

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 1547

**Primjena koncepta jezgrenih funkcija u
dubinskoj analizi teksta**

Dajana Buršić

Zagreb, srpanj 2005.

*Zahvaljujem se sjajnoj mentorici i voditeljici
prof. dr.sc. Bojani Dalbelo-Bašić
i kolegi Domagoju Tomincu na uvijek spremnoj pomoći*



SADRŽAJ

1	UVOD	3
2	JEZGRENE METODE	4
2.1	Uvod u jezgrene metode.....	4
2.2	Redukcija dimenzionalnosti jezgrenim metodama	7
2.2.1	Jezgrena metoda glavnih komponenata.....	7
2.2.2	Jezgrena kanonska korelacijska analiza	12
2.3	Jezgrene funkcije za klasifikaciju teksta	16
2.4	Istraživanje višejezičnih korpusa.....	22
2.5	Metoda potpornih vektora	24
2.5.1	Upotreba metode potpornih vektora	33
2.5.2	Biblioteka funkcija LIBSVM te OSU SVM	35
2.5.3	Klasifikacija u više klase.....	36
3	REZULTATI.....	37
3.1	Ulagni podaci	37
3.1.1	Croatia Weekly.....	37
3.1.2	Baza novinskih članaka Vjesnik	38
3.1.3	Baza novinskih članaka Vjesnik AMNv1.0.....	39
3.2	Rezultati usporedbe PCA i CCA	41
4	KONCEPTNI VEKTORI.....	48
4.1	Konceptni vektori i objektivna funkcija	48
4.2	Spherical k-means algoritam	49
4.3	Eksperimentalni rezultati.....	51
4.3.1	Objektivna funkcija je rastuća.....	51
4.3.2	Konceptni vektori su lokalni, rijetki i teže ortonormalnosti	52

4.3.3	Struktura unutar i između grupa	57
5	UPORABA RAZLIČITIH METRIKA.....	59
5.1	Euklidska mjera udaljenosti	59
5.2	Mahalanobisova udaljenost	60
5.3	Učenje funkcije udaljenosti	62
5.3.1	Učenje metrike	62
5.3.2	Jezgreno učenje	66
5.4	Jezgreni k-means algoritam.....	67
5.5	Odabir jezgrene funkcije	69
6	KONCEPTNO INDEKSIRANJE.....	73
6.1	Eksperimentalni rezultati.....	74
6.1.1	Nadzirana i nenadzirana redukcija dimenzionalnosti.....	74
6.1.2	Utjecaj morfološke normalizacije na hrvatski jezik.....	76
6.1.3	Mjere uspješnosti za različite morfološke normalizacije, sa i bez kategorije tema dana	79
7	ZAKLJUČAK.....	85
8	LITERATURA	86
9	PRIMJER RADA METODE GLAVNIH KOMPONENTA I KANONSKE KORELACIJSKE ANALIZE.....	89

1 UVOD

Osnovni zadatak ovog diplomskog rada bio je proučiti primjenu koncepta jezgrenih funkcija u dubinskoj analizi teksta. Sastoji se od 2 dijela. U prvom dijelu proučava se mogućnost uvođenja koncepta jezgrenih funkcija u metode za dubinsku analizu teksta, posebno metode za redukciju dimenzionalnosti. Opisani su neki od osnovnih algoritma za rad u jezrenom području – algoritmi za redukciju dimenzionalnosti i klasifikaciju. Algoritmi za redukciju dimenzionalnosti temelje se na općenitom problemu svojstvenih vrijednosti i za njihovo razumijevanje potrebno je osnovno znanje linearne algebre te se mogu lako riješiti ili aproksimirati koristeći dobro poznate tehnikе numeričke algebre. Nadalje ovakvi problemi mogu se riješiti u dualnoj reprezentaciji tj. potrebne su nam samo informacije o skalarном produktu između točaka. Općeniti problemi svojstvenih vrijednosti obuhvaćaju pronalaženje k smjerova u preslikanom prostoru koji sadrže maksimalnu vrijednost varijance u podacima (metoda glavnih komponenata) ili pronalaze korelacije između različitih reprezentacija istih podataka (kanonska korelacijska analiza). Algoritme redukcije dimenzionalnosti slijedi klasifikacija podataka jednom od najpoznatijih metoda, metodom potpornih vektora.

