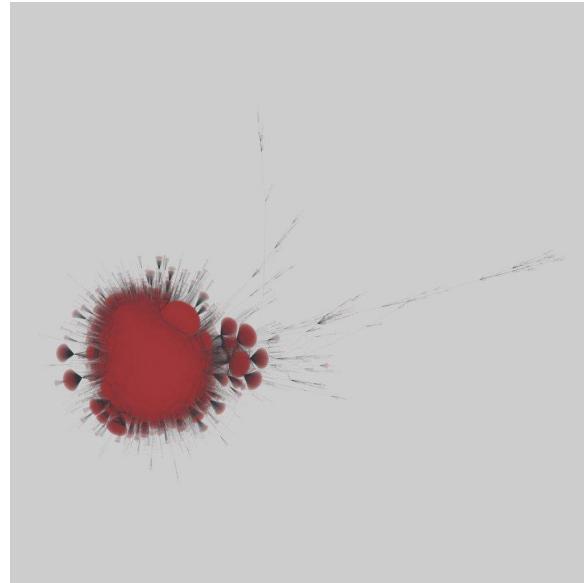


# Network visualisaion

Can hairy ball be a good visualisation?

Discussions for class on visualisations

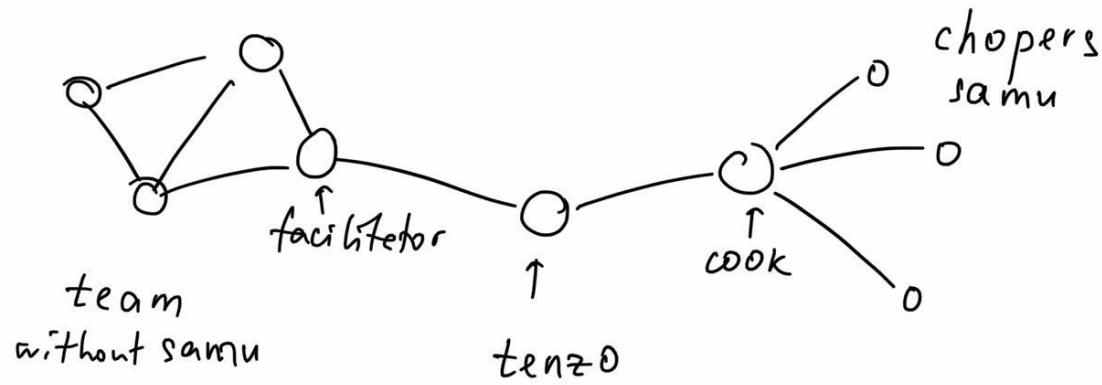
Linux kernel mailing list. A bipartite network of contributions by users to threads on the Linux kernel mailing list. 379554 nodes, 1565683 edges. [https://networks.skewed.de/net/lkml\\_thread](https://networks.skewed.de/net/lkml_thread)



# Hands-on session

1. How to generate your own network?
2. How to create networks using models?
3. How to load your network using network-dataset?
4. How to study your data using the network representation of data?
5. What if the data is not in the network format? Preparing network format.
6. What are ways to create statistics and geometry on your newly represented data as a network?

# Analysis of communities



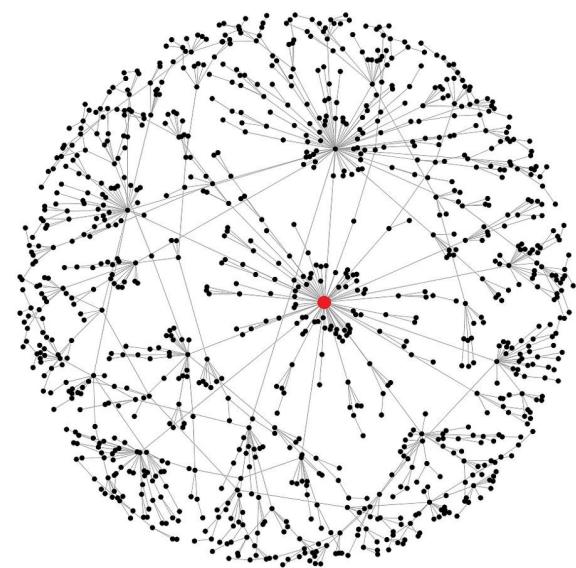
# How to construct networks from data?

