

Informationsåtervinning på webben

Sökmotorernas framtid



KUNGL. INGENJÖRSVETENSKAPSAKADEMIEN
Royal Swedish Academy of Engineering Sciences

Seminarium 2004-09-02

1

Informationsåtervinning på webben

Sökmotorernas framtid

- Ge inspiration till
 - forskning
 - att skapa nya affärsmöjligheter
 - smart användning av sökverktyg i den egna organisationen
- Belysa sökmotorer ur ett tekniskt, juridiskt och affärsmässigt perspektiv.



KUNGL. INGENJÖRSVETENSKAPSAKADEMIEN
Royal Swedish Academy of Engineering Sciences

2

Informationsåtervinning på webben

Sökmotorernas framtid - **Program**

- Sökmotorer och tekniken bakom, *Bo Kågström*
- Sökmotorer i framtiden, *Jussi Karlgren*
- Integritet, säkerhet och manipulation, *Nicklas Lundblad*
- Affärsperspektiv och applikationer
- Web4Health, *Jacob Palme*
- SiteSeeker och Euroling, *Hercules Dalianis*
- Askology och QuickAsk, *Erik Sneiders*
- Öppen diskussion, *Sture Hägglund* moderator

3

Sökmotorer och tekniken bakom

Ranking av webbsidor med länkanalys -
Googles PageRank och liknande metoder

Bo Kågström

Dept. of Computing Science and HPC2N

Umeå University

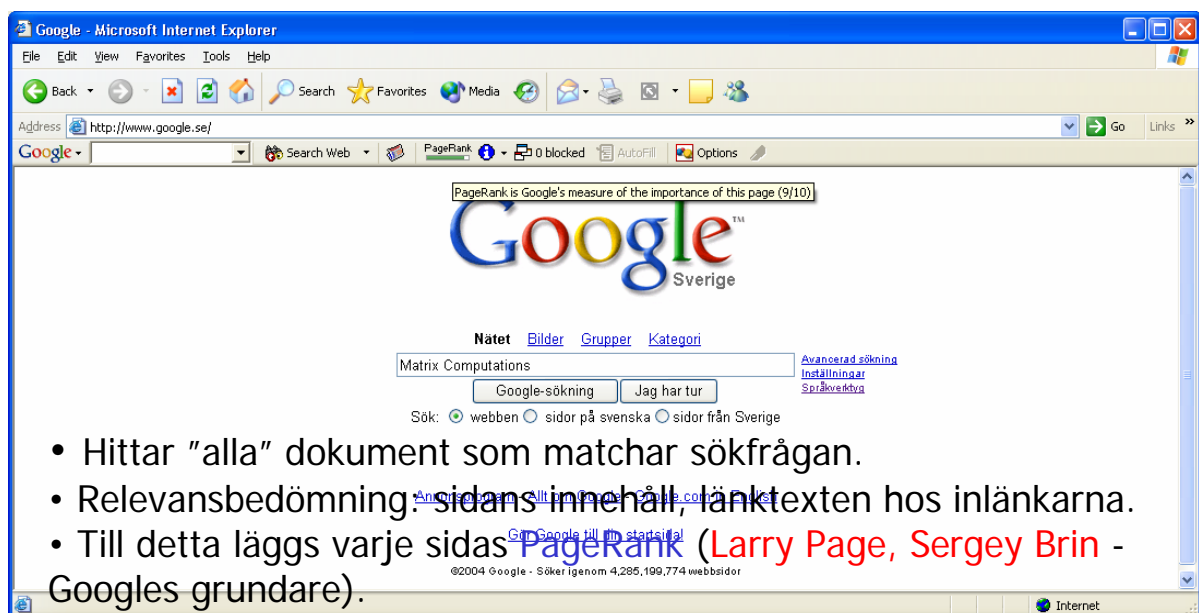
bokg@cs.umu.se

Lite bakgrund

- Webb-sökning:
 - Angelägen om "breda förfrågningar" (*broad-topic queries*), t.ex. "web-browsers".
 - *Överflödsproblematik*: # sidor som hittas och bedöms som relevanta ("träffar") är alldeles för stort!
 - Behövs en mekanism för att rangordna dessa sidor.
- Hypotes: Om sidan *j* har en länk till sidan *i*, så ger den auktoritet till *i*.
 - Hur används *länkinformationen* för att rangordna "träffar"?
 - Önskar *relevanta* och *auktoritativa* sidor.

5

Hur går en Googlesökning till?

