

Link Analysis

Leonid E. Zhukov

School of Data Analysis and Artificial Intelligence
Department of Computer Science
National Research University Higher School of Economics

Network Science



NATIONAL RESEARCH
UNIVERSITY

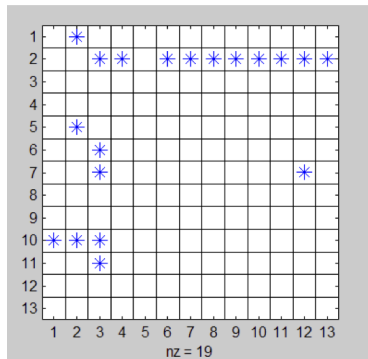
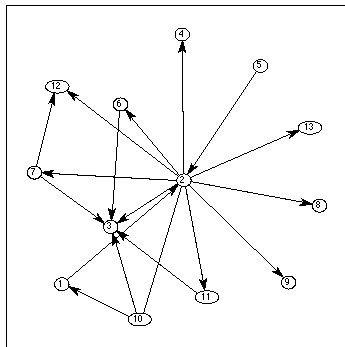
Lecture outline

- 1 Graph-theoretic definitions
- 2 Web page ranking algorithms
 - Pagerank
 - HITS
- 3 The Web as a graph
- 4 PageRank beyond the web

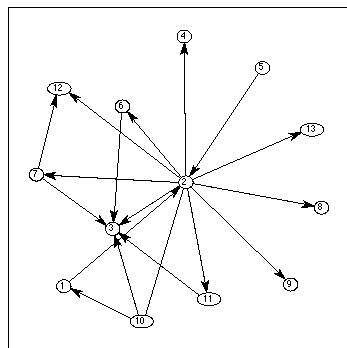
Graph theory

Graph $G(E, V)$, $|V| = n$, $|E| = m$

Adjacency matrix $\mathbf{A}^{n \times n}$, A_{ij} , edge $i \rightarrow j$

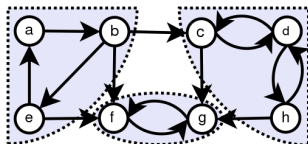


Graph is directed, matrix is non-symmetric: $\mathbf{A}^T \neq \mathbf{A}$, $A_{ij} \neq A_{ji}$



- sinks: zero out degree nodes, $k_{out}(i) = 0$, absorbing nodes
- sources: zero in degree nodes, $k_{in}(i) = 0$

- Graph is **strongly connected** if every vertex is reachable from every other vertex.
- **Strongly connected components** are partitions of the graph into subgraphs that are strongly connected



- In strongly connected graphs there is a path in each direction between any two pairs of vertices

image from Wikipedia

- A directed graph is **aperiodic** if the greatest common divisor of the lengths of its cycles is one (there is no integer $k > 1$ that divides the length of every cycle of the graph)

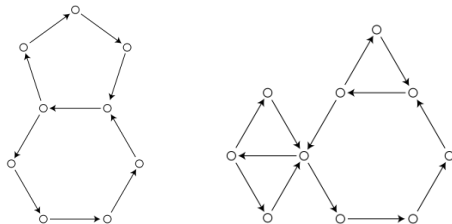
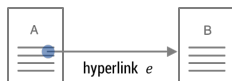


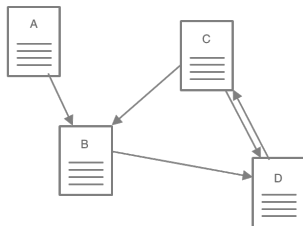
image from Wikipedia

Web as a graph

- Hyperlinks - implicit endorsements

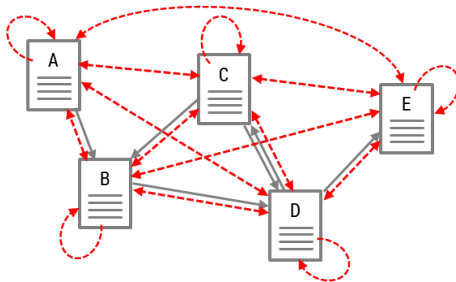


- Web graph - graph of endorsements (sometimes reciprocal)



PageRank

"PageRank can be thought of as a model of user behavior. We assume there is a "random surfer" who is given a web page at random and keeps clicking on links, never hitting "back" but eventually gets bored and starts on another random page. The **probability** that the random surfer visits a page is its **PageRank**."



