



Web Search (Google)

MA 723

Department of Mathematics
North Carolina State University
Raleigh, NC

4/28/2005

Newsweek

March 29, 2004 : \$3.95

newsweek.msnbc.com

The Next Frontiers

The New Age of Google

The Search Giant Has Changed
Our Lives. Can Anybody
Catch These Guys? **By Steven Levy**

Google founders Larry Page and Sergey Brin



Using Link Structure

Indexing

- First index key terms on each page
Robots crawl the web — software does indexing
- File structure: Terms \longrightarrow Pages (similar to book index)
 $Term_1 \rightarrow P_i, P_j, \dots$
 $Term_2 \rightarrow P_k, P_l, \dots$
 \vdots



Using Link Structure

Indexing

- First index key terms on each page
Robots crawl the web — software does indexing
- File structure: Terms \longrightarrow Pages (similar to book index)
 $Term_1 \rightarrow P_i, P_j, \dots$
 $Term_2 \rightarrow P_k, P_l, \dots$
 \vdots

Importance Rankings

- Attach an “importance rank” r_i to each page: $P_i \hookrightarrow r_i$



Using Link Structure

Indexing

- First index key terms on each page
Robots crawl the web — software does indexing
- File structure: Terms \longrightarrow Pages (similar to book index)
 $Term_1 \rightarrow P_i, P_j, \dots$
 $Term_2 \rightarrow P_k, P_l, \dots$
 \vdots

Importance Rankings

- Attach an “importance rank” r_i to each page: $P_i \hookrightarrow r_i$
— r_i based on link structure (i.e., query independent)



Using Link Structure

Indexing

- First index key terms on each page
Robots crawl the web — software does indexing
- File structure: Terms \longrightarrow Pages (similar to book index)
 $Term_1 \rightarrow P_i, P_j, \dots$
 $Term_2 \rightarrow P_k, P_l, \dots$
 \vdots

Importance Rankings

- Attach an “importance rank” r_i to each page: $P_i \hookrightarrow r_i$
 - r_i based on link structure (i.e., query independent)
 - r_i computed prior to any query



Using Link Structure

Indexing

- First index key terms on each page
Robots crawl the web — software does indexing
- File structure: Terms \longrightarrow Pages (similar to book index)
 $Term_1 \rightarrow P_i, P_j, \dots$
 $Term_2 \rightarrow P_k, P_l, \dots$
 \vdots

Importance Rankings

- Attach an “importance rank” r_i to each page: $P_i \hookrightarrow r_i$
 - r_i based on link structure (i.e., query independent)
 - r_i computed prior to any query

Direct Query Matching

- Query = $(Term_1, Term_2) \longrightarrow (P_i, r_i), (P_j, r_j), (P_k, r_k), \dots$



Using Link Structure

Indexing

- First index key terms on each page
Robots crawl the web — software does indexing
- File structure: Terms \longrightarrow Pages (similar to book index)
 $Term_1 \rightarrow P_i, P_j, \dots$
 $Term_2 \rightarrow P_k, P_l, \dots$
 \vdots

Importance Rankings

- Attach an “importance rank” r_i to each page: $P_i \hookrightarrow r_i$
 - r_i based on link structure (i.e., query independent)
 - r_i computed prior to any query

Direct Query Matching

- Query = $(Term_1, Term_2) \longrightarrow (P_i, r_i), (P_j, r_j), (P_k, r_k), \dots$

Return P_i, P_j, P_k, \dots in order of ranks r_i, r_j, r_k, \dots



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent

- Depends primarily on link structure of web



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent

- Depends primarily on link structure of web
 - Off-line calculations
 - No computation at query time



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent

- Depends primarily on link structure of web
 - Off-line calculations
 - No computation at query time

$r(P)$ Depends On Ranks Of Pages Pointing To P

- Importance is not number of in-links or out-links



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent

- Depends primarily on link structure of web
 - Off-line calculations
 - No computation at query time

$r(P)$ Depends On Ranks Of Pages Pointing To P

- Importance is not number of in-links or out-links
 - One link to P from Yahoo! is important



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent

- Depends primarily on link structure of web
 - Off-line calculations
 - No computation at query time

$r(P)$ Depends On Ranks Of Pages Pointing To P

- Importance is not number of in-links or out-links
 - One link to P from Yahoo! is important
 - Many links to P from me is not



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent

- Depends primarily on link structure of web
 - Off-line calculations
 - No computation at query time

$r(P)$ Depends On Ranks Of Pages Pointing To P

- Importance is not number of in-links or out-links
 - One link to P from Yahoo! is important
 - Many links to P from me is not

PageRank Shares The Vote

- Yahoo! casts many “votes” \implies value of vote from Y is diluted



Google's PageRank

(Lawrence Page & Sergey Brin 1998)

PageRank $r(P)$ Is Not Query Dependent

- Depends primarily on link structure of web
 - Off-line calculations
 - No computation at query time

$r(P)$ Depends On Ranks Of Pages Pointing To P

- Importance is not number of in-links or out-links
 - One link to P from Yahoo! is important
 - Many links to P from me is not

PageRank Shares The Vote

- Yahoo! casts many “votes” \implies value of vote from Y is diluted
 - If Yahoo! “votes” for n pages
 - then P receives only $r(Y)/n$ credit from Y



