# Link Analysis

I. Makarov & L.E. Zhukov

Moscow Institute of Physics and Technology

**Network Science**
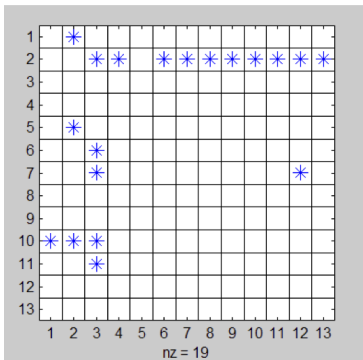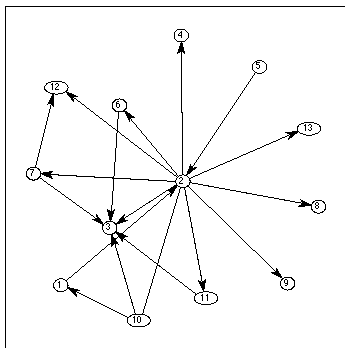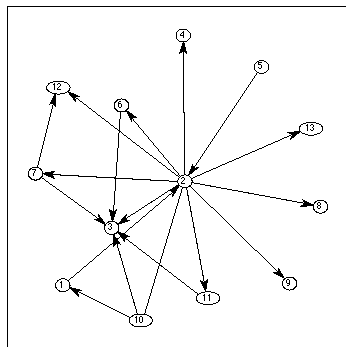
# Lecture outline

# Graph theory

Graph $G(E, V)$, $|V| = n$, $|E| = m$
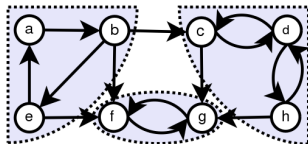Adjacency matrix $A^{n \times n}$, $A_{ij}$, edge $i \to j$



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

# Graph theory



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

# Graph theory

- Graph is **strongly connected** if every vertex is reachable form 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 is each direction between any two pairs of vertices

image from Wikipedia

# Graph theory

- 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)
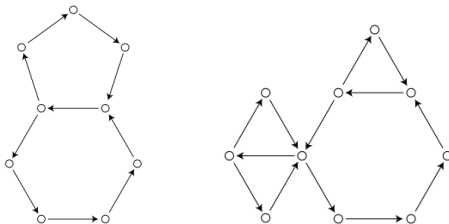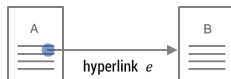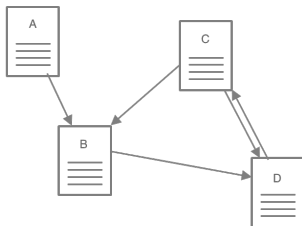


image from Wikipedia

# Web as a graph

- Hyperlinks - implicit endorsements



- Web graph - graph of endorsements (sometimes reciprocal)

# Random walk

- Random walk on a directed graph:

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

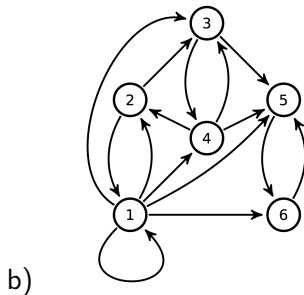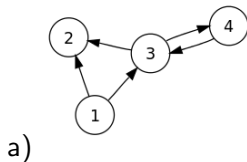$$D_{ii} = diag\{d_i^{out}\}$$

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

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

- Power iterations

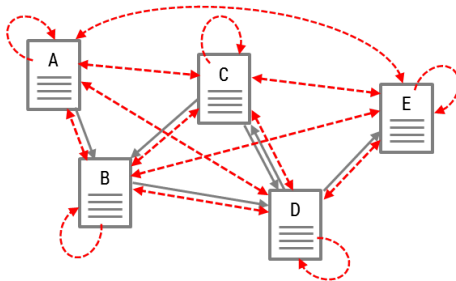$$p^{t+1} \leftarrow P^T p^t$$

a)

b)

# 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

# PageRank formulation

- Power iterations:

$$p \leftarrow \alpha P^T p + (1 - \alpha)\frac{e}{n}, \quad \alpha \text{ - teleportation coefficient}$$

- Sparse linear system:

$$(I - \alpha P^T)p = (1 - \alpha)\frac{e}{n}$$

- Eigenvalue problem ($\lambda = 1$):

$$\left( \alpha P^T + (1 - \alpha)E \right) p = \lambda p$$

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

# 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

$$\exists \lim_{t \to \infty} \bar{\mathsf{p}}^t = \bar{\pi}$$

and can be found as a left eigenvector

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

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

Oscar Perron, 1907, Georg Frobenius,1912.

# PageRank variations

- Power iterations

$$p \leftarrow \alpha P^T p + (1 - \alpha)v, \quad v \text{ - teleportation vector}$$

$$P' = \alpha P + (1 - \alpha)ev^T$$

$$p \leftarrow P'^T p, \ \|p\| = 1$$

- Topic specific PageRank

  v - set of pages on specific topics

- TrustRank

  v - set of trusted pages

- Personalized PageRank

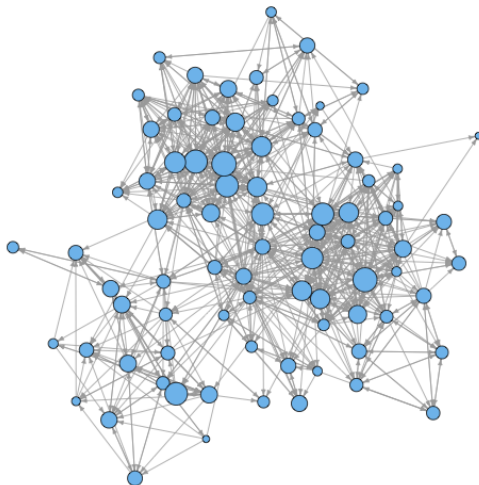  v - set of personal preference pages

Bow tie structure of the web





Andrei Broder et al, 1999

# PageRank

# 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, $a_i$
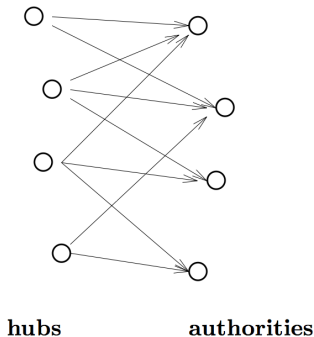- hubs, contains links to authorities, $h_i$

Mutual recursion

- Good authorities reffered 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**

Jon Kleinberg, 1999

# HITS

System of linear equations

$$
\begin{aligned}
a &= \alpha A^T h \\
h &= \beta A a
\end{aligned}
$$

Symmetric eigenvalue problem

$$
\begin{aligned}
(A^T A)a &= \lambda a \\
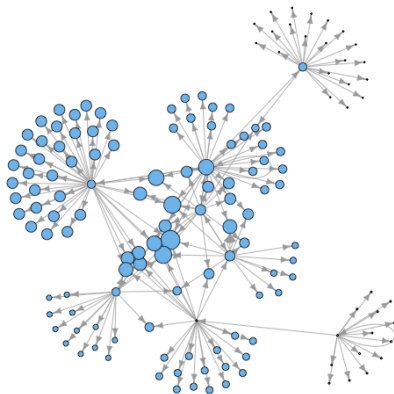(A A^T)h &= \lambda h
\end{aligned}
$$

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

# Hubs and Authorities

Hubs

Authorities

# References

- 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