Sergey Brin and Larry Page, 1998

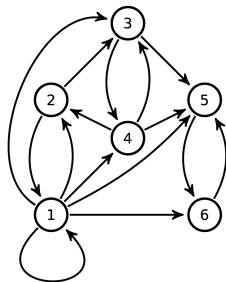
- Random walk on a directed graph

$$p_i^{t+1} = \sum_{j \in N(i)} \frac{p_j^t}{d_j^{\text{out}}} = \sum_j \frac{A_{ji}}{d_j^{\text{out}}} p_j^t$$

$$\mathbf{D}_{ii} = \text{diag}\{d_i^{\text{out}}\}$$

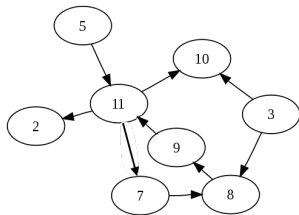
$$\mathbf{p}^{t+1} = (\mathbf{D}^{-1}\mathbf{A})^T \mathbf{p}^t$$

$$\mathbf{p}^{t+1} = \mathbf{P}^T \mathbf{p}^t$$



Ranking on directed graph

- Absorbing nodes
- Source nodes
- Cycles



Perron-Frobenius Theorem

Perron-Frobenius theorem (Fundamental Theorem of Markov Chains)

If matrix is

- stochastic (non-negative and rows sum up to one, describes Markov chain)
- irreducible (strongly connected graph)
- aperiodic

then

$$\exists \lim_{t \rightarrow \infty} \bar{\mathbf{p}}^t = \bar{\pi}$$

and can be found as a left eigenvector

$$\bar{\pi} \mathbf{P} = \bar{\pi}, \quad \text{where } \|\bar{\pi}\|_1 = 1$$

$\bar{\pi}$ - stationary distribution of Markov chain, row vector

Oscar Perron, 1907, Georg Frobenius, 1912.

PageRank

Transition matrix:

$$\mathbf{P} = \mathbf{D}^{-1}\mathbf{A}$$

Stochastic matrix:

$$\mathbf{P}' = \mathbf{P} + \frac{\mathbf{s}\mathbf{e}^T}{n}$$

PageRank matrix:

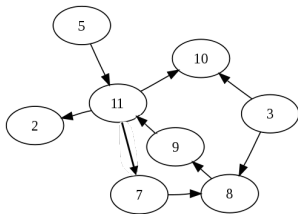
$$\mathbf{P}'' = \alpha\mathbf{P}' + (1 - \alpha)\frac{\mathbf{e}\mathbf{e}^T}{n}$$

Eigenvalue problem (choose solution with $\lambda = 1$):

$$\mathbf{P}''^T \mathbf{p} = \lambda \mathbf{p}$$

Notations:

\mathbf{e} - unit column vector, \mathbf{s} - absorbing nodes indicator vector (column)



- Eigenvalue problem ($\lambda = 1$, $\|\mathbf{p}\|_1 = \mathbf{p}^T \mathbf{e} = 1$):

$$\left(\alpha \mathbf{P}' + (1 - \alpha) \frac{\mathbf{e}\mathbf{e}^T}{n} \right)^T \mathbf{p} = \lambda \mathbf{p}$$

$$\mathbf{p} = \alpha \mathbf{P}'^T \mathbf{p} + (1 - \alpha) \frac{\mathbf{e}}{n}$$

- Power iterations:

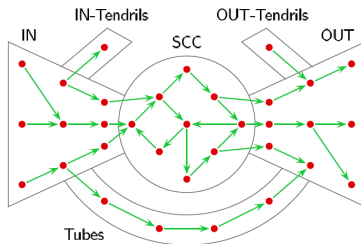
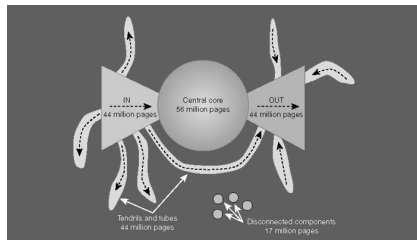
$$\mathbf{p} \leftarrow \alpha \mathbf{P}'^T \mathbf{p} + (1 - \alpha) \frac{\mathbf{e}}{n}$$