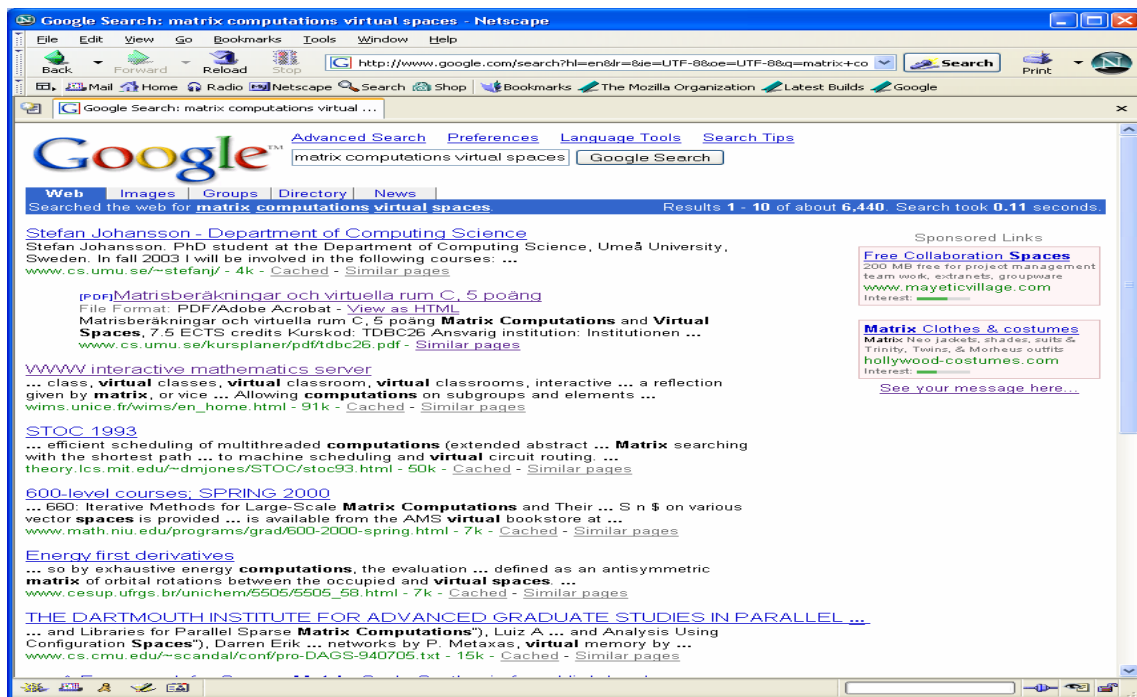


- Hittar "alla" dokument som matchar sökfrågan.
- Relevansbedömning: sidans innehåll, länktexten hos inlänkarna.
- Till detta läggs varje sidas *PageRank* (Larry Page, Sergey Brin - Googles grundare).

• Relevanta dokument rangordnas och listas utifrån sina PageRank-värden.

6

Exempel på sökresultat



Innehåll

- Sökning på webben – lite bakgrund
- Googles PageRank-algoritm
 - Definition av PageRank
 - PageRank är en "dominerande" egenvektor
 - Rank-sinks och -sources – modifierad definition
 - Beräkning av PageRank – världens största matrisberäkning?
 - Varför fungerar det? Konvergensgenskaper
- HITS-algoritmen: Hypertext Internet Topic Search
 - Auktoriteter och hubbar är "dominerande" singulära vektorer

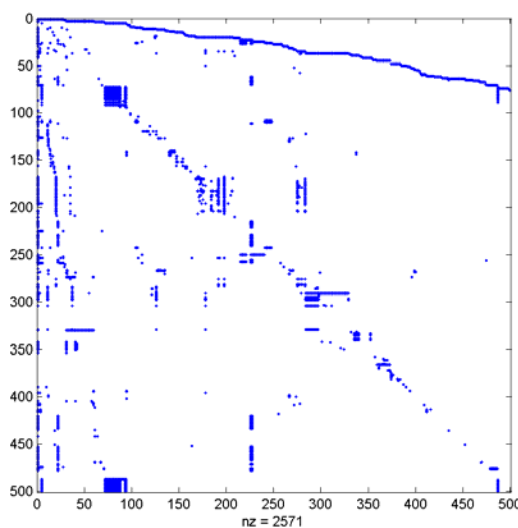
Webb-grafen & Webb-matrisen

- WWW kan representeras som en graf med "sites" (t.ex. hemsidor) som noder and länkar som kanter.
- Webbmatrisen A (adjacency or connectivity) representerar länkstrukturen mellan sidor:
 $A(i,j) = 1$ om sida i pekar på sida j
 $A(i,j) = 0$ annars
- A är en gles (sparse) av storlek $n \times n$,
 $n > 4$ miljarder ($4 \cdot 10^9$)!

