

Matematika delovanja spletnih iskalnikov

Rok Požar

Univerza na Primorskem

Fakulteta za matematiko, naravoslovje in informacijske tehnologije

Famnitovi izleti v matematično vesolje

Januar 2015





Informacije zbrane na papirusovih zvitkih

povzetki kot značke na dokumentih,
povzetki kot ustna sporočila: grške igre (5. st. pr. n. št.)

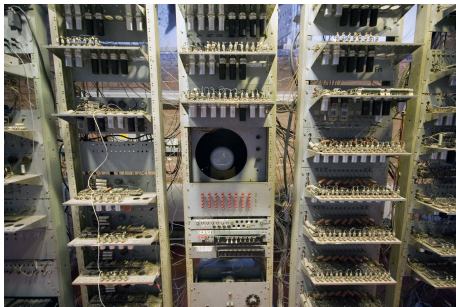
Prvi knjižnični katalog

knjižnica v Aleksandriji (2. st. pr. n. št.),
skoraj 500 000 zvitkov razvrščenih po vsebini
Pinakes: tablice z naslovom in bibliografskimi podatki o avtorju



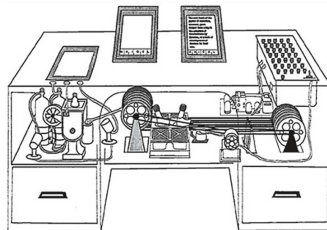
Tisk s premičnimi črkami (1450)

hierarhičen sistem klasifikacije: skupine s podobno tematiko,
nasveti knjižničarja,
urejanje po naslovu in avtorju



Digitalni računalnik (1940-1950)

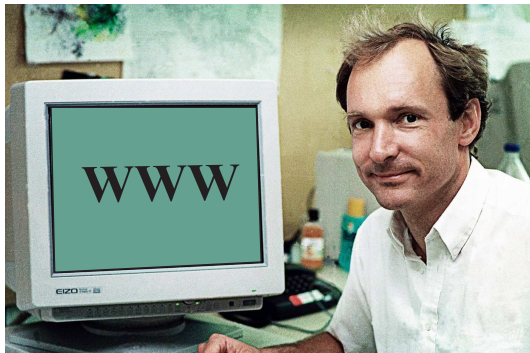
prvi računalniški iskalni sistemi,
avtomatično povpraševanje po knjigah in člankih (Cornel SMART sistem)



Futuristični sestav *Memex*

esej: „As we may think“ (1945),
pristop kazal ni ustrezen,

„človeško razmišljanje ne deluje tako, s tekoče misli preskoči na naslednjo“

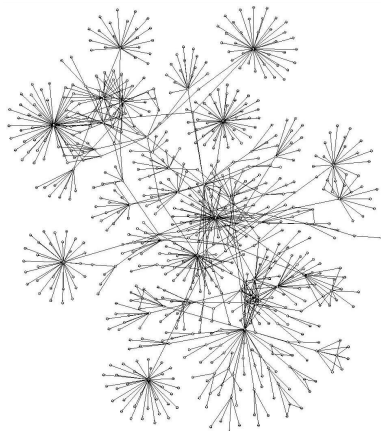


Ustanovitelj svetovnega spleta (World Wide Web, 1989)

povezovanje dokumentov na podlagi hiperpovezav in hiperteksta,
revolucija shranjevanja in dostopnosti informacij

Spletno iskanje kot iskanje igle v kupu sena

pregled nad spletom je temeljil na kazalnih: razvrščanje po temah (Yahoo)