1. Decide what are links and nodes in your data?
2. Preprocess data (first build correlation from data)

Examples: social networks, flights data in the practical part

# How to construct networks from data?

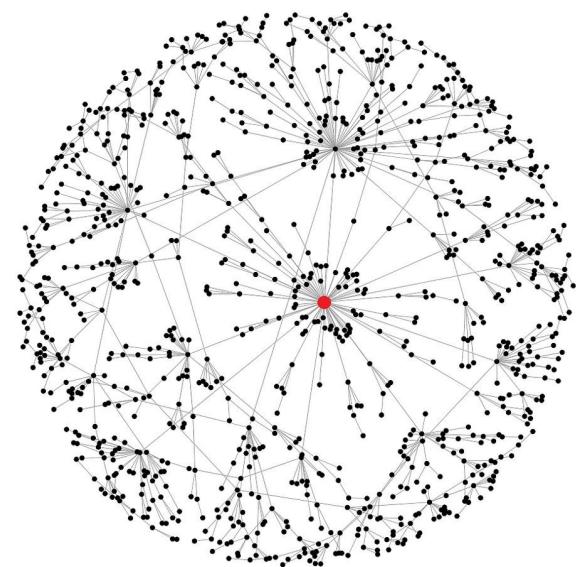
1. Directly build correspondence between links and edges with data (social networks, flights data)
2. Preprocess data (first build correlation from data)



# How to construct networks from data?

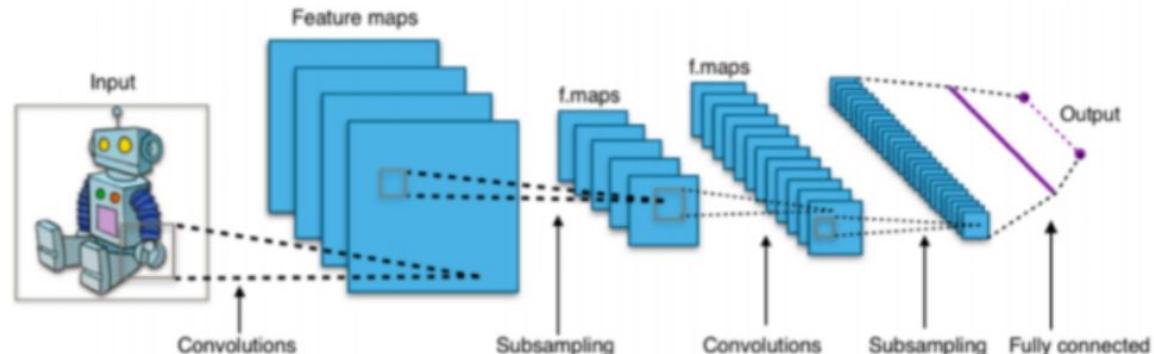
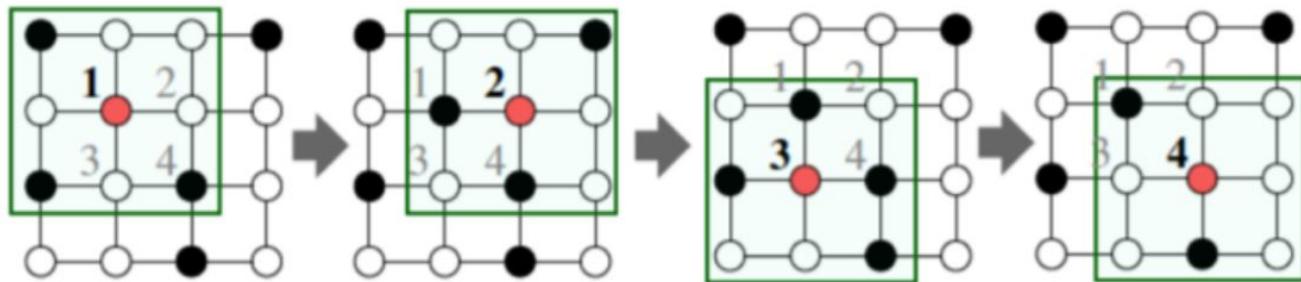
1. Directly build correspondence between links and edges with data (social networks, flights data)
2. Preprocess data (first build correlation from data)

Remember that networks do not give one-to-one  
Correspondence of your data.  
Hence do not generalize



# Algorithms on graphs (examples like GNNs)

The most fundamental part of GNN is a Graph. In computer science, a **graph** is a data structure consisting of two components: **nodes** (vertices) and **edges**. A graph  $G$  can be defined as  $G = (V, E)$ , where  $V$  is the set of nodes, and  $E$  are the edges between them  
<https://neptune.ai/blog/graph-neural-network-and-some-of-gnn-applications>



# Network resources

<http://networkrepository.com/networks.php>

<http://networksciencebook.com/chapter/3#advanced-b>



Network Repository. An Interactive *Scientific* Network Data Repository.  
THE FIRST SCIENTIFIC NETWORK DATA REPOSITORY WITH INTERACTIVE VISUAL ANALYTICS.  
NEW [GraphVis: interactive visual graph mining and machine learning](#)

The first interactive data and network data repository with real-time visual analytics. Network repository is not only the first interactive repository, but also the *largest network repository* with thousands of donations in 30+ domains (from biological to social network data). This large comprehensive collection of network graph data is useful for making significant research findings as well as benchmark network data sets for a wide variety of applications and domains (e.g., network science, bioinformatics, machine learning, data mining, physics, and social science) and includes relational, attributed, heterogeneous, streaming, spatial, and time series network data as well as non-relational machine learning data. All graph data sets are easily downloaded into a standard consistent format. We also have built a multi-level interactive graph analytics engine that allows users to visualize the structure of the network data as well as macro-level graph data statistics as well as important micro-level network properties of the nodes and edges. Check out [GraphVis](#): the interactive visual network mining and machine learning tool.