9

Webb-matris: Harvard 500x500

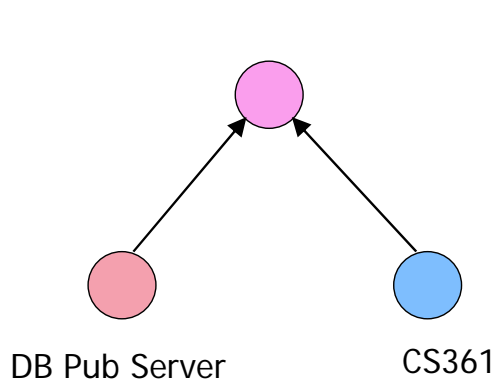
$G(i,j) = 1$
om url{i}
länkar till
url{j}.



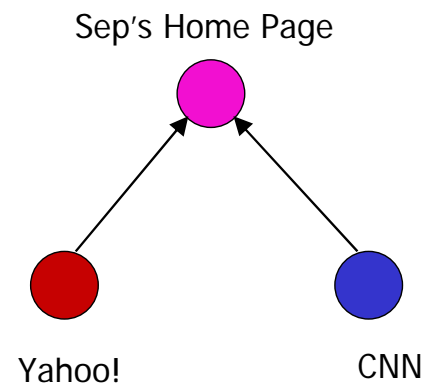
Skapad med `[U,G] = surfer('http://www.harvard.edu',500)`
U cell-array med besökta URL:er

10

Inlänkar räknas!



Länkad av två
"mindre viktiga" sidor



Länkad av två
"viktiga" sidor

11

Definition av PageRank

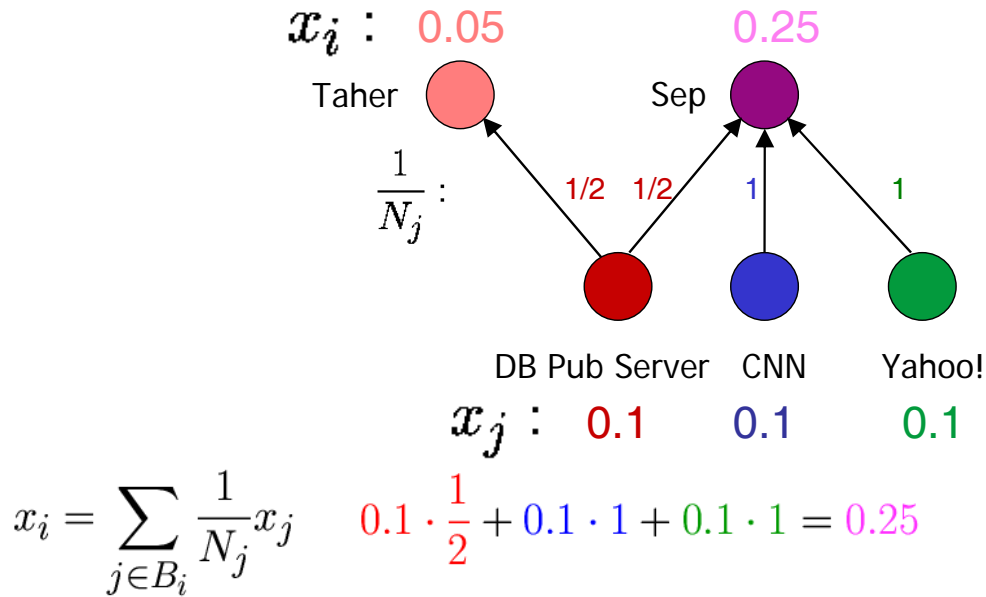
- En sidas "betydelse" (vikt) ges av "vikten" hos de sidor som pekar på den.

$$x_i = \sum_{j \in B_i} \frac{1}{N_j} x_j$$

importance of page i pages j that link to page i importance of page j number of out-links from page j