PageRank

The Definition

$$r(P) = \sum_{P \in \mathcal{B}_P} \frac{r(P)}{|P|}$$

$\mathcal{B}_P = \{\text{all pages pointing to } P\}$

$|P| = \text{number of out links from } P$



PageRank

The Definition

$$r(P) = \sum_{P \in \mathcal{B}_P} \frac{r(P)}{|P|}$$

$\mathcal{B}_P = \{\text{all pages pointing to } P\}$

$|P| = \text{number of out links from } P$

Successive Refinement

Start with $r_0(P_i) = 1/n$ for all pages P_1, P_2, \dots, P_n



PageRank

The Definition

$$r(P) = \sum_{P \in \mathcal{B}_P} \frac{r(P)}{|P|}$$

$\mathcal{B}_P = \{\text{all pages pointing to } P\}$

$|P| = \text{number of out links from } P$

Successive Refinement

Start with $r_0(P_i) = 1/n$ for all pages P_1, P_2, \dots, P_n

Iteratively refine rankings for each page



PageRank

The Definition

$$r(P) = \sum_{P \in \mathcal{B}_P} \frac{r(P)}{|P|}$$

$\mathcal{B}_P = \{\text{all pages pointing to } P\}$

$|P| = \text{number of out links from } P$

Successive Refinement

Start with $r_0(P_i) = 1/n$ for all pages P_1, P_2, \dots, P_n

Iteratively refine rankings for each page

$$r_1(P_i) = \sum_{P \in \mathcal{B}_{P_i}} \frac{r_0(P)}{|P|}$$



PageRank

The Definition

$$r(P) = \sum_{P \in \mathcal{B}_P} \frac{r(P)}{|P|}$$

$\mathcal{B}_P = \{\text{all pages pointing to } P\}$

$|P| = \text{number of out links from } P$

Successive Refinement

Start with $r_0(P_i) = 1/n$ for all pages P_1, P_2, \dots, P_n

Iteratively refine rankings for each page

$$r_1(P_i) = \sum_{P \in \mathcal{B}_{P_i}} \frac{r_0(P)}{|P|}$$

$$r_2(P_i) = \sum_{P \in \mathcal{B}_{P_i}} \frac{r_1(P)}{|P|}$$



PageRank

The Definition

$$r(P) = \sum_{P \in \mathcal{B}_P} \frac{r(P)}{|P|}$$

$\mathcal{B}_P = \{\text{all pages pointing to } P\}$

$|P| = \text{number of out links from } P$

Successive Refinement

Start with $r_0(P_i) = 1/n$ for all pages P_1, P_2, \dots, P_n

Iteratively refine rankings for each page

$$r_1(P_i) = \sum_{P \in \mathcal{B}_{P_i}} \frac{r_0(P)}{|P|}$$

$$r_2(P_i) = \sum_{P \in \mathcal{B}_{P_i}} \frac{r_1(P)}{|P|}$$

\vdots

$$r_{j+1}(P_i) = \sum_{P \in \mathcal{B}_{P_i}} \frac{r_j(P)}{|P|}$$



In Matrix Notation

After Step j

$$\boldsymbol{\pi}_j^T = [r_j(P_1), r_j(P_2), \dots, r_j(P_n)]$$



In Matrix Notation

After Step j

$$\boldsymbol{\pi}_j^T = [r_j(P_1), r_j(P_2), \dots, r_j(P_n)]$$

$$\boldsymbol{\pi}_{j+1}^T = \boldsymbol{\pi}_j^T \mathbf{P} \quad \text{where} \quad p_{ij} = \begin{cases} 1/|P_i| & \text{if } i \rightarrow j \\ 0 & \text{otherwise} \end{cases}$$



In Matrix Notation

After Step j

$$\pi_j^T = [r_j(P_1), r_j(P_2), \dots, r_j(P_n)]$$

$$\pi_{j+1}^T = \pi_j^T \mathbf{P} \quad \text{where} \quad p_{ij} = \begin{cases} 1/|P_i| & \text{if } i \rightarrow j \\ 0 & \text{otherwise} \end{cases}$$

$$\text{PageRank} = \lim_{j \rightarrow \infty} \pi_j^T = \pi^T \quad (\text{provided limit exists})$$



In Matrix Notation

After Step j

$$\pi_j^T = [r_j(P_1), r_j(P_2), \dots, r_j(P_n)]$$

$$\pi_{j+1}^T = \pi_j^T \mathbf{P} \quad \text{where} \quad p_{ij} = \begin{cases} 1/|P_i| & \text{if } i \rightarrow j \\ 0 & \text{otherwise} \end{cases}$$

$$\text{PageRank} = \lim_{j \rightarrow \infty} \pi_j^T = \pi^T \quad (\text{provided limit exists})$$

Maybe It's A Markov Chain?

$$\text{If } \mathbf{P} = [p_{ij}] \text{ is a stochastic matrix} \quad (p_{ij} \geq 0 \text{ and } \sum_j p_{ij} = 1)$$



In Matrix Notation

After Step j

$$\pi_j^T = [r_j(P_1), r_j(P_2), \dots, r_j(P_n)]$$

$$\pi_{j+1}^T = \pi_j^T \mathbf{P} \quad \text{where} \quad p_{ij} = \begin{cases} 1/|P_i| & \text{if } i \rightarrow j \\ 0 & \text{otherwise} \end{cases}$$

$$\text{PageRank} = \lim_{j \rightarrow \infty} \pi_j^T = \pi^T \quad (\text{provided limit exists})$$

Maybe It's A Markov Chain?

If $\mathbf{P} = [p_{ij}]$ is a stochastic matrix ($p_{ij} \geq 0$ and $\sum_j p_{ij} = 1$)

Each π_j^T is a probability vector ($\pi_i \geq 0$ and $\sum_i \pi_i = 1$)



In Matrix Notation

After Step j

$$\pi_j^T = [r_j(P_1), r_j(P_2), \dots, r_j(P_n)]$$

$$\pi_{j+1}^T = \pi_j^T \mathbf{P} \quad \text{where} \quad p_{ij} = \begin{cases} 1/|P_i| & \text{if } i \rightarrow j \\ 0 & \text{otherwise} \end{cases}$$

$$\text{PageRank} = \lim_{j \rightarrow \infty} \pi_j^T = \pi^T \quad (\text{provided limit exists})$$

Maybe It's A Markov Chain?

If $\mathbf{P} = [p_{ij}]$ is a stochastic matrix ($p_{ij} \geq 0$ and $\sum_j p_{ij} = 1$)

Each π_j^T is a probability vector ($\pi_i \geq 0$ and $\sum_i \pi_i = 1$)

$\pi_{j+1}^T = \pi_j^T \mathbf{P}$ is random walk on the graph defined by links



In Matrix Notation

After Step j

$$\pi_j^T = [r_j(P_1), r_j(P_2), \dots, r_j(P_n)]$$

$$\pi_{j+1}^T = \pi_j^T \mathbf{P} \quad \text{where} \quad p_{ij} = \begin{cases} 1/|P_i| & \text{if } i \rightarrow j \\ 0 & \text{otherwise} \end{cases}$$

$$\text{PageRank} = \lim_{j \rightarrow \infty} \pi_j^T = \pi^T \quad (\text{provided limit exists})$$

Maybe It's A Markov Chain?

If $\mathbf{P} = [p_{ij}]$ is a stochastic matrix ($p_{ij} \geq 0$ and $\sum_j p_{ij} = 1$)

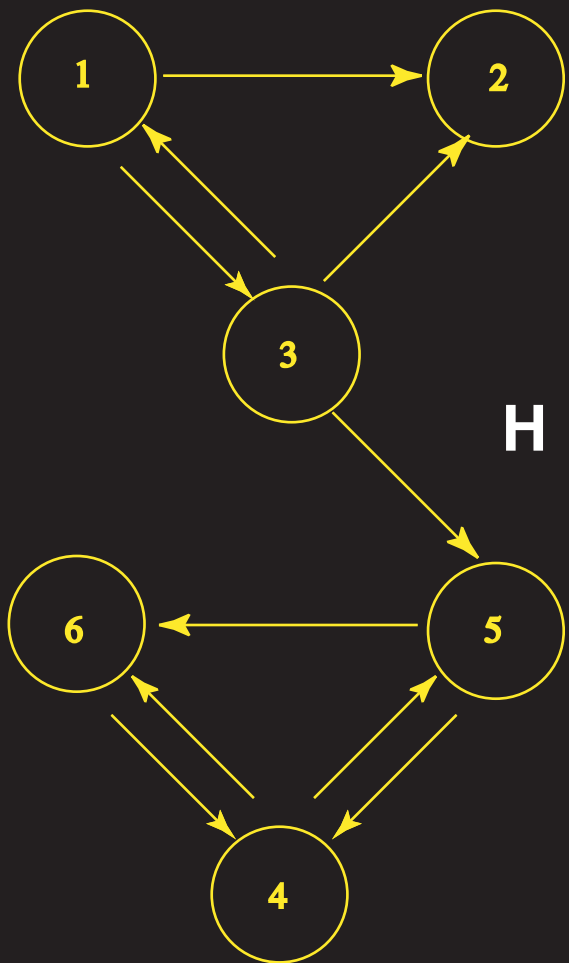
Each π_j^T is a probability vector ($\pi_i \geq 0$ and $\sum_i \pi_i = 1$)

