M1 CompuPhys 2022-2023
Introduction to Python Language

# *Network Analysis*

Célestin Coquidé
PostDoc at UTINAM (Observatory building office n°13)

# Table of contents :

# Lectures and Practical Works

Lectures
26/10 1.5 h
28/10 1.5 h

Practicing
17/11 3h[1]
24/11 3h[2]

[1]: Handling of different packages + Exercises

[2]: Preparation for a graded homework

# I Introduction

# Data science

# Data science

- A new discipline

- Few amount of book with a global and unified sight

# Evolutions in data

# Evolutions in data

# Evolutions in data

# Multidisciplinary research field



Source: Palmer, Shelly. Data Science for the C-Suite.
New York: Digital Living Press, 2015. Print.

# Different types of data

- Linear data
  ex: Correlation matrix between proteins

- Sequential data (implicit time)
  ex: Geographical itineraries

- Temporal data
  ex: Wall Street market

# Different types of data

- Linear data
    ex: Correlation matrix between proteins

- Sequential data (implicit time)
    ex: Geographical itineraries

- Temporal data
    ex: Wall Street market



(a) Input

(b) $P_x$ for $x \in \{xac, yac, bc\}$

(c) VON network

Example of sequential data and its analysis through variable-order Markov chains

# Different types of data

- Linear data
    ex: Correlation matrix between proteins

- Sequential data (implicit time)
    ex: Geographical itineraries

- Temporal data
    ex: Wall Street market

Network representation



Example of sequential data and its analysis through variable-order Markov chains

# II Network Representations

# Graph theory



Kaliningrad (2016)

First paper on graph theory :

**The Seven Bridges of Königsberg**
Euler 1736

Is there a path connecting each territory passing by every bridge (but just once) ?

# Graph theory



2
3
B
4
A
C
1 7
5
D
6

Kaliningrad (2016)

First paper on graph theory :

**The Seven Bridges of Königsberg**
Euler 1736

Is there a path connecting each
territory passing by every bridge
(but just once) ?
Answer : **No**



k=3

2
B

3        4

5

k=5  A        C  k=3

7        6

1        D

k=3

G=(V,E)
V : set of vertices (nodes)
E : set of edges (links)

k : vertex degree

# Network representation



Path = sequence of links

$L_1$(A,C) = [(A,D), (D,C)]          Length 2
$L_2$(A,C) = [(A,C)]          Length 1  = *shortest path*
$L_3$(A,C) = [(A,B), (B,C)]          Length 2
$L_4$(A,C) = [(A,B), (B,A), (A,C)]          Length 3
...          ...

# Network representation



No
Eulerian path

Eulerian path = Sequence of links without duplicates representing a closed loop containing all nodes

# A trivial model: Random networks

(Erdös-Rényi random graph model from '59)

**Random Network**



**Bell Curve Distribution of Node Linkages**



Number of Nodes

— Typical node

Number of Links

k

Model 1 :
**G=G(n,M)**
Graphs with n vertices and M links

Model 2 :
**G=G(n,p)**
Graphs with n vertices and linking probabilit
p.
Number of expected links  = pn(n-1)/2

$$n \rightarrow \infty$$

**Properties studied :**
Connectedness
Graph diameter
Subgraphs
…

# More complex model: Scale-free network

(Barabasi-Albert model from '99)

## Scale-Free Network



## Power Law Distribution of Node Linkages



Preferential attachment :
**G=G(m$_0$,$\vec{p}$ )**
Initial graphs with m$_0$ vertices and
preferential attachment vector $\vec{p}$

$$p_i = \frac{k_i}{\sum_{j=1}^{m_0} k_j}$$

Where indexes i and j denote vertex "i" and "j" and $k$ is a number of link.