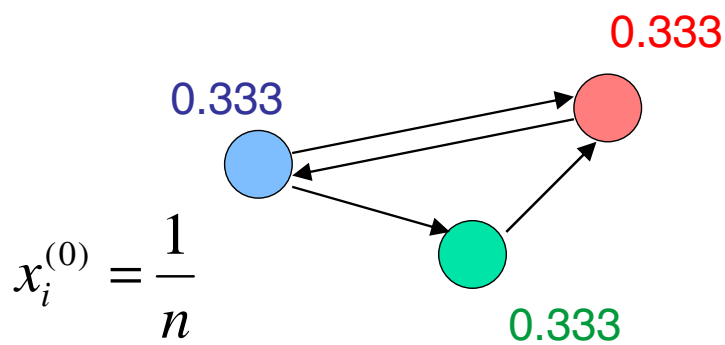
12

Definition av PageRank - exempel



13

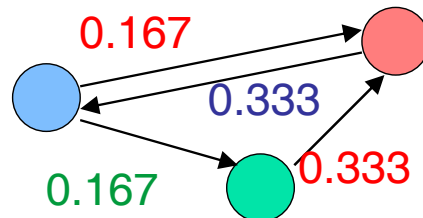
PageRank-diagram



Initialisera alla noder till samma rang (vikt)

14

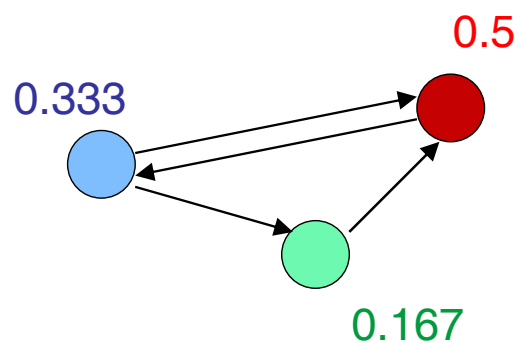
PageRank-diagram



Propagera rangerna över länkarna
(multiplicera med länkvikter)

15

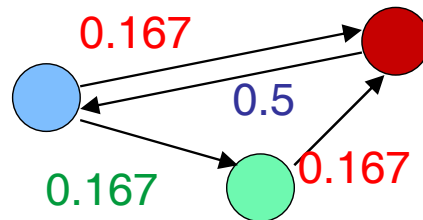
PageRank-diagram



$$x_i^{(1)} = \sum_{j \in B_i} \frac{1}{N_j} x_j^{(0)}$$

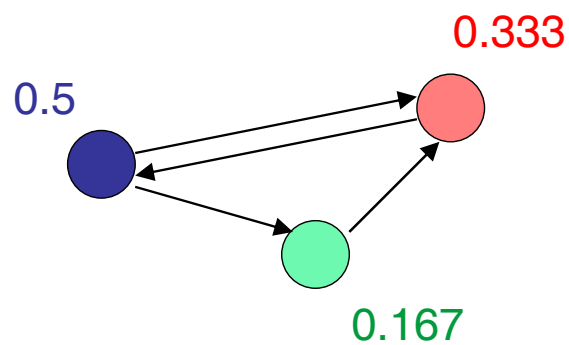
16

PageRank-diagram



17

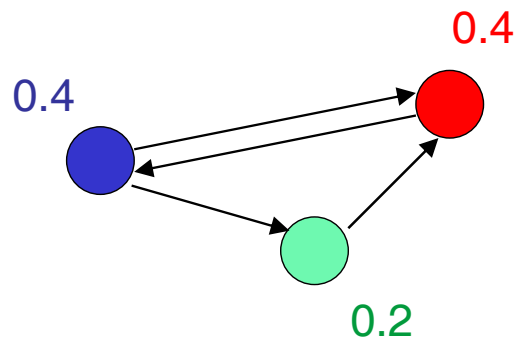
PageRank-diagram



$$x_i^{(2)} = \sum_{j \in B_i} \frac{1}{N_j} x_j^{(1)}$$

18

PageRank-diagram



Efter ett tag ...

$$x_i = \sum_{j \in B_i} \frac{1}{N_j} x_j$$

19

Förenklad beräkning av PageRank

- Initialisering: $x_i^{(0)} = \frac{1}{n}$

- Upprepa tills konvergens:

$$x_i^{(k+1)} = \sum_{j \in B_i} \frac{1}{N_j} x_j^{(k)}$$

importance of page i

importance of page j

pages j that link to page i

number of outlinks from page j

20

Matrisnotation – rätt abstraktion

$$x_i = \sum_{j \in B_i} \frac{1}{N_j} x_j$$

The diagram shows the equation $x_i = \sum_{j \in B_i} \frac{1}{N_j} x_j$ being represented as a matrix multiplication $\mathbf{x} = \mathbf{P}^T \mathbf{x}$. The vector \mathbf{x} is a column vector with values $.1, .3, .2, .3, .1, .1$. The matrix \mathbf{P}^T is a square matrix with a red row containing values $0.2, 0.3, 0, 0.1, .4, 0.1$. The vector on the right is a column vector with values $.1, .3, .2, .3, .1, .1$. Red arrows point from the equation to the matrix and vector, and a blue arrow points from the vector \mathbf{x} to the matrix.

