



---

# Social and Technological Network Analysis

## Lecture 5: Web Search and Random Walks Dr. Cecilia Mascolo

# In This Lecture



- 
- We describe the concept of search in a network.
  - We describe powerful techniques to enhance these searches.

# Search



- When searching “Computer Laboratory” on Google the first link is for the department’s page.
- How does Google know this is the best answer?
  - Information retrieval problem: synonyms (jump/leap), polysemy (Leopard), etc
  - Now with the web: diversity in authoring introduces issues of common criteria for ranking documents
  - The web offers abundance of information: whom do we trust as source?
- Still one issue: static content versus real time
  - World trade center query on 11/9/01
  - Twitter helps solving these issues these days

# Automate the Search



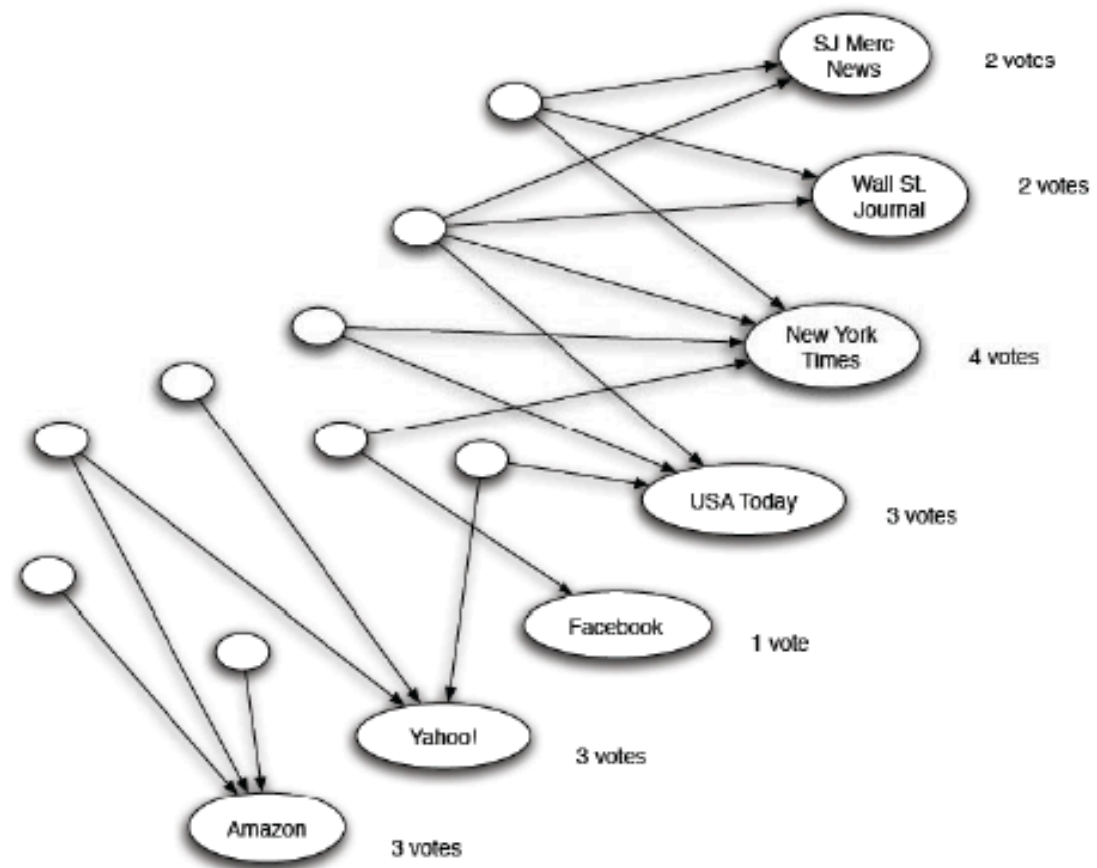
- We could collect a large sample of pages relevant to “computer laboratory” and collect their votes through their links.
- The pages receiving more in-links are ranked first.
- But if we use **the network structure** more deeply we can improve results.

