

**Laboratorij za tehnologije znanja (KTLab)**

Zavod za elektroniku, mikroelektroniku, računalne i inteligentne sustave

Fakultet elektrotehnike i računarstva

Sveučilište u Zagrebu

## **Interni dokument**

**© 2009 KTLab**

Niti jedan dio ovog dokumenta ne smije se fotokopirati,  
umnožavati niti prevoditi na drugi jezik  
bez prethodnog pismenog odobrenja.

SVEU ILIŠTE U ZAGREBU  
FAKULTET ELEKTROTEHNIKE I RA UNARSTVA

DIPLOMSKI ZADATAK br. 1614

**PRIMJENA TEORIJE GRAFOVA U DUBINSKOJ  
ANALIZI TEKSTA**

Frane Šari

Zagreb, rujan 2006.

# Sadržaj

SADRŽAJ .....	1
1. SAŽETAK .....	3
2. UVOD .....	3
3. ANALIZA VEZA U GRAFOVIMA .....	7
3.1. PageRank.....	7
3.2. Metoda vorišta i autoriteta .....	9
3.2.1. Nalaženje stranica autoriteta .....	10
3.3. Sličnost me u vrhovima dva grafa .....	13
4. NALAŽENJE SINONIMA U RJEČNIKU .....	15
4.1. Primjena matrice sličnosti .....	15
4.2. Homografi i polisemi .....	17
4.3. Problem zaustavnih riječi.....	18
5. PRIMJENA NA INDEKSIRANJE DOKUMENATA.....	20
5.1. Nalaženje sličnih dokumenata .....	21
5.2. Osnovni sustav .....	24
6. REZULTATI .....	26
6.1. Nalaženje sinonima.....	26
6.1.1. Nalaženje sinonima u engleskom jeziku .....	26
6.1.2. Evaluacija .....	27
6.1.3. Usporedba s implementacijom u Blondelovom linku.....	27
6.1.4. Nalaženje sinonima u hrvatskom jeziku .....	29
6.2. Automatsko indeksiranje dokumenata .....	31
6.2.1. Osnovna metoda .....	33
6.2.2. Korištenje matrice sličnosti .....	34
7. ZAKLJUČAK .....	36
8. LITERATURA .....	37
9. DODATAK .....	40

<b>9.1. Popis tablica.....</b>	<b>40</b>
<b>9.2. Popis slika .....</b>	<b>41</b>

INTERNI DOKUMENT

# 1. Sažetak

Posljednjih desetak godina web-tražilice poele su primjenjivati metode analize grafova radi rangiranja rezultata pretrage. Budu i da su se te metode pokazale uspješnima, njihova primjena na razna područja sve više se istražuje.

U Blondelovom članku [3] opisana je mjera sličnosti među vrhovima dva grafa, te primjena te metode na nalaženje sinonima u rječniku. Nalaženje sinonima jedan je od alata koji se može koristiti prilikom dubinske analize teksta. U ovom radu proučena je primjena mjere sličnosti na automatsko indeksiranje dokumenata.

Uspješnost metode nalaženja sinonima ocijenjena je korištenjem engleskog i hrvatskog rječnika. Automatsko indeksiranje korištenjem mjere sličnosti evaluirano je na skupu dokumenata indeksiranim pojmovima iz hijerarhijskog pojmovnika EUROVOC [5].

U radu se koriste neki od rezultata projekta „Text mining system – Sustav za automatsko indeksiranje, kategorizaciju i semantičko pretraživanje teksta“ koji je nastao kao suradnja tima Fakulteta elektrotehnike i računarstva i Filozofskog fakulteta Sveučilišta u Zagrebu, Hrvatske informacijsko-dokumentacijske referalne agencije i Zajedničkog istraživačkog centra Europske komisije u Ispri, Italija.

## 2. Uvod

Porastom popularnosti Interneta, a posebno web-stranica, broj dostupnih dokumenata poveao se gotovo geometrijskom progresijom. Uz porast koli ine dokumenata pojavio se i problem nalaženja informacija, tj. dokumenata koji sadrže korisne informacije. Iako su prve web tražilice (npr. Excite<sup>1</sup>, Hotbot<sup>2</sup>, Altavista<sup>3</sup>) omogu ile jednostavnu i brzu pretragu tog skupa, relevantni dokumenti esto se nisu nalazili me u prvim rezultatima pretrage. Nalaženje poretka prona enih dokumenata prema tome koliko su oni relevantni za korisni ki upit ostao je veliki problem. Najve i korak u rješavanju tog problema na inila je web tražilica Google<sup>4</sup> kada je prije desetak godina po ela primjenjivati metode teorije grafova na analizu veza me u web stranicama, tj. algoritam PageRank [6]. Od tada se pokušavaju sli ne metode primijeniti na mnoga druga podru je, ne nužno vezana uz pretragu web stranica.

Podru je u kojem se tek po inju primjenjivati sli ne metode jest dubinska analiza teksta. Dubinska analiza teksta podrazumijeva pronalaženje i sintezu informacija iz isklju ivo tekstnih podataka. Jedan od alata koji se koristi prilikom analize teksta jest nalaženje sinonima neke rije i. Ponekad je korisno svesti razne sinonime neke rije i u jednu kategoriju. Sinonimi se mogu koristiti i za transformaciju upita koji korisnik zadaje na na in da rije i iz upita budu sli nije rije ima korištenim u dokumentu [14, 13]. Blondel [4] opisuje metodu nalaženja sinonima u engleskom jeziku korištenjem engleskog rije nika. U ovom radu ta metoda primijenjena je na nalaženje sinonima u hrvatskom jeziku.

---

<sup>1</sup> <http://www.excite.com>

<sup>2</sup> <http://www.hotbot.com>

<sup>3</sup> <http://www.altavista.com>

<sup>4</sup> <http://www.google.com>

Prilikom pretrage web stranica korisnik zadaje upit slobodnim odabirom rije i. Ponekad je teško odabrati prave rije i ijim zadavanjem bi tražilica vratila rezultate koje korisnik traži. U slučaju da se pretražuju znatno manje kolekcije dokumenata (od svih web stranica na Internetu) može se svaki dokument unaprijed označiti sa skupom pojmova relevantnim za taj dokument. Sami pojmovi kojima se označavaju dokumenti odabire se iz kontroliranog rječnika ili pojmovnika. Proces dodjeljivanja relevantnih pojmova svakom dokumentu zovemo indeksiranje. Indeksirane dokumente korisnik može pretraživati tako da koristi samo pojmove iz kontroliranog rječnika. Ukoliko su pojmovi hijerarhijski organizirani, lako se mogu naći i specifičiji ili općenitiji rezultati pretrage promjenom pojmova upita na specifičnije ili općenitije pojmove.

Jedan od problema kojima se bavi dubinska analiza teksta jest automatska kategorizacija dokumenata, tj. svrstavanje dokumenata u jednu ili više zadanih kategorija. Indeksiranje dokumenata možemo shvatiti kao kategorizaciju dokumenata u niz kategorija koje određeni pojmovi kojima se indeksiraju dokumenti. U Blondelovom članku [3] opisana je mjera složenosti me u vrhovima dva grafa. U ovom radu pokušava se odgovoriti na pitanje kako se može koristiti složenost na metoda u svrhu automatskog indeksiranja dokumenata korištenjem pojmovnika EUROVOC [5].

U literaturi su opisane razne automatske metode nalaženja sinonima. Metode koje su u svojim radovima opisali Wu [24] i Curran [7] koriste samo velike količine tekstova u jednom jeziku. Poput Blondelove metode [4] koja je implementirana u ovom radu i metoda ArcRank [11] koristi rječnik za nalaženje sinonima.

Već su razvijeni sustavi za automatsko indeksiranje dokumenata korištenjem kontroliranog pojmovnika koji ne zahtijevaju eksperta, tj. uvijek koji bi ručno zadao pravila dodjele pojmova dokumentima. Sustav sveučilišta u Berkeleyu [18] koristi složenost naziva pojmova i raznih dijelova dokumenta za indeksiranje dokumenta. Sustav CONDORCET [23] analizom sažetaka i

naslova dokumenata određuje kojim pojmovima se indeksira dokument. Pouliquen [19] razvio je sustav za indeksiranje dokumenata korištenjem pojmovnika EUROVOC. Ovaj sustav, koristeći skup indeksiranih dokumenata, svakom pojmu iz kontroliranog rječnika dodjeljuje niz ključnih riječi i težina. Prilikom indeksiranja novog dokumenta traže se pojmovi čije riječi najbolje odgovaraju riječi ima u dokumentu. U literaturi se metode analize grafova nisu koristile za potrebe indeksiranja dokumenata.

INTERNI DOKUMENT

### 3. Analiza veza u grafovima

Značajni skok u kvaliteti rangiranja web stranica dogodio se kada je tražilica Google uvela metodu *PageRank* [6]. Google je počeo rangirati stranice ne samo prema analizi pojavljivanja upita na traženoj stranici, već i korištenjem mjere PageRank koja svakoj stranici dodjeljuje važnost neovisno o upitu korisnika. PageRank u potpunosti ovisi o vezama među različitim web stranicama.

Metoda koristi autoriteta tako što koristi strukturu povezanosti web stranica. Tom metodom određuju se mjerodavnosti (koliko je neka stranica dobar autoritet) stranica s obzirom na dani korisnički upit. Generalizacijom ove metode dolazimo do pojma sličnosti među vrhovima u grafu, tj. do matrice sličnosti.

Ove metode kao relevantnu karakteristiku web stranica uzimaju u obzir **veze** (hiperlinkove) koje određena stranica sadrži. Pri tome se gledaju samo veze koje ta stranica ima prema drugim stranicama, a ignoriraju se veze koje pokazuju na drugi dio iste stranice. Nadalje, ako stranica *a* ima vezu koja pokazuje na stranicu *b*, onda ćemo tu vezu zvat **izlaznom vezom** stranice *a* i **ulaznom vezom** stranice *b*.

#### 3.1. PageRank

Uspjeh ove metode može se pripisati iskorištavanju strukture veza među web stranicama radi ocjenjivanja važnosti neke stranice. PageRank vrijednost neke stranice zapravo predstavlja vjerojatnost da će se korisnik nakon odabira slučajne veze u svom web-pregledniku zadržati na toj stranici. Pretpostavka je da će korisnik ako naiđe na stranicu koja nema izlaznih veza nastaviti sa pregledavanjem neke druge, slučajne, stranice na Internetu. Još jedna pretpostavka je da korisnik nastavlja postupak uz određenu vjerojatnost

– svakom korisniku nakon nekog vremena dosadi pregledavanje Interneta. Kada prestane odabirati veze, korisnik se zadržava na zadnjoj stranici koju je posjetio.

Vrijednost PR (PageRank) neke stranice može se izraziti sljedećom formulom:

$$PR(p_i) = \frac{1-d}{N} + d \sum_{p_j \in M(p_i)} \frac{PR(p_j)}{L(p_j)}. \quad (3.1)$$

- $PR(p)$  - vrijednost PageRank stranice  $p$ .  
 $N$  - ukupni broj stranica.  
 $L(p)$  - broj veza na stranici  $p$ .  
 $M(p)$  - skup stranica koje veze pokazuju na  $p$ .  
 $d$  - vjerojatnost da će korisnik nastaviti pregledavati Internet, tj. da će odabrati jednu od veza na stranici  $p$ .

Početni PageRank svake stranice je jednak. Iterativnim računom brzo se dolazi do konvergencije i PageRank se stabilizira. Izraz 3.1 može se napisati ovako:

$$R = \begin{bmatrix} (1-d)/N \\ (1-d)/N \\ \vdots \\ (1-d)/N \end{bmatrix} + d \cdot \frac{A}{\|A\|_1} \cdot R. \quad (3.2)$$

$A$  je matrica susjedstva s elementima  $a_{ij} = \begin{cases} 1, & \text{postoji veza od } p_j \text{ prema } p_i \\ 0, & \text{inače} \end{cases}$

$\|A\|_1$  je 1-norma matrice  $A$ .

Vektor  $R = \begin{bmatrix} PR(p_1) \\ PR(p_2) \\ \vdots \\ PR(p_N) \end{bmatrix}$  sadrži vrijednosti PageRank svake stranice.

Vektor  $R$  može se računati iterativno tako da se u jednadžbi 3.2 na desnoj strani koristi stara vrijednost vektora  $R$ , a na lijevoj strani se računa nova vrijednost. Do konvergencije dolazi nakon malog broja iteracija, što je jedna od prednosti PageRank metode.

### 3.2. Metoda vorišta i autoriteta

Uvjerivši se u učinkovitost PageRank metode, znanstvenici su sve više proučavali nove načine analize veza među web stranicama. Još jedan od razloga zašto je bilo potrebno razviti i neku novu metodu jest taj što je PageRank metoda patentirana. Kleinberg [13] predlaže novu metodu koja za dani korisnikov upit (pretragu) ocjenjuje koliko su stranice u rezultatima dobre vorišta (engl. *hubs*) i koliko su stranice dobri autoriteti (engl. *authorities*). Na primjer, za pretragu „ZEMRIS“ dobar autoritet bila bi stranica [www.fer.hr](http://www.fer.hr), a dobro vorište stranica koja ima vezu na [www.fer.hr](http://www.fer.hr), npr. [www.carnet.hr/clanice/punopravne](http://www.carnet.hr/clanice/punopravne). Za danu pretragu korisniku tražilice prvo bi se predstavile stranice koje su dobri autoriteti.

Za razliku od PageRank metode, gdje je vrijednost PageRank unaprijed izračunata za svaku stranicu neovisno o upitu, metoda vorišta i autoriteta dinamički, tj. nakon što korisnik pošalje upit, računa se za neki reducirani skup stranica  $S$  koliko je svaka stranica dobro vorište ili autoritet.

Skup  $S$  trebao bi biti relativno malen, ali u isto vrijeme i sadržavati veći broj stranica koje su dobri autoriteti za zadani upit. Skup  $S$  gradi se na sljedeći način:

1. Na početku se skup stranica  $S_0$  koji se sastoji od prvih  $t$  rezultata pretrage koriste i pretraživači koji uzima u obzir samo tekst stranice (npr. AltaVista<sup>5</sup>).

---

<sup>5</sup> <http://www.altavista.com>

2.  $S \leftarrow S_0$ .
3. Skup  $S$  proširi se skupom svih stranica na koje pokazuju veze stranica iz skupa  $S_0$ .
4. Skup  $S$  proširi se skupom svih stranica koje pokazuju na stranice iz skupa  $S_0$ . Uzima se u obzir najviše  $d$  ulaznih veza svake stranice iz skupa  $S_0$ .

Prema metodi koju opisuje Kleinberg [13], odabiru se vrijednosti  $t = 200$  i  $d = 50$ . Tvrdi se da tako izgrađen skup  $S$  sadrži većinu autoriteta za zadani upit, te je i dalje relativno malen, a stranice skupa  $S$  međusobno su dobro povezane. Te karakteristike važne su za postupak opisan u sljedećem poglavlju.

### 3.2.1. Nalaženje stranica autoriteta

Glavna ideja Kleinbergovog algoritma jest ta da su stranice dobri autoriteti ne samo ako imaju puno ulaznih veza, već i ako se stranice koje pokazuju na njih preklapaju. Postavljaju se sljedeće uzajamno podupirive pretpostavke:

stranica je dobar autoritet ako su stranice koje pokazuju na nju dobra  
vorišta

stranica je dobro vorište ako pokazuje na dobre autoritete

Svatom vrhu  $i$ , tj. stranici  $p_i$ , pridružit ćemo nenegativne vrijednosti  $a_i$  ( $a$  dolazi od engl. *authority*) i  $h_i$  ( $h$  dolazi od engl. *hub*) koje označavaju koliko je stranica  $p_i$  dobar autoritet ili dobro vorište. Veće vrijednosti znače da je stranica bolji autoritet, odnosno bolje vorište.

Koristeći skup  $S$  možemo konstruirati usmjereni graf  $G = (V, E)$  čiji su vrhovi stranice iz skupa  $S$ , a bridovi veze između u tim stranicama. Formalno,  $V = \{i : p_i \in S\}$ , te  $E = \{(i, j) : p_i, p_j \in V \wedge p_i \in M(p_j)\}$ . Sada možemo numerički zapisati dvije pretpostavke:

$$\begin{aligned}
 a_j &\leftarrow \sum_{i:(i,j) \in E} h_i, \\
 h_j &\leftarrow \sum_{i:(j,i) \in E} a_i.
 \end{aligned}
 \tag{3.3}$$

Ako s  $B$  označimo matricu susjedstva grafa  $G$ , tj. matricu u kojoj je element u retku  $i$  i stupcu  $j$  jednak broju bridova između vrhova  $i$  i  $j$ , onda relacije 2.3 možemo zapisati na sljedeći način:

$$\begin{bmatrix} h \\ a \end{bmatrix}_{k+1} = \begin{bmatrix} 0 & B \\ B^T & 0 \end{bmatrix} \begin{bmatrix} h \\ a \end{bmatrix}_k, \quad k = 0, 1, \dots$$

Skraćeno zapisano,

$$\begin{aligned}
 x_{k+1} &= M x_k, \quad k = 0, 1, \dots \\
 x_k &= \begin{bmatrix} h \\ a \end{bmatrix}_k, \quad M = \begin{bmatrix} 0 & B \\ B^T & 0 \end{bmatrix}.
 \end{aligned}
 \tag{3.4}$$

Zanimaju nas jedino relativni rezultati, pa stoga promatramo normirani niz gdje vektore  $x_k$  normiramo Euklidskom normom  $\| \cdot \|_2$ :

$$z_{k+1} = \frac{M z_k}{\|M z_k\|_2}, \quad k = 0, 1, \dots$$

Idealno bi bilo kada bi mogli definirati vrijednosti  $h$  i  $a$  kao limes vektora  $z_k$  kada  $k \rightarrow \infty$ . S takvom definicijom postoje dva problema. Prvi problem je taj što limes u velikom broju slučajeva i ne postoji, već alternira između vrijednosti  $z_{parni} = \lim_{k \rightarrow \infty} z_{2k}$  i  $z_{neparni} = \lim_{k \rightarrow \infty} z_{2k+1}$ . Drugi problem jest odabir početne vrijednosti  $z_0$ . Naime, vrijednosti  $z_{parni}$  i  $z_{neparni}$  ovise o  $z_0$ , pa ćemo ih označavati sa  $z_{parni}(z)$  i  $z_{neparni}(z)$ .

U Kleinbergovom [13] članku odabire se  $z_0 = \mathbf{1}$ , gdje  $\mathbf{1}$  predstavlja vektor ili matricu sa svim elementima jednakim 1. Blondel [3] obrazlaže odabir  $z_0 = \mathbf{1}$  time što se tada mogu jednostavno izraziti  $z_{parni}(\mathbf{1})$  i  $z_{neparni}(\mathbf{1})$ , a  $z_{parni}(\mathbf{1})$  je u tom slučaju najveći vektor od svih  $z_{parni}(z)$  i  $z_{neparni}(z)$  gledano po 1-normi.

Formalni dokaz za postojanje  $z_{parni}(1)$  ne erno iznositi jer je detaljno opisan u Kleinbergovom i Blondelovom lanku, ali erno iznijeti skicu. Vektori u nizu  $z_k$  imaju smjerove (ne nužno i magnitude)

$$x_k, Mx_k, M^2x_k, \dots$$

ili

$$\begin{bmatrix} h \\ a \end{bmatrix}_k, \begin{bmatrix} 0 & B \\ B^T & 0 \end{bmatrix} \begin{bmatrix} h \\ a \end{bmatrix}_k, \begin{bmatrix} BB^T & 0 \\ 0 & B^TB \end{bmatrix} \begin{bmatrix} h \\ a \end{bmatrix}_k, \dots$$

iz ega vidimo da je:

$$\begin{bmatrix} h \\ a \end{bmatrix}_{k+2} = t \begin{bmatrix} BB^T & 0 \\ 0 & B^TB \end{bmatrix} \begin{bmatrix} h \\ a \end{bmatrix}_k.$$

Sa  $t$  ozna avamo faktor koji uzima u obzir magnitudu, tj. normalizaciju vektora. Raspišemo li jednadžbe, dobijemo:

$$h_{k+2} = tBB^T h_k, \quad a_{k+2} = tB^TB a_k,$$

te,

$$h_{2k} = (BB^T)^k t \mathbf{1}, \quad a_{2k} = (B^TB)^k t \mathbf{1}.$$

Ako je neka matrica  $M$  simetri na, a svojstveni vektor  $w$  koji odgovara dominantnoj (po apsolutnoj vrijednosti najve ojoj) svojstvenoj vrijednosti matrice  $M$  nije okomit na neki vektor  $v$ , onda  $\lim_{k \rightarrow \infty} M^k v = w$  [10]. Tako er, ako  $M$  nema negativnih vrijednosti, nema ih ni  $w$ .

Ako znamo da je  $v = t \mathbf{1}$  i da su matrice  $BB^T$  i  $B^TB$  simetri ne i nemaju negativnih elemenata (jer ni matrica susjedstva  $B$  nema negativnih elemenata) onda vrijedi  $w \cdot t \mathbf{1} > 0$  (jer je suma komponenti vektora  $w$  pozitivna). Zato vrijedi i da su  $\lim_{k \rightarrow \infty} (BB^T)^k t \mathbf{1}$  i  $\lim_{k \rightarrow \infty} (B^TB)^k t \mathbf{1}$  dominantni svojstveni vektori matrica  $BB^T$  i  $B^TB$ .

Za ra unanje egzaktnih vrijednosti  $a_i$  i  $h_i$  svake stranice mogu se na i dominantni svojstveni vektori matrica  $BB^T$  i  $B^TB$ . Budu i da za potrebe rangiranja web-stranica nisu potrebne egzaktne vrijednosti, ve samo aproksimacije, a vrijednosti  $a_i$  i  $h_i$  konvergiraju brzo, iterativni postupak eš e se koristi.

### 3.3. Sli nostme u vrhovima dva grafa

U daljnjem tekstu opisujemo kako se u Blondelovom lanku generalizira metoda vorišta i autoriteta.

Neka je  $G = (V, E)$  graf kojemu se ra unale vrijednosti  $a_i$  i  $h_i$ , tj. mjere koliko je neka stranica dobar autoritet ili dobro vorište za odre eni upit. Vrijednost  $a_i$  vrha  $i$  možemo poistovjetiti s mjerom sli nosti vrha  $i$  i vrha *autoritet* sljede eg grafa:

$$\text{čvorište} \rightarrow \text{autoritet.}$$

Sli no, vrijednost  $h_i$  možemo izjednati s mjerom sli nosti vrha  $i$  i vrha *vorš*. Za ra unanje sli nosti grafa  $G$  s grafom  $\text{čvorište} \rightarrow \text{autoritet}$  može se koristiti metoda vorišta i autoriteta.

Ako vrijednosti  $h_i$  i  $a_i$  ozna imo redom s  $x_{i1}$  i  $x_{i2}$  onda izraz 3.3 možemo zapisati ovako:

$$\begin{aligned} x_{i1} &\leftarrow \sum_{j: (i,j) \in E} x_{j2}, \\ x_{i2} &\leftarrow \sum_{j: (j,i) \in E} x_{j1}. \end{aligned} \quad (3.5)$$

Postavlja se pitanje mogu li se relacije 3.5 poop iti na složenije grafove od grafa oblika  $1 \rightarrow 2$ . Prvo emo pokušati pokazati kako bi se ra unala matrica sli nosti s grafom oblika  $1 \rightarrow 2 \rightarrow 3$ . Za ra unanje sli nosti vrhova grafa  $G$  i vrhova 1 i 3 grafa  $1 \rightarrow 2 \rightarrow 3$  relacije e biti sli ne relaciji 3.5. Nešto druk ija e biti relacija za ra unanje sli nosti s vrhom 2:

$$\begin{aligned} x_{i1} &\leftarrow \sum_{j: (i,j) \in E} x_{j2}, \\ x_{i2} &\leftarrow \sum_{j: (j,i) \in E} x_{j1} + \sum_{j: (i,j) \in E} x_{j3}, \end{aligned} \quad (3.6)$$

$$x_{i3} \leftarrow \sum_{j: (j,i) \in E} x_{j2}.$$

Po uzoru na izraz 3.4 pišemo:

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{k+1} = \begin{bmatrix} 0 & B & 0 \\ B^T & 0 & B \\ 0 & B^T & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_k, \quad k = 0, 1, \dots \quad (3.7)$$

Sva promatranja vezana uz grani ne vrijednosti koja su vrijedila za graf  $1 \rightarrow 2$  vrijede i za graf  $1 \rightarrow 2 \rightarrow 3$ . Za po etnu vrijednost odabire se  $[1 \ 1 \ 1]^T$ . Sli nost grafa  $G$  s grafom  $1 \rightarrow 2 \rightarrow 3$  definiramo kao grani nu vrijednost za parne elemente niza iz relacije 3.7.

U op em slu aju promatramo grafove  $G_A = (V_A, E_A)$  i  $G_B = (V_B, E_B)$ , uz  $n_A = |V_A|$  i  $n_B = |V_B|$ , te ra unamo matricu sli nosti  $X = (x_{ij})$  dimenzija  $n_B \times n_A$  na sljede i na in:

$$x_{ij} \leftarrow \sum_{\substack{r: (r,i) \in E_B \\ s: (s,j) \in E_A}} x_{rs} + \sum_{\substack{r: (i,r) \in E_B \\ s: (j,s) \in E_A}} x_{rs}, \quad \begin{array}{l} i \in \{1, \dots, n_B\} \\ j \in \{1, \dots, n_A\} \end{array} \quad (3.8)$$

Operacija nalaženja matrice sli nosti može se promatrati kao produkt grafova  $G_A$  i  $G_B$  iji je rezultat graf  $G_{A \times B} = (V_{A \times B}, E_{A \times B})$ , gdje je

$$\begin{aligned} V_{A \times B} &= \{(i, j) : i \in V_A \wedge j \in V_B\}, \\ E_{A \times B} &= \{((i_1, j_1), (i_2, j_2)) : (i_1, i_2) \in E_A \wedge (j_1, j_2) \in E_B\}. \end{aligned}$$

U takvom grafu vrijednost svakog vrha jednaka je sumi vrijednosti svih susjednih vrhova. Vrijednost vrha  $(i, j)$  novonastalog grafa jednaka je elementu matrice sli nosti u retku  $i$  i stupcu  $j$ .

Promatraju i izraz 3.7 vidimo da se izraz 3.8 može napisati u kra em obliku:

$$X_{k+1} = BX_k A^T + B^T X_k A, \quad (3.9)$$

gdje su  $A$  i  $B$  matrice susjedstva grafova  $G_A$  i  $G_B$ .

Izraz u 3.9 konvergira za parne  $k$  uz  $X_0 = \mathbf{1}$  (prisjetimo se,  $\mathbf{1}$  nije jedini na matrica ve matrica ili vektor  $i$ ji su svi elementi jednaki 1). Dokaz ove tvrdnje može se na  $i$  u Blondelovom lanku [3],.

INTERNI DOKUMENT

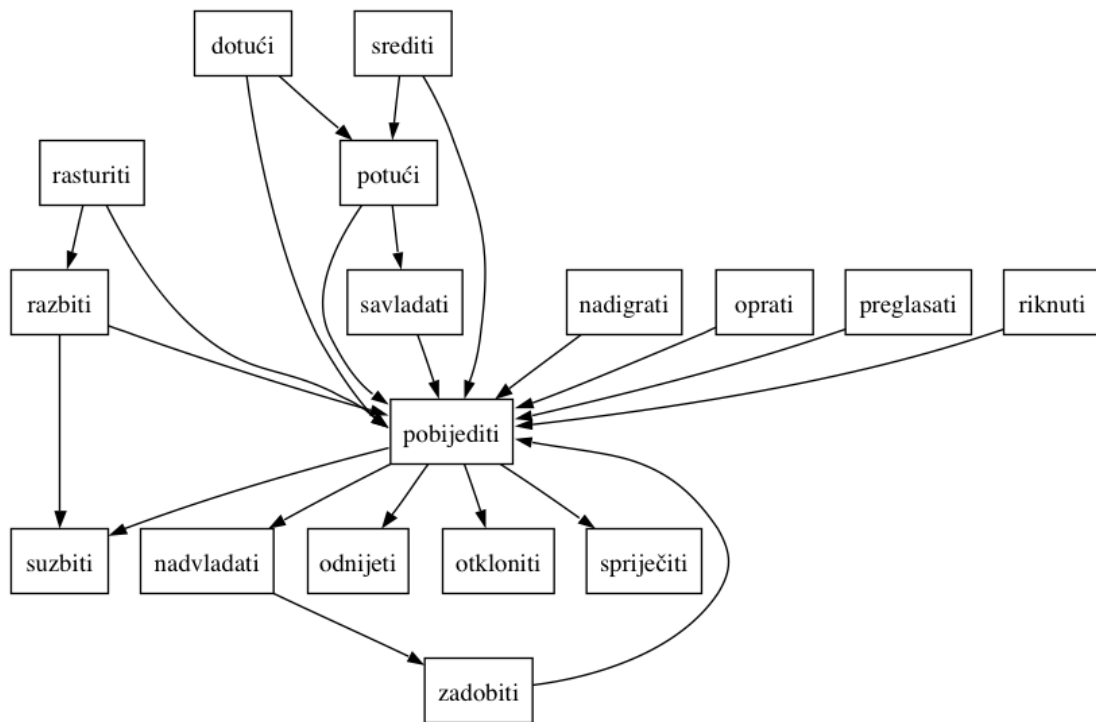
## 4. Nalaženje sinonima u rječniku

### 4.1. Primjena matrice sličnosti

Možemo se prisjetiti da se u metodi *voršta* i *autoriteta* izdvajao onaj podgraf grafa povezanosti svih web stranica koji je sadržavao prvih 200 rezultata korisničke pretrage. Zatim se taj podgraf proširio koristeći i ulazne i izlazne veze od po etnih 200 stranica. Na kraju se tražila sličnost vrhova dobivenog podgrafova s vrhom 2. Ta sličnost predstavljala je mjeru autoritativnosti web stranice.

Metoda nalaženja sinonima u rječniku koju opisuje Blondel [4] koristi englesko-engleski rječnik. Glavna pretpostavka u njegovom radu jest da sinonimi isto imaju zajedničke riječi u svom opisu, te da se sinonimi isto koriste u opisu istih riječi.

Prvi korak algoritma nalaženja sinonima jest konstruiranje **grafa rječnika**  $G_R$  čiji su vrhovi sve riječi u rječniku. Postoji usmjereni brid od riječi  $a$  do riječi  $b$  u tom grafu samo ako se riječ  $b$  pojavljuje u opisu riječi  $a$ . Po uzoru na postupak koji opisuje Kleinberg [13] drugi korak algoritma jest konstrukcija grafa  $G$ , podgrafova grafa  $G_R$  za zadanu riječ  $w$  (upit). Graf  $G$  sadrži samo one vrhove iz grafa  $G_R$  koji su prvi susjedi vrha  $w$  i sam vrh  $w$ . Graf  $G$  nazivamo **graf susjedstva** riječi  $w$ .



**Slika 1.** Graf susjedstva rije i *pobijediti*.

Rije i u grafu susjedstva trebamo poredati po važnosti, tj. za svaku rije trebamo procijeniti koliko je ta rije dobar sinonim rije i  $w$ . Razumna ideja jest da rangiramo vrhove prema tome koliko su oni dobra vorišta ili autoriteti. Drugim rije ima, mogli bismo izraziti sličnost grafa  $G$  sa grafom  $1 \rightarrow 2$ , te koristiti sličnost sa vrhom 1 ili sličnost sa vrhom 2 grafa  $1 \rightarrow 2$ . Time bi postupak bio gotovo identičan postupku rangiranja web stranica metodom vorišta i autoriteta. Ipak, rezultati dobiveni takvom metodom nisu zadovoljavajući. Razlog tome je što pravi sinonimi u ovako konstruiranom grafu susjedstva rije i  $w$  ne odgovaraju ni pojmu vorišta ni pojmu autoriteta. Na primjeru grafa susjedstva rije i *pobijediti* (slika 1), rije *suzbiti* ima relativno visok rezultat kao autoritet, a rije i *rasturiti*, *dotući*, *srediti* relativno visok rezultat kao vorišta. Ipak, o čemu je da te rije i nisu sinonimi jer nemaju ni jednu izlaznu ili pak nemaju ni jednu ulaznu vezu.

U Blondelovom članku [4] predlaže se umjesto mjerenja sličnosti grafa  $G$  sa vrhovima grafa  $1 \rightarrow 2$  mjerenje sličnosti grafa  $G$  sa vrhom 2 grafa  $1 \rightarrow 2 \rightarrow 3$ . Naime, po etnoj pretpostavka bila je da se više sinonima nalazi

u opisu zajedničke riječi i da sinonimi u svom opisu imaju zajedničke riječi. Mjere sličnosti vrhova grafa  $G$  s vrhom 2 zapravo mjerimo koliko strukture tih riječi odgovaraju strukturi tipičnih sinonima. Usporedimo li međusobno te vrijednosti, redovito će riječ  $w$  biti najbolji sinonim (tj. riječ  $e$  sama sebi biti vrlo dobar sinonim).

## 4.2. Homografi i polisemi

Homografi su riječi koje se isto pišu, ali imaju različite značenja. U rječnicima se obično oni homografi koji su različite vrste riječi smatraju zasebnim pojmovima, tj. odvojeno su definirani. Npr. riječ „lije ni ki“ može biti pridjev (koji se odnosi na liječnike) ili prilog (kao liječnik, na način liječnika). Polisemi su riječi koje se isto pišu, imaju različita značenja, ali dijele istu osnovu. Npr. riječ „brod“ je imenica muškog roda, ali može predstavljati plovilo ili uzdužni dio crkve. Svi polisemi obično su navedeni u definiciji jednog pojma. U rječnicima je uz svaku definiciju pojma navedena i vrsta riječi (imenica, zamjenica, itd.).

Za pojavnicu  $B$  koja se nalazi u opisu riječi  $A$ , ne možemo uvijek odrediti o kojoj vrsti riječi se radi. To se događa u slučajevima ako postoje riječi i druge vrste koje se pišu kao i riječ  $B$  (tj. homografi su). Prilikom konstrukcije grafa riječnika u idealnom slučaju povezali bi smo  $A$  s  $B$  samo ako su  $A$  i  $B$  iste vrste. Nema smisla za sinonim imenice kao rezultat dati glagol.

Problem možemo riješiti tako da provedemo pojednostavljenja navedena u Blondelovom članku [4]:

opisi svih homografa i polisema spojeni su u jedan opis i promatraju se kao jedna riječ ;

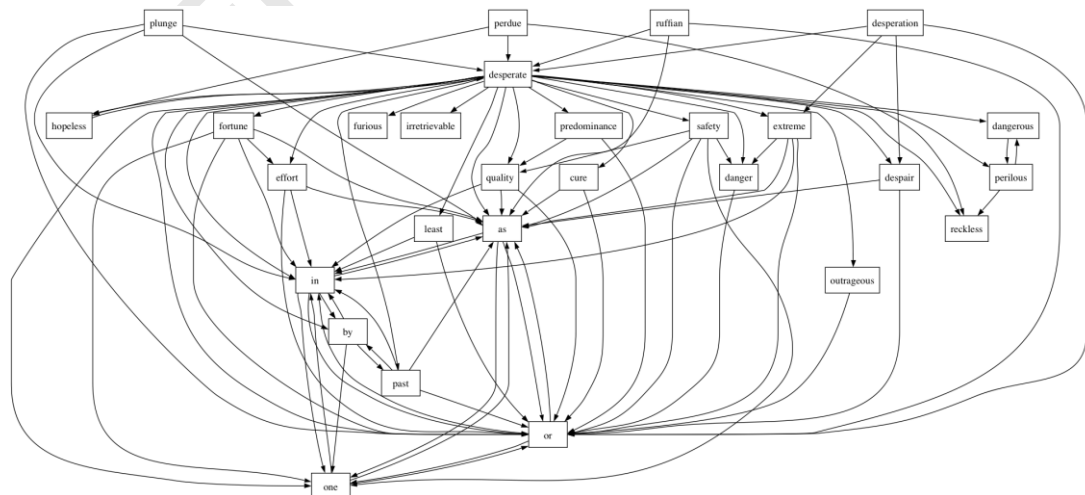
reći ćemo da je vrsta novonastale riječi skup vrsta svih riječi od kojih je nastala;

kako nema smisla tražiti sinonime među riječima ima različite vrste, prilikom konstrukcije grafa povezanosti riječi spojiti ćemo dvije riječi i samo ako

njihove vrste imaju neprazan presjek (tj. riječi od kojih se sastoje imaju barem jednu zajedničku vrstu).

### 4.3. Problem zaustavnih riječi

Prilikom konstrukcije grafa susjedstva riječi često se događa da riječ  $w$  ima vezu prema nekoj zaustavnoj riječi (npr. engleske riječi *of*, *and*, *by*, itd.). To samo po sebi ne bi trebao biti problem, jer je pretpostavka da će takve zaustavne riječi imati samo ulazne veze, a ne i izlazne veze, budući da je mala vjerojatnost da određena zaustavna riječ ima vezu sa sinonimom kojeg tražimo. Zbog provjere da sinonimi moraju biti riječi iste vrste ili da im se barem vrste moraju preklapati, nameće se zaključak da ne postoje veze od zaustavnih riječi ili prema njima. Međutim, takve se veze ipak javljaju. Primjerice, riječ *by* može biti prijedlog, pridjev ili prilog, riječ *in* može biti prijedlog, prilog ili imenica, itd. Riječi *in* i *by* međusobno se spominjati u opisima i zato će biti vrlo slične vrhu 2 grafa 1 → 2 → 3. Stoga se te zaustavne riječi initi kao idealni sinonimi bilo koje riječi (vidi sliku 2).



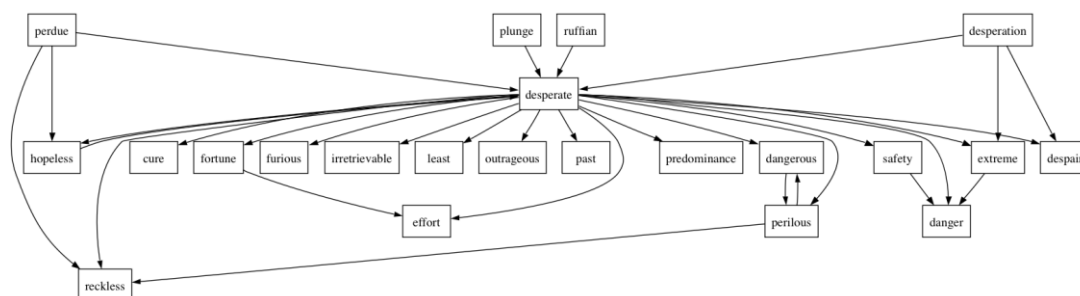
**Slika 2.** Graf susjedstva riječi *desperate* bez filtriranja zaustavnih riječi. *On*, *of*, *by*, *and* i *one* imaju jako puno ulaznih i izlaznih veza, te znatno pogoršavaju rad metode.

Ovaj problem može se jednostavno riješiti na na in da se iz grafa izbace sve rije i koje se pojavljuju u više od  $k$  opisa. U slu aju engleskog rje nika odabran je  $k = 1000$ .

Zaustavna rije	Broj pojavljivanja
of	67337
a	47188
the	43402
or	40835
to	31441
in	23600
as	22245
and	16491
an	13705
by	12061
one	11185
with	10753
which	10326
see	8422
is	8350

**Tablica 1.** Popis naj eš ih zaustavnih rije i u engleskom rje niku.

Nakon što su izba ene zaustavne rije i, grafovi susjedstva postaju znatno jednostavniji i upotrebljiviji (slika 3).



**Slika 3.** Graf susjedstva rije i *desperate* znatno je jednostavniji nakon što su izba ene zaustavne rije i.

Negativan utjecaj zaustavne riječi ovisi o načinu sastavljanja riječi, a ponajviše o samom jeziku, pa se o tome i ovisi izbor *k*.

INTERNI DOKUMENT

## 5. Primjena na indeksiranje dokumenata

Indeksiranje dokumenata jest postupak dodjeljivanja jedne ili više oznaka ili pojmova svakom dokumentu. Skup pojmova kojima se dokumenti označuju, tj. indeksiraju, unaprijed je određen. Ako su dokumenti kvalitetno indeksirani, pretraživanje dokumenata može biti znatno lakše. Umjesto da se pretraga vrši nad cijelim skupom dokumenata, može se ograničiti samo na manji dio koji je označen pojmom vezanim uz pretragu. Indeksiranje se može vršiti ručno, poluautomatski ili automatski. U prvom slučaju uvijek sam odabire pojmove za koje misli da najbolje opisuju dokument. U slučaju automatskog indeksiranja računalni program sam odabire sve pojmove koji odgovaraju nekom dokumentu, a kod poluautomatskog indeksiranja računalno samo predlaže skraćenu listu pojmova koja koristi uvijek da brže indeksira dokument. Samo indeksiranje može se svesti na problem klasifikacije, tj. svrstavanje dokumenta u jednu ili više kategorija.

Uspješnost metoda koje poboljšavaju kvalitetu pretrage kod Internet tražilica uvelike je ovisila o strukturi grafa koji opisuje veze između web stranica. U mnogim područjima koja nisu vezana uz pretragu web-a pokušavaju se iskoristiti slične metode tako da se pokuša izgraditi struktura slična web-u. Na primjer, za automatsko generiranje sažetaka radova razvijen je algoritam LexRank [9] po uzoru na algoritam PageRank. Dokumenti su razdijeljeni u manje segmente, a slični segmenti su međusobno povezani. Metodom LexRank mjeri se korisnost svake rečenice za uključivanje u sažetak. Druga primjena je rješavanje višeznačnosti riječi uz korištenje PageRank algoritma, koju je opisao Agirre [1]. Sve riječi koje su se pojavljivale uz višeznačnicu u korpusu ušle su u graf. Dvije riječi bile su povezane ako su se pojavile zajedno u grafu. Na kraju, prikazana je i primjena u slučaju traženja sinonima zadane riječi. Ono što je svojstveno svakoj od navedenih primjena jest inovativan način na koji su se konstruirali grafovi.

Treba primijetiti da se ni u jednoj primjeni nije koristilo znanje naučeno na primjerima. Metoda za rangiranje stranica nije iskoristila podatke nekih

korisnika da im ta stranica za određeni upit odgovara ili ne odgovara. Sli no tome, metoda nalaženja sinonima ne može unaprijed iskoristiti primjere za u enje te pokušati iz njih generalizirati nau eno. To nas navodi na sljede i zaklju ak: ako bi za svrhu klasifikacije konstruirali graf koji na neki na in ima ugra ene informacije dobivene od rezultata klasifikacije dobivenih na primjerima za u enje, te informacije ne bismo mogli iskoristiti tijekom rada na skupu za testiranje. Zato ostaje otvoreno pitanje kako iskoristiti podatke za u enje ako primjenjujemo metode analize grafova.

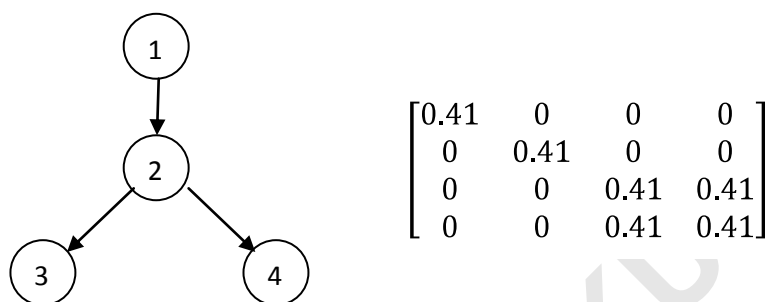
### *5.1. Nalaženje sli nih dokumenata*

Jedan od na ina za iskorištavanje rezultata indeksiranja dobivenih na primjerima za u enje jest da novi dokument usporedimo s ve indeksiranim dokumentima i iskoristimo pojmove kojima su indeksirani najslji niji dokumenti. Time svodimo problem na problem traženja sli nih dokumenata. Odabir strukture grafa kojim bi to ostvarili otvoren je problem.

Kako sada na raspolaganju imamo operaciju nalaženja sli nosti me u vrhovima dva grafa, logi no rješenje bilo bi da konstruiramo grafove na na in da usporedbom sli nosti vrhova dva grafa dobijemo mjeru sli nosti me u dokumentima koje predstavljaju ti vrhovi. Barem jedan od tih grafova trebao bi imati ugra ene informacije o nekim svojstvima tog dokumenta. Svojstvo dokumenta mogu biti, primjerice, rije i koje se pojavljuju u dokumentu. Kako tekstni dokumenti mogu biti vrlo veliki, a njihov broj može biti tako er velik, bilo bi teško analizirati dobiveni graf radi njegove veli ine i vremena potrebnog za analizu. Jedno od svojstava dokumenta jesu pojavljivanja pojmova kojima pokušavamo indeksirati te iste dokumente, tako da njih možemo koristiti pri izgradnji grafa.

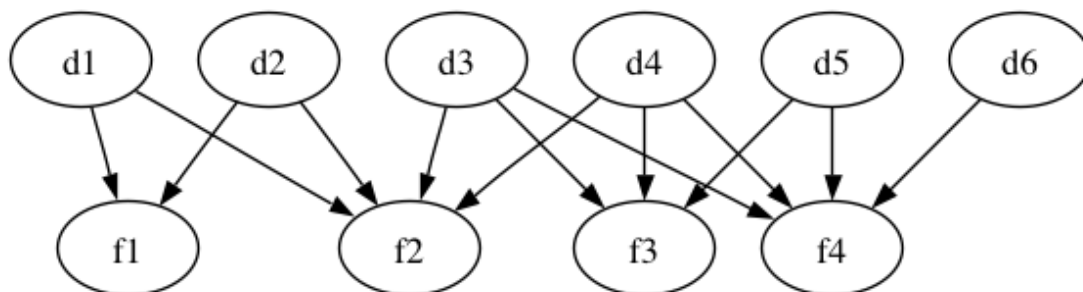
Operacija nalaženja sli nosti me u grafovima funkcionira i na primjeru jednog grafa, tj. kad tražimo samosli nost nekog grafa [3]. Na slici 4 vidimo primjer matrice samosli nosti. Izgleda da je dovoljno da izgradimo graf kojemu e

neki vorovi predstavljati dokumente, a drugi vorovi svojstva tih dokumenata. Kao mogu a svojstva dokumenata možemo uzeti sve pojmove kojima indeksiramo dokumente. Ako se u tekstu dokumenta pojavljuju pojmovi kojima ina e indeksiramo dokumente, onda možemo u grafu povezati dokument s vrhovima koji predstavljaju te pojmove. Sli nost me u dokumentima nalazit emo ra unaju i sli nost me u vrhovima koji predstavljaju te dokumente.



**Slika 4.** Matrica sli nosti grafa sa samim sobom.

Jedan graf koji zadovoljava ta svojstva je graf prikazan na slici 5.



**Slika 5.** Graf u kojemu je svaki od šest dokumenata spojen sa vrhovima koji predstavljaju svojstva tih dokumenata. Vidimo da su dokumenti  $d1$  i  $d2$ , te  $d3$  i  $d4$  sli ni.

Matricu sli nosti ra unamo iterativnim postupkom, a njena veli ina može biti npr. 5000 (broj dokumenata) s 10000 (broj svojstava). Matrica susjedstva je tada veli ine 15000 s 15000. Zato trebamo pojednostaviti izraze za ra unanje matrice sli nosti za ovaj slu aj. Ako s  $A$  ozna imo matricu susjedstva i uvrstimo u izraz 3.9 za ra unanje matrice sli nosti, dobijemo sljede i izraz:

$$X_{k+1} = AX_kA^T + A^T X_k A. \quad (5.1)$$

Matrica  $A$  ima sljedeći oblik:

$$A = \begin{bmatrix} \mathbf{0} & B \\ \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad (5.2)$$

gdje je  $B$  matrica dimenzija  $M \times N$ ,  $M$  broj dokumenata, a  $N$  broj svojstvenih vrhova.  $\mathbf{0}$  označava nul-matricu ili vektor. Uvrstimo 5.2 u 5.1, te uvrstimo vrijednost za  $X_{k=0} = \mathbf{1}_{M \times N}$ .  $\mathbf{1}_{M \times N}$  je matrica dimenzija  $M \times N$  u kojoj su svi elementi jednaki 1. Dobijemo sljedeći izraz:

$$X_1 = \begin{bmatrix} \mathbf{0} & B \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{1}_{M \times M} & \mathbf{1}_{M \times N} \\ \mathbf{1}_{N \times M} & \mathbf{1}_{N \times N} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ B^T & \mathbf{0} \end{bmatrix} + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ B^T & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{1}_{M \times M} & \mathbf{1}_{M \times N} \\ \mathbf{1}_{N \times M} & \mathbf{1}_{N \times N} \end{bmatrix} \begin{bmatrix} \mathbf{0} & B \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

Nakon sreivanja:

$$X_1 = \begin{bmatrix} B\mathbf{1}_{N \times N}B^T & \mathbf{0} \\ \mathbf{0} & B^T\mathbf{1}_{M \times M}B \end{bmatrix}. \quad (5.3)$$

Radujemo  $X_2$  uvrštavanjem 5.3 u 5.1:

$$X_2 = \begin{bmatrix} BB^T\mathbf{1}_{M \times M}BB^T & \mathbf{0} \\ \mathbf{0} & B^TB\mathbf{1}_{N \times N}B^TB \end{bmatrix}.$$

Sada je očit da je:

$$X_{2k} = \begin{bmatrix} (BB^T)^k\mathbf{1}_{M \times M}(BB^T)^k & \mathbf{0} \\ \mathbf{0} & (B^TB)^k\mathbf{1}_{N \times N}(B^TB)^k \end{bmatrix}.$$

Nas ne zanima sličnost me u dokumentima i svojstvenim vrhovima, već samo me u dokumentima. Zato nam je potreban samo gornji lijevi dio matrice. Dobivamo  $S$ , matricu sličnosti me u dokumentima:

$$S_k = (BB^T)^k\mathbf{1}_{M \times M}(BB^T)^k.$$

Bez dokaza navodimo da je  $\lim_k (S_k)^k = ww^T$ , gdje je  $w$  dominantni svojstveni vektor matrice  $BB^T$ . Takvu matricu koristimo za usporedbu dva dokumenta. Znači:

$$\lim_{k \rightarrow \infty} X_{2k} = \begin{bmatrix} ww^T & \mathbf{0} \\ \mathbf{0} & qq^T \end{bmatrix}, \quad (5.4)$$

gdje je  $w$  dominantni svojstveni vektor matrice  $BB^T$ , a  $q$  dominantni svojstveni vektor matrice  $B^TB$ . Budući da ćemo koristiti samo gornji lijevi dio matrice, tj. relaciju sličnosti među dokumentima, trebamo izraziti samo vektor  $w$ .

Sada možemo iskoristiti metodu poput metode traženja najbližih susjeda [8], gdje se nekom dokumentu pronađu i određenu metriku, k najbližih susjeda. Dokument se svrstava u kategoriju u koju je svrstano najviše od tih k dokumenata. Da bi ocijenili rad metode koristimo sljedeći vrlo jednostavan sustav.

## 5.2. Osnovni sustav

Neka je  $P = \{p_1, p_2, \dots, p_N\}$  skup pojmova kojim indeksiramo dokumente, a  $D = \{d_1, d_2, \dots, d_N\}$  skup dokumenata. Nađimo matricu  $A_{M \times N} = (a_{ij})$ , u kojoj je element  $a_{ij}$  jednak broju pojavljivanja pojma  $p_j$  u dokumentu  $d_i$ . Neka su  $a_{it_1}, a_{it_2}, \dots, a_{it_k}$  k najvećih elemenata  $i$ -tog retka matrice  $A$ , gdje je  $k$  neka unaprijed utvrđena konstanta. Tada ćemo indeksirati dokument  $d_i$  pojmovima  $p_{t_1}, p_{t_2}, \dots, p_{t_k}$ . Očekujemo da ćemo postići poboljšanje ako umjesto broja ponavljanja pojma u dokumentu za element matrice  $a_{ij}$  upotrijebimo TFIDF mjeru [12] na elemente matrice  $A$ .

TFIDF računamo na sljedeći način:

$$TFIDF = TF \cdot IDF,$$

gdje je:

$TF$  = broj pojavljivanja pojma  $p$  u dokumentu  $d$

$$IDF = \log\left(\frac{\text{broj dokumenata}}{\text{broj dokumenata koji sadrže pojam } p}\right)$$

INTERNI DOKUMENT

## 6. Rezultati

### 6.1. Nalaženje sinonima

U ovom poglavlju prikazani su rezultati metode nalaženja sinonima, te rezultati eksperimenata s automatskim indeksiranjem deskriptora. Radi provjere valjanosti algoritma i boljeg shvaćanja rada, prvi cilj pri izradi programa je bio pokušaj repliciranja rezultata iz Blondelovog rada [4]. Nakon što je provjerena ispravnost implementacije algoritma, krenulo se u obavljanje eksperimenata s hrvatskim jezikom.

#### 6.1.1. Nalaženje sinonima u engleskom jeziku

Na temelju projekta „Project Gutenberg Text of Webster's Unabridged Dictionary“ koji je temeljen na rječniku „1913 US Webster's Unabridged Dictionary“ nastao je „Online Plain Text English Dictionary“ [16], rječnik koji koristimo pri nalaženju sinonima. Taj rječnik sastoji se od 27 HTML datoteka formatiranih na vrlo jednostavan način: svaka definicija u rječniku pojavljuje se u zasebnom redu datoteke. Definicija, vrsta i opis riječi i omeđeni su uvijek istim HTML tagovima. Homografi (riječi koje se jednako pišu, a imaju različita značenja) i polisemije (slučaj kada više oblika imaju istu osnovu, ali različita i nesrodna značenja) imaju zaseban unos, tj. navedeni su u odvojenim redovima. Rječnik ima ukupno oko 176000 unosa, odnosno riječi, ako brojimo zasebno homografije i polisemije.

Jedan od problema u ovom rječniku je taj što definicija svake riječi počinje velikim početnim slovom. Radi lakšeg stvaranja grafa svih riječi u rječniku, sve su riječi i opisi riječi prebačeni u mala slova. Određene kategorije riječi izbačene su iz rječnika:

definicije prefiksa ili sufiksa. Npr. „Aero-“ ili „-escence“,

fraze i drugi višerje ni unosi jer je teško odrediti kada određeni niz riječi u opisu predstavlja više zasebnih riječi ili jedan unos, te sve definicije koje sadrže znamenke.

Nakon spajanja opisa kako je opisano u poglavlju 4.2, broj unosa smanjio se na 111617 riječi (za usporedbu, Blondel je dobio 112169 riječi). Svaki spojeni unos označen je jednom ili više od ukupno 309 vrsta riječi. Ukupan broj pojava korišten u opisima svih riječi je otprilike 1957000.

### 6.1.2. Evaluacija

Za evaluaciju automatskog nalaženja sinonima obično se koriste dva pristupa. Jedan pristup evaluaciji jest automatska usporedba s rječnikom riječi nikom sinonima [24], npr. s rječnikom WordNet [1]. Problem kod tog pristupa jest što rječnikom izraženim riječima često nedostaju mnogi sinonimi. Zato je preciznost riječnika koji se ocjenjuje često biti preniska. Drugi pristup je ljudska evaluacija više različitih metoda (npr. [4]). Jasno je da su mane ovog pristupa vremenska zahtijevnost, te nekonzistentnost. Ponekad se koristi i hibridni pristup, tj. uz automatsku evaluaciju s rječnikom riječi nikom vrši se ljudska evaluacija nad podskupom riječi koji nije pokriven u rječnikom riječi nikom [17].

U Blondelovom članku [4] uspoređuje se pet metoda nalaženja sinonima: mjera udaljenosti, ArcRank [11], rječnik WordNet [15] i rječnik iz programa Microsoft Word 97. Svakom metodom nađeni su sinonimi za engleske riječi *disappear*, *sugar*, *parallelogram*, *science*. Deset osoba ocijenilo je svaku listu sinonima ocjenom 0-10. Lako se uočava propust ove metode: ispitane su samo četiri riječi iz rječnika od preko 100000 riječi.

### 6.1.3. Usporedba s implementacijom u Blondelovom lanku

Ponovno ćemo naći i sinonime za četiri riječi koje su navedene u Blondelovom radu [3]. To su riječi **disappear**, **parallelogram**, **sugar** i **science**. Kao razlog zašto su baš ove četiri riječi odabrane navodi se prva riječ ima razne sinonime, druga ih nema, treća ima mnogo značenja, a četvrta je vrlo česta i teško je uopće reći što je njen sinonim.

<i>Blondelov lank</i>	<i>Ovaj rad</i>
vanish	vanish
pass	pass
die	die
wear	faint
faint	cease
fade	dissipate
sail	evanesce
light	wear
dissipate	light
ease	ship

**Tablica 2.** Sinonimi riječi *disappear*.

<i>Blondelov lank</i>	<i>Ovaj rad</i>
square	square
rhomb	rhomb
parallel	parallel
figure	equal
prism	opposite
equal	figure
opposite	parallelepiped
angles	altitude
quadrilateral	quadrilateral
rectangle	rhomboid

**Tablica 3.** Sinonimi riječi *parallelogram*

<i>B bnde bv anak</i>	<i>Ovaj rad</i>
cane	cane
starch	sucrose
sucrose	starch
milk	molasses
sweet	juice
dextrose	milk
molasses	sweet
juice	dextrose
glucose	glucose
lactose	lactose

**Tablica 4.** Sinonimi rije i *sugar*.

<i>B bnde bv anak</i>	<i>Ovaj rad</i>
art	art
branch	branch
law	study
study	practice
practice	learning
natural	natural
knowledge	application
learning	skill
theory	theory
principle	history

**Tablica 5.** Sinonimi rije i *science*

Na primjeru tablice 2 do tablice 5 vidimo usporedbu implementacije algoritma u ovom radu i u Blondelovom radu. Za svaku rije na eno je deset najbolje rangiranih sinonima, s time da se sama rije za koju su se tražili sinonimi nije uzimala u obzir. Postoje male razlike u dobivenim listama sinonima, ali one se mogu pripisati malo druk ijoj konstrukciji grafa rije i, te injenicom da se u Blondelovom radu [4] transformiralo 309 kategorija rije i u kombinaciju 5 kategorija (detalji nisu navedeni). Unato manjim razlikama, možemo re i da su implementacije vrlo sli ne.

U inkovitost algoritma isprobana je i na nekim drugim rije ima. Primije en je jedan sustavni nedostatak ove metode: u listi sinonima esto se pojavljuju rije i suprotnog zna enja. Neobi no je što u svom radu Blondel [4] uop e ne spominje taj problem. U lanku je ak i prikazan graf susjedstva rije i *likely*, ali ispod slike piše da graf nije potpun (ru no su odstranjeni neki vrhovi). Primjerice, me u prvih 10 sinonima rije i *quick* i *likely* pojavljuju se rije i *slow* i *unlikely* koje imaju suprotno zna enje.

#### 6.1.4. Nalaženje sinonima u hrvatskom jeziku

Prilikom primjene opisane metode na rije i iz hrvatskog jezika javlja se veliki problem u izboru hrvatskog rje nika. Naime, metoda zahtijeva korištenje hrvatskog rje nika koji za svaku definiciju rije i ima opis te rije i na istom jeziku. Za razliku od engleskog jezika, besplatnih i javno dostupnih materijala za hrvatski jezik ima vrlo malo. Javno je dostupan hrvatski rje nik<sup>6</sup> koji ne sadrži ni opise ni vrste rije i. Dostupan je i englesko-hrvatsko-engleski rje nik<sup>7</sup>, pa bi mogli probati koristiti taj rje nik za prijevod hrvatske rije na englesku, a zatim za prijevod liste sinonima s engleskog jezika na hrvatski. Ta metoda bila bi vrlo neprecizna jer se u prijevodu svake hrvatske rije i nalazi nekoliko engleskih i obratno. Prijevodom samo u jednom smjeru (hrvatski u engleski) ve znatno opada preciznost.

U ovom trenutku jedini dostupan hrvatsko hrvatski rje nik u elektronskom izdanju je „Veliki rje nik hrvatskoga jezika“ Vladimira Anića [2], pa su svi eksperimenti ra eni korištenjem izdanja tog rje nika na CD-u. Rje nik se sastoji od približno 70000 natuknica. Homografi se nalaze u zasebnim natuknicama, a svi polisemi u jednoj. U rje niku je korištena 21 vrsta rije i<sup>8</sup>.

Za razliku od engleskog rje nika koji je korišten, u hrvatskom rje niku možemo razlikovati izme u rije i pisanih velikim i malim slovom. Opisi rije i nikada nisu pune re enice i ne po inju velikim slovom samo zbog po etka re enice. Iz svake rije i uklonjeni su naglasci – npr. rije nàst njen st zamijenjena je s rije i nastanjenost. Opisi rije i sadrže gramati ke natuknice, etimologije, frazeologije, te sintagme. Te dodatne informacije su uklonjene.

---

<sup>6</sup> <http://cvs.linux.hr/spell/ispell/>

<sup>7</sup> <http://web.math.hr/~igaly/EHrjecnik.htm>

<sup>8</sup> Vrsta rije i koja se spominje u tekstu ne odgovara nužno jednoj od 10 vrsta rije i u hrvatskom jeziku. U rje niku se ponekad pod vrstom rije i spominje „srednji rod“, što je specifi nije od vrste rije i „imenica“ u gramatici.

Nakon pretprocesiranja rije nika, tj. izbacivanja višerje njih ulaza, micanja natuknica koje su samo prefiksi i sufiksi, te spajanja svih homografa dobijemo oko 65000 rije i opisanih s oko 705000 rije i. Vidimo da engleski rije nik ima skoro dvostruko više rije i sa skoro tri puta dužim opisima. U slučaju hrvatskog jezika problem zaustavne rije i nije toliko izražen. Jedan od razloga je taj što ve ina zaustavnih rije i pripada vrsti rije i koja nije ista kao i vrsta rije i ije sinonime tražimo.

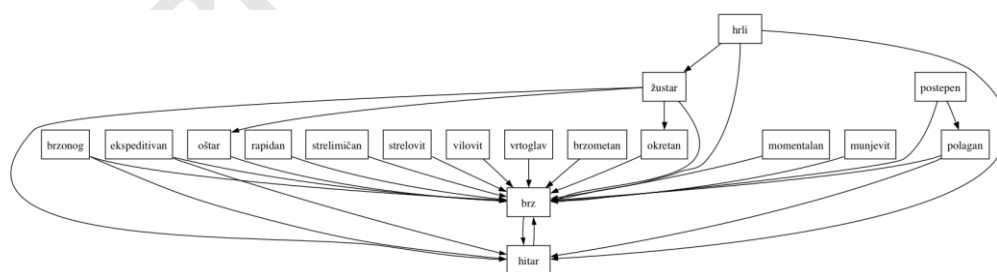
Automatsku evaluaciju u ovom trenutku gotovo je nemoguće napraviti jer ne postoji javno dostupan rije nik sinonima u elektronskom obliku. Zato ćemo samo navesti nekoliko primjera i proučiti njihove karakteristike.

INTERNI DOKUMENT

<i>brz</i>	<i>jama</i>	<i>nestati</i>	<i>ukrasti</i>	<i>ko ta</i>
hitar	rupa	umrijeti	zdipiti	kolotur
žustar	spilja	za i	maznuti	velosiped
polagan	vapnenica	pro i	odnijeti	zup anik
hrli	budža	izvjetriti	drpnuti	to ak
brzonog	rova a	zamaknuti	gepiti	kotur
ekspeditivan	golubinka	oti i	š apiti	elunik
okretan	kre ana	utonuti	otu iti	bucanj
oštar	fojba	utišati	dignuti	razmjer
postepen	graba	iš eznuti	mrknuti	senzor

**Tablica 6.** Deset predloženih sinonima za odabrane rije i

Za rije *brz* predloženi su i sinonim *polagan* i *postepen*, što nikako nije dobro. Ako pogledamo graf susjedstva na slici 6, vidimo da rije i *polagan* i *postepen* spominju rije *brz* u svom opisu. Kako rije *postepen* spominje i rije *polagan*, vrh rije i *polagan* postaje sli an vrhu 2 grafa  $1 \rightarrow 2 \rightarrow 3$ .



**Slika 6.** Graf susjedstva rije i *brz*.

Kod rije i *ko ta* primje uju se rije i koje nikako ne mogu biti sinonimi, jer je teško prona i deset sinonima za rije *ko ta* . Svi ostali prona eni sinonimi dobro su odabrani.

## 6.2. Automatsko indeksiranje dokumenata

Za ulazni skup korišten je skup dokumenata indeksiran pojmovnikom EUROVOC [5]. Pojmovnik EUROVOC koristi se prilikom indeksiranja službenih dokumenata institucija Republike Hrvatske koje provodi Hrvatska informacijsko-dokumentacijska referalna agencija (HIDRA). EUROVOC je hijerarhijski pojmovnik koji se sastoji od oko 6500 deskriptora (pojmovi) svrstanih u 21 područje i 127 podpodručja. Preveden je na 20 službenih jezika Europske Unije, a preveden je i na hrvatski jezik. HIDRA je dodala oko 3700 pojmova specifičnih za Hrvatsku. Primjer dijela EUROVOCa može se vidjeti na slici 7.



**Slika 7.** Mali dio pojmovnika EUROVOC. Pojmovi kojima se indeksiraju dokumenti označeni su tamnom bojom.

Skup dokumenata na kojem je isprobavano indeksiranje sastoji se od 4556 službenih dokumenata Republike Hrvatske iz Narodnih novina. Dokumente je indeksirala HIDRA ručnim postupkom. U skupu je oko 3300

dokumenata indeksirano jednim pojmom, oko 1200 sa dva pojma, te svega 9 dokumenata sa tri pojma.

Za pronalaženje pojavljivanja pojmova u dokumentu korištena je metoda koju u svom radu opisuju Šari i suradnici [21]. Pojmovi se mogu pojavljivati u bilo kojem padežu zahvaljujući korištenju lematizacijske metode „Automatska morfološka normalizacija“ koju je razvio Šnajder [22].

Korištene su tri mjere za evaluaciju indeksiranja: preciznost, odaziv i F1. Definiramo ih na sljedeći način:

$$\text{preciznost} = \frac{\text{broj točnih pojmova koje je odabrao program}}{\text{ukupni broj pojmova koje je odabrao program}},$$

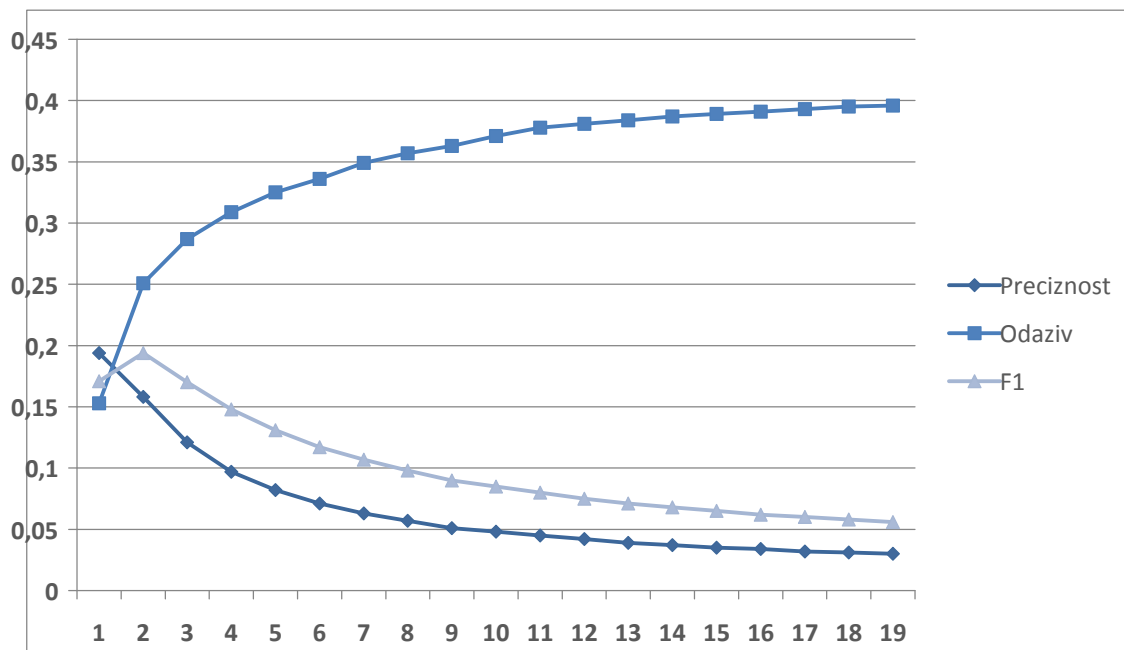
$$\text{odaziv} = \frac{\text{broj točnih pojmova koje je odabrao program}}{\text{ukupni broj pojmova koje su odabrali indeksatori}},$$

a F1 je mjera koja objedinjuje preciznost i odaziv u jednu i njihovu harmonijsku sredinu:

$$F1 = \frac{2}{\frac{1}{\text{preciznost}} + \frac{1}{\text{odaziv}}}.$$

### 6.2.1. Osnovna metoda

Za osnovnu metodu skup nije bilo potrebno dijeliti u skup za učenje i skup za testiranje jer ta metoda ne koristi već klasificirane primjere. Ovisno o  $k$ , broju odabranih pojmova po dokumentu dobiveni su rezultati za osnovnu metodu, prikazani na slici 8.



**Slika 8.** Rezultati osnovne metode ovisno o broju izabranih pojmova po dokumentu.

Za odabrani  $k = 2$  postiže se  $F1 = 0.194$ . Ako dodatno odbacimo sve pojmove s TFIDF vrijednoš u manjom od 0.3, dobijemo  $F1 = 0.218$ . Bez TFIDF mjere najbolji rezultat je  $F1 = 0.143$ .

### 6.2.2. Korištenje matrice sli nosti

Skup od 4556 dokumenata podijeljen je u omjeru na na in da je 70% slu ajno odabranih dokumenata pripalo skupu za u enje, a 30% skupu za testiranje. Za svaki dokument na eno je  $k$  njemu najslji njih dokumenata. Taj dokument indeksiran je svim pojmovima kojima je indeksirano preko 50% od  $k$  dokumenata. Sli nost dva dokumenta izra unata je metodom opisanom u poglavlju 5.1.

Najbolji rezultat dobio se korištenjem  $k = 2$  i iznosi  $F1 = 0.03$ . Tako loši rezultati bili su izvan svakog o ekivanja, pa je valjalo prona i uzrok za to. Numeri ki podatci su više puta provjereni i isprobani na znatno manjem

primjeru. Glavni razlog za loše performanse na n je u činjenici da operacija ranjanja matrice sli nosti ne daje intuitivne (o ekivane) rezultate.

Ime metode impliciralo bi da se tom operacijom stvarno ranja sli nost me u vrhovima dva grafa. Jedno od o ekivanja takvog razmišljanja bilo je da e nakon ranjanja sli nosti grafa sa samim sobom svakom vrhu najslji niji vrh biti upravo on sam. Iako se u Blondelovom radu [3] dokazuje tvrdnja da je najve i element u matrici sli nosti kod ranjanja samosli nosti upravo element na dijagonali, to ne zna i da su svi elementi na dijagonali najve i u svom retku (ili stupcu).

Ovakav rezultat mogao se predvidjeti malo boljim razmatranjem jednadžbe 5.4. Primijetimo da za ovako konstruiran graf cijela matrica sli nosti ovisi o dominantnim svojstvenim vektorima  $BB^T$  i  $B^TB$ . Drugim rije ima, cijela matrica sli nosti ovisi o  $M + N$  realnih vrijednosti ( $M$  i  $N$  su brojevi dokumenata, odnosno pojmova u pojmovniku), a o ekivali smo red veli ine  $M^2/2$  korisnih informacija iz te matrice. Jasno je da ovako definirana operacija sli nosti me u vrhovima grafova ne daje o ekivane rezultate. U prakti noj primjeni takva operacija jedino je bila korištena kod metode vorišta i autoriteta (iz koje je i generalizirana), te kod metode nalaženja sinonima.

## 7. Zaključak

U radu je proučena metoda nalaženja sličnosti među vrhovima dva grafa koja je nastala generalizacijom metode vorišta i autoriteta, poznate i u inkovite metode za rangiranje web-stranica. U inkovitost metoda koje koriste analizu grafova ovisi o dobro izgrađenoj strukturi grafa, u kojoj su veze među vrhovima najviše nositelj informacija.

Metoda nalaženja sinonima opisana u Blondelovom radu [3] iskorištava specifičnu strukturu riječnika, tj. gradi veze među riječima u riječniku na način koji podsjeća na Kleinbergovu [13] metodu rangiranja web stranica. Blondel je također uspješno modificirao Kleinbergovu metodu tako da daje bolje rezultate pri rangiranju sinonima. Ta metoda implementirana je u ovom radu i evaluirana na primjeru engleskog riječnika i na primjeru hrvatskog riječnika. Zanimljivo je i neke nedostatke, Blondelova modifikacija Kleinbergove metode pokazala se kao dobrom metodom nalaženja sinonima.

U radu je operacija traženja sličnosti među vrhovima istog grafa iskorištena za nalaženje sličnosti među dokumentima indeksiranim pojmovima pojmovnika EUROVOC. Izražena je i vrlo jednostavna osnovna metoda kako bi se bolje ocijenila u inkovitost metode temeljene na nalaženju matrice samosličnosti jednog grafa. Primjena nove metode na navedenom skupu dokumenata nije pokazala zadovoljavajuće rezultate. Ova metoda nalaženja sličnosti među grafovima morala bi se primijeniti na drukčije izgrađene grafove. Način kako bi se to postiglo radi bolje klasifikacije tekstualnih dokumenata ostaje otvoren problem.

## 8. Literatura

- [1] Agirre E; Martinez D; Lopez de Lacalle, O; Soroa A: *Two graph-based algorithms for state-of-the-art WSD*, Proc. of the 2006 Conference on Empirical Methods in Natural Language Processing (EMNLP 2006), pp. 585–593, Sydney, 2006.
- [2] Ani , V: *Veliki rječnik hrvatskoga jezika*, Novi Liber, Zagreb, 2004.
- [3] Blondel, V. D; Gajardo, A; Heymans, M; Senellart, P. P; Van Dooren, P: *A measure of similarity between graph vertices: Applications to synonym extraction and web searching*, SIAM Review, Vol. 46, No. 4. pp. 647-666, 2004.
- [4] Blondel, V. D; Senellart, P. P: *Automatic extraction of synonyms in a dictionary*, Technical Report 89, Université catholique de Louvain, Louvain-la-neuve, Belgium, 2001.
- [5] Bratani , M., ur: *Pojmovnik EUROVOC, drugo izdanje*, HIDRA, Zagreb, 2000.
- [6] Brin, S; Page, L: *The anatomy of a large-scale hypertextual Web search engine*, Proc. of the 7th International World Wide Web Conference, 1998.
- [7] Curran, J: *Ensemble Methods for Automatic Thesaurus Extraction*, In Proc. of the First International Workshop on Multiple Classifier Systems, pp. 1-15., 2000.
- [8] Dasarathy, B. V., ur: *Nearest Neighbor (NN) Norms: NN Pattern Classification Technique*, IEEE Computer Society Press, Los Alamitos, 1991.

- [9] Erkan, G; Radev, D. R: *Graph-based lexical centrality as salience in text summarization*, Journal of Artificial Intelligence Research, Vol 22, pp. 457-479, 2004.
- [10] Ivanšić, I: *Numerička matematika*, Element, Zagreb, 2002.
- [11] Jannink, J; Wiederhold G; *Thesaurus entry extraction from an online dictionary*, Proc. of Fusion 1999. Sunnyvale CA, 1999.
- [12] Joachims, T: *A probabilistic analysis of the Rocchio algorithm TFIDF for text categorization*, Proceedings of the International Conference on Machine Learning, Nashville, 1997.
- [13] Kleinberg, J. M: *Authoritative sources in a hyperlinked environment*, Journal of the ACM, 46:5, pp. 604-632, 1999.
- [14] Mandala, R; Tokunaga, T; Tanaka, H: *Combining Multiple Evidence from Different Type of Thesaurus for Query Expansion*, In Proc. of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval, Berkeley, 1999.
- [15] Miller, G: *WordNet: A lexical database for English*, Communications of the ACM 38(11), pp. 39-41, 1995.
- [16] *The Online Plain Text English Dictionary*, <http://msowww.anu.edu.au/~ralph/OPTED/>, 2000.  
[1. 6. 2006.]
- [17] van der Plas, L; Tiedemann, J: *Finding synonyms using automatic word alignment and measures of distributional similarity*, Proc. of the COLING/ACL, pp. 866-873, Sydney, 2006

- [18] Plaunt, C; Norgard, B: *An association-based method for automatic indexing with a controlled vocabulary*, Journal of the American Society for Information Science, Vol. 49(10), 1998.
- [19] Pouliquen, B; Steinberger, R; Ignat, C: *Automatic annotation of multilingual text collections with a conceptual thesaurus*, Proc. of In: Proceedings of the Workshop Ontologies and Information Extraction at the Summer School The Semantic Web and Language Technology - Its Potential and Practicalities (EUROLAN'2003), Bucharest, 2003.
- [20] Radev, D; Qi, H; Zheng, Z; Goldensohn, S; Zhang, Z; Fang, W; Prager, J: *Mining the Web for Answers to Natural Language Questions*, In the 10th International ACM Conference on Information and Knowledge Management, Atlanta, 2001.
- [21] Šari , F; Šnajder, J; Dalbelo-Baši , B; Ekli , H: *Enhanced thesaurus terms extraction for document indexing*, Proc. of the 27th International Conference ITI, Cavtat, 2005.
- [22] Šnajder, J: *Rule-based automatic acquisition of large-coverage morphological lexicons for information retrieval*, Tech. Report, MZOŠ 2003-082, ZEMRIS, FER, University of Zagreb, 2005.
- [23] van der Vet, P. E; Mars, N. J. I: *Condorcet Query Engine: A Query Engine for Coordinated Index Terms*, Journal of the American Society for Information Science Vol. 50(6), pp. 485-492, 1999.
- [24] Wu, H; Zhou, M: *Optimizing synonym extraction using monolingual and bilingual resources*, Proc. of the Second International Workshop on Paraphrasing: Paraphrase Acquisition and Applications (IWP2003), Sapporo, Japan, 2003.

## 9. Dodatak

### 9.1. Popis tablica

Tablica 1. Popis naj eš ih zaustavnih rije i u engleskom rje niku.....	19
Tablica 2. Sinonimi rije i <i>dissapear</i> .....	28
Tablica 3. Sinonimi rije i <i>parallelogram</i> .....	28
Tablica 4. Sinonimi rije i <i>sugar</i> .....	28
Tablica 5. Sinonimi rije i <i>science</i> .....	28
Tablica 6. Deset predloženih sinonima za odabrane rije i.....	31

INTERNI DOKUMENT

## 9.2. Popis slika

Slika 1. Graf susjedstva rije i <i>pobijediti</i> .....	16
Slika 2. Graf susjedstva rije i <i>desperate</i> bez filtriranja zaustavnih rije i. <i>On, of, by, and</i> i <i>one</i> imaju jako puno ulaznih i izlaznih veza, te znatno pogoršavaju rad metode. ....	18
Slika 3. Graf susjedstva rije i <i>desperate</i> znatno je jednostavniji nakon što su izba ene zaustavne rije i. ....	19
Slika 5. Graf u kojemu je svaki od šest dokumenata spojen sa vrhovima koji predstavljaju svojstva tih dokumenata. Vidimo da su dokumenti <i>d1</i> i <i>d2</i> , te <i>d3</i> i <i>d4</i> sli ni. ....	22
Slika 4. Matrica sli nosti grafa sa samim sobom.....	22
Slika 6. Graf susjedstva rije i <i>brz</i> . ....	31
Slika 7. Mali dio pojmovnika EUROVOC. Pojmovi kojima se indeksiraju dokumenti ozna eni su tamnom bojom.....	32
Slika 8. Rezultati osnovne metode ovisno o broju izabranih pojmova po dokumentu. ....	34
ova po dokumentu. ....	34