21

Sökt: egenvektor svarande mot största egenvärdet = 1

Hitta \mathbf{x} som uppfyller:

$$\mathbf{x} = \mathbf{P}^T \mathbf{x}$$

The diagram shows the equation $\mathbf{x} = \mathbf{P}^T \mathbf{x}$ being represented as a matrix multiplication. The vector \mathbf{x} is a column vector with values $.1, .3, .2, .3, .1, .1$. The matrix \mathbf{P}^T is a square matrix with a red row containing values $0.2, 0.3, 0, 0.1, .4, 0.1$. The vector on the right is a column vector with values $.1, .3, .2, .3, .1, .1$. Red arrows point from the equation to the matrix and vector.

22

Tillämpa potensmetoden (Power Method)

- Initialisering:

$$\mathbf{x}^{(0)} = \begin{bmatrix} 1 & \dots & 1 \\ n & & n \end{bmatrix}^T$$

- Upprepa tills konvergens:

$$\mathbf{x}^{(k+1)} = \mathbf{P}^T \mathbf{x}^{(k)}$$

Stoppkriterium: $\text{norm}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}) < \text{tolerance}$

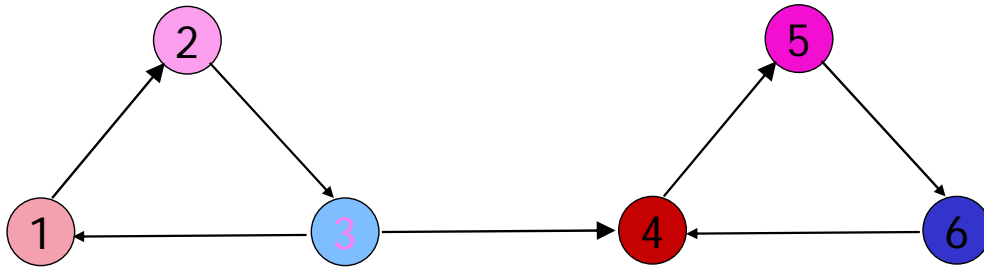
23

Random Walk på webben

- **Definitionen av PageRank** kan ses som slumpvandringar (random walks) på grafer.
- Surfa från sida till sida genom att **slumpmässigt välja en utlänk från en sida** för att komma till nästa.
- Kan leda till "dead ends" hos **sidor som saknar utlänkar (dangling pages)**, eller **cykler** kring klickar av **sammanshängande sidor (loops)**.

24

Loop som en rang-ansamlare



Rank Sink: Loopen 4 -> 5 -> 6 ackumulerar rang men kommer aldrig att distribuera någon rang (inga utgående länkar).

25

Rank Sink - problematik

- Alla egenvektorer till webbmatriosen P^T i sista exemplet har nollor i de tre första komponenterna.



$$P = \left[\begin{array}{ccc|ccc} 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{array} \right]$$

PageRank för sidorna 1, 2 and 3 are 0!

Botemedel: introducera **artificiella länkar** (rank sources).

26

PageRank med Rank Sources

$$\text{PageRank}(i) \equiv x_i = \sum_{j \in B_i} \frac{1}{N_j} x_j + s_i \quad s = \begin{bmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{bmatrix}$$

Alla PageRank är **skilda från noll från början!**

- En "random surfer" kan följa vilken utlänk som helst från en sida med samma sannolikhet (förenklade definitionen).
- **Då och då**, blir hon "less" and **hoppas till en slumpvis sida på Webben** (ny definition med "rank sources").

27

Googles PageRank-matris

- Periodiskt, väljs en slumpvis sida på webben för att överkomma "dangling pages" och loopar.
- $A = c P^T + (1 - c) E^T$
 - C = bråkdel av tiden som en "surfare" (random walk) följer en länk (t.ex., $c = 0.85$)
 - $1 - C$ = bråkdel av tiden som en godtycklig sida väljs
 - E är $n \times n$ med $E(i,j) = 1/n$ ($n = \#$ länkar i Webbmatriken)
- A är "tät" (**dense**), **rang-1**-modifiering av en **gles** matris - de flesta $A(i, j) = (1 - c) / n$.