# Example: Query “newspaper”

## Authorities



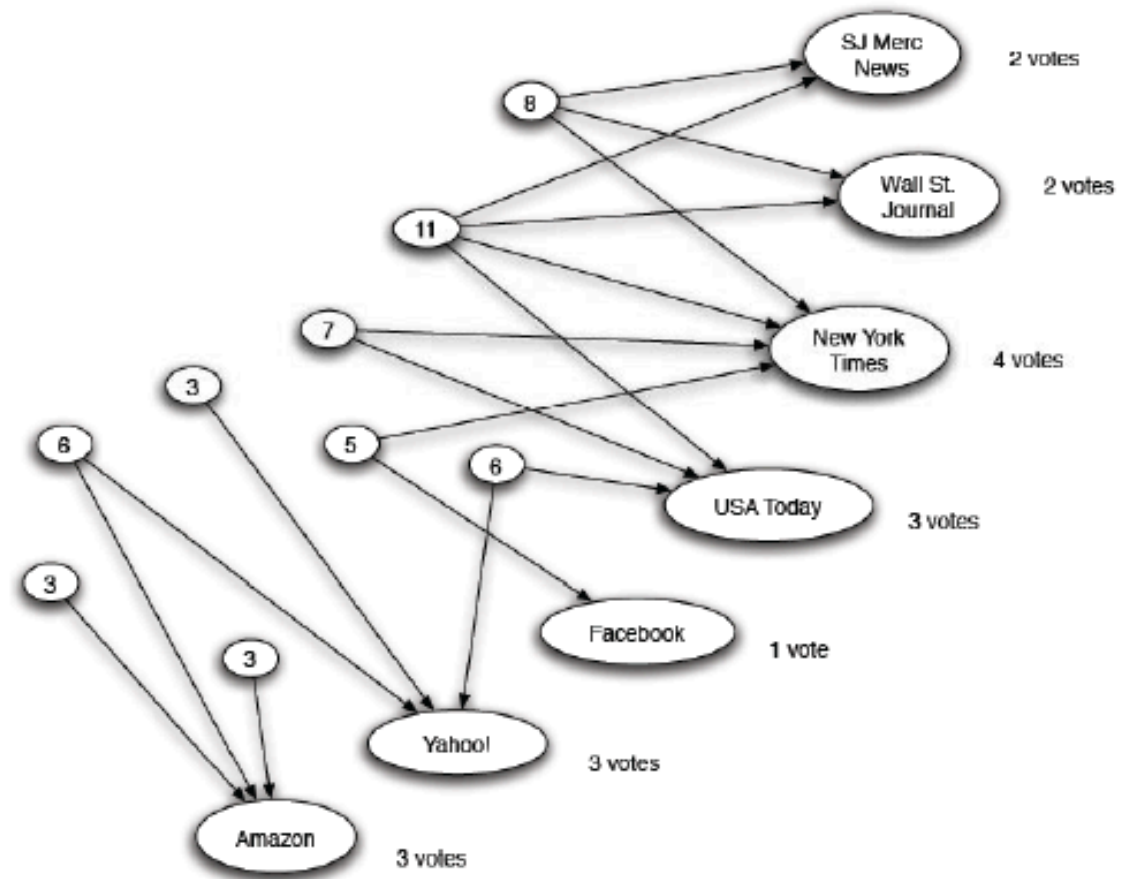
- Links are seen as votes.
- **Authorities** are established: the highly endorsed pages