U drugom dijelu ispitana su svojstva i mogućnosti primjene metrika različitih od Euklidske na postupke grupiranja i klasificiranja podataka u dubinskoj analizi teksta za model vektorskog prostora. U zadnjem dijelu kao redukcija dimenzionalnosti odabrana je metoda konceptnog indeksiranja te je na temelju nje i metode potpornih vektora izrađen klasifikator. Klasifikacija je testirana na testnim skupovima dokumenata na hrvatskom jeziku.

2 JEZGRENE METODE

2.1 Uvod u jezgrene metode

Metode raspoznavanja uzorka primjenjuju se u mnogo područja gdje je potrebno klasificirati uzorke u određen broj klasa. Primjeri klasičnih problema raspoznavanja uzorka su prepoznavanje slova, prepoznavanje govora, medicinska dijagnostika... Novije primjene raspoznavanja uzorka su u području dubinske analize podataka te u dubinskoj analizi teksta (eng. text mining), tj. radu na sintezi novih informacija iz podataka koji dolaze isključivo u tekstualnom obliku.

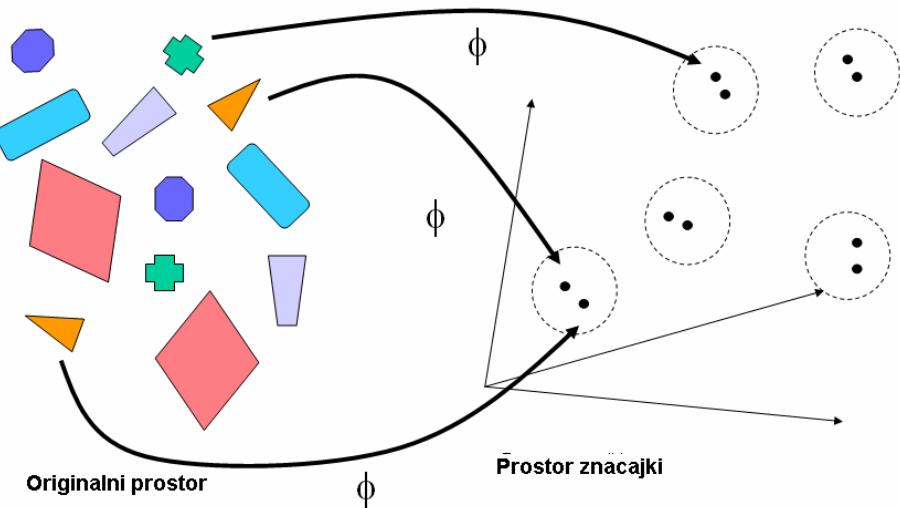
Evoluciju kroz koju su prolazili algoritmi raspoznavanja uzorka možemo podijeliti u 3 etape. Sve je počelo 1960-ih otkrićem linearnih algoritama, među kojima je bitno spomenuti perceptron. Takvi algoritmi mogli su pronaći samo linearne relacije među uzorcima. Uslijedila je „nelinearna revolucija“ sredinom 1980-ih koja je dovela do razvoja raspoznavanja uzorka i njegove primjene na do tada nezamislivim područjima. Iako je riješeno najvažnije pitanje postavljano do tada, korišteni su heuristički algoritmi i nedorađena statistička analiza, a algoritmi gradijentnog spusta i pohlepna heuristika doveli su do novih problema, među ostalim lokalnog minimuma i pretreniranosti. Treća etapa započinje sredinom 1990-ih uvođenjem novog pristupa poznatog kao jezgrene metode. Ovakav pristup omogućuje analizu nelinearnih relacija među uzorcima učinkovitošću prije rezerviranom samo za linearno odvojive uzorke, te nadalje, napredak u statističkoj analizi omogućio je analizu u višedimenzionalnom prostoru značajki. Prevladani su tipični problemi 3 etape, pretreniranost i lokalni minimum. Sa svih pogleda, računskog, statističkog i konceptualnog, ovi algoritmi predstavljaju novu stepenicu u razvoju raspoznavanja uzorka.

Prvi put su uvedeni u metodi potpornih vektora, klasifikacijskom algoritmu koji je opravdao sve gore navedene prednosti jezgrenih metoda. Zatim su se proširili i na druga područja uz klasifikaciju uzorka, ovdje možemo spomenuti još redukciju dimenzionalnosti i dvije metode kojima ćemo se baviti: jezgrena metoda glavnih komponenata i jezgrena kanonska korelacijska analiza.