$\pi_{j+1}^T = \pi_j^T \mathbf{P}$ is random walk on the graph defined by links

$\pi^T = \lim_{j \rightarrow \infty} \pi_j^T =$ steady-state probability distribution



Tiny Web

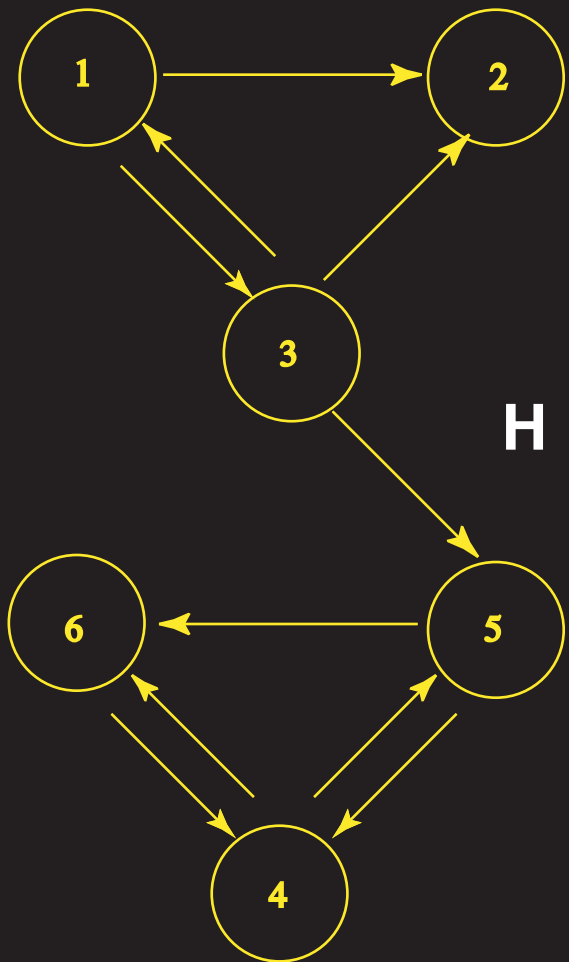


H =

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \end{pmatrix}$$



Tiny Web

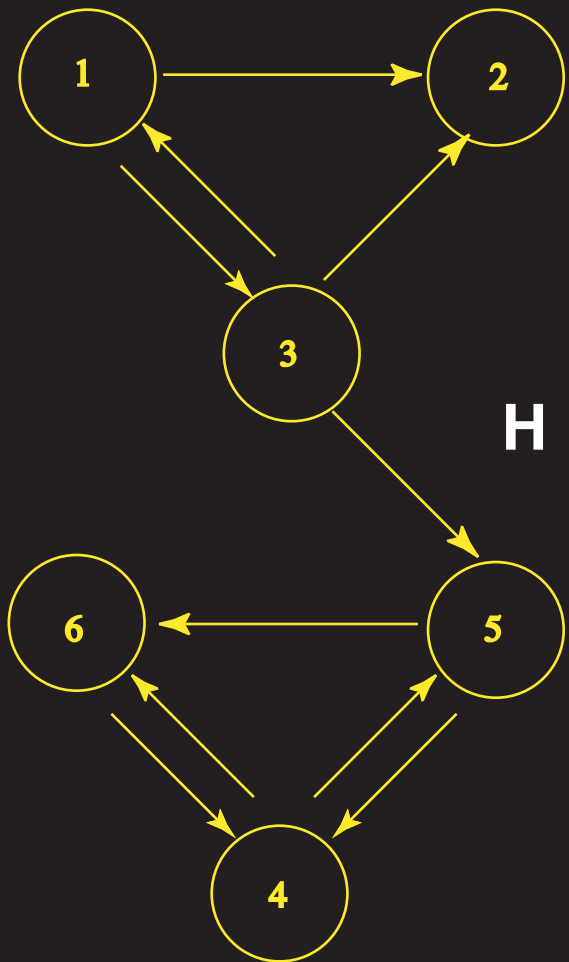


H =

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ 0 & 1/2 & 1/2 & 0 & 0 & 0 \end{pmatrix}$$