# A Refinement: Hubs



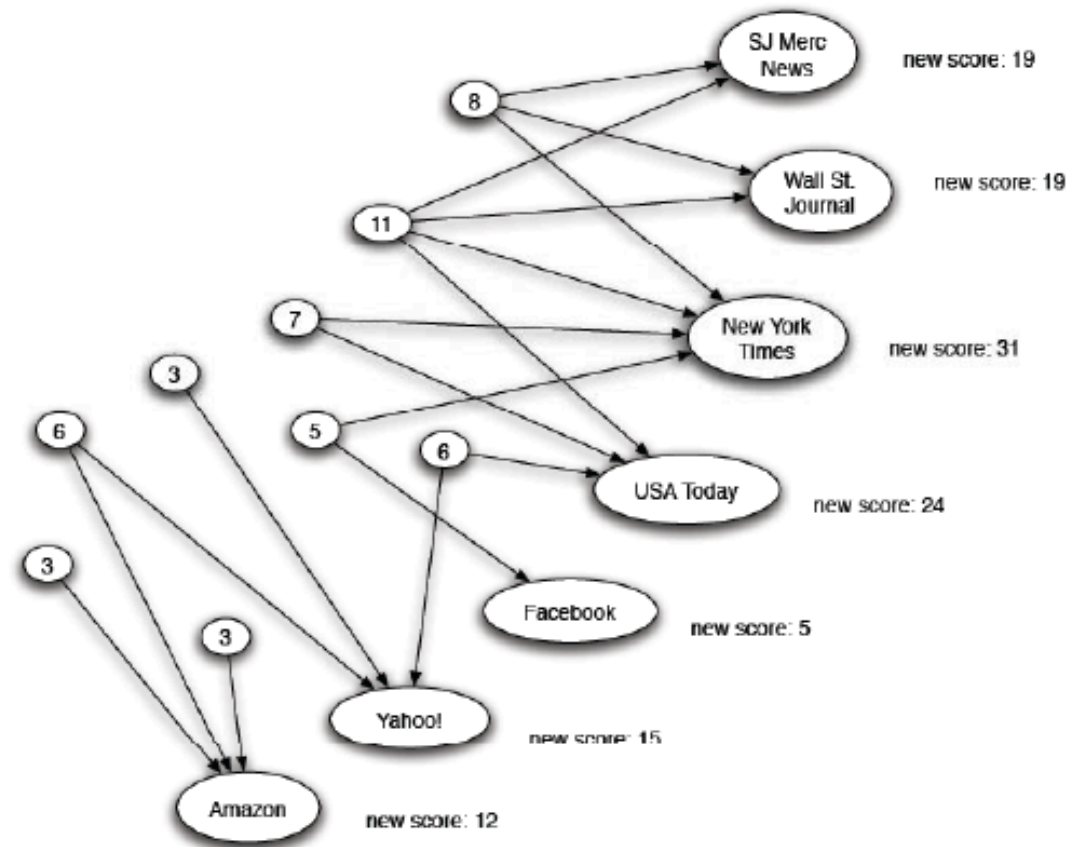
- Numbers are reported back on the source page and aggregate.
- Hubs are high value lists



# Principle of Repeated Improvement



- And we are now reweighting the authorities
- When do we stop?



# Repeating and Normalizing



- The process can be repeated
- Normalization:
  - Each authority score is divided by the sum of all authority scores
  - Each hub score is divided by the sum of all hub scores



# More Formally: does the process converge?

---



- Each page has an authority  $a_i$  and a hub  $h_i$  score
- Initially  $a_i = h_i = 1$

