

Univerza v Ljubljani
Fakulteta za *matematiko in fiziko*



Katzova središčnost in Googlov PageRank

POROČILO PROJEKTA PRI PREDMETU FINANČNI PRAKTIKUM

Anamari Oštarijaš, Tina Ražić

Ljubljana, december 2018

Kazalo

1	Opis projekta	3
2	Katzova središčnost	3
2.1	Matematično ozadje	3
2.2	Psevdokoda	3
3	Googlov PageRank	3
3.1	Potenčna metoda	6
3.2	Algoritem za Googlov PageRank	6
3.3	Analiza algoritma	7
3.3.1	Izbira parametra α	7
4	Primerjava algoritmov	9

1 Opis projekta

Kompleksna omrežja lahko analiziramo z uporabo različnih kvantitetnih merjenj, imenujemo jih tudi mere središčnosti, ki intuitivno zajamejo pomembnost določenih vozlišč. V projektu bova implementirali Googlov PageRank in Katzovo središčnost z uporabo potenčne metode. Na različnih grafih (tudi socialnih omrežjih) bova analizirali in primerjali, kako merjenji razvrstita vozlišča po pomembnosti.

2 Katzova središčnost

2.1 Matematično ozadje

Katzova središčnost izmeri vpliv igralca v omrežju tako, da upošteva direktne sosedo igralca in vse druge igralce, ki so posredno povezani s tem igralcem preko njegovih direktnih sosedov.

Naj bo naše omrežje graf z n vozlišči oziroma igralci. Vsaka povezava v grafu dobi utež α in z α^d izračunamo težo povezave vozlišča z drugim vozliščem, pri čemer je d število povezav med njima. Naš graf predstavimo z matriko sosednosti A , torej element matrike a_{ij} ima vrednost 1, če je vozlišče i povezano z vozliščem j , in 0, če nista povezana. Potence matrike A nam povejo, če je vozlišče povezano s drugimi indirektnimi vozlišči preko sosedov. Na primer, če je v matriki A^3 element $a_{2,5} = 1$, pomeni, da sta vozlišče 2 in vozlišče 5 povezana s tremi povezavami preko sosedov prve stopnje in sosedov druge stopnje.

Označimo s $C_{Katz}(i)$ Katzovo središčnost vozlišča i . Potem lahko izračunamo središčnost na sledeči način:

$$C_{Katz}(i) = \sum_{k=1}^{\infty} \sum_{j=1}^n \alpha^k (A^k)_{ij}.$$

Pri izbiri α moramo upoštevati zgornjo omejitev

$$\alpha < \frac{1}{|\lambda_{max}|}.$$

2.2 Psevdokoda

3 Googlov PageRank

Nemogoče je definirati splošno mer pomembnosti, ki bi bila sprejemljiva za vse uporabnike iskalnika. Google uporablja pagerank za mero kakovosti spletnih strani. Temelji na predpostavki, da število linkov (povezav) do in iz strani daje informacijo o pomembnosti strani.

Naj bodo vse spletne strani urejene s števili od 1 do n in naj bo i neka spletna stran. Potem O_i določa množico strani, s katerimi je i povezana, tako da i vsebuje link do strani v množici O_i (*outlink*). Število outlinkov označimo z $N_i = \|O_i\|$. Množica *inlinkov*, označena z I_i , je množica strani, ki imajo outlink do i (strani v I_i vsebujejo linke do i). Splošno, več ko ima strani inlinkov, pomembnejša je.

Vseeno pa bi bilo s takim sistemom preprosto manipulirati (če bi nekdo želel, da

njegovo spletno stran vidi čim več ljudi, bi ustvaril veliko število nepomembnih spletnih strani, ki bi vsebovale linke do njegove spletne strani). Da bi tako manipulacijo preprečili, definiramo rang vozlišča i tako, da če ima visoko rangirana stran j outlink do i , to doda pomembnosti i na sledeč način: rang strani i je utežena vsota rangov strani, ki imajo outlink do i . Obteženost je taka, da je rang strani j razdeljen enakomerno med njenimi outlinki. Z enačbo:

$$r_i = \sum_{j \in I_1} \frac{r_j}{N_j}.$$

Ta definicija je rekurzivna, zato pageranki ne morejo biti izračunani direktno. Uporabimo iteracijo. Najprej ugibamo začetni rangni vektor r^0 . Potem iteriramo:

$$r_i^{(k+1)} = \sum_{j \in I_1} \frac{r_j^{(k)}}{N_j}$$

Raje zapišimo problem z matrikami. Naj bo Q kvadratna matrika dimenzije n . Definiramo:

$$Q_{ij} = \begin{cases} 1/N_j & , \text{ če obstaja link od } j \text{ do } i \\ 0 & , \text{ sicer} \end{cases}$$

Torej ima vrstica i neničelne elemente na mestih inlinkov i . Podobno ima stolpec j neničelne elemente enake N_j na mestih outlinkov j , vsota vseh elementov v stolpcu je enaka 1. Definicija je ekvivalentna skalarnem produktu vrstice i in vektorja r , ki vsebuje range vseh strani.