28

Perron-Frobenius teorem

- $A = c P^T + (1 - c) E^T$ är övergångsmatrix hos en Markovkedja (**transition probability matrix**)

$0 \leq A(i,j) \leq 1$, alla kolumnsummor = 1

- PF: A :s **största egenvärde** = 1 svarande till en entydig egenvektor x med $x_i \geq 0$

$$Ax = x$$

$$\sum_{i=1}^n x_i = 1 \quad \Longrightarrow$$

x är Markovkedjans tillståndsvektor (**state vector of the Markov chain**)

29

Potensmetoden tillämpad på A – världens största matrisberäkning?

- Initialisering:

$$\mathbf{x}^{(0)} = \left[\frac{1}{n} \quad \dots \quad \frac{1}{n} \right]^T$$

- Upprepa tills konvergens:

$$\mathbf{x}^{(k+1)} = \mathbf{A}\mathbf{x}^{(k)}$$

Beräkning av $y = Ax$ där $A = c P^T + (1 - c)/n e e^T$
- A beräknas ej explicit
- Utnyttja A :s struktur

Elementen i x är **Googles PageRank!**

30

Varför fungerar det?

- Antag att $n \times n$ matrisen \mathbf{A} har n egenvektorer \mathbf{u}_i .

$$\mathbf{A}\mathbf{u}_i = \lambda_i \mathbf{u}_i$$

- Då kan en godtycklig n -dimensionell vektor skrivas som en linjärkombination av egenvektorerna till \mathbf{A} .

$$\mathbf{x}^{(0)} = \mathbf{u}_1 + \alpha_2 \mathbf{u}_2 + \dots + \alpha_n \mathbf{u}_n$$

$$\lambda_1 = 1; \quad \lambda_1 > |\lambda_2| \geq \dots$$

$$\begin{array}{ccccc} \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 & \mathbf{u}_4 & \mathbf{u}_5 \\ 1 & \alpha_2 & \alpha_3 & \alpha_4 & \alpha_5 \end{array}$$

31

Konvergenssegenskaper

$$\mathbf{x}^{(k)} = \mathbf{u}_1 + \alpha_2 \lambda_2^k \mathbf{u}_2 + \dots + \alpha_n \lambda_n^k \mathbf{u}_n$$

$$\lambda_1 = 1; \quad \lambda_1 > |\lambda_2| \geq \dots$$

$$\begin{array}{ccccc} \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 & \mathbf{u}_4 & \mathbf{u}_5 \\ 1 & \alpha_2 \lambda_2^k & \alpha_3 \lambda_3^k & \alpha_4 \lambda_4^k & \alpha_5 \lambda_5^k \end{array}$$

- Konvergensthastigheten** bestäms av $1/|\lambda_2|$ dvs beloppet av kvoten mellan det största och näst största egenvärdet.
- Ju mindre λ_2 , desto snabbare konvergerar potensmetoden ($\text{abs}(\lambda_2) \leq c, 0 \leq c \leq 1$)

32

Är potensmetoden (PM) bästa valet?

- Traditionellt:
 - A , $n \ll 4$ miljarder, ofta tät matris.
 - Risk för att λ_2 är nära $\lambda_1 \rightarrow$ potensmetoden långsam!.
 - För detta problem:
 - A , enormt stor, kolumnstokastisk, ofta tät, rank-1 modifiering av en gles matris, där λ_2 är liten¹
 - \rightarrow Potensmetoden fungerar mycket bra!!
- ¹ – Se Haveliwala T.H. and Kamvar S.D., “The Second Eigenvalue of the Google Matrix” dbpubs.stanford.edu/pub/2003-20.
- Det pågår forskning med att “snabba upp” PM.
 - Andra metoder kan fungera lika bra eller bättre för beräkning av PageRank för mer begränsade domäner.

33

PageRank - sammanfattning

- Sökning av webbsidor är huvudtillämpningen – används i fulltext-sökmotorn Google.
- PageRank är en global ranking av alla webbsidor, oberoende av dess innehåll, enbart baserad på dess plats i Webbgraf-strukturen (länk-baserad) - beräknas om ca 1 gång/månad, tar 1-2 veckor!?!)
- Ranking används för att rangordna sidorna så att mer centrala webbsidor ges preferens.
- Bakåtlänkar från ”viktiga” sidor är mer signifikanta än bakåtlänkar från ”oviktiga” sidor (rekursiv definition av PageRank).

34