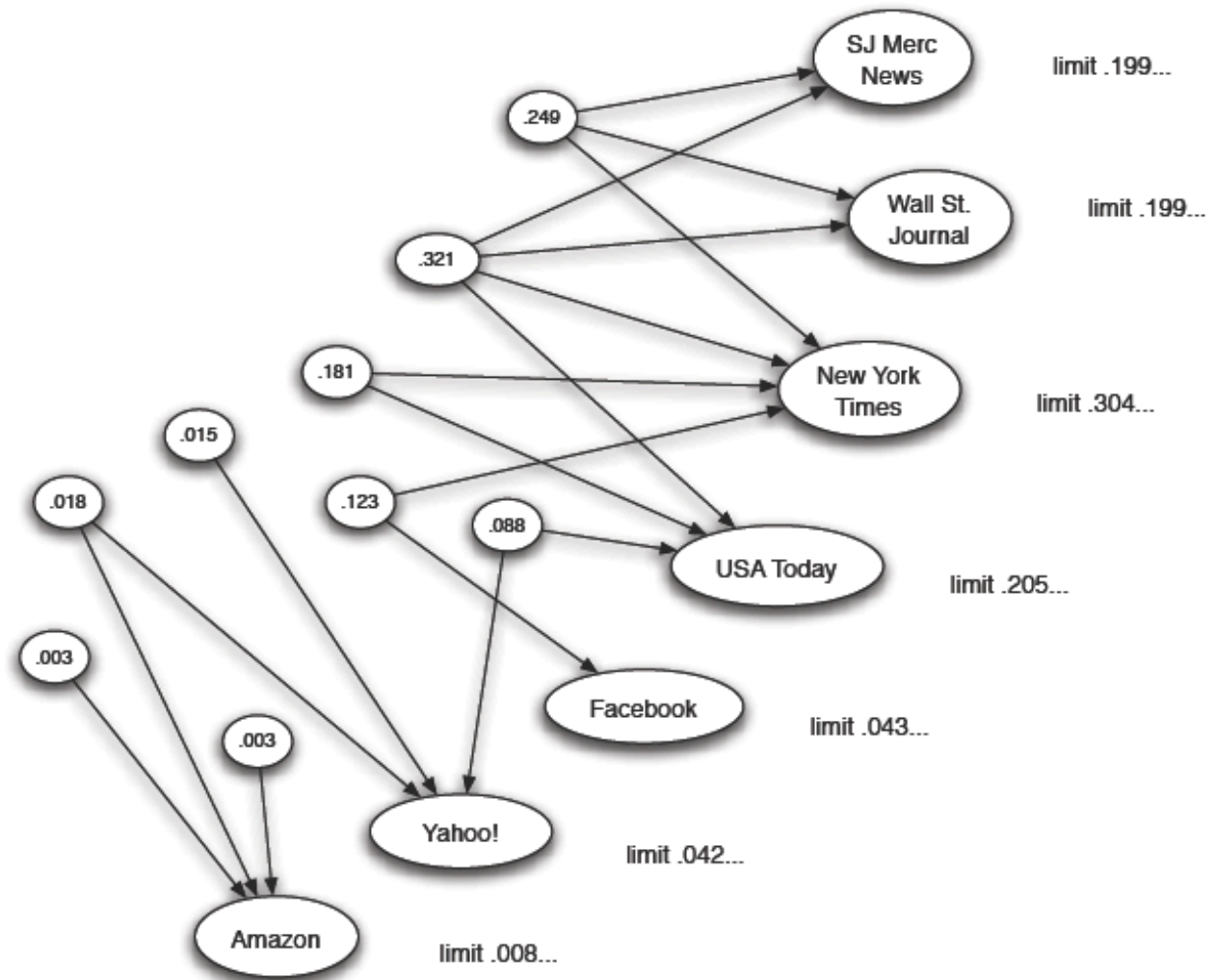
- At each step 
$$a_i = \sum_{j \rightarrow i} h_j$$

$$h_j = \sum_{j \rightarrow i} a_i$$

- Normalize 
$$\sum a_i = 1$$

$$\sum h_j = 1$$

# The process converges





# Let's prove it

---

- Graph seen as matrix  $M_{ij}$  [where  $m_{ij}=1$  if  $i \rightarrow j$ ]
- Vectors  $a=(a_1, a_2, \dots, a_n)$ ,  $h=(h_1, h_2, \dots, h_n)$

$$h_i = \sum_{i \rightarrow j} a_j \Leftrightarrow h_i = \sum_j M_{ij} a_j$$

- So  $h=Ma$  and  $a=M^T h$

# Algorithm



- Initially  $a=h=1^n$
- Step
  - $h=Ma$ ,  $a=M^T h$
  - Normalize
- Then  $a=M^T (Ma)=(M^T M)a$
- And  $h=M(M^T h)=(MM^T)h$
  