Enačbo sedaj lahko zapišemo v matrični obliki:

$$\lambda r = Qr, \quad \lambda = 1,$$

Tako dobimo, da je r lastni vektor matrike Q z lastno vrednostjo $\lambda = 1$. Sedaj preprosto vidimo, da je iteracija ekvivalentna

$$r^{(k+1)} = Qr^{(k)}, \quad k = 0, 1, \dots,$$

kar je potenčna metoda za izračun lastnega vektorja.

Problem: Ni jasno, da je pagerank dobro definiran, saj ne vemo ali obstaja lastna vrednost enaka 1. Pomagamo si s teorijo Markovske verige.

Obstaja interpretacija pageranka z naključnimi sprehodi. Predpostavimo, da uporabnik na spletni strani izbere naslednjo stran med outlinki z enako verjetnostjo. Markovska veriga je slučajni proces v katerem je naslednje stanje določeno le s trenutnim stanjem, proces nima spomina. Prehodna matrika Markovske verige je Q^T .

Slučajni uporabnik nikoli ne sme obtičati. Z drugimi besedami, naš model slučajnega sprehoda ne sme imeti strani brez outlinkov (taka stran bi imela stolpec iz ničel v Q). Torej je model prilagojen tako, da so ničelni stolpci nadomeščeni s konstantnimi vrednostmi na vseh mestih. To pomeni, da je enaka verjetnost, da pridemo na katero koli drugo spletno stran.

Definiramo vektorje:

$$d_j = \begin{cases} 1 & , \text{če } N_j = 0 \\ 0 & , \text{sicer} \end{cases}$$

za $j = 1, \dots, n$ in

$$e = [1 \dots 1]^T \in \mathbb{R}^n.$$

Prilagojena matrika je definirana s $P = Q + \frac{1}{n}ed^T$. S to prilagoditvijo je matrika P desna stohastična matrika; vsi elementi so nenegativni in elementi vsakega stolpca se seštevajo v 1.

Desna stohastična matrika P zadošča

$$e^T P = e^T$$

Želimo definirati pagerank vektor kot enoličen lastni vektor matrike P z lastno vrednostjo 1:

$$Pr = r$$

Lastni vektor prehodne matrike je stacionarna verjetnostna porazdelitev Markovske verige. Element na mestu i (r_i) je verjetnost, da po velikem številu korakov slučajni uporabnik pristane na strani i . Za zagotovitev enoličnosti mora biti matrika ireducibilna (slučajni uporabnik se lahko v nekem delu grafa zagodzi, v tem primeru ima matrika več lastnih vrednosti enakih 1).

Edinstvenost največje lastne vrednosti ireducibilne, pozitivne, desne stohastične matrike je zagotovljena s *Perron-Frobeniusovim izrekom*, največja singularna vrednost bo enaka 1, pripadajoč slastni vektor je pozitiven in je edini lastni vektor, ki je nenegativen.

Zaradi velikosti spleta je matrika linkov P reducibilna, torej pagerank lastni vektor ni dobro definiran. Da si zagotovimo ireducibilnost, umetno dodamo linke iz vsake spletne strani do vseh drugih. To lahko storimo, če vzamemo konveksno kombinacijo P in matrike ranga 1:

$$A = \alpha P + (1 - \alpha) \frac{1}{n} ee^T,$$

za nek α , ki zadošča $0 \leq \alpha \leq 1$. Matrika A je desna stohastična. Razlaga naključnega sprehoda dodatnega rang-1 izraza je, da bo uporabnik na vsakem časovnem koraku skočil na naključno stran z verjetnostjo $1 - \alpha$. Za konvergenco numeričnega algoritma za lastno vrednost je pomembno, da vemo, kako so s prilagoditvijo ranga-1 spremenjene lastne vrednosti.

Izrek: Naj bodo lastne vrednosti desne stohastične matrike P enake $1, \lambda_2, \lambda_3, \dots, \lambda_n$. Potem so lastne vrednosti $A = \alpha P + (1 - \alpha) \frac{1}{n} ee^T$ enake $1, \alpha\lambda_2, \alpha\lambda_3, \dots, \alpha\lambda_n$.

Ta izrek nam pove, da tudi če ima P več lastnih vrednosti enakih 1, kar je po navadi res, bo druga največja lastna vrednost matrike A enaka α . Namesto prilagoditve lahko definiramo:

$$A = \alpha P + (1 - \alpha) ve^T,$$

Kjer je v nenegativen vektor z normo 1, ki je lahko izbran tako, da je iskanje pristransko do strani določene vrste. Zato ga včasih imenujemo *personaliziran vektor*. Uporaben je tudi v izogib manipulacijam proti t.i. *link farms*.

3.1 Potenčna metoda

Želimo rešiti problem lastne vrednosti:

$$Ar = r,$$

Kjer je r normaliziran $\|r\|_1 = 1$. Iskani lastni vektor označimo s t_1 . Prepostavimo, da imamo dan začetni približek $r(0)$.

Algoritem:

za $k = 1, 2, \dots$ do konvergence

$$q^{(k)} = Ar^{(k-1)}$$

$$r^{(k)} = q^{(k)} / \|q^{(k)}\|_1$$

Konvergenca je odvisna od porazdelitve lastnih vrednosti. Če je druga največja lastna vrednost blizu 1, bo iteracija zelo počasna. To na srečo ne velja za Google matriko. Vektor normaliziramo, da ne bi z iteracijami postali preveliki ali premajhni, posledično nereprezentativni s števili s plavajočo vejico. To v resnici ni potrebno, saj se v primeru desnih stohastičnih matrik temu izognemo.

3.2 Algoritem za Googlov PageRank

```

1 def Qmatrix(matrika):
2     '''za dano matriko sosednosti grafa vrne desno stohastično
3     matriko Q, kjer so nekateri stolpci se vedno ničelni'''
4     sums = np.sum(matrix, axis=0) #tabela vsot stolpcev
5     matrike
6     return np.nan_to_num(matrix/sums) #vsak stolpec matrike
7     delimo z njegovo vsoto, ce je vsota 0, vrne 0 namesto nan
8
9 def dvector(matrika):
10    '''za dano matriko sosednosti vrne transponiran vektor d, kjer
11    d[j] = 1 ; ce #(outlinkov od j) = 0 in d[j] = 0 ; sicer'''
12    Q = Qmatrix(matrix) #izracunamo matriko Q
13    sums = np.sum(matrix, axis=0) #tabela vsot stolpcev matrike
14    d_bool = (sums == 0) #bool tabela, True ce je vsota
15    stolpca enaka 0 in False sicer
16    d = d_bool*1 #False v 0 in True v 1
17    return d
18
19 def evector(n):
20    '''vrne vektor v iz samih enk velikosti n, kjer n = st.
21    stolpcev v matriki'''
22    return np.ones((n,1))
23
24 def Pmatrix(matrix):
25    '''vrne desno stohastično matriko P'''
26    n = matrix.shape[1]
27    Q = Qmatrix(matrix)

```

```

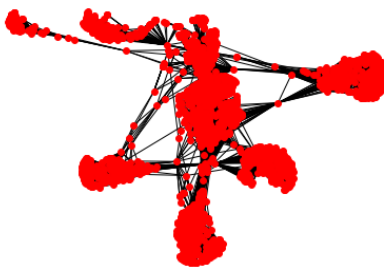
22     e = evector(n)
23     d = dvector(matrix)
24     return Q + 1/n*e.dot(d)
25
26 def Amatrix(matrix, alpha):
27     '''vrne ireducibilno desno stohastično matriko A, kjer je 0 <=
28     alpha <= 1'''
29     n = matrix.shape[1]
30     e = evector(n)
31     P = Pmatrix(matrix)
32     return alpha*P + (1-alpha)*1/n*e.dot(e.transpose())
33
34 def pagerank2(matrix, max_num_of_steps, tolerance, alpha):
35     '''vrne pagerank vektor matrike sosednosti grafa, 0<=alpha<=1
36     ,,,
37
38     start = timer()
39     A = Amatrix(matrix, alpha)                #desno stohastična matrika
40     dane matrike
41     r = A[:,0]                                #ugibanje pagerank vektorja s prvim
42     stolpcem matrike Q
43     diff = 1000
44     k = 0
45     while diff > tolerance and k < max_num_of_steps:
46         r, q = A.dot(r), r                    #potencna
47         metoda
48         #r = q/np.linalg.norm(q, ord=1)       #ni potrebno pri Google
49         matriki
50         diff = np.linalg.norm(q-r, ord=1)     #razlika med prejsnjim
51         in novim vektorjem r
52         k+=1
53     print('proces končan po {} iteracijah'.format(k))
54     end = timer()
55     print('trajanje procesa: {} sekund'.format(end-start))
56     return r

```

3.3 Analiza algoritma

Analizirali sva učinkovitost algoritma na podlagi večih parametrov. Najbolj naju je zanimalo, kako izbira parametra α , toleranca in število korakov vplivajo na hitrost algoritma in natančnost rezultata.

Uporabljali sva *Facebook graf*, ki vsebuje 4039 vozlišč in 88234 povezav.



3.3.1 Izbira parametra α

Vemo, da večji ko je α , več podatkov o omrežju izgubimo. Torej želimo, da je ta parameter čim večji.

Poglejmo najprej kako velikost vpliva na čas izračuna pagerank vektorja.

```
1 M_fb = nx.adjacency_matrix(G_fb)
2 r_1 = pagerank2(M_fb,10000000,0.0000001, 0.98)
3 print(r_1)
4
5 proces končan po 575 iteracijah
6 trajanje procesa: 11.7118154744594 sekund
7 [[ 4.79508113e-03]
8 [ 2.17910466e-04]
9 [ 1.52316417e-04]
10 ...,
11 [ 6.77736016e-05]
12 [ 1.25649865e-04]
13 [ 2.72346363e-04]]
```

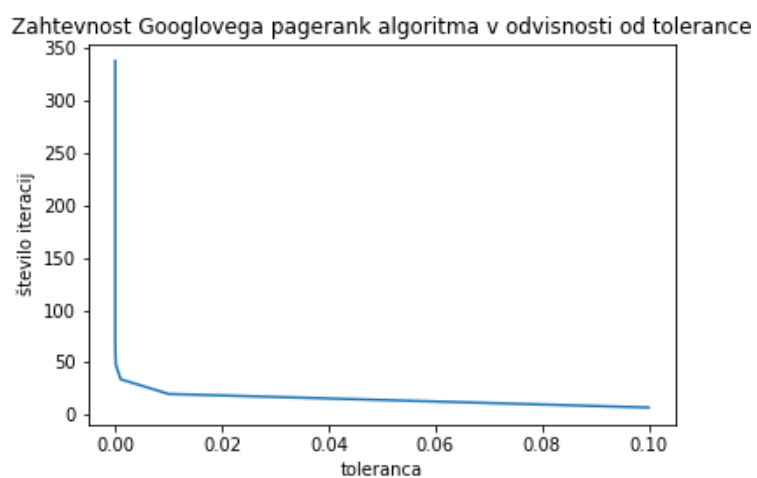
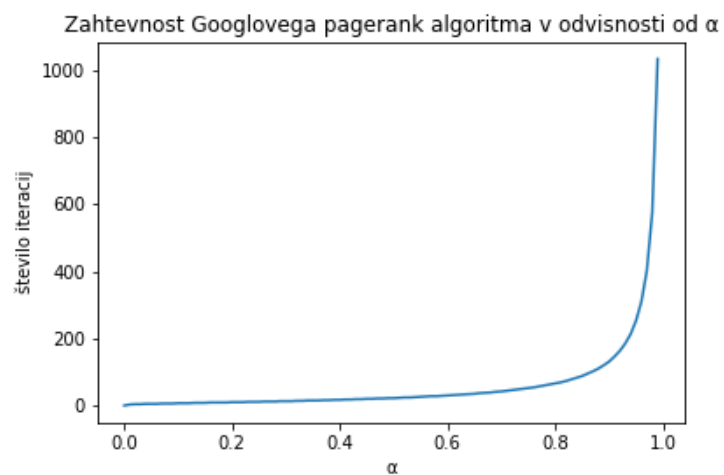
Izračun je bil končan po 575 iteracijah in 11,71 sekundah.

```
1 r_2 = pagerank2(M_fb,10000000,0.0000001, 0.999)
2 print(r_2)
3
4 proces končan po 4719 iteracijah
5 trajanje procesa: 61.73710875091092 sekund
6 [[ 2.37207123e-03]
7 [ 1.15387662e-04]
8 [ 6.99886732e-05]
9 ...,
10 [ 2.37766591e-05]
11 [ 4.70770334e-05]
12 [ 1.04002744e-04]]
```

S povečanjem α smo povročili, da se je čas izračuna povečal na 61,74 sekund, torej za več kot 5 kratnik prejšnjega.

```
1 r_4 = pagerank2(M_fb,10000000,0.0000001, 0.9999)
2 print(r_4)
3
4 proces končan po 8104 iteracijah
5 trajanje procesa: 105.30294071864228 sekund
6 [[ 1.99799361e-03]
7 [ 9.77897213e-05]
8 [ 5.77554403e-05]
9 ...,
10 [ 1.28541957e-05]
11 [ 2.56615989e-05]
12 [ 5.75079410e-05]]
```


Večji ko je α , več iteracij potrebujemo, da pridemo do željene tolerance.
 Razlika med r_1 in r_2 v drugi normi je 0,00717920758147, med r_2 in r_4 pa 0,00371765512728, torej lahko z višjim α pričakujemo bolj natančne rezultate.



4 Primerjava algoritmov