[GET NETWORK DATA](#)   [COMPARE GRAPH DATA](#)   [VISUALIZE NETWORKS](#)

# Network resources

<http://networkrepository.com/networks.php>

<http://networksciencebook.com/chapter/3#advanced-b>

## Spatial Networks

Marc Barthélémy\*

Institut de Physique Théorique, CEA, IPHT CNRS, URA 2306 F-91191 Gif-sur-Yvette France and  
Centre d'Analyse et de Mathématique Sociales (CAMS), UMR 8557 CNRS-EHESS  
Ecole des Hautes Etudes en Sciences Sociales, 54 bd. Raspail, F-75270 Paris Cedex 06, France.

Complex systems are very often organized under the form of networks where nodes and edges are embedded in space. Transportation and mobility networks, Internet, mobile phone networks, power grids, social and contact networks, neural networks, are all examples where space is relevant and where topology alone does not contain all the information. Characterizing and understanding the structure and the evolution of spatial networks is thus crucial for many different fields ranging from urbanism to epidemiology. An important consequence of space on networks is that there is a cost associated to the length of edges which in turn has dramatic effects on the topological structure of these networks. We will expose thoroughly the current state of our understanding of how the spatial constraints affect the structure and properties of these networks. We will review the most recent empirical observations and the most important models of spatial networks. We will also discuss various processes which take place on these spatial networks, such as phase transitions, random walks, synchronization, navigation, resilience, and disease spread.

## Contents

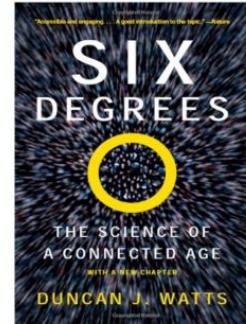
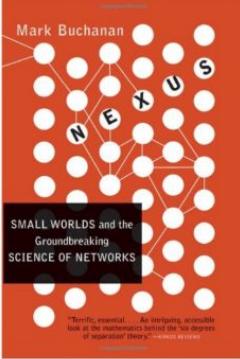
I. Networks and space	2	1. Erdos-Renyi graph	39
A. Introduction	2	2. Planar Erdos-Renyi graph	39
B. Quantitative geography and networks	2	3. The hidden variable model for spatial networks	40
C. What this review is (not) about	2	4. The Waxman model	41
II. Characterizing spatial networks	3	C. Spatial small worlds	41
A. Generalities on planar networks	3	1. The Watts-Strogatz model	41
1. Spatial and planar networks	3	2. Spatial generalizations	42
2. Classical results for planar networks	3	D. Spatial growth models	43
3. Voronoi tessellation	4	1. Generalities	43
		2. Preferential attachment and distance selection	43
		3. Growth and local optimization	46
		E. Optimal networks	49

How Everything Is Connected to  
Everything Else and What It Means for  
Business, Science, and Everyday Life

Linked



"Linked could alter the way we think about all of the networks that affect our lives." —The New York Times  
Albert-László Barabási  
With a New Afterword



Connected

The Surprising Power of Our Social Networks  
and How They Shape Our Lives

Copyrighted Material

# Jupyter notebooks and resources

<https://classroom.google.com/u/0/w/NTQ0NTIwOTczMTQ2/t/all>

# GNNs

Why GNNs?

For more resources:

[https://towardsdatascience.com/an-introduction-to-graph-neural-network-gnn-for-a  
nalysing-structured-data-afce79f4cfdc](https://towardsdatascience.com/an-introduction-to-graph-neural-network-gnn-for-analysing-structured-data-afce79f4cfdc)

# GNNs

## Traditional Graph Analysis Methods

Traditional methods are mostly algorithm-based, such as :

1. searching algorithms, e.g. BFS, DFS
2. shortest path algorithms, e.g. Dijkstra's algorithm, Nearest Neighbour
3. spanning-tree algorithms, e.g. Prim's algorithm
4. clustering methods, e.g. Highly Connected Components, k-mean

# GNNs

Graph Neural Network, as how it is called, is a neural network that can directly be applied to graphs. It provides a convenient way for node level, edge level, and graph level prediction task.

There are mainly three types of graph neural networks in the literature:

Recurrent Graph Neural Network

Spatial Convolutional Network

Spectral Convolutional Network

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

Adjacency matrix (A)

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

Feature matrix (X)