- In  $2k$  steps
  - $a=(M^T M)^k a$  and  $h=(MM^T)^k h$

# Eigenvalues and eigenvectors



- **Definition:**
- *Let  $Ax = \lambda x$  for scalar  $\lambda$ , vector  $x$  and matrix  $A$*
- *Then  $x$  is an eigenvector and  $\lambda$  is its eigenvalue*
  
- If  $A$  is symmetric ( $A^T = A$ ) then  $A$  has  $n$  orthogonal unit eigenvectors  $w_1, \dots, w_n$  that form a basis (ie linear independent) with eigenvalues  $\lambda_1, \dots, \lambda_n$  ( $|\lambda_i| \leq |\lambda_{i+1}|$ )
  - Tip:  $M^T M$  and  $M M^T$  are symmetric

So...



- Given the linear independence of the  $w_i$  we can write  $x = \sum p_i w_i$
- And  $x$  has coordinates  $[p_1 \dots p_n]$
- If  $\lambda_1, \dots, \lambda_n$  ( $|\lambda_i| \leq |\lambda_{i+1}|$ ) then
- $A^k x = (\lambda_1^k p_1 w_1 + \lambda_2^k p_2 w_2 + \dots + \lambda_n^k p_n w_n) = \sum \lambda_i^k p_i w_i$
- Let's divide by  $\lambda_1^k$   
 $A^k x / \lambda_1^k = \lambda_1^k p_1 w_1 / \lambda_1^k + \lambda_2^k p_2 w_2 / \lambda_1^k + \dots + \lambda_n^k p_n w_n / \lambda_1^k$
- When  $k \rightarrow \infty$ ,  $A^k x / \lambda_1^k \rightarrow p_1 w_1$

# Convergence



- $(M^T M)^k a$  and  $(M M^T)^k h$
- So for  $k \rightarrow \infty$  these sequences converge

# PageRank



- We have seen hubs and authorities
  - Hubs can “collect” links to important authorities who do not point to each others
  - There are other models: better for the web, where one prominent can endorse another.
- The **PageRank** model is based on transferrable importance.





# PageRank Concepts

---

- Pages pass endorsements on outgoing links as fractions which depend on out-degree
- Initial PageRank value of each node in a network of  $n$  nodes:  $1/n$ .
- Choose a number of steps  $k$ .
- **[Basic] Update rule:** each page divides its pagerank equally over the outgoing links and passes an equal share to the pointed pages. Each page's new rank is the sum of received pageranks.