Tiny Web

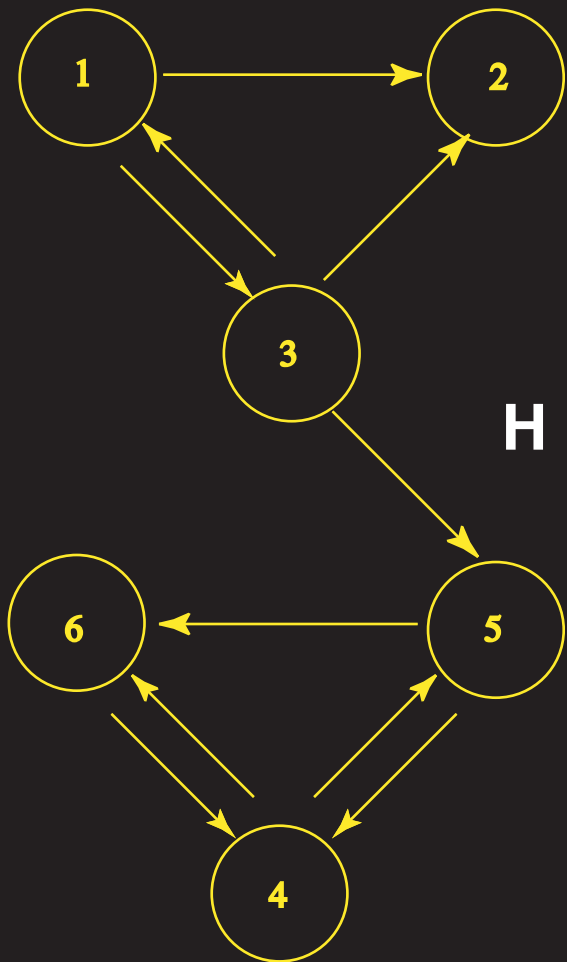


H =

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix}$$



Tiny Web

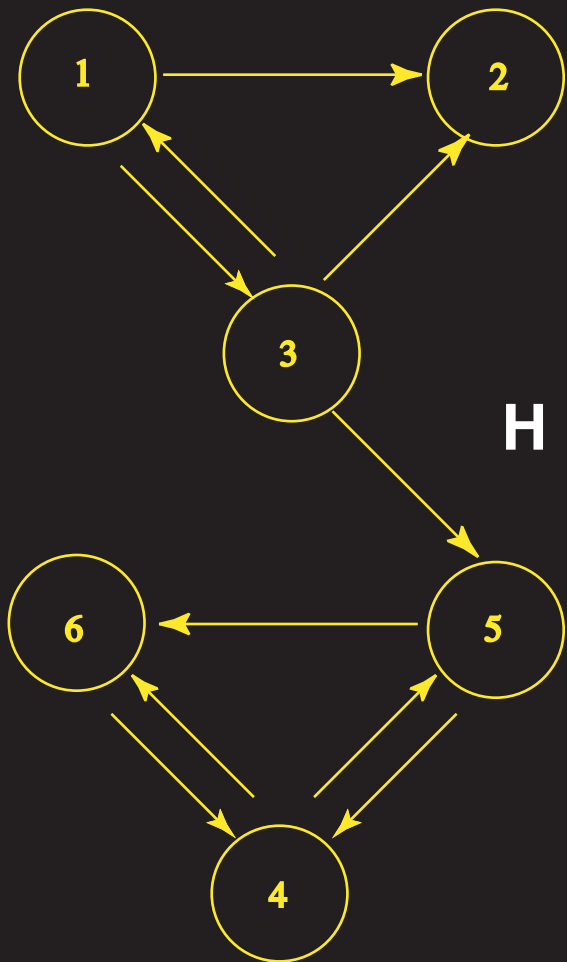


H =

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$



Tiny Web

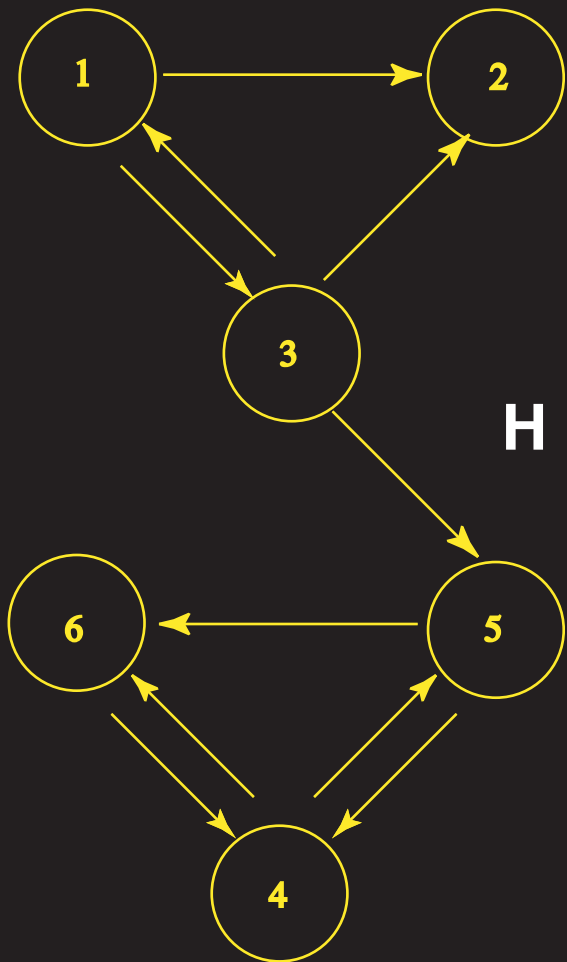


H =

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$



Tiny Web

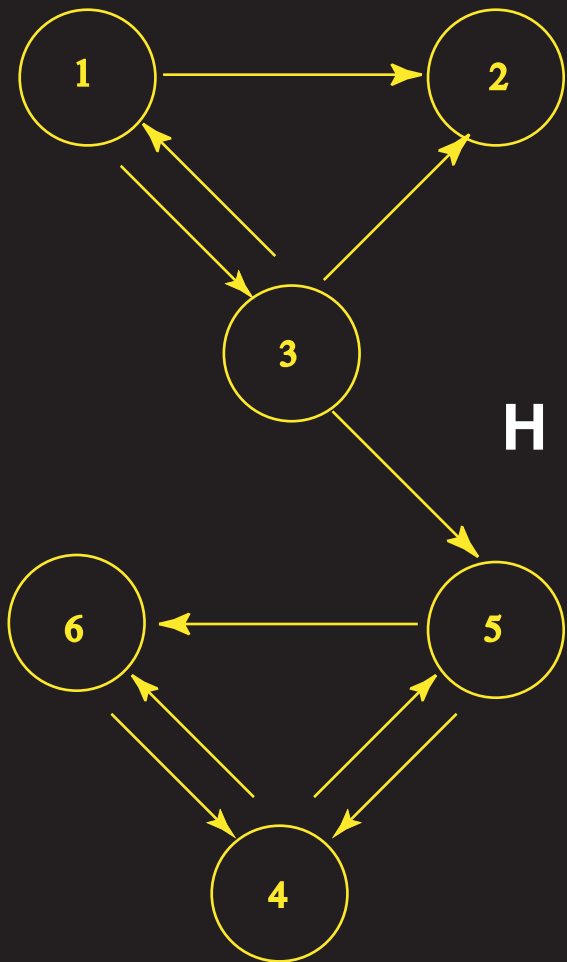


H =

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ 0 & 0 & 0 & 1/2 & 0 & 1/2 \end{pmatrix}$$



Tiny Web

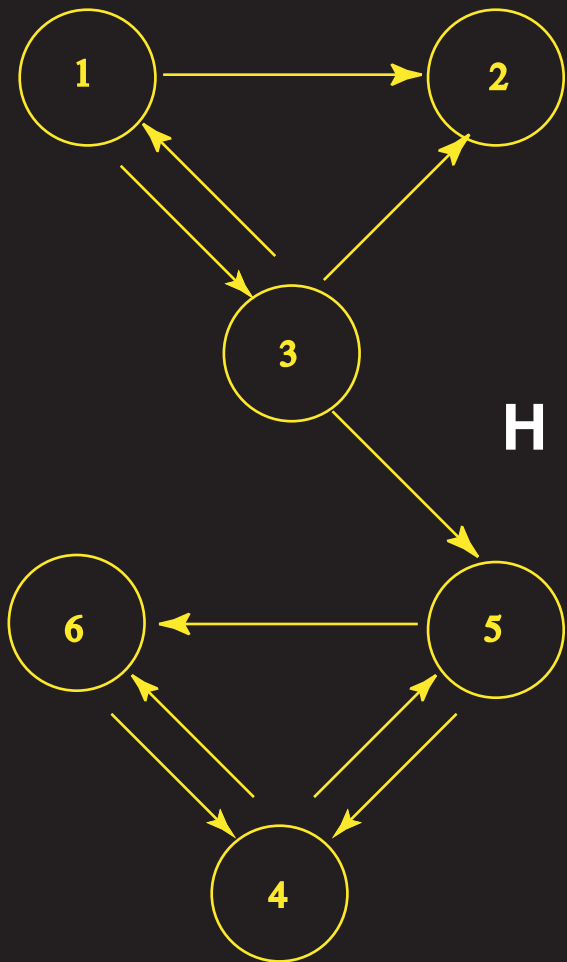


H =

$$\begin{matrix} & P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ P_1 & 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ P_2 & 0 & 0 & 0 & 0 & 0 & 0 \\ P_3 & 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ P_4 & 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ P_5 & 0 & 0 & 0 & 1/2 & 0 & 1/2 \\ P_6 & 0 & 0 & 0 & 1 & 0 & 0 \end{matrix}$$