**Properties studied :**
Connectedness
Graph diameter
Subgraphs
…

# Random graphs

(Erdös-Rényi random graph model from '59)

**VS**

# (Real) complex networks

(A.-L. Barabási, R. Albert '99)

**Random Network**

**Scale-Free Network**

**No singular node**

**Hubs**

**Bell Curve Distribution of Node Linkages**

— Typical node

Number of Nodes

Number of Links

$k$

**Power Law Distribution of Node Linkages**

Number of Nodes

Number of Links

$k$

$k^{-\mu}$

Number of Nodes (log scale)

Number of Links (log scale)

$k$

$\mu \sim 2 - 3$

# Random graphs

**VS**

# (Real) complex networks

## Random Network, Accidental Node Failure



Before

After

# Random graphs    VS    (Real) complex networks



## Scale-Free Network, Accidental Node Failure

Hub

Before

Failed node

After

## Scale-Free Network, Attack on Hubs

Hub

Before

Attacked hub

After

# Internet network in year 2000



| | | | | | |
|---|---|---|---|---|---|
| 🟥 Switzerland | 🟨 Spain | 🟩 Japan | 🟦 Russian Federation | 🟦 UK | ⬛ Unknown |
| 🟧 Germany | 🟩 Italy | 🟩 Netherlands | 🟦 Sweden | 🟪 USA | |

# Real scale-free networks

## Examples of Scale-Free Networks

| NETWORK | NODES | LINKS |
|---|---|---|
| Cellular metabolism | Molecules involved in burning food for energy | Participation in the same biochemical reaction |
| Hollywood | Actors | Appearance in the same movie |
| Internet | Routers | Optical and other physical connections |
| Protein regulatory network | Proteins that help to regulate a cell's activities | Interactions among proteins |
| Research collaborations | Scientists | Co-authorship of papers |
| Sexual relationships | People | Sexual contact |
| World Wide Web | Web pages | URLs |

# Real scale-free networks

## Examples of Scale-Free Networks

| NETWORK | NODES | LINKS |
|---|---|---|
| Cellular metabolism | Molecules involved in burning food for energy | Participation in the same biochemical reaction |
| Hollywood | Actors | Appearance in the same movie |
| Internet | Routers | Optical and other physical connections |
| Protein regulatory network | Proteins that help to regulate a cell's activities | Interactions among proteins |
| Research collaborations | Scientists | Co-authorship of papers |
| Sexual relationships | People | Sexual contact |
| World Wide Web | Web pages | URLs |

## Construction property

"The rich get richer"
or
Matthew effect (Mt 25:29 NT)
or
Preferential attachment

**Type of links:**
Directed ⟶
Undirected ⟶
Weighted(D/U) ⟶ 0.9

# Scale-free networks

## Examples of Scale-Free Networks

| NETWORK | NODES | LINKS |
|---|---|---|
| Cellular metabolism | Molecules involved in burning food for energy | Participation in the same biochemical reaction |
| Hollywood | Actors | Appearance in the same movie |
| Internet | Routers | Optical and other physical connections |
| Protein regulatory network | Proteins that help to regulate a cell's activities | Interactions among proteins |
| Research collaborations | Scientists | Co-authorship of papers |
| Sexual relationships | People | Sexual contact |
| World Wide Web | Web pages | URLs |

## Construction property

"The rich get richer"
or
Matthew effect (Mt 25:29 NT)
or
Preferential attachment

**What about directionality ?**
In fact any network can be considered as directed.

Now there is a difference between
**node outdegree**
And
**node indegree**



FIG. 2 (Color online) Distribution $w_{in,out}(k)$ of number of in-going (a) and outgoing (b) links $k$ for $N = 3282257$ Wikipedia English articles (Aug 2009) of Fig. 1 with total number of links $N_\ell = 71012307$. The straight dashed fit line shows the slope with $\mu_{in} = 2.09 \pm 0.04$ (a) and $\mu_{out} = 2.76 \pm 0.06$ (b). After (Zhirov et al., 2010).

Bipartite networks (ex: election network)
Multi-layer networks (ex: Wikipedia network)
Higher-order networks (Sequential data repr.)
Temporal networks (Time series)
Trees



$k=3$

$k=2$

$k=1$

$i$

Multi-layer network

Bipartite network and its projection

# III Centrality Measures

# What is a Centrality measure?

Number of Link?
Efficiency of path?
Number of path?