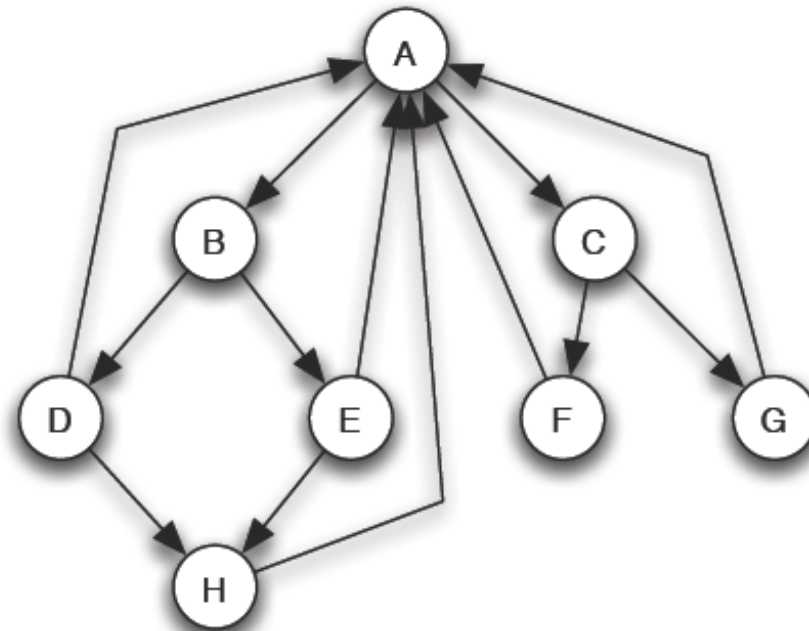
# Example



- All pages start with PageRank=  $1/8$

Step	A	B	C	D	E	F	G	H
1	$1/2$	$1/16$	$1/16$	$1/16$	$1/16$	$1/16$	$1/16$	$1/8$
2	$3/16$	$1/4$	$1/4$	$1/32$	$1/32$	$1/32$	$1/32$	$1/16$

A becomes important and B,C benefit too at step 2

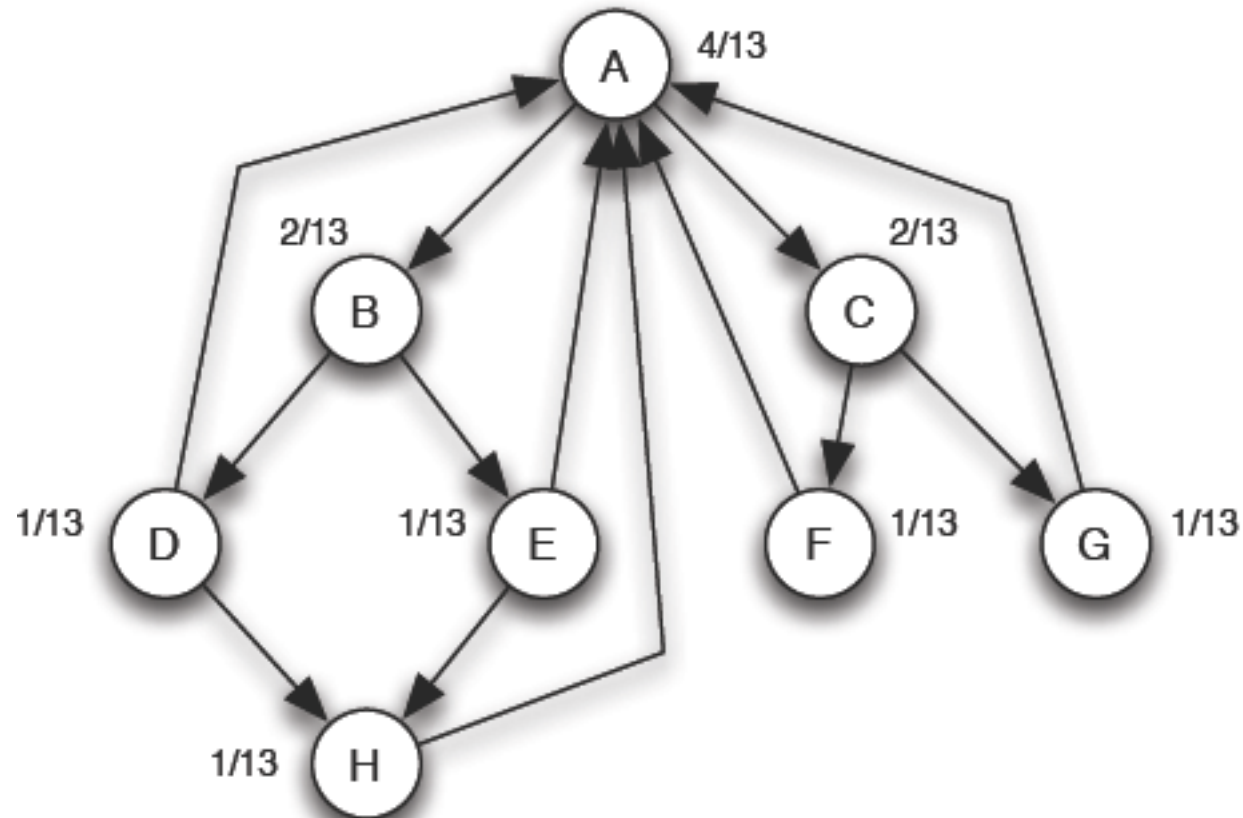


# Convergence



- Except for some special cases, PageRank values of all nodes converge to limiting values when the number of steps goes to infinity.
- The convergence case is one where the PageRank of each page does not change anymore, i.e., they regenerate themselves.