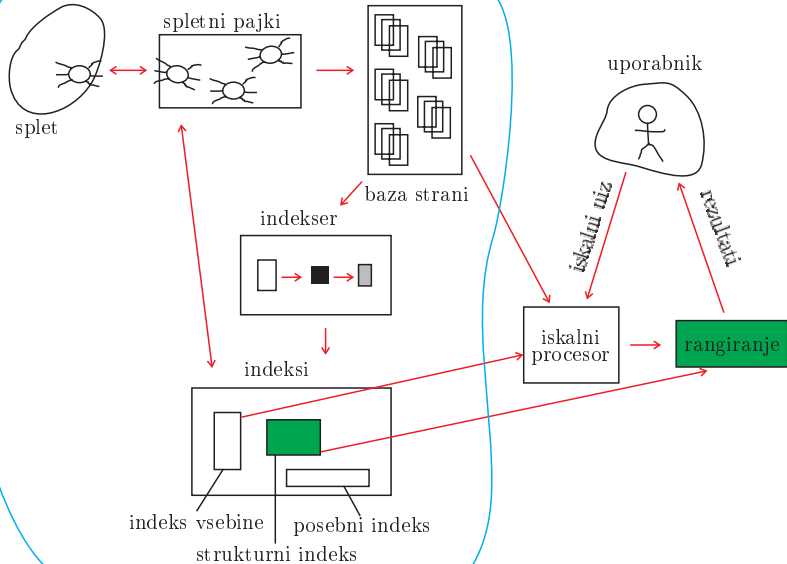


Študij strukturnih lastnosti svetovnega spleta (1998)

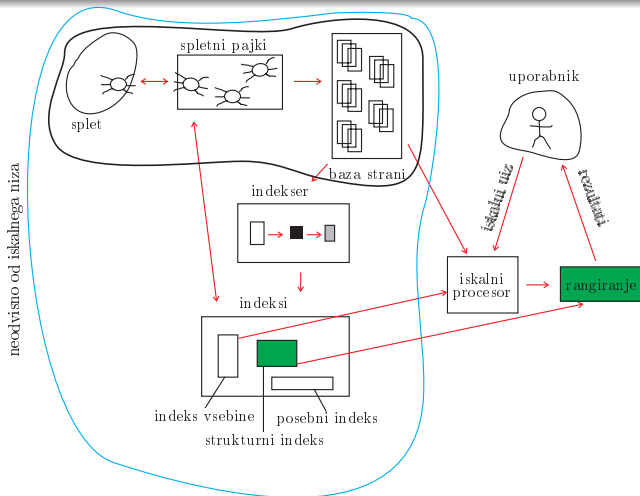
spletno iskanje strmo napreduje (Google)

Kako deluje spletni iskalnik?

neodvisno od iskalnega niza



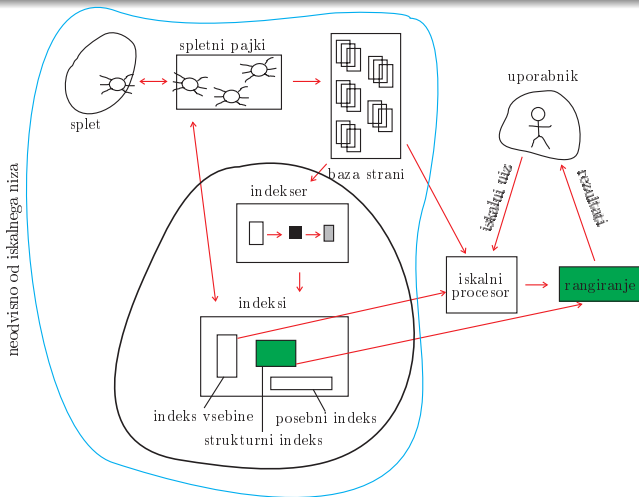
Kopiranje podatkov s spleta



Spletni pajki

iščejo in zapisujejo vsebino spleta v svojo bazo, strani ostanejo v bazi do indeksiranja, popularne strani so v bazi shranjene dlje

Indeksiranje spletnih strani



Indekser

izluči le važne informacije,
naredi skrčen opis strani,
opis shrani v različne indekse