Osnovni koncept jezgrenih metoda je preslikavanje uzorka u novi pogodniji prostor u kojem se zatim koriste algoritmi temeljeni na linearnoj algebri, geometriji i statistici za pronalaženje relacija među preslikanim uzorcima. Vidjet ćemo da je bit jezgrenih metoda

računski prečac za predstavljanje uzoraka u prostoru značajki koji zovemo jezgrena funkcija.

Zadani je skup $S = \{(x_1, y_1), (x_2, y_2) \dots, (x_i, y_i)\}$ točaka $x_i \subseteq \mathbb{R}^N$ zajedno s odgovarajućim oznakama $y_i \subseteq \mathbb{R}$. U dalnjem tekstu ℓ će označavati broj zadanih točaka, dok je N njihova dimenzionalnost. Za taj skup definirano je preslikavanje $\phi: x \in \mathbb{R}^N \rightarrow \phi(x) \in F$ na prostor značajki takvo da je



Slika 1. Preslikavanje točaka

rješavanje problema u novom prostoru jednostavnije npr. linearno. Ponekad nije potrebno eksplicitno znanje o preslikavanju ϕ , već je skalarni produkt u novom prostoru odabran kao potrebna mjera sličnosti između točaka. Jezgrena funkcija je ta mjera sličnosti i definirana je kao:

$$k(\mathbf{x}, \mathbf{z}) = \langle \phi(x), \phi(z) \rangle.$$

Moguće je, dok za sada izgleda nemoguće, koristiti prostor značajki sa eksponencijalnim ili beskonačnim brojem dimenzija. Krenimo sa jednostavnim primjerom jezgrene funkcije da shvatimo ključnu ideju ovakvog pristupa.

Za dvodimenzionalni ulazni prostor $X \subseteq \mathbb{R}^2$ zadano je preslikavanje

$$\phi: x = (x_1, x_2) \rightarrow \phi(x) = (x_1^2, x_2^2, \sqrt{2}x_1x_2) \in F = \mathbb{R}^3.$$

podaci se preslikavaju iz dvodimenzionalnog prostora u trodimenzionalni na taj način da linearne relacije u prostoru značajki odgovaraju kvadratnim relacijama na ulaznom prostoru. Skalarni produkt u prostoru značajki se može izračunati kao

$$\begin{aligned}
\langle \phi(x), \phi(z) \rangle &= \left\langle (x_1^2, x_2^2, \sqrt{2}x_1x_2), (z_1^2, z_2^2, \sqrt{2}z_1z_2) \right\rangle \\
&= x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 x_2 z_1 z_2 \\
&= (x_1 z_1 + x_2 z_2)^2 = \langle x, z \rangle^2
\end{aligned}$$

Dakle $k(x, z) = \langle x, z \rangle^2$ je jezgrena funkcija. Vidimo da možemo izračunati skalarni produkt između projekcija bez eksplisitnog znanja njihovih koordinata. Međutim nisu sve funkcije jezgrene funkcije.

Gram matrica je definirana kao $\ell \times \ell$ matrica sa elementima $G_{ij} = \langle \phi(x_i), \phi(x_j) \rangle = k(x_i, x_j)$. Češće se ta matrica naziva jezgrenom matricom. Definirajmo još i matricu \mathbf{X} dimenzija $\ell \times N$

$$\mathbf{X} = [\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_N)]'.$$

Jezgrene matrice definirane su elementima koji su skalarni produkt točaka u prostoru značajki. Iz gornje definicije proizlazi da takve matrice moraju biti simetrične i pozitivne definitne tj. sve svojstvene vrijednosti matrice su veće od nule. Tako funkcija $k(x, z)$ zadovoljava svojstvo pozitivnosti i definitnosti ako je simetrična funkcija za koju i konstruirana jezgrena matrica zadovoljava ovo svojstvo. Prema Merceru jezgrene funkcije moraju zadovoljavati sljedeći uvjet:

$$\int_{X \times X} k(x, z) f(x) f(z) dx dz \geq 0 \quad \forall f \in L_2(X)$$

Napomenimo još jedanput da rad u jezgrenom području znači da ne možemo eksplisitno prikazati preslikane točke. Npr. slika ulazne točke \mathbf{x} je $\phi(\mathbf{x})$, ali nemamo pristupa komponentama ovog vektora već samo vrijednosti skalarnog produkta između slike točaka. Unatoč ovakvoj smetnji začuđujuće je koliko se informacija može dobiti o $\phi(\mathbf{x})$ koristeći jezgrene funkcije.