# Example of Equilibrium

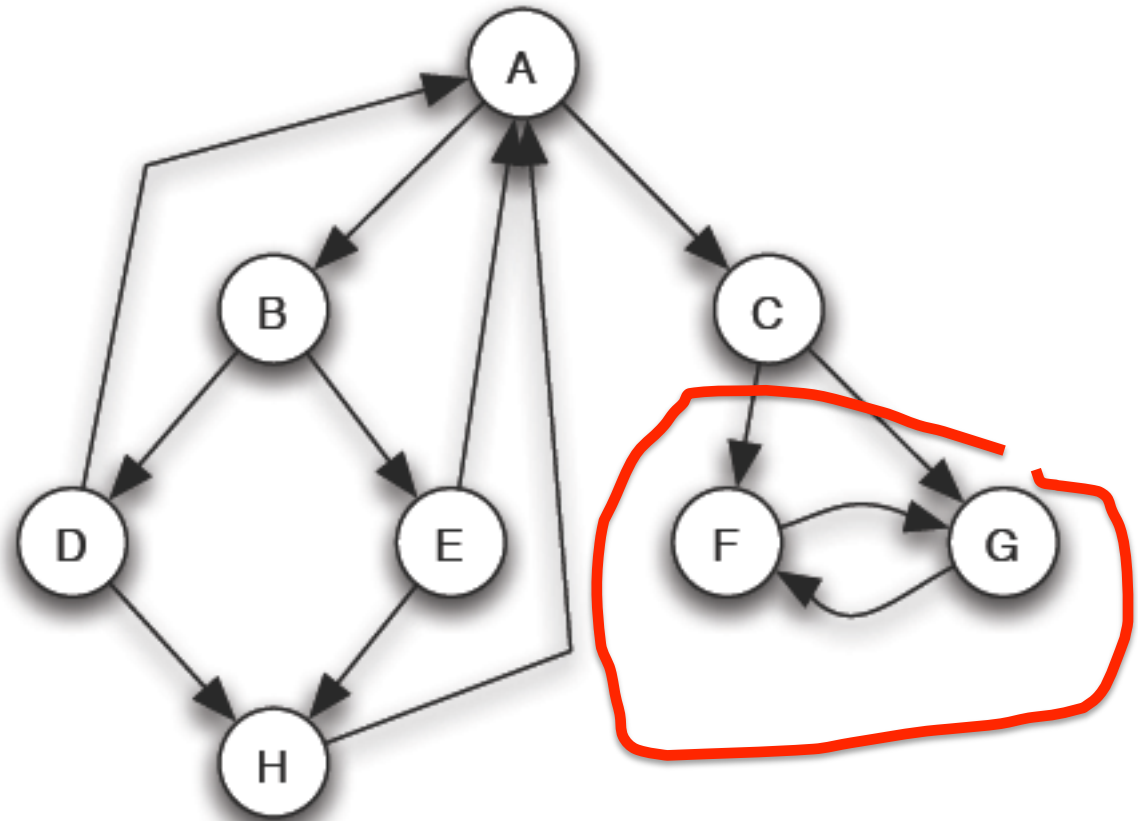


# Problems with the basic PageRank

## Dead ends



- F,G converge to  $\frac{1}{2}$  and all the other nodes to 0



# Solution: The REAL PageRank



- **[Scaled] Update Rule:**
  - Apply basic update rule. Then, scale down all values by scaling factor  $s$  [chosen between 0 and 1].
  - **[Total network PageRank value changes from 1 to  $s$ ]**
  - Divide  $1-s$  residual units of PageRank equally over all nodes:  $(1-s)/n$  each.
- It can be proven that values converge again.
- Scaling factor usually chosen between 0.8 and 0.9

# Search Ranking is very important to business

---



- A change in results in the search pages might mean loss of business
  - I.e., not appearing on first page.
- Ranking algorithms are kept very secret and changed continuously.