Indeks vsebine

tekstna vsebina, shranjena v obliki obrnjenih seznamov

abak – 3, 117, 3961

⋮

abonent – 3, 5, 19, 101, 117, 367, 3961

otrok – 3, 31, 56, 94, 367, 673, 909, 11114

⋮

zrak – 344521

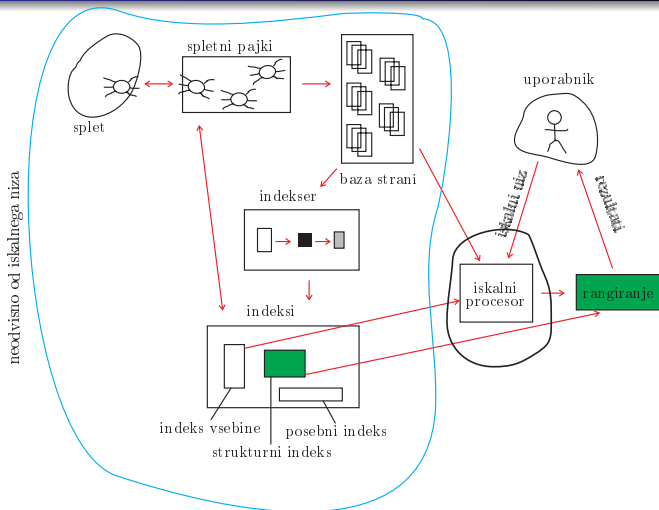
Strukturni indeks

struktura hiperpovezav med spletnimi stranmi, shranjena v skrčeni oblik, pajki ga včasih uporabljajo za iskanje novih strani

Posebni indeks

slike, pdf dokumenti,...

Iskanje relevantnih strani



Iskalni procesor

pretvori iskalni niz v sebi razumljiv jezik,
poveže se z indeksom vsebine in obrnjenim seznamom,
poišče **relevantne strani** (tiste, na katerih se pojavi iskalni niz)

Ocena glede na vsebino

upošteva, kje se nahaja iskalni niz (naslov, povzetek ali jedro),
upošteva relativne frekvence pojavitev ključnih besed,
izračuna se na podlagi indeksa vsebine in obrnjene seznama,
odvisna od iskalnega niza

Ocena glede na popularnost

upošteva analizo strukturnih lastnosti hiperpovezav,
izračuna se iz strukturnega indeksa,
ponavadi **neodvisna od iskalnega niza**

Skupna ocena = vsebinska ocena + ocena popularnosti

rezultat je urejen seznam relevantnih strani,
bolj „verjetne“ strani so višje na seznamu

Rangiranje strani po popularnosti

Google

Splet Slike Videoposnetki Zemljevidi Knjige Več ▾ Orodja za iskanje

Približno 289.000.000 rez. (0,22 sek.)

Michael Jordan - Wikipedija, prosta enciklopedija ✓
sl.wikipedia.org/wiki/Michael_Jordan ▾

Michael Jordan je nekdanji košarkar v ligi NBA, poslovnež in večinski lastnik NBA kluba Charlotte Bobcats. V njegovi biografiji na straneh lige NBA je zapisano, ...

Michael Jordan - Wikipedia, the free encyclopedia ✓
en.wikipedia.org/wiki/Michael_Jordan ▾ [Prevedi to stran](#)

Michael Jeffrey Jordan (born February 17, 1963), also known by his initials, MJ, is an American former professional basketball player, entrepreneur, and ...
[Jeffrey Jordan](#) - [Marcus Jordan](#) - [1984 - Air Jordan](#)

Michael Jordan NBA Stats | Basketball-Reference.com ✓
www.basketball-reference.com ▾ [Players](#) ▾ [J](#) ▾ [Prevedi to stran](#)

Michael Jordan - Career stats, game logs, biographical info, awards, and achievements for the NBA and NCAA.

Michael Jordan Stats, Bio - ESPN ✓
espn.go.com/nba/player/_id/1035/michael-jordan ▾ [Prevedi to stran](#)

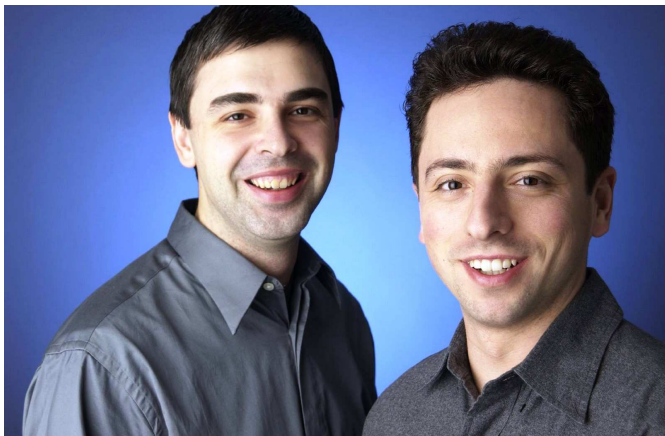
Get the latest news, career stats and more about guard **Michael Jordan** on ESPN.com.

Michael Jordan - Biography - Basketball Player - Biography ... ✓
www.biography.com/people/michael-jordan-9358066 ▾ [Prevedi to stran](#)

Follow the career of former basketball star **Michael Jordan**, from his college career to being the Chicago Bulls' MVP, to his multiple retirements, ...

Michael Jordan - Forbes ✓
www.forbes.com/profile/michael-jordan/ ▾ [Prevedi to stran](#)

Basketball's greatest player realized his dream of becoming an NBA owner this year when he became the majority shareholder in the Charlotte Bobcats. **Jordan** ...



PageRank algoritem

Larry Page in Sergey Brin

članek *The anatomy of large-scale hypertextual Web search engine*

Spletne strani, vhodne in izhodne povezave

poveži eno stran z drugo stranjo



Spletni graf

vozlišča = spletne strani

usmerjene povezave = hiperpovezave

Kaj vpliva na oceno popularnosti strani?

Hiperpovezave kot glasovi na volitvah

stran i kaže na stran $j \Leftrightarrow$ obstaja usmerjena povezava od i do j

Več glasov večja popularnost

število vseh strani, ki kažejo na njo

Pomemben status glasovalca

popularnost strani, ki kaže na njo

Glas volilca se sorazmerno uteži glede na število oddanih glasov

število izhodnih povezav strani, ki kaže na njo

Popularnost strani je določena s popularnostjo strani, ki kažejo na njo.

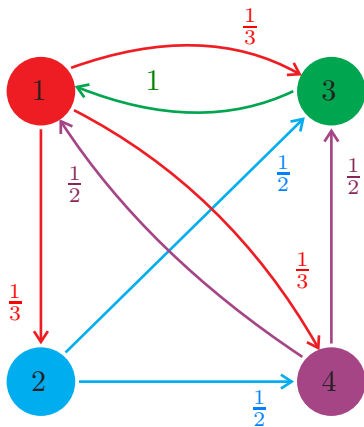
Popularnost strani je določena s popularnostjo strani, ki kažejo na njo:

$$P_1 = P_3 + \frac{1}{2}P_4$$

$$P_2 = \frac{1}{3}P_1$$

$$P_3 = \frac{1}{3}P_1 + \frac{1}{2}P_2 + \frac{1}{2}P_4$$

$$P_4 = \frac{1}{3}P_1 + \frac{1}{2}P_2$$



Popularnost strani je določena s popularnostjo strani, ki kažejo na njo:

$$P_1 = P_3 + \frac{1}{2}P_4$$

$$P_i = \sum_{j \in V_i} \frac{1}{N_j} P_j$$

pomembnost strani i

strani j, ki kažejo na stran i

število izhodnih povezav iz vozlišča j

pomembnost strani j

Prepišimo v primerno matematično obliko:

$$P_1 = 0P_1 + 0P_2 + P_3 + \frac{1}{2}P_4$$

$$P_2 = \frac{1}{3}P_1 + 0P_2 + 0P_3 + 0P_4$$

$$P_3 = \frac{1}{3}P_1 + \frac{1}{2}P_2 + 0P_3 + \frac{1}{2}P_4$$

$$P_4 = \frac{1}{3}P_1 + \frac{1}{2}P_2 + 0P_3 + 0P_4$$

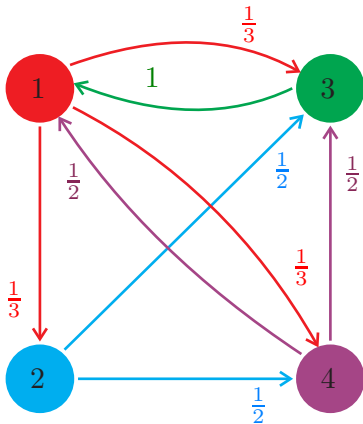
Ocena popularnosti kot lastni vektor matrike

In še malo bolj primerno obliko:

$$\begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \\ P_4 \end{bmatrix}$$

ali

$$x = Hx$$



element j v stolpcu i = verjetnost prehoda iz vozlišča i v vozlišče j

Uporaba iterativne metode

Iščemo: x , da velja $x = Hx$

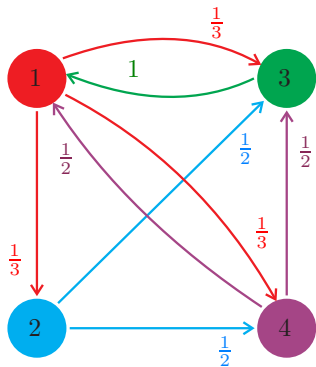
Potenčna metoda:

Vzemi začetni približek $x^{(0)}$

Ponavljaj

$$x^{(k+1)} := Hx^{(k)}$$

dokler „ne konvergira“



$$H = \begin{bmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{bmatrix}$$

$$x^{(0)} = \begin{bmatrix} 0,25 \\ 0,25 \\ 0,25 \\ 0,25 \end{bmatrix} \quad x^{(1)} = \begin{bmatrix} 0,37 \\ 0,08 \\ 0,33 \\ 0,20 \end{bmatrix} \quad x^{(2)} = \begin{bmatrix} 0,43 \\ 0,12 \\ 0,27 \\ 0,16 \end{bmatrix}$$

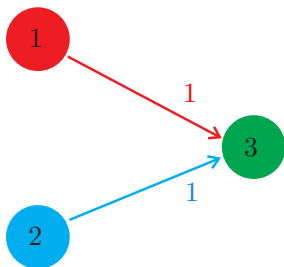
$$x^{(3)} = \begin{bmatrix} 0,35 \\ 0,14 \\ 0,29 \\ 0,20 \end{bmatrix} \quad x^{(4)} = \begin{bmatrix} 0,39 \\ 0,11 \\ 0,29 \\ 0,19 \end{bmatrix} \quad x^{(5)} = \begin{bmatrix} 0,39 \\ 0,13 \\ 0,28 \\ 0,19 \end{bmatrix}$$

$$x^{(6)} = \begin{bmatrix} 0,38 \\ 0,13 \\ 0,29 \\ 0,19 \end{bmatrix} \quad x^{(7)} = \begin{bmatrix} 0,38 \\ 0,12 \\ 0,29 \\ 0,19 \end{bmatrix} \quad x^{(8)} = \begin{bmatrix} 0,38 \\ 0,12 \\ 0,29 \\ 0,19 \end{bmatrix}$$

Ali konvergira k čemu smiselnemu v kontekstu ocene popularnosti?

Izolirana vozlišča brez izhodnih povezav (H ni stohastična):

pdf dokumenti, slike, tabele

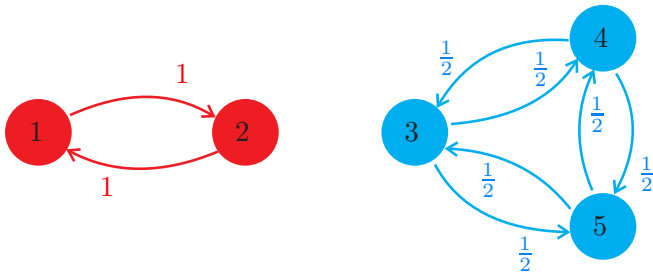


$$\mathbf{H} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix} \quad x^{(0)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \quad x^{(1)} = \begin{bmatrix} 0 \\ 0 \\ \frac{2}{3} \end{bmatrix} \quad x^{(2)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Ali je rešitev enolična?

Graf ima več komponent (H je razcepna):

večina spletni strani kaže samo na prgišče ostalih strani



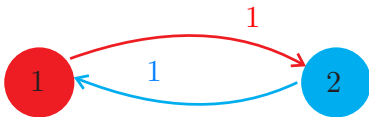
$$\mathbf{H} = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

Ali je konvergenca odvisna od izbire začetnega približka?

Cikli (H je periodična):

medsebojne združbe, ki glasujejo samo med sabo

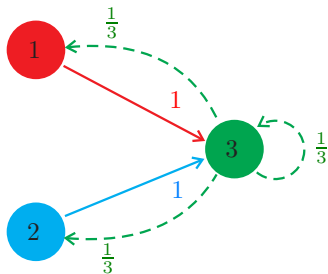


$$\mathbf{H} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

$$x^{(0)} = \begin{bmatrix} 0,5 \\ 0,5 \end{bmatrix} \quad x^{(1)} = \begin{bmatrix} 0,5 \\ 0,5 \end{bmatrix}$$

$$x^{(0)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad x^{(1)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad x^{(2)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Popravimo ničelne stolpce



postavi elemente 3. stolpca na $\frac{1}{3}$

$$\mathbf{H} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 1 & 0 \end{bmatrix} \quad \mathbf{S} = \begin{bmatrix} 0 & 0 & \frac{1}{3} \\ 0 & 0 & \frac{1}{3} \\ 1 & 1 & \frac{1}{3} \end{bmatrix}$$

Zamenjaj matriko \mathbf{H} z matriko $\mathbf{S} = \mathbf{H} + \begin{bmatrix} \frac{a_1}{n} & \dots & \frac{a_n}{n} \\ \vdots & \ddots & \vdots \\ \frac{a_1}{n} & \dots & \frac{a_n}{n} \end{bmatrix}$

$a_i = 1$, če vozlišče i nima izhodnih povezav, in 0 sicer

Zamenjaj matriko \mathbf{S} z matriko $\mathbf{G} = \alpha \mathbf{S} + (1 - \alpha) \begin{bmatrix} \frac{1}{n} & \dots & \frac{1}{n} \\ \vdots & \ddots & \vdots \\ \frac{1}{n} & \dots & \frac{1}{n} \end{bmatrix}$

„teleportacijski“ koeficient $0 \leq \alpha < 1$ (običajno $\alpha = 0,85$)

G je stohastična, nerazcepna in aperiodična

enoličen dominanten lastni vektor:

$$\mathbf{x} = \mathbf{G}\mathbf{x}, \quad \mathbf{x} \geq 0, \quad \sum_{i=1}^n x_i = 1$$

x_i je ocena popularnosti za spletno stran i

G ima vse elemente neničelne

velikosti več deset milijard

Toda matrika H ima večino elementov ničelnih

$$x^{(k+1)} = \mathbf{G}x^{(k)} = \alpha \mathbf{H}x^{(k)} + \alpha \begin{bmatrix} \frac{a_1}{n} & \dots & \frac{a_n}{n} \\ \vdots & \ddots & \vdots \\ \frac{a_1}{n} & \dots & \frac{a_n}{n} \end{bmatrix} x^{(k)} + (1-\alpha) \begin{bmatrix} \frac{1}{n} & \dots & \frac{1}{n} \\ \vdots & \ddots & \vdots \\ \frac{1}{n} & \dots & \frac{1}{n} \end{bmatrix} x^{(k)}$$

Implementacija

niti **G** niti **S** nista eksplicitno shranjeni,
pri računanju se matrika **G** ne spreminja,
v praksi za izračun zadostuje 50-100 iteracij

- S. Brin and L. Page. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1-7):107–117, 1998.
- A. Borodin, G.O. Roberts, J.S. Rosenthal, P. Tsaparas. *Link Analysis Ranking Algorithms Theory And Experiments*, 2004
<http://cs-people.bu.edu/evimaria/cs565/lar.pdf>
- J. Khoury. How is it made? Google Search Engine, 2010.
<http://aix1.uottawa.ca/jkhoury/>
- A. N. Langville and C. D. Meyer. *Google's PageRank and Beyond: The Science of Search Engine Rankings*. Princeton University Press, 2006.
- R. S. Wills. *Google's PageRank: The Math Behind the Search Engine*, 2006.
http://www.cems.uvm.edu/tlakoba/AppliedUGMath/other_Google/Wills.pdf

WHEN IN DOUBT..


Hvala za pozornost!