Degree (K)

$$K(i) = \sum_{j \in V} x_{ij}$$

$$x_{ij} = \begin{cases} 1 & \text{if} \quad x_{ij} \in E \\ 0 & \text{else} \end{cases}$$

Closeness

$$C(i) = (\sum_{j \in V} dist(i,j))^{-1}$$

where $dist(i,j)$ is the shortest path distance

Betweeness

$$B(i) = \sum_{s,t \in V} \frac{\sigma_{st}(i)}{\sigma_{st}}$$

where $\sigma_{st}$ is the number of shortest path between $s, t$ and $\sigma_{st}(i)$ the number of shortest path going through $i$

# What is a Centrality measure?

The method you use depends
on what you are investigating

For directed network such that
the World Wide Web, the
eigenvector centrality is used

Ex: Google's research engine
algorithm

DATA → Matrix → Network → Centrality

# Modeling Random Walk on a network



Dangling node

Sink

**Adjacency matrix :**

$$A_{ij} = \begin{cases} 1 & \text{if } j \to i \\ 0 & \text{otherwise} \end{cases}$$

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

# The Stochastic Matrix



**Stochastic matrix :**

$$S_{ij} = \begin{cases} \dfrac{1}{N} & \text{if } j \text{ is a dangling node} \\[2em] \dfrac{A_{ij}}{\sum_{i=1}^{N} A_{ij}} & \text{otherwise} \end{cases}$$

$$S = \begin{pmatrix} 0 & 0 & 1/8 & 0 & 0 & 0 & 0 & 0 \\ 1/2 & 0 & 1/8 & 1/3 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 1/8 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/2 & 1/8 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1/8 & 1/3 & 0 & 0 & 1/2 & 0 \\ 0 & 0 & 1/8 & 1/3 & 1/3 & 0 & 0 & 1/2 \\ 0 & 0 & 1/8 & 0 & 1/3 & 0 & 0 & 1/2 \\ 0 & 0 & 1/8 & 0 & 1/3 & 1 & 1/2 & 0 \end{pmatrix}$$

# The Google Matrix



Dangling node

Sink

**Google matrix** :

$$G_{ij} = \alpha S_{ij} + (1 - \alpha)/N$$

$\alpha \in [0, 1]$ is the <u>damping factor</u>

$$G = \begin{pmatrix} 1/40 & 1/40 & 1/8 & 1/40 & 1/40 & 1/40 & 1/40 & 1/40 \\ 17/40 & 1/40 & 1/8 & 7/24 & 1/40 & 1/40 & 1/40 & 1/40 \\ 17/40 & 17/40 & 1/8 & 1/40 & 1/40 & 1/40 & 1/40 & 1/40 \\ 1/40 & 17/40 & 1/8 & 1/40 & 1/40 & 1/40 & 1/40 & 1/40 \\ 1/40 & 1/40 & 1/8 & 7/24 & 1/40 & 1/40 & 17/40 & 1/40 \\ 1/40 & 1/40 & 1/8 & 7/24 & 7/24 & 1/40 & 1/40 & 17/40 \\ 1/40 & 1/40 & 1/8 & 1/40 & 7/24 & 1/40 & 1/40 & 17/40 \\ 1/40 & 1/40 & 1/8 & 1/40 & 7/24 & 33/40 & 17/40 & 1/40 \end{pmatrix}$$

# The Google matrix



Dangling node

Sink

$$A = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{pmatrix}$$
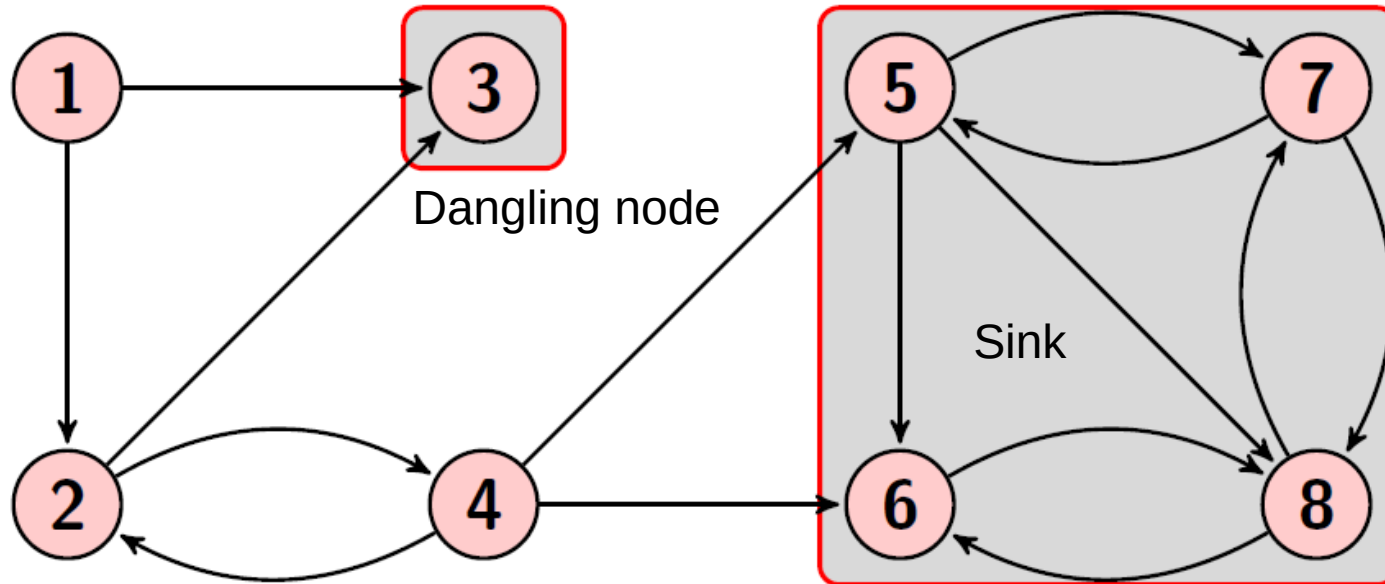
$$S = \begin{pmatrix} 0 & 0 & 1/8 & 0 & 0 & 0 & 0 & 0 \\ 1/2 & 0 & 1/8 & 1/3 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 1/8 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1/2 & 1/8 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1/8 & 1/3 & 0 & 0 & 1/2 & 0 \\ 0 & 0 & 1/8 & 1/3 & 1/3 & 0 & 0 & 1/2 \\ 0 & 0 & 1/8 & 0 & 1/3 & 0 & 0 & 1/2 \\ 0 & 0 & 1/8 & 0 & 1/3 & 1 & 1/2 & 0 \end{pmatrix}$$
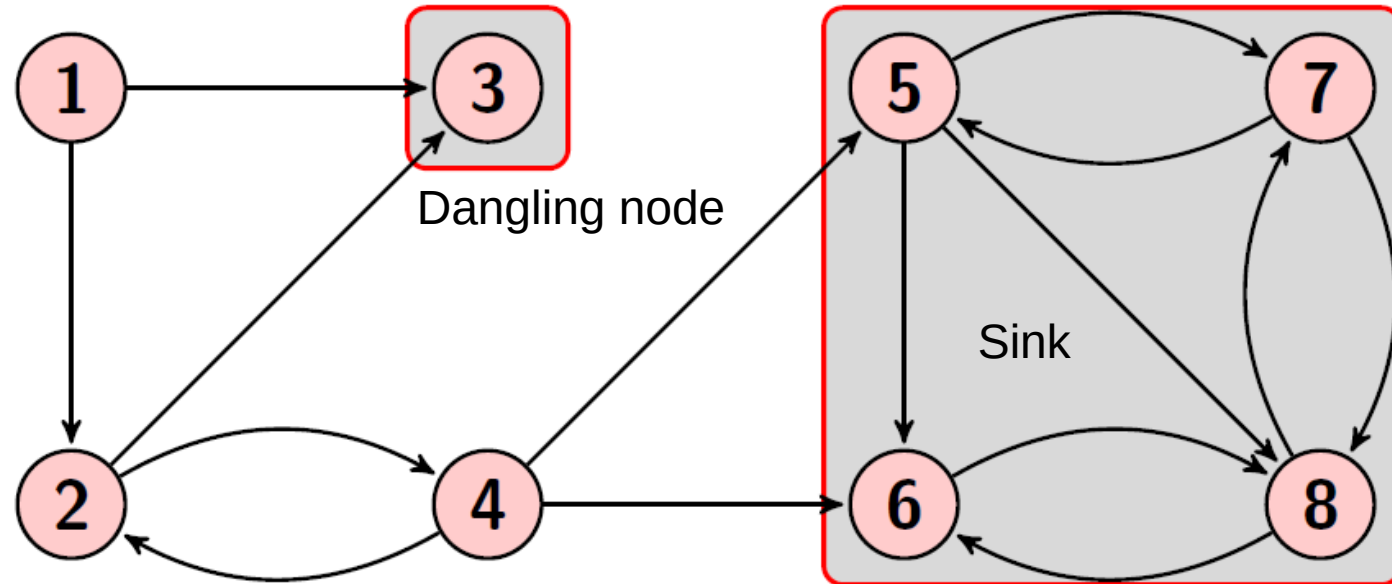
$$\alpha = 0.8$$

# Fundamentals of Google research engine

$$G = \alpha S + \frac{1-\alpha}{N} e e^T \qquad e^T = (1, 1, \ldots, 1)$$

If the spectrum of the stochastic matrix $S$ is $\{1, \lambda_1, \lambda_2, ..., \lambda_N\}$, then the spectrum of the Google matrix $G = \alpha S + (1-\alpha)\mathbf{e}\mathbf{v}^T$ is $\{1, \alpha\lambda_1, \alpha\lambda_2, ..., \alpha\lambda_N\}$, where $\mathbf{v}^T$ is a probability vector.

## PageRank probability vector P

$$P = GP$$

*P* is the eigenvector associated to the largest eigenvalue, ie 1

For very large network (such as WWW), no way to directly diagonalize *G* → **Powermethod**

$$P = \lim_{n \to \infty} G^n v_0 \qquad \forall v_0$$

## Interpretation

$$P_i = \sum_{j \in B_i} \frac{P_j}{k_{out}(j)}$$

The more a node is pointed by important node, the more it is important

**Algortithm at the hearth of** Google **ΤM search engine.**
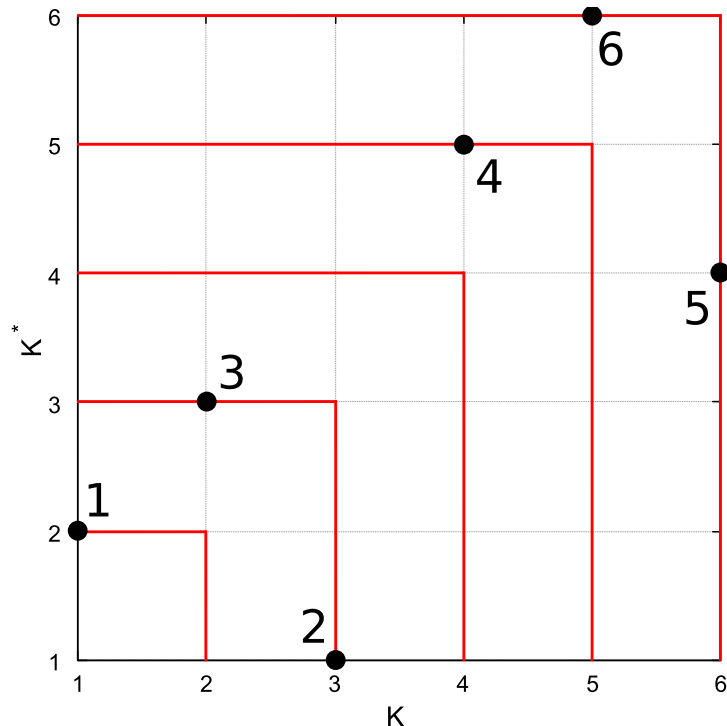
(Brin & Page cofounders)

**PageRank**
Measures the <u>importance</u> of a node as the property to be pointed by other important nodes (influent node).

**CheiRank**
Same as PageRank but considering the adjacency matrix of the <u>inverted network</u>.
Measures the importance of a node as the property to point to other important nodes = <u>communicability</u>



$$P(K) \sim \frac{1}{K^\beta}$$

$$\beta = \frac{1}{\mu - 1}$$



P : PageRank vector
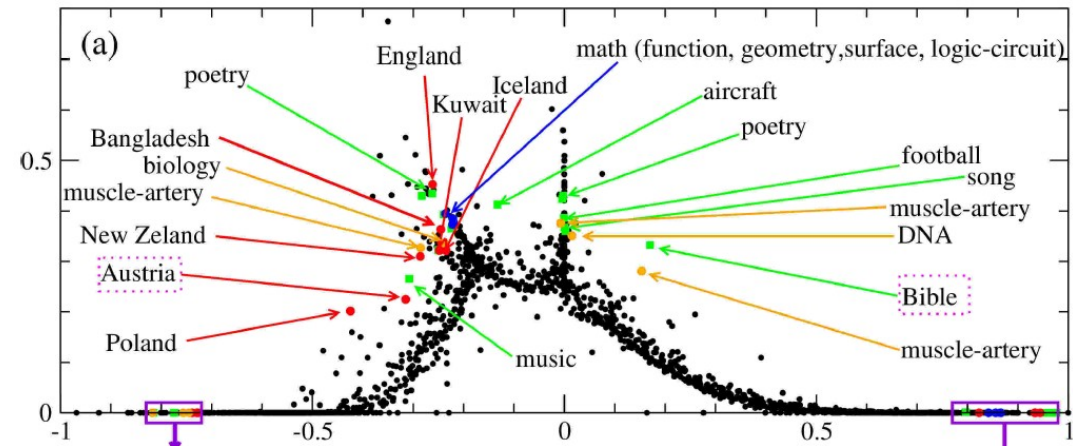K : PageRank
P* : CheiRank vector
K* : CheiRank

**2DRank**
Mix *PageRank* and CheiRank, measures the propension of a node to be pointed and to point towards other nodes.

# Plenty of applications in scientific research

Non exhaustive list of applications :
- Linux Kernel networks
- Wikipedia networks
- Twitter network
- World trade network
- Brain neural network
- DNA sequences
- Networks of Game Go moves
- Opinion formation on directed networks
- Contagion on networks
- ...



Spectrum of Google matrix of English Wikipedia



Crisis contagion on World trade network



DNA sequence analysis

# IV Data and Network analysis with Python

# Networkx and Pandas



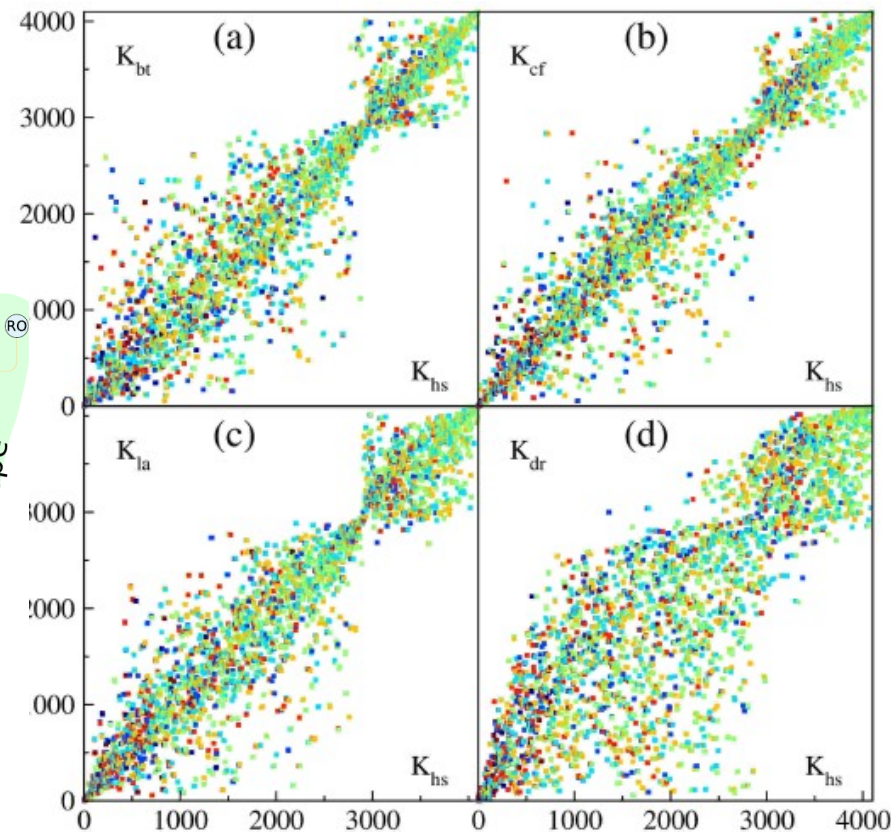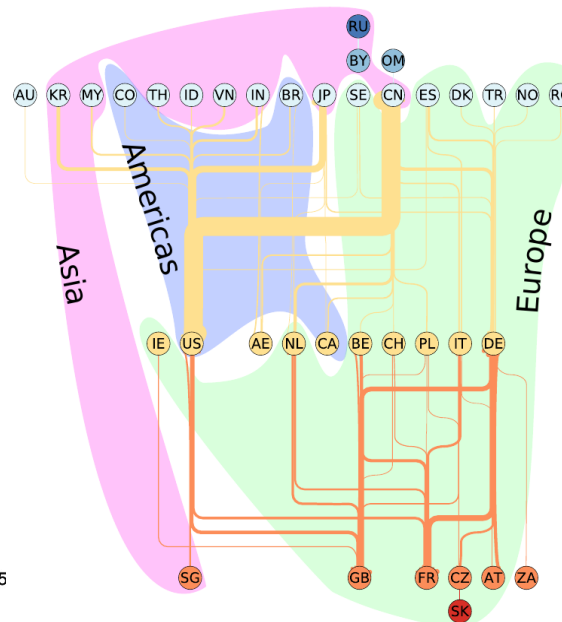Creation and modification of network objects
    From generative models
    From data

Algorithm ready to use
    Clustering
    Centrality

Network visualization

Documentation and guide available at https://networkx.org/

Installing with `pip install networkx`

Using with `import networkx as nx`

# Networkx and Pandas



Data manipulation tools
　　CSV to Python Dictionary

Data visualization and analysis tools

Documentations and guide available at https://pandas.pydata.org/

Installing with **pip install pandas**

Using with **import pandas as pd**

# Numpy



Scientific computing for Python

Linear algebra

Documentations and guide available at http://numpy.org/

Installing with `pip install numpy`

Using with `import numpy as np`

# Pyplot



Powerful Matplotlib object t generate graphics

PDF, SVG, JPG, GIF and other output formats

Documentations and guide available at https://matplotlib.org/

Installing with **pip install matplotlib**

Using with **import matplotlib.pyplot as plt**

Create a python script taking as input: *links.dat* and *nodes.dat*

Generate outputs A and B

A: Network representation in PDF
B: File consisting in the list of nodes with:  degree, closeness and betweeness

In case of B, try you own function for degree, betweeness and
 closeness, before using networkx's ones