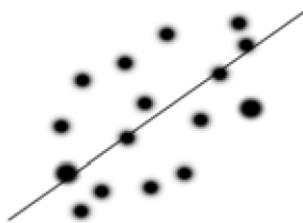
2.2 Redukcija dimenzionalnosti jezgrenim metodama

Sljedeća poglavlja analiziraju dvije metode redukcije dimenzionalnosti, tj. nastoje smisleno smanjiti dimenzionalnost N koje smo definirali u skupu S .

Analizirane redukcije dimenzionalnosti temelje se na općim problemima svojstvenih vektora. Proučavaju se jednostavnom linearnom algebrrom, te je također izračunavanje ili njihova aproksimacija učinkovita. Nadalje, takvi problemi se mogu rješavati i u jezgrenom prostoru značajki koristeći dualnu reprezentaciju, tj. potrebni su nam samo podatci o skalarnom produktu između podatkovnih vektora.

2.2.1 Jezgrena metoda glavnih komponenata¹

Za početak ćemo prepostaviti da su podaci centrirani tj. $\sum_i x_i = 0$ ². Cilj nam je pronaći smisleno preslikavanje podataka i time smanjiti dimenzionalnost. Ipak, problem je postavljen tako da nemamo pristupa ikakvima oznakama y_i . Da imamo, okrenuli bi se na PLS³. I bez oznaka, cilj nam je pronaći potprostor najvećih varijanci gdje broj novih dimenzija unaprijed odredimo. Varijanca je mjera rasipanja podataka, a računa se kao srednje kvadratno odstupanje od aritmetičke sredine. Može se dogoditi da je zanimljivi signal zarobljen u smjeru male varijance, i u tom slučaju metoda glavnih komponenata ili PCA nije dobra tehnika. Međutim, obično vrijedi da smjerovi male varijance pretstavljaju samo šum, te odšumljavanjem možemo samo popraviti reprezentaciju podataka.



¹ Eng. Kernel principal component analysis KPCA

² Centriranje je translacija osi.

³ Eng. Partial Last Squares

Slika 2. Smjer najveće varijance podataka

Zapišimo kovarijacijsku matricu **C**

$$C_{st} = \frac{1}{l} \sum_i^l \phi(x_i)_s \phi(x_i)_t^T, s, t = 1, \dots, N$$

te razmotrimo izračunavanje varijance projekcije na normalizirani vektor **w** (napomenimo da su podaci centrirani)

$$\frac{1}{l} \sum_{i=1}^l (P_w(\phi(x_i)))^2 = E[w' \phi(x) \phi(x)' w] = w' E[\phi(x) \phi(x)'] w = w' X' X w = w' C w,$$

$$\text{gdje se } E \text{ koristi kao empirička srednja vrijednost ili } E[f(x)] = \frac{1}{l} \sum_{i=1}^l f(x).$$

Smjer koji maksimizira varijancu dobiva se rješavanjem sljedećeg problema

$$\max_w \quad w' C w \quad \text{uz} \quad \|w\|_2 = 1.$$

Rješenje ovog problema dano je svojstvenim vektorom matrice **C** kojem je pridružena najveća svojstvena vrijednost. Svojstveni vektor koji odgovara najvećoj svojstvenoj vrijednosti je zapravo smjer u kojem su podaci najviše rastegnuti. Svojstvene vrijednosti ove matrice predstavljaju količinu varijance u smjerovima pridruženih svojstvenih vektora. Sljedeći smjer je ortogonalan na njega, i odgovara vrijednosti najveće varijance u ortogonalnom potprostoru itd. Dakle, za smanjenje dimenzionalnosti projiciramo podatke na svojstvene - smjerove pridružene najvećim varijancama.

PCA uzima podskup k svojstvenih vektora kao osnovnih osi, dobivenih na skupu za učenje, te projicira podatke u prostor razapet tim osima. Nove koordinate nazivaju se glavne koordinate⁴, a svojstveni vektori glavnim osima⁵.