# Examples of Google Bombs



Tue, Jan 27 2009 15:05 CET | by Rene Beekman | 1422 V

Google™ who is a failure? Search

oogle™ провал

Web Groups

[President of the United States - George W. Bush](#) · Government  
Article from Encarta Encyclopedia provides an overview of Bush's life.  
[www.whitehouse.gov/president/](#) - 21k - [Cached](#) - [Similar pages](#)

[Historians vs. George W. Bush](#) · Politics  
Of 415 historians who expressed a view of President Bush's administration to this point  
success or **failure**, 338 classified it as a **failure** and 77 as a ...  
[hnn.us/articles/5019.html](#) - 38k - [Cached](#) - [Similar pages](#)

[Heart failure - Wikipedia, the free encyclopedia](#)  
Congestive heart **failure** (CHF), congestive cardiac **failure** (CCF) or just heart **failure**, i  
condition that can result from any structural or functional ...  
[en.wikipedia.org/wiki/Heart\\_failure](#) - 146k - [Cached](#) - [Similar pages](#)

[ал или как да им спретнем един Googlebomb](#)  
2009 ... **провал** или как да направим Googlebomb на едно тотално  
пено правителство.  
[net/vile/failure/](#) - 91k - [Cached](#) - [Similar pages](#)

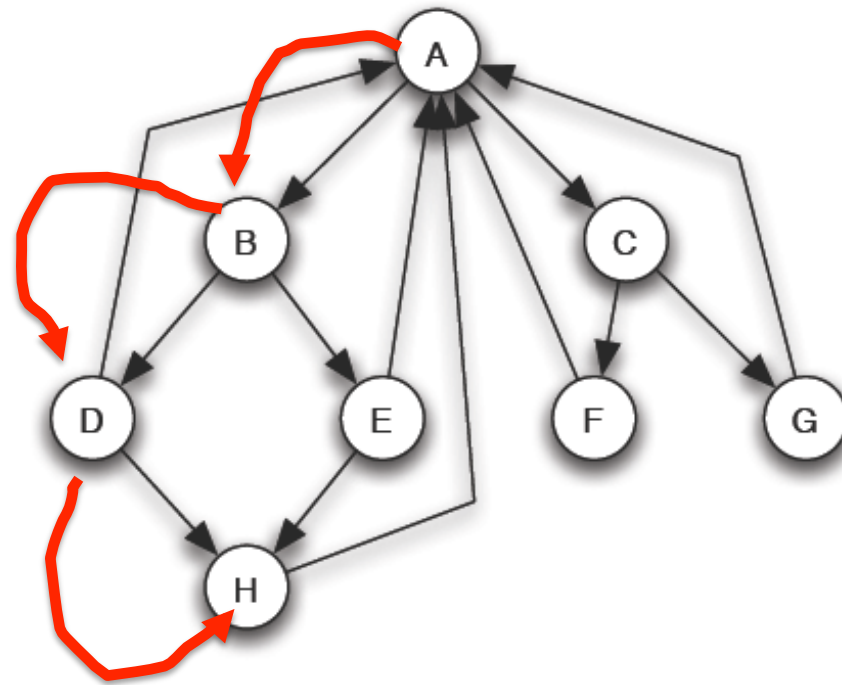
1 of 1



# Random Walks



- Starting from a node, follow one outgoing link with random probability



# PageRank as Random Walk



- The probability of being at a page  $X$  after  $k$  steps of a random walk is precisely the PageRank of  $X$  after  $k$  applications of the Basic PageRank Update Rule
- Scaled Update Rule equivalent: follow a random outgoing link with probability  $s$  while with probability  $1-s$  jump to a random node in the network.

# References

---



- Chapter 14