- Sparse linear system:

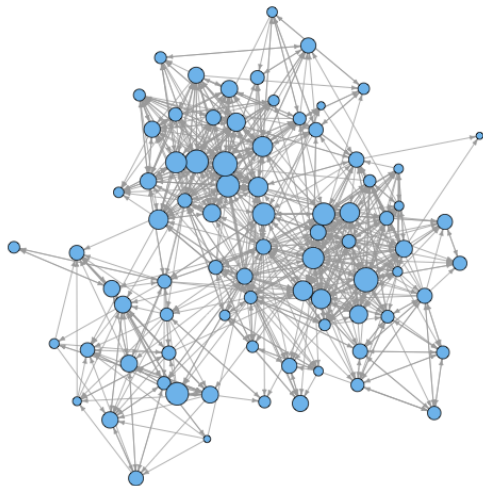
$$(\mathbf{I} - \alpha \mathbf{P}'^T) \mathbf{p} = (1 - \alpha) \frac{\mathbf{e}}{n}$$

Graph structure of the web

Bow tie structure of the web



Andrei Broder et al, 1999



PageRank beyond the Web

1. GeneRank
2. ProteinRank
3. FoodRank
4. SportsRank
5. HostRank
6. TrustRank
7. BadRank
8. ObjectRank
9. ItemRank
10. ArticleRank
11. BookRank
12. FutureRank
13. TimedPageRank
14. SocialPageRank
15. DiffusionRank
16. ImpressionRank
17. TweetRank
18. TwitterRank
19. ReversePageRank
20. PageTrust
21. PopRank
22. CiteRank
23. FactRank
24. InvestorRank
25. ImageRank
26. VisualRank
27. QueryRank
28. BookmarkRank
29. StoryRank
30. PerturbationRank
31. ChemicalRank
32. RoadRank
33. PaperRank
34. Etc...

Hubs and Authorities (HITS)

Citation networks. Reviews vs original research (authoritative) papers

- authorities, contain useful information, a_i
- hubs, contains links to authorities, h_i

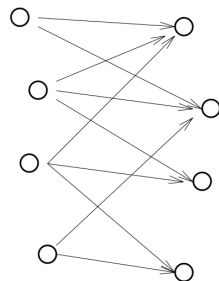
Mutual recursion

- Good authorities referred by good hubs

$$a_i \leftarrow \sum_j A_{ji} h_j$$

- Good hubs point to good authorities

$$h_i \leftarrow \sum_j A_{ij} a_j$$



hubs

authorities

System of linear equations

$$\mathbf{a} = \alpha \mathbf{A}^T \mathbf{h}$$

$$\mathbf{h} = \beta \mathbf{A} \mathbf{a}$$

Symmetric eigenvalue problem

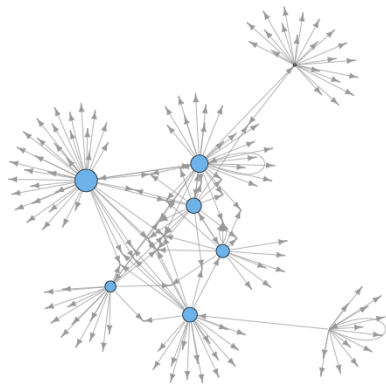
$$(\mathbf{A}^T \mathbf{A}) \mathbf{a} = \lambda \mathbf{a}$$

$$(\mathbf{A} \mathbf{A}^T) \mathbf{h} = \lambda \mathbf{h}$$

where eigenvalue $\lambda = (\alpha\beta)^{-1}$

Hubs and Authorities

Hubs



Authorities



- The PageRank Citation Ranknig: Bringing Order to the Web. S. Brin, L. Page, R. Motwany, T. Winograd, Stanford Digital Library Technologies Project, 1998
- Authoritative Sources in a Hyperlinked Environment. Jon M. Kleinberg, Proc. 9th ACM-SIAM Symposium on Discrete Algorithms,
- Graph structure in the Web, Andrei Broder et all. Procs of the 9th international World Wide Web conference on Computer networks, 2000
- A Survey of Eigenvector Methods of Web Information Retrieval. Amy N. Langville and Carl D. Meyer, 2004
- PageRank beyond the Web. David F. Gleich, arXiv:1407.5107, 2014