Projekcije su dane sa:

$$y_i = U_k^T x_i \quad \forall i$$

gdje je U_k $m \times k$ podmatrica sa prvih k svojstvenih vektora u stupcima, a algoritam:

⁴ Eng. principal coordinates

⁵ Eng. principal axes

Ulaz podaci $S = \{x_1, x_2, \dots, x_l\} \subset \Re^n$

$$\mathbf{C} = \frac{1}{l} \mathbf{X}' \mathbf{X} \text{ kovarijacijska matrica}$$

$$[\mathbf{U}, \Lambda] = \text{eig}(\ell \mathbf{C}) \text{ dekompozicija}$$

$$x_i^* = \mathbf{U}_k' x_i, i = 1, \dots, \ell \text{ preslikavanje}$$

Izlaz $S^* = \{x_1^*, x_2^*, \dots, x_k^*\}, \text{ dimenzija } k$

Algoritam metode glavnih komponenata

Kako je već napomenuto svojstvene vrijednosti predstavljaju količinu varijance uhvaćene pridruženim svojstvenim vektorom. Čim je veća dimenzionalnost k veća je obuhvaćena varijanca. Nadalje, ortogonalna projekcija P_{U_k} u potprostor razapet sa prvih k svojstvenih vektora matrice \mathbf{C} je k -dimenzionalna ortogonalna projekcija koja minimalizira prosječnu kvadratnu udaljenosti između podatkovnih točaka i njihovih slika. Odnosno \mathbf{U}_k rješava optimizacijski problem

$$\min \sum_{i=1}^l \|P_{U_k}^\perp(\phi(x_i))\|_2^2.$$

Dokaz se može pronaći u (Shawe-Taylor et al., 2004.)

Posljedice Singular Value Decomposition.

Vidjet ćemo kako možemo naučiti nešto o kovarijacijskoj matrici \mathbf{C} koristeći jezgenu matricu $\mathbf{K} = \mathbf{XX}'$. Varijance koje su nam potrebne dane su u kovarijacijskoj matrici, ali se također mogu izračunati koristeći jezgenu matricu. Njihova veza postat će očita ako razmotrimo njihove dekompozicije

$$\ell \mathbf{C} = \mathbf{X}' \mathbf{X} = \mathbf{U} \Lambda_k \mathbf{U}' \quad \text{i} \quad \mathbf{K} = \mathbf{XX}' = \mathbf{V} \Lambda_k \mathbf{V}',$$

gdje su redci \mathbf{u}_i ortonormalne matrice \mathbf{U} svojstveni vektori $\ell \mathbf{C}$, a stupci \mathbf{v}_i ortonormalne matrice \mathbf{V} svojstveni vektori matrice \mathbf{K} . Sad razmotrimo svojstveni vektor i pridruženu svojstvenu vrijednost \mathbf{v}, λ matrice \mathbf{K} . Dakle,

$$\ell \mathbf{C}(\mathbf{X}' \mathbf{v}) = \mathbf{X}' \mathbf{X} \mathbf{X}' \mathbf{v} = \mathbf{X}' \mathbf{K} \mathbf{v} = \lambda \mathbf{X}' \mathbf{v}$$

možemo zaključiti da je $(\mathbf{X}' \mathbf{v})$, λ svojstveni vektor i svojstvena vrijednost za $\ell \mathbf{C}$. Nadalje norma $\mathbf{X}' \mathbf{v}$ je

$$\| \mathbf{X}' \mathbf{v} \|^2 = \mathbf{v}' \mathbf{X} \mathbf{X}' \mathbf{v} = \lambda$$

Iz gornjeg dobivamo normalizirani svojstveni vektor matrice $\ell \mathbf{C}$, $\mathbf{u} = \lambda^{-1/2} \mathbf{X}' \mathbf{v}$. Također postoji simetrija:

$$\lambda^{-1/2} \mathbf{X} \mathbf{u} = \lambda^{-1} \mathbf{X} \mathbf{X}' \mathbf{v} = \mathbf{v}$$

te možemo zapisati

$$\mathbf{u} = \lambda^{-1/2} \mathbf{X}' \mathbf{v} \quad \text{i} \quad \mathbf{v} = \lambda^{-1/2} \mathbf{X} \mathbf{u}$$

Dekompozicija matrice \mathbf{X} slijedi iz sljedeće dvije jednakosti

$$\mathbf{X}^* \mathbf{X}^* = (\mathbf{X} - \mathbf{v} \mathbf{v}' \mathbf{X}) (\mathbf{X} - \mathbf{v} \mathbf{v}' \mathbf{X})' = \mathbf{X} \mathbf{X}' - \lambda \mathbf{v} \mathbf{v}', \text{ i}$$

$$\mathbf{X}^* \mathbf{X}^* = (\mathbf{X} - \mathbf{v} \mathbf{v}' \mathbf{X})' (\mathbf{X} - \mathbf{v} \mathbf{v}' \mathbