

Simultano klasteriranje dokumenata i riječi

Petra Sočo, Jelena Zaninović

22. prosinca 2020.

- Klasteriranje skupa podataka je grupiranje elemenata na temelju sličnosti.
- Konkretni problem: klasteriranje dokumenata i riječi
- Dualnost između ta dva procesa \rightarrow ima ih smisla pokušati simultano provesti.
- Ideja: Riječi i dokumente organizirati kao bipartitni graf \rightarrow problem particioniranja bipartitnog grafa.
- "Algoritam": dekompozicija singularnih vrijednosti pripadajuće matrice i provođenje k-means algoritma na jednodimenzionalnom skupu podataka.

Model bipartitnog grafa

- Težinski, neusmjereni graf je uređeni par $G = (\mathcal{V}, E)$ zadan skupom vrhova $\mathcal{V} = \{1, 2, \dots, |\mathcal{V}|\}$ i bridova $\{i, j\}$, gdje svaki brid ima težinu E_{ij} .
- Definiramo **matricu susjedstva** M na sljedeći način:

$$M_{ij} = \begin{cases} E_{ij} & \exists \text{ brid } \{i, j\}, \\ 0 & \text{inače.} \end{cases}$$

- Za proizvoljnu 2-particiju \mathcal{V}_1 i \mathcal{V}_2 skupa \mathcal{V} definiramo rez particije:

$$\text{cut}(\mathcal{V}_1, \mathcal{V}_2) := \sum_{i \in \mathcal{V}_1, j \in \mathcal{V}_2} M_{ij}.$$

Prethodna definicija u slučaju k-particije:

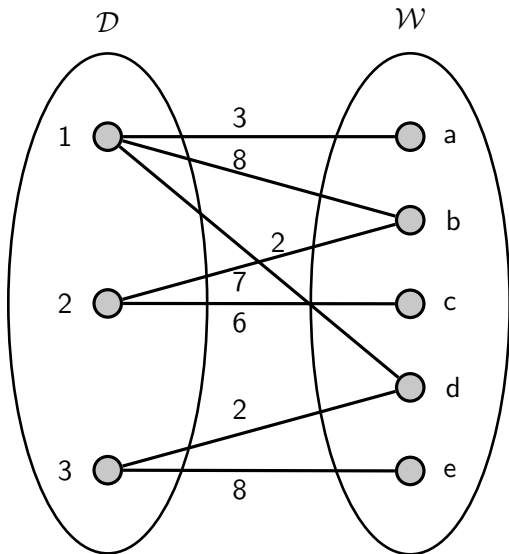
$$\text{cut}(\mathcal{V}_1, \dots, \mathcal{V}_k) := \sum_{i < j} M_{ij}.$$

- Neusmjereni, bipartitni graf je uređena trojka $G = (\mathcal{D}, \mathcal{W}, E)$ gdje su $\mathcal{D} = \{d_1, \dots, d_n\}$ i $\mathcal{W} = \{w_1, \dots, w_m\}$ dva skupa vrhova (dokumenti i riječi redom), a $E = \{\{d_i, w_j\} : d_i \in \mathcal{D}, w_j \in \mathcal{W}\}$ je skup bridova.
- Brid $\{d_i, w_j\}$ postoji ako se riječ w_j pojavljuje u dokumentu d_i .
- Jačina te veze izražena je pozitivnom težinom E_{ij} koju postavljamo na brid koji ih povezuje. Matrica susjedstva bipartitnog grafa M ima oblik

$$M = \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix},$$

gdje je A dimenzije $m \times n$ i vrijedi $A_{ij} = E_{ij}$.

- Prvih m redaka i stupaca odgovara riječima, a zadnjih n odgovara dokumentima.



- Klasteriranje riječi generira klastere dokumenata i obratno:

$$\{\mathcal{D}_1, \dots, \mathcal{D}_k\} \longleftrightarrow \{\mathcal{W}_1, \dots, \mathcal{W}_k\}$$

- Riječ w_i pripada klasteru \mathcal{W}_m ako je njena veza s klasterom \mathcal{D}_m jača od veza s ostalim klasterima. Klastere riječi \mathcal{W}_m možemo karakterizirati s:

$$\mathcal{W}_m = \{w_i : \sum_{j \in \mathcal{D}_m} A_{ij} \geq \sum_{j \in \mathcal{D}_l} A_{ij}, \forall l = 1, \dots, k\}, m = 1, \dots, k.$$

- Slično, iz klastera riječi $\mathcal{W}_1, \dots, \mathcal{W}_k$ može se dobiti:

$$\mathcal{D}_m = \{d_j : \sum_{i \in \mathcal{W}_m} A_{ij} \geq \sum_{i \in \mathcal{W}_l} A_{ij}, \forall l = 1, \dots, k\}, m = 1, \dots, k.$$

- Tražimo particiju grafa u kojoj će težine koje povezuju različite klastere biti što manje moguće \rightarrow minimizacija reza:

$$cut(\mathcal{W}_1 \cup \mathcal{D}_1, \dots, \mathcal{W}_k \cup \mathcal{D}_k) = \min_{\mathcal{V}_1, \dots, \mathcal{V}_k} cut(\mathcal{V}_1, \dots, \mathcal{V}_k),$$

gdje je $\mathcal{V}_1, \dots, \mathcal{V}_k$ neka k-particija grafa.

Uloga svojstvenog problema u particioniranju grafa

- Naći optimalnu 2-particiju grafa \longrightarrow naći podskupove \mathcal{V}_1^* i \mathcal{V}_2^* takve da je $cut(\mathcal{V}_1, \mathcal{V}_2)$ minimiziran.
- Neka je graf $G = (\mathcal{V}, E)$ zadan s n vrhova i m bridova kojima su pridružene težine E_{ij}
- Definiramo prvo $n \times m$ matricu incidencije I_G : stupac matrice I_G pridružen bridu $\{i, j\}$ ima na i -tom i j -tom mjestu $\sqrt{E_{ij}}$, $-\sqrt{E_{ij}}$ redom, dok su na ostalim mjestima nule.
- Laplaceova matrica L grafa G dana je s:

$$L_{ij} = \begin{cases} \sum_k E_{ik} & i = j, \\ -E_{ij} & i \neq j, \exists \text{ brid } \{i, j\}, \\ 0 & \text{inače.} \end{cases}$$

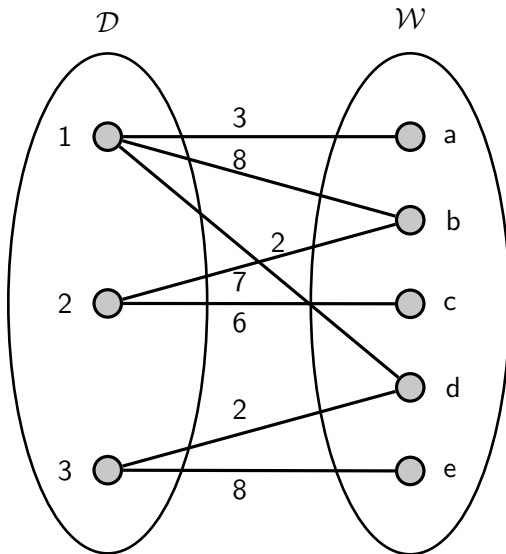
Teorem

Laplaceova matrica $L = L_G$ grafa G ima sljedeća svojstva:

1. $L = D - M$, gdje je M matrica incidencije, a D dijagonalna matrica takva da $D_{ii} = \sum_k E_{ik}$.
2. $L = |_G|_G^T$.
3. L je simetrična pozitivno semi-definitna matrica.
4. Neka je $e = [1, \dots, 1]^T$. Tada je $Le = 0$.
5. Ako graf G ima c komponenti povezanosti, tada L ima c svojstvenih vrijednosti jednakih 0.
6. Za proizvoljni vektor x , $x^T L x = \sum_{\{i,j\} \in E} E_{ij} (x_i - x_j)^2$.
7. Za proizvoljni vektor x i skalare α i β vrijedi

$$(\alpha x + \beta e)^T L (\alpha x + \beta e) = \alpha^2 x^T L x.$$

- 3. tvrdnja \implies sve svojstvene vrijednosti matrice L su realne i nenegativne
- 4. tvrdnja \implies $(0, e)$ je jedan svojstveni par matrice L .



$$A = \begin{bmatrix} 3 & 0 & 0 \\ 8 & 7 & 0 \\ 0 & 6 & 0 \\ 2 & 0 & 2 \\ 0 & 0 & 8 \end{bmatrix}, \quad M = \begin{bmatrix} 0 & A \\ A^T & 0 \end{bmatrix},$$

$$L = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 & -3 & 0 & 0 \\ 0 & 15 & 0 & 0 & 0 & -8 & -7 & 0 \\ 0 & 0 & 6 & 0 & 0 & 0 & -6 & 0 \\ 0 & 0 & 0 & 4 & 0 & -2 & 0 & -2 \\ 0 & 0 & 0 & 0 & 8 & 0 & 0 & -8 \\ -3 & -8 & 0 & -2 & 0 & 13 & 0 & 0 \\ 0 & -7 & -6 & 0 & 0 & 0 & 13 & 0 \\ 0 & 0 & 0 & -2 & -8 & 0 & 0 & 10 \end{bmatrix}, \quad D = \text{diag}(L).$$

- Nadalje pretpostavljamo da graf G ima jednu komponentu povezanosti.
- Za proizvoljnu 2-particiju $\mathcal{V}_1, \mathcal{V}_2$ definiramo vektor p kao vektor particije na sljedeći način

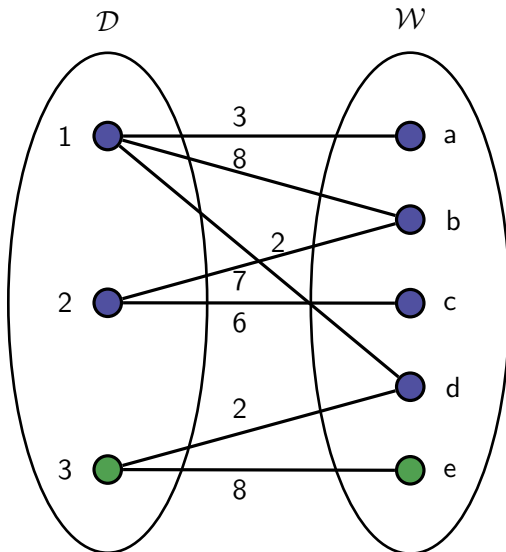
$$p_i := \begin{cases} +1 & i \in \mathcal{V}_1 \\ -1 & i \in \mathcal{V}_2. \end{cases}$$

Teorem

Neka je L Laplaceova matrica grafa G i p vektor particije. Tada za Rayleighov koeficijent vrijedi

$$\frac{p^T L p}{p^T p} = \frac{4}{n} \text{cut}(\mathcal{V}_1, \mathcal{V}_2).$$

Primjer particije



- Želimo "balansirane" particije!
- Svakom vrhu pridružujemo težinu $w(i)$ i definiramo dijagonalnu $n \times n$ matricu W sastavljenu od tih težina.
- Za podskup vrhova \mathcal{V}_I definiramo težinu podskupa kao:

$$w(\mathcal{V}_I) = \sum_{i \in \mathcal{V}_I} w(i) = \sum_{i \in \mathcal{V}_I} W_{ii}.$$

- 2-particija je "balansirana" ako su težine podskupova približno jednake \rightarrow nova funkcija cilja:

$$Q(\mathcal{V}_1, \mathcal{V}_2) = \frac{cut(\mathcal{V}_1, \mathcal{V}_2)}{w(\mathcal{V}_1)} + \frac{cut(\mathcal{V}_1, \mathcal{V}_2)}{w(\mathcal{V}_2)} \quad (1)$$

Lema

Neka su L i W Laplaceova i matrica težina vrhova nekog grafa G . Tada generalizirani vektor particije q definiran s

$$q_i := \begin{cases} +\sqrt{\frac{\eta_2}{\eta_1}} & i \in \mathcal{V}_1 \\ -\sqrt{\frac{\eta_1}{\eta_2}} & i \in \mathcal{V}_2. \end{cases}$$

zadovoljava $q^T W e = 0$ i $q^T W q = w(\mathcal{V})$. Gdje je $\eta_1 = w(\mathcal{V}_1)$ i $\eta_2 = w(\mathcal{V}_2)$.

Teorem

$$\frac{q^T L q}{q^T W q} = \frac{\text{cut}(\mathcal{V}_1, \mathcal{V}_2)}{w(\mathcal{V}_1)} + \frac{\text{cut}(\mathcal{V}_1, \mathcal{V}_2)}{w(\mathcal{V}_2)}$$

- Desna strana izraza u prethodnom teoremu je upravo funkcija koju želimo minimizirati (1) pa pažnju sada možemo usmjeriti na traženje "optimalnog" generaliziranog vektora particije q koji je dan prethodnom lemom.
- Sljedeći teorem omogućuje da diskretni problem nalaženja optimalnog generaliziranog vektora q zamijenimo kontinuiranim.

Teorem

Minimum problema

$$\min_{q \neq 0} \frac{q^T L q}{q^T W q}, \text{ uz uvjet } q^T W e = 0$$

se postiže za svojstveni vektor druge najmanje svojstvene vrijednosti generaliziranog svojstvenog problema

$$Lz = \lambda Wz. \quad (2)$$

Kako pridružiti težine vrhovima?

- $w(i) = 1$ za svaki vrh i pa funkcija iz (1) ima sljedeći oblik

$$Q(\mathcal{V}_1, \mathcal{V}_2) = \frac{cut(\mathcal{V}_1, \mathcal{V}_2)}{|\mathcal{V}_1|} + \frac{cut(\mathcal{V}_1, \mathcal{V}_2)}{|\mathcal{V}_2|}.$$

- Takvim odabirom težina ćemo za vrhove dobiti 2-particije koje će biti "balansirane" u smislu veličina.
- Drugi pristup je da težinu vrha mjerimo pomoću težina koje "izlaze" iz njega, tj. stavit ćemo $w(i) = \sum_k E_{ik}$.
- U tom slučaju će matrica W svojstvenog problema (2) biti jednaka matrici D iz pethodnog teorema.

Veza s dekompozicijom singularnih vrijednosti

- Tražimo drugi najmanji svojstveni vektor generaliziranog svojstvenog problema $Lz = \lambda Dz$.
- U slučaju bipartitnog grafa:

$$L = \begin{bmatrix} D_1 & -A \\ -A^T & D_2 \end{bmatrix}, \quad D = \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix}$$

gdje je $D_1(i, i) = \sum_j A_{i,j}$, $D_2(j, j) = \sum_i A_{i,j}$.

$$\Rightarrow \begin{bmatrix} D_1 & -A \\ -A^T & D_2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \lambda \begin{bmatrix} D_1 & 0 \\ 0 & D_2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \quad (3)$$

- Pretpostavimo da su D_1 i D_2 regularne

$$D_1^{1/2}x - D_1^{-1/2}Ay = \lambda D_1^{1/2}x,$$

$$-D_2^{-1/2}A^T x + D_2^{1/2}y = \lambda D_2^{1/2}y.$$

- Neka su $u = D_1^{1/2}x$ i $v = D_2^{1/2}y$

$$D_1^{-1/2}AD_2^{-1/2}v = (1 - \lambda)u,$$

$$D_2^{-1/2}A^T D_1^{-1/2}u = (1 - \lambda)v,$$

- Prethodne jednačbe definiraju dekompoziciju singularnih vrijednosti (SVD) normalizirane matrice $A_n = D_1^{-1/2} A D_2^{-1/2}$.
- Svojstveni vektor druge (najmanje) svojstvene vrijednosti iz (3) \rightarrow lijevi i desni singularni vektor koji pripadaju drugoj (najvećoj) singularnoj vrijednosti od A_n

$$A_n v_2 = \sigma_2 u_2, \quad A_n^T u_2 = \sigma_2 v_2 \quad (4)$$

gdje je $\sigma_2 = 1 - \lambda_2$.

- Desni singularni vektor $v_2 \rightsquigarrow$ biparticija dokumenata; lijevi singularni vektor $u_2 \rightsquigarrow$ biparticija riječi.

Biparticioniranje

- Drugi svojstveni vektor matrice L je dan s

$$z_2 = \begin{bmatrix} D_1^{-1/2} u_2 \\ D_2^{-1/2} v_2 \end{bmatrix}. \quad (5)$$

- Tražimo bi-modalne vrijednosti m_1 i m_2 koje ćemo pridružiti riječima i dokumentima kako bismo direktno čitali rješenje.
- Jedan od pristupa je tražiti m_j takve da minimiziraju funkciju

$$\sum_{j=1}^2 \sum_{z_2(i) \in m_j} (z_2(i) - m_j)^2.$$

- Prethodni izraz je upravo funkcija cilja kakvu minimizira klasični k-means algoritam.

Algoritam biparticije

1. Za dani A , izračunati $A_n = D_1^{-1/2} A D_2^{-1/2}$,
2. Izračunati druge singularne vektore u_2 i v_2 matrice A_n i definirati z_2 kao $u(5)$,
3. Provesti *k-means* algoritam na jednodimenzionalnim podacima z_2 .

k-particioniranje

- Iskoristimo $l = \lceil \log_2 k \rceil$ singularnih vektora u_2, \dots, u_{l+1} i v_2, \dots, v_{l+1} i definiramo

$$Z = \begin{bmatrix} D_1^{-1/2} U \\ D_2^{-1/2} V \end{bmatrix}. \quad (6)$$

gdje je $U = [u_2, \dots, u_{l+1}]$ i $V = [v_2, \dots, v_{l+1}]$.

- Tražimo l -dimenzionalne točke m_j , $j = 1, \dots, k$ koje ćemo pridružiti dokumentima i riječima. Stoga, minimiziramo funkciju

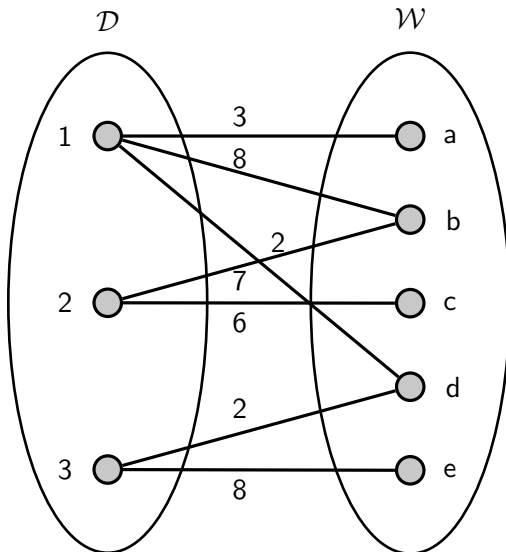
$$\sum_{j=1}^k \sum_{Z(i) \in m_j} \|Z(i) - m_j\|^2,$$

gdje je redak $Z(i)$, $i = 1, \dots, d + w$ dokument ili riječ.

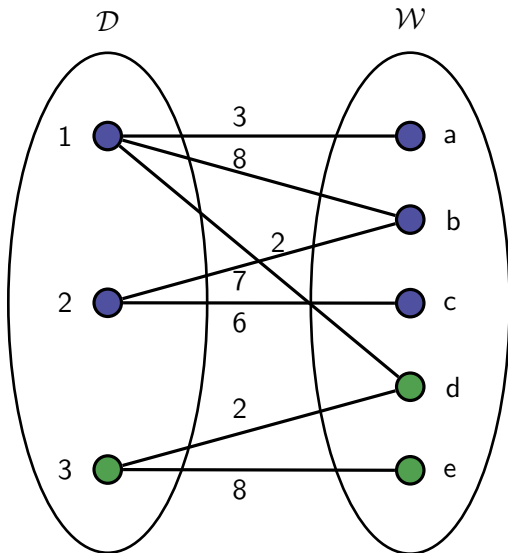
Algoritam k -particije

1. Za dani A , izračunati $A_n = D_1^{-1/2} A D_2^{-1/2}$,
2. Izračunati $l = \lceil \log_2 k \rceil$ singularnih vektora u_2, \dots, u_{l+1} i v_2, \dots, v_{l+1} matrice A_n i definirati matricu Z kao u (6),
3. Provesti *k-means* algoritam na l -dimenzionalnim podacima Z .

Testni primjer



Rezultat



Priprema podataka

- Manji skupovi dokumenata pripremljeni su manualno, a veći su preuzeti iz "TMDB 5000 Movie" (www.kaggle.com/tmdb/tmdb-movie-metadata) i "BBC News Summary" (www.kaggle.com/pariza/bbc-news-summary)
- Tekstovi su obrađeni koristeći MATLAB-ov *Text Analytics Toolbox*:
 - tekstovi → niz riječi (*tokenizedDocument*);
 - odstraniti smo najčešće riječi (*RemoveStopWords*);
 - izbrisati interpunkcijske znakove (*removePunctuation*);
 - izvaditi korijene riječi (*normalizeWords*).

Primjer obrađenog dokumenta:

*childhood hour other seen other saw bring passion common spring
same sourc taken sorrow awaken heart joi same tone lovd lovd alon
childhood dawn ...*

- Tako uređeni dokumenti su kasnije reprezentirani *bag-of-words* modelom, a pritom je zanemaren poredak riječi i sačuvan podatak o broju pojavljivanja svake riječi.
- Primjerice, iz modela ćemo znati da se u gornjem tekstu "childhood" javlja 2 puta, ali ne znamo da se "hour" javlja prije "spring".
- Dobivene podatke ćemo analizirati mjerom *term frequency-inverse document frequency*.

- Term frequency-inverse document frequency je mjera kojom određujemo težinu riječi u skupu dokumenata.
- Jedna varijanta formule za riječ r , dokument d i skup dokumenata D je

$$tf(r, d) = f_{r,d} , \quad idf(r, D) = \ln \frac{N}{n_r}$$

$$tfidf(r, d, D) = tf(r, d) \times idf(r, D)$$

gdje je $f_{r,d}$ "broj ponavljanja r u d ", $N = card(D)$ i $n_r = \{d \in D : r \in d\}$.

- "težina" riječi u dokumentu raste s brojem pojavljivanja u tom dokumentu, a opada s brojem dokumenata u kojima se pojavljuje.

Primjer 1

	<i>recepti</i>	<i>lijekovi</i>	<i>politika</i>
\mathcal{D}_0 :	50	0	0
\mathcal{D}_1 :	0	50	0
\mathcal{D}_2 :	0	0	50
\mathcal{W}_0 :	minut, add, cook, heat, salt		
\mathcal{W}_1 :	medicin, take, doctor, effect, side		
\mathcal{W}_2 :	mr, govern, blair, labour, ministr		

Primjer 2

	<i>sport</i>	<i>business</i>	<i>tech</i>
\mathcal{D}_0 :	1	48	1
\mathcal{D}_1 :	49	0	1
\mathcal{D}_2 :	0	2	48
\mathcal{W}_0 :	year, growth, new, compani		
\mathcal{W}_1 :	olymp, world, athlet, athen		
\mathcal{W}_2 :	mobil, peopl, gadget, user		

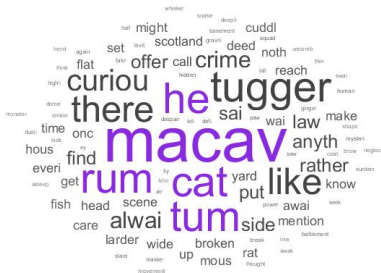
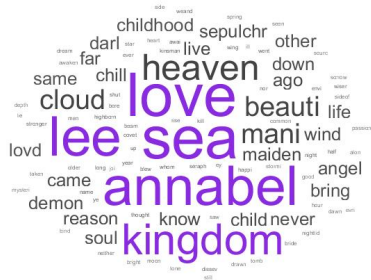
Primjer 3a

	<i>Christmas</i>	<i>E.A.Poe</i>	<i>Cats</i>
\mathcal{D}_0 :	3	0	0
\mathcal{D}_1 :	0	3	0
\mathcal{D}_2 :	0	0	2
\mathcal{W}_0 :	snow, let, christma, merri, dai		
\mathcal{W}_1 :	sea, love, annabel, lee, kingdom		
\mathcal{W}_2 :	macav, cat, rum, tum, curious		

Primjer 3b

	<i>Christmas</i>	<i>E.A.Poe</i>	<i>Cats</i>
\mathcal{D}_0 :	3	1	0
\mathcal{D}_1 :	0	2	0
\mathcal{D}_2 :	0	0	2
\mathcal{W}_0 :	snow, let, christma, merri, dai		
\mathcal{W}_1 :	sea, love, annabel, lee, kingdom		
\mathcal{W}_2 :	macav, cat, rum, tum, curious		

Primjer 3b



Primjer 3c

	<i>Christmas</i>	<i>E.A.Poe</i>	<i>Cats</i>
\mathcal{D}_0 :	3	2	0
\mathcal{D}_1 :	0	1	0
\mathcal{D}_2 :	0	0	2
\mathcal{W}_0 :	snow, christma, love, sea, annabel		
\mathcal{W}_1 :	childhood, other, same, lovd, hour		
\mathcal{W}_2 :	macav, cat, rum, tum, curious		

Primjer 4

	<i>sport</i>	<i>tech</i>
\mathcal{D}_0 :	172	400
\mathcal{D}_1 :	328	0
\mathcal{W}_0 :	olymp, world, athlet, athen	
\mathcal{W}_1 :	year, new, mobil, peopl	

	<i>sport</i>	<i>tech</i>
\mathcal{D}_0 :	4	0
\mathcal{D}_1 :	23	400
\mathcal{D}_2 :	471	0
\mathcal{D}_3 :	2	0
\mathcal{W}_0 :	harrier, ac, stade, treviso	
\mathcal{W}_1 :	peopl, game, mobil, phon	
\mathcal{W}_2 :	game, win, plai, against	
\mathcal{W}_3 :	republ, ireland, franc, faro	

Primjer 4



Primjer 5

	<i>sport</i>	<i>entertainment</i>	<i>politics</i>	<i>business</i>	<i>tech</i>
\mathcal{D}_0 :	4	19	0	0	1
\mathcal{D}_1 :	0	1	0	0	0
\mathcal{D}_2 :	0	2	0	0	0
\mathcal{D}_3 :	45	23	48	50	49
\mathcal{D}_4 :	1	5	2	0	0
\mathcal{W}_0 : best, award, film, book, year, winner					
\mathcal{W}_1 : la, fenic, viotti, opera, director, includ					
\mathcal{W}_2 : christian, andersen, han, prize, author, booker					
\mathcal{W}_3 : mr, year, new, govern, peopl, world					
\mathcal{W}_4 : film, famili, border, ballet, white, year					

Primjer 5