Tiny Web



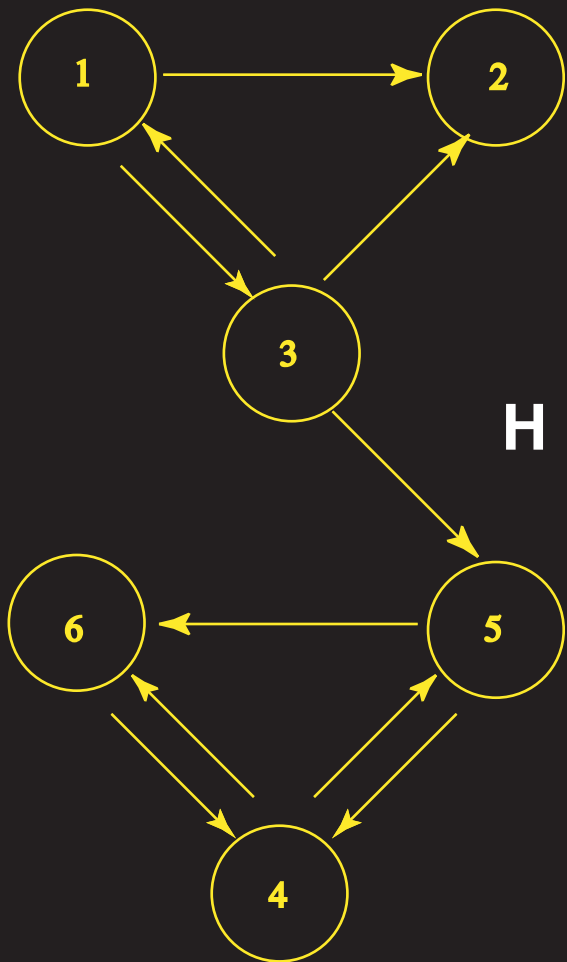
$\mathbf{H} =$

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ 0 & 0 & 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

- Dead end page (nothing to click on) — a “dangling node”



Tiny Web



$\mathbf{H} =$

$$\begin{matrix} P_1 \\ P_2 \\ P_3 \\ P_4 \\ P_5 \\ P_6 \end{matrix} \begin{pmatrix} P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ 0 & 0 & 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

- Dead end page (nothing to click on) — a “dangling node”

✓ π^T not well defined



The Fix

- Replace zero rows with $(1/n)\mathbf{e}^T = (1/n, 1/n, \dots, 1/n)$

$$\mathbf{S} = \begin{matrix} & P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ P_1 & 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ P_2 & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} & \mathbf{1/6} \\ P_3 & 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ P_4 & 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ P_5 & 0 & 0 & 0 & 1/2 & 0 & 1/2 \\ P_6 & 0 & 0 & 0 & 1 & 0 & 0 \end{matrix}$$



The Fix

- Replace zero rows with $(1/n)\mathbf{e}^T = (1/n, 1/n, \dots, 1/n)$

$$\mathbf{S} = \begin{matrix} & P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ P_1 & \begin{pmatrix} 0 & 1/2 & 1/2 & 0 & 0 & 0 \end{pmatrix} \\ P_2 & \begin{pmatrix} 1/6 & 1/6 & 1/6 & 1/6 & 1/6 & 1/6 \end{pmatrix} \\ P_3 & \begin{pmatrix} 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \end{pmatrix} \\ P_4 & \begin{pmatrix} 0 & 0 & 0 & 0 & 1/2 & 1/2 \end{pmatrix} \\ P_5 & \begin{pmatrix} 0 & 0 & 0 & 1/2 & 0 & 1/2 \end{pmatrix} \\ P_6 & \begin{pmatrix} 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix} \end{matrix}$$

$$\mathbf{S} = \mathbf{H} + \frac{\mathbf{a}\mathbf{e}^T}{6} = \mathbf{H} + \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \frac{1}{6} (1 \ 1 \ 1 \ 1 \ 1 \ 1)$$



Another Problem

- S is reducible

$$\mathbf{S} = \begin{array}{c|ccc|ccc} & P_1 & P_2 & P_3 & P_4 & P_5 & P_6 \\ \hline P_1 & 0 & 1/2 & 1/2 & 0 & 0 & 0 \\ P_2 & 1/6 & 1/6 & 1/6 & 1/6 & 1/6 & 1/6 \\ P_3 & 1/3 & 1/3 & 0 & 0 & 1/3 & 0 \\ \hline P_4 & 0 & 0 & 0 & 0 & 1/2 & 1/2 \\ P_5 & 0 & 0 & 0 & 1/2 & 0 & 1/2 \\ P_6 & 0 & 0 & 0 & 1 & 0 & 0 \end{array}$$

✓ π^T not well defined



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$

✓ $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ won't convergence



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$

✓ $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ won't convergence

Convergence Requirement

Markov chain must be irreducible and aperiodic



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$

✓ $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ won't convergence

Convergence Requirement

Markov chain must be irreducible and aperiodic

- This means \mathbf{P} must be a primitive matrix

✓ No eigenvalues other than $\lambda = 1$ on unit circle



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$

✓ $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ won't convergence

Convergence Requirement

Markov chain must be irreducible and aperiodic

- This means \mathbf{P} must be a primitive matrix
 - ✓ No eigenvalues other than $\lambda = 1$ on unit circle
 - ✓ $\mathbf{P}^k > 0$ for some k



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$

✓ $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ won't convergence

Convergence Requirement

Markov chain must be irreducible and aperiodic

- This means \mathbf{P} must be a primitive matrix

✓ No eigenvalues other than $\lambda = 1$ on unit circle

✓ $\mathbf{P}^k > \mathbf{0}$ for some k

The Google Fixes

- $\mathbf{P} = \alpha \mathbf{S} + (1 - \alpha) \mathbf{e} \mathbf{e}^T / n$ $\alpha \approx .85$



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$

✓ $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ won't convergence

Convergence Requirement

Markov chain must be irreducible and aperiodic

- This means \mathbf{P} must be a primitive matrix

✓ No eigenvalues other than $\lambda = 1$ on unit circle

✓ $\mathbf{P}^k > 0$ for some k

The Google Fixes

- $\mathbf{P} = \alpha \mathbf{S} + (1 - \alpha) \mathbf{e} \mathbf{e}^T / n$ $\alpha \approx .85$

- $\mathbf{P} = \alpha \mathbf{S} + (1 - \alpha) \mathbf{e} \mathbf{v}^T$ $\mathbf{v}^T =$ positive probability vector



Yet More Problems

Could get trapped into a cycle $(P_i \rightarrow P_j \rightarrow P_i)$

✓ $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ won't convergence

Convergence Requirement

Markov chain must be irreducible and aperiodic

- This means \mathbf{P} must be a primitive matrix

✓ No eigenvalues other than $\lambda = 1$ on unit circle

✓ $\mathbf{P}^k > 0$ for some k

The Google Fixes

- $\mathbf{P} = \alpha \mathbf{S} + (1 - \alpha) \mathbf{e} \mathbf{e}^T / n$ $\alpha \approx .85$

- $\mathbf{P} = \alpha \mathbf{S} + (1 - \alpha) \mathbf{e} \mathbf{v}^T$ $\mathbf{v}^T =$ positive probability vector

- $\mathbf{P} = \alpha \mathbf{H} + (\alpha \mathbf{a} + (1 - \alpha) \mathbf{e}) \mathbf{v}^T$

THE WALL STREET JOURNAL

© 2003 Dow Jones & Company. All Rights Reserved

WEDNESDAY, FEBRUARY 26, 2003 - VOL. CCXLI NO. 39 - ★★★ \$1.00

WSJ.com

What's News—

Business and Finance

World-Wide

NEWS CORP. and Liberty are no longer working together on a joint offer to take control of Hughes, with News Corp. proceeding on its own and Liberty considering an independent bid. The move threatens to cloud the process of finding a new owner for the GM unit.

(Article on Page A3)

The SEC signaled it may file civil charges against Morgan Stanley, alleging it doled out IPO shares based partly on investors' commitments to buy more stock.

(Article on Page C1)

Ahold's problems deepened as U.S. authorities opened inquiries into accounting at the Dutch company's U.S. Foodservice unit.

Fleming said the SEC upgraded to a formal investigation an inquiry into the food wholesaler's trade practices with suppliers.

(Articles on Page A2)

Consumer confidence fell to its lowest level since 1993, hurt by energy costs, the terrorism threat and a stagnant job market.

(Article on Page A3)

The industrials rebounded on

BUSH IS PREPARING to present Congress a huge bill for Iraq costs.

The total could run to \$95 billion depending on the length of the possible war and occupation. As horse-trading began at the U.N. to win support for a war resolution, the president again made clear he intends to act with or without the world body's imprimatur. Arms inspectors said Baghdad provided new data, including a report of a possible biological bomb. Gen. Franks assumed command of the war-operations center in Qatar. Allied warplanes are aggressively taking out missile sites that could threaten the allied troop buildup. (Column 4 and Pages A4 and A6)

Turkey's parliament debated legislation to let the U.S. deploy 62,000 to open a northern front. Kurdish soldiers lined roads in a show of force as U.S. officials traveled into Iraq's north for an opposition conference.

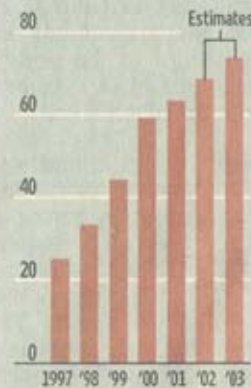
Powell said North Korea hasn't restarted a reactor and plutonium-processing facility at Yongbyon, hinting such forbearance might constitute an overture. But saber rattling continued a day after a missile test timed for the inauguration in Seoul. Pyongyang accused U.S. spy planes of violating its airspace and told its army to prepare for U.S. attack. (Page A14)

The FBI came under withering bipartisan criticism in a Senate Judiciary report in which Sen. Specter

Web Master

As the Web spreads...

Total Internet users, by household, in millions



Sources: Forrester Research; Nielsen NetRatings

Google's U.S. presence expands

Top search engines, in millions of unique visitors¹



¹Including visitors from home and work, in January 2003

Top shopping-referral sites, in millions of referrals²



²Number of people the sites send to major online stores, including only visitors from home, for Q4 2002

Bush to Seek up to \$95 Billion To Cover Costs of War on Iraq

By GREG JAFFE
And JOHN D. MCKINNON

WASHINGTON—The Bush administration is preparing supplemental spending requests totaling as much as \$95 billion for a war with Iraq, its aftermath and new expenses to fight terrorism, officials said.

The total could be as low as \$60 billion because Pentagon budget planners don't know how long a military conflict will last, whether U.S. allies will contribute more than token sums to the effort and what damage Saddam Hussein might do

to his own country to retaliate against conquering forces.

Budget planners also are awaiting the outcome of an intense internal debate over whether to include \$13 billion in the requests to Congress that the Pentagon says it needs to fund the broader war on terrorism, as well as for stepped up homeland security. The White House Office of Management and Budget argues that the money might not be necessary. President Bush, Defense Secretary Donald Rumsfeld and budget director Mitchell Daniels Jr. met yesterday to discuss the matter but didn't reach a final agreement. Mr. Rumsfeld plans to continue pressing his

Cat and Mouse

As Google Becomes Web's Gatekeeper, Sites Fight to Get In

Search Engine Punishes Firms That Try to Game System; Outlawing the 'Link Farms'

Exoticleatherwear Gets Cut Off

By MICHAEL TOTT
And MYLENE MANGALINDAN

Joy Holman sells provocative leather clothing on the Web. She wants what nearly everyone doing business online wants: more exposure on Google.

So from the time she launched exoticleatherwear.com last May, she tried all sorts of tricks to get her site to show up among the first listings when a user of Google Inc.'s popular search engine typed in "women's leatherwear" or "leather apparel." She buried hidden words in her Web pages intended to fool Google's computers. She signed up with a service that promised to have hundreds of sites link to her online store—thereby boosting a crucial measure in Google's system of ranking sites.

The techniques worked—



Web Sites Fight for Prime Real Estate on Google

Continued From First Page

advertising that tried to capitalize on Google's formula for ranking sites. In effect, SearchKing was offering its clients a chance to boost their own Google rankings by buying ads on more-popular sites. SearchKing filed suit against the search company in federal court in Oklahoma, claiming that Google "purposefully devalued" SearchKing and its customers, damaging its reputation and hurting its advertising sales.

Google won't comment on the case. In court filings, the company said SearchKing "engaged in behavior that would lower the quality of Google search results" and alter the company's ranking system.

Google, a closely held company founded by Stanford University graduate students Sergey Brin and Larry Page, says Web companies that want to rank high should concentrate on improving their Web pages rather than gaming its system. "When people try to take scoring into their own hands, that turns into a worse experience for users," says Matt Cutts, a Google software engineer.

Coding Trickery

Efforts to outfox the search engines have been around since search engines first became popular in the early 1990s. Early tricks included stuffing thousands of widely used search terms in hidden coding, called "metatags." The coding fools a search engine into identifying a site with popular words and phrases that may not actually appear on the site.

Another gimmick was hiding words or terms against a same-color background. The hidden coding deceived search engines that relied heavily on the number of times a word or phrase appeared in ranking a site. But Google's system, based on links, wasn't fooled.

Mr. Brin, 29, one of Google's two founders and now its president of technology, boasted to a San Francisco search-engine conference in 2000 that Google wasn't worried about having its results clogged with irrelevant results because its search methods couldn't be manipulated.

That didn't stop search optimizers from finding other ways to outfox the system. Attempts to manipulate Google's results even became a sport, called Google-hacking. Hacksters would try to

creating Web sites that were nothing more than collections of links to the clients' site, called "link farms." Since Google ranks a site largely by how many links or "votes" it gets, the link farms could boost a site's popularity.

In a similar technique, called a link exchange, a group of unrelated sites would agree to all link to each other, thereby fooling Google into thinking the sites have a multitude of votes. Many sites also found they could buy links to themselves to boost their rankings.

Ms. Holman, the leatherwear retailer, discovered the consequences of trying to fool Google. The 42-year-old hospital laboratory technician, who learned computer skills by troubleshooting her hospital's

'The big search engines determine the laws of how commerce runs,' says Mr. Massa.

equipment, operates her online apparel store as a side business that she hopes can someday replace her day job.

When she launched her Exotic Leather Wear store from her home in Mesa, Ariz., she quickly learned the importance of appearing near the top of search-engine results, especially on Google. She boned up on search techniques, visiting online discussion groups dedicated to search engines and reading what material she could find on the Web.

At first, Ms. Holman limited herself to modest changes, such as loading her page with hidden metatag coding that would help steer a search toward her site when a user entered words such as "haltertops" or "leather miniskirts." Since Google doesn't give much weight to metatags in determining its rankings, the efforts had little effect on her search results.

She then received an e-mail advertisement from AutomatedLinks.com, a Wirral, England, company that promised to send traffic "through the roof" by linking more than 2,000 Web sites to hers. Aside from attracting customers, the links were designed to improve her site's search engine rankings by taking

In theory, when Google encounters the AutomatedLinks code, it treats it as a legitimate referral to the other sites and counts them in totting up the sites' popularity.

Shortly after Ms. Holman signed up with AutomatedLinks in July, she read on an online discussion group that Google objected to such link arrangements. She says she immediately stripped the code from her Web pages. For a while her site gradually worked its way up in Google search results, and business steadily improved because links to her site still remained on the sites of other AutomatedLinks customers. Then, sometime in November, her site was suddenly no longer appearing among the top results. Her orders plunged as much as 80%.

Ms. Holman, who e-mailed Google and AutomatedLinks, says she has been unable to get answers. But in the last few months, other AutomatedLinks customers say they have seen their sites apparently penalized by Google. Graham McLeay, who runs a small chauffeur service north of London, saw revenue cut in half during the two months he believes his site was penalized by Google.

The high-stakes fight between Google and the optimizers can leave some Web-site owners confused. "I don't know how people are supposed to judge what is right and wrong," says Mr. McLeay.

AutomatedLinks didn't respond to requests for comment. Google declined to comment on the case. But Mr. Cutts, the Google engineer, warns that the rules are clear and that it's better to follow them rather than try to get a problem fixed after a site has been penalized. "We want to return the most relevant pages we can," Mr. Cutts says. "The best way for a site owner to do that is follow our guidelines."

Crackdown

Google has been stepping up its enforcement since 2001. It warned Webmasters that using trickery could get their sites kicked out of the Google index and it provided a list of forbidden activities, including hiding text and "link schemes," such as the link farms. Google also warned against "cloaking"—showing a search engine a page that's designed to score well while giving visitors a different, more attractive page—or creating multiple Web addresses that take visitors to a single site.

To stay one step ahead of the Web

homa City-based SearchKing, an online directory for hundreds of small, specialty Web sites. SearchKing also sells advertising links designed both to deliver traffic to an advertiser and boost its rankings in Google and other search results.

Bob Massa, SearchKing's chief executive, last August launched the PR Ad Network as a way to capitalize on Google's page-ranking system, known as PageRank. PageRank rates Web sites on a scale of one to 10 based on their popularity, and the rankings can be viewed by Web users if they install special Google software. PR Ad Network sells ads that are priced according to a site's PageRank, with higher-ranked sites commanding higher prices. When a site buys an advertising link on a highly ranked site, the ad buyer could see its ratings improve because of the greater weight Google gives to that link.

Shortly after publicizing the ad network, Mr. Massa discovered that his site suddenly dropped in Google's rankings. What's more, sites that participated in the separate SearchKing directory also had their Google rankings lowered. He filed a lawsuit in Oklahoma City federal court, claiming Google was punishing him for trying to profit from the company's page-ranking system.

A Google spokesman won't comment on the case. In its court filings, Google said it demoted pages on the SearchKing site because of SearchKing's attempts to manipulate search results. The company has asked for the suit to be dismissed, arguing that the PageRank represents its opinion of the value of a Web site and as such is protected by the First Amendment.

"The big search engines determine the laws of how commerce runs," says Mr. Massa, who is persisting with the lawsuit even though the sites have had their page rankings partly restored. "Someone needs to demand accountability."

Google is taking steps that many say could satisfy businesses trying to boost their rankings. Google has long sold sponsored links that show up on the top of many search-results pages, separate from the main listings. Last year, the company expanded its paid-listings program, so that there are now more slots where sites can pay for a prominent place in the results. Many sites now are turning to advertising instead of tactics to optimize their rankings.

Home Depot Amid First

By CHAD TERHUN

ATLANTA—Home Depot's fiscal fourth-quarter earnings fell 3.4% on disappointing sales.

Speaking to investors and analysts, the company's chief executive, Bob Nardelli, said Home Depot is prepared to meet dissatisfied customers and a competitive challenge from competitors with remodeled stores, inventory and improved customer service.

The nation's largest home improvement retailer said net income for the quarter ended Feb. 2 decreased to 30 cents a share, from \$71.30 cents a share, a year earlier. Revenue rose 2% to \$13.21 billion from \$13.45 billion a year earlier. Home Depot's first quarterly sales decline in its 24-year history. Home Depot's latest quarter was a week earlier than a year earlier. Using comparable periods, the company said quarterly sales increased 5% and net income rose 10%.

Same-store sales, or sales from stores open at least a year, declined 1% in the quarter. Home Depot said stock prices last month offset a disastrous quarter and helped the retailer avoid a 10% decline in its stock price. Analysts estimate that same-store sales rose as much as 10%. In 4 p.m. trading on the New York Stock Exchange composite trading index, Home Depot shares rose 66 cents to \$48.75.

Fiat Patria Is Set to Be

By ALESSANDRA GAIARDI

ROME—Umberto Agnelli, the former Fiat SpA chairman, is stepping into the driver's seat as the conglomerate works on an 11th-hour effort to salvage the car unit.

Mr. Agnelli, the 68-year-old Fiat patriarch Gianni Agnelli's son, was widely expected to take over from current chairman Luca Cordero di Montezemolo last month, but he was replaced by Agnelli last month. But he was widely expected to take over from current chairman Luca Cordero di Montezemolo last month. But he was widely expected to take over from current chairman Luca Cordero di Montezemolo last month.



Back To Tiny Web

The Google Matrix

$$\mathbf{P} = \alpha \mathbf{H} + (\alpha \mathbf{a} + (1 - \alpha) \mathbf{e}) \mathbf{v}^T \quad (\text{with } \alpha = .9 \text{ and } \mathbf{v} = \mathbf{e})$$

$$= \begin{bmatrix} 1/60 & 7/15 & 7/15 & 1/60 & 1/60 & 1/60 \\ 1/6 & 1/6 & 1/6 & 1/6 & 1/6 & 1/6 \\ 19/60 & 19/60 & 1/60 & 1/60 & 19/60 & 1/60 \\ 1/60 & 1/60 & 1/60 & 1/60 & 7/15 & 7/15 \\ 1/60 & 1/60 & 1/60 & 7/15 & 1/60 & 7/15 \\ 1/60 & 1/60 & 1/60 & 11/12 & 1/60 & 1/60 \end{bmatrix}$$

The PageRank Vector

$$\pi_{j+1}^T = \pi_j^T \mathbf{P} \rightarrow \pi^T$$

$$\pi^T = \begin{pmatrix} & \mathbf{1} & \mathbf{2} & \mathbf{3} & \mathbf{4} & \mathbf{5} & \mathbf{6} \\ \mathbf{.03721} & \mathbf{.05396} & \mathbf{.04151} & \mathbf{.3751} & \mathbf{.206} & \mathbf{.2862} \end{pmatrix}$$



Computing π^T

A Big Problem

$$\text{Solve } \pi^T = \pi^T \mathbf{P}$$

(eigenvector problem)



Computing π^T

A Big Problem

$$\text{Solve } \pi^T = \pi^T \mathbf{P}$$

(eigenvector problem)

$$\pi^T (\mathbf{I} - \mathbf{P}) = \mathbf{0}$$

(too big for direct solves)

Google's PageRank is an eigenvector of a matrix of order 2.7 billion.

One of the reasons why Google is such an effective search engine is the PageRank™ algorithm, developed by Google's founders, Larry Page and Sergey Brin, when they were graduate students at Stanford University. PageRank is determined entirely by the link structure of the Web. It is recomputed about once a month and does not involve any of the actual content of Web pages or of any individual query. Then, for any particular query, Google finds the pages on the Web that match that query and lists those pages in the order of their PageRank.

Imagine surfing the Web, going from page to page by randomly choosing an outgoing link from one page to get to the next. This can lead to dead ends at pages with no outgoing links, or cycles around cliques of interconnected pages. So, a certain fraction of the time, simply choose a random page from anywhere on the Web. This theoretical random walk of the Web is a *Markov chain* or *Markov process*. The limiting probability that a dedicated random surfer visits any particular page is its PageRank. A page has high rank if it has links to and from other pages with high rank.

Let W be the set of Web pages that can be reached by following a chain of hyperlinks starting from a page at Google and let n be the number of pages in W . The set W actually varies with time, but in May 2002, n was about 2.7 billion. Let G be the n -by- n connectivity matrix of W that is, $G_{ij} = 1$ if there is a hyperlink from page i to page j and 0

BY CLEVE MOLER

It tells us that the largest eigenvalue of A is equal to one and that the corresponding eigenvector, which satisfies the equation

$$x = Ax,$$

exists and is unique to within a scaling factor. When this scaling factor is chosen so that

$$\sum_i x_i = 1$$

then x is the state vector of the Markov chain. The elements of x are Google's PageRank.

If the matrix were small enough to fit in MATLAB, one way to compute the eigenvector x would be to start with a good approximate solution, such as the PageRanks from the previous month, and simply repeat the assignment statement

$$x = Ax$$

until successive vectors agree to within specified tolerance. This is known as the power method and is about the only possible approach for very large n . I'm not sure how Google actually computes PageRank, but one step of the power method would require one pass over a database of Web pages, updating weighted reference counts generated by the hyperlinks between pages.



Computing π^T

A Big Problem

Solve $\pi^T = \pi^T \mathbf{P}$ (eigenvector problem)

$\pi^T (\mathbf{I} - \mathbf{P}) = \mathbf{0}$ (too big for direct solves)

Start with $\pi_0^T = \mathbf{e}/n$ and iterate $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ (power method)



Computing π^T

A Big Problem

Solve $\pi^T = \pi^T \mathbf{P}$ (eigenvector problem)

$\pi^T (\mathbf{I} - \mathbf{P}) = \mathbf{0}$ (too big for direct solves)

Start with $\pi_0^T = \mathbf{e}/n$ and iterate $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ (power method)

Convergence Time

Measured in days



Computing π^T

A Big Problem

Solve $\pi^T = \pi^T \mathbf{P}$ (eigenvector problem)

$\pi^T (\mathbf{I} - \mathbf{P}) = \mathbf{0}$ (too big for direct solves)

Start with $\pi_0^T = \mathbf{e}/n$ and iterate $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ (power method)

Convergence Time

Measured in days

A Bigger Problem — Updating

Pages & links are added, deleted, changed continuously



Computing π^T

A Big Problem

Solve $\pi^T = \pi^T \mathbf{P}$ (eigenvector problem)

$\pi^T (\mathbf{I} - \mathbf{P}) = \mathbf{0}$ (too big for direct solves)

Start with $\pi_0^T = \mathbf{e}/n$ and iterate $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ (power method)

Convergence Time

Measured in days

A Bigger Problem — Updating

Pages & links are added, deleted, changed continuously

Google says just start from scratch every 3 to 4 weeks



Computing π^T

A Big Problem

Solve $\pi^T = \pi^T \mathbf{P}$ (eigenvector problem)

$\pi^T (\mathbf{I} - \mathbf{P}) = \mathbf{0}$ (too big for direct solves)

Start with $\pi_0^T = \mathbf{e}/n$ and iterate $\pi_{j+1}^T = \pi_j^T \mathbf{P}$ (power method)

Convergence Time

Measured in days

A Bigger Problem — Updating

Pages & links are added, deleted, changed continuously

Google says just start from scratch every 3 to 4 weeks

Prior results don't help to restart



Conclusions

✦ Elegant Blend of NA, LA, Graph Theory, MC, & CS ✦



Conclusions

- ✦ Elegant Blend of NA, LA, Graph Theory, MC, & CS ✦
- ✦ Google Now Uses Many Other “Metrics” to augment PR ✦



Conclusions

- ✦ Elegant Blend of NA, LA, Graph Theory, MC, & CS ✦
- ✦ Google Now Uses Many Other “Metrics” to augment PR ✦
- ✦ Search Is Opening New Areas Ripe For Innovative Ideas ✦
- ✦ Exciting Work That Is Changing The World ✦

✦ Thanks For Your Attention ✦