

**WWW, WIKIPEDIA AND OSNS**

AP

# NETWORKS OF HUMANS

Theme: no one controls the evolution of the network, which is self-organizing

What is represented (self, news, opinion, concept) and its **lifecycle** determines the structure and the research questions

look at how they connect and when

Direction of communication is important

```
1 import networkx as nx
2
3 eu_DG = nx.DiGraph()
```

# GETTING DATA

# WWW

- a Networkx digraph will represent connectivity
- a companion dictionary maps vertices to URLs of the relative pages
- source: a *scrape* of the 2005 “.eu” domain

# TWITTER

- supported by the Twython module
- requires Twitter registration/API token
- alternative platforms exist, e.g. Tweepy (+NLTK)
- interesting: the network of mentions as a voting system

# WIKIPEDIA

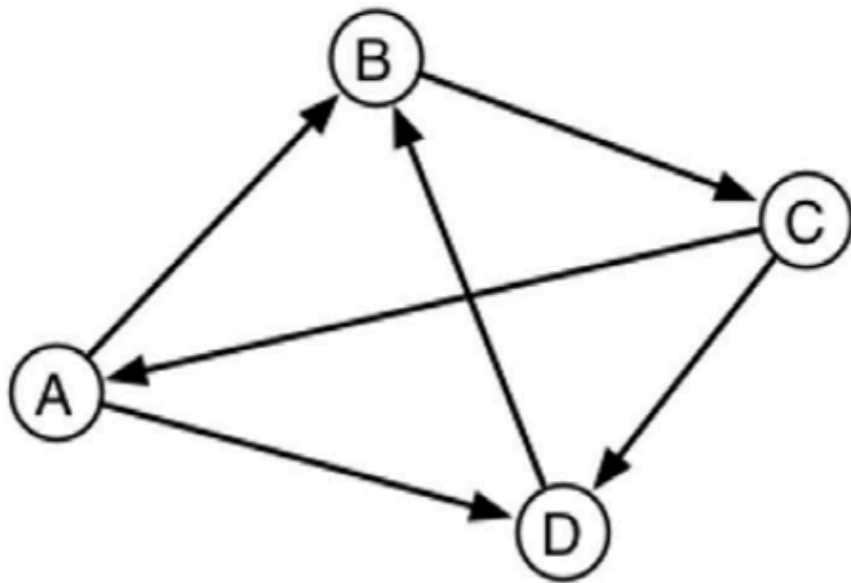
- a network of concepts (lemmas/lemmata) maintained by humans (and some bot)
- time-stamped evolution of the network is available [\[here\]](#)
- contrary to *curated* taxonomies, e.g., [Linnaeus, 1735], this is not a tree

a **directed acyclic graph** is the reference model

# **RANKING ALGORITHMS: PAGERANK**

# PAGERANK IDEA

Rank vertices (page) on the basis of their *importance* for network *navigation*  
Importance is captured by the relative value in the dominant Eigenvector of a new matrix  $P$  that represents *navigation*



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

**Fig. 4.1** A simple oriented graph with its adjacency matrix.



# VARIABLES USED

$A$ : the directed adjacency matrix (admits *dangling* ends)

$K0^{-1}$ : 0 everywhere but  $\frac{1}{k_i^{\text{out}}}$  on the main diagonal:  $(K0^{-1})_{ii} = \frac{1}{k_i^{\text{out}}}$

$$N = A \times K0^{-1}$$

Matrix  $E$ : it models *escape* from the navigation and *forces* irreducibility in the sense of Perron-Frobenius

Set  $E_{ii} = \frac{1}{n}$ , where  $n = |V|$ ; set 0 everywhere else

(current versions put different values on the main diagonal of  $E$ )

$$P = \alpha N + (1 - \alpha)E$$

Experimentally, Page and Brin set  $\alpha = 0.85$

I.e., we assign a  $1 - \alpha = 0.15$  probability to *jump out* of a navigation *path* and into an arbitrary *restart* node.

Compute the dominant eigenvector  $\mathbf{r}_1$  of  $P$

Perron-Frobenius applies so  $\mathbf{r}_1$  will be a 'proper' vector of real, positive values

Compute  $\mathbf{r}_1$  by iterating the  $\mathbf{r}_1^{(t+1)} = P\mathbf{r}_1^{(t)}$  eq. until it converges to a fixpoint value.

# **RANKING ALGORITHMS: HITS**

# HITS IDEA

Hyperlink-Induced Topic Search by [Kleinberg, 1999]

Sees importance of a node in a more nuanced way:

Pages that are important for consultation, e.g., train schedules, have *authority* and tend to be *terminal*

Well-connected *hub* pages that facilitate navigation, e.g., Time Out, are useful but not authoritative per se

1. authority score  $\mathbf{au(i)}$

2. hub score  $\mathbf{h(i)}$

# HITS AS MUTUAL RECURSION

Hub-iness influences authority which in turns influences hub-iness:

$$au(i) \propto \sum_{j \rightarrow i} h(j)$$

page  $i$  is authoritative proportionally to the sum of the hub-iness of the pages that link to it.

$$h(i) \propto \sum_{i \rightarrow j} au(j)$$

page  $i$  is hub proportionally to the sum of the authoritativeness of pages that it links to.

# COMPUTING HITS SCORES

We could start with assigning 1 everywhere and hoping that mutual recursion will converge to stable  $au$  and  $h$  values.

As with Von Mises' method, we normalise vectors to 1 at each iteration.

# LINEAR ALGEBRA DERIVATIONS

$$\mathbf{h} \propto AA^T \mathbf{h} = \lambda_h AA^T \mathbf{h}$$

$$\mathbf{au} \propto A^T A \mathbf{au} = \lambda_{au} A^T A \mathbf{au}$$

I.e., we can find  $\mathbf{h}$  and  $\mathbf{au}$  separately by solving the eigenvalue problem for the matrices  $AA^T$  and  $A^T A$



# MAIN RESULT

For *primitive* matrices (i.e., connected networks, no dead-ends/sinks)

$$\mathbf{h} \propto AA^T \mathbf{h} = \lambda_h AA^T \mathbf{h}$$

$$\mathbf{au} \propto A^T A \mathbf{au} = \lambda_{au} A^T A \mathbf{au}$$

- convergence is assured;
- dominant  $\lambda$  is unique and
- values for  $\mathbf{h}$  and  $\mathbf{au}$  will be all positive, as desired.

(negative values have no interpretation here)

# COMMUNITY DETECTION

# FINDING SOCIAL STRUCTURES

this is an example of Provost-Fawcett's problems

- 4: Clustering
- 5: co-occurrence grouping

For homogeneous networks, eg., country-to-county of Ch. 2

Community: nodes that are closely connected with each other by *strenght* or *density*

Resolution limit: communities with less than  $\sqrt{|V|}$  members cannot be properly identified.

# GIVAN-NEWMAN

1. Rank edges by their *help to connectivity*
2. remove the top-ranking edge
3. repeat until loss of connection
4. now-isolated areas are called communities

Hyp: Betweenness centrality captures *help to connectivity*

# MODULARITY

# AS AN OPTIMIZATION PROB.

**Instance:** an adj. matrix  $A$ , a small integer  $g$

**Solution:** a partition of  $V$  into  $g$  groups

**Measure:** maximise  $Q$ : the overall modularity measure

**Interpretation:** how likely is a random walker to leave the community?

# THE Q FACTOR

Let  $E_{g \times g}$  be the cross-group matrix and  $f_i$  the sum of col.  $i$

Electrical conductance:

$$Q = \sum_{i=1}^g e_{ii} - f_i^2$$

Complexity: NP-complete

Even random networks might exhibit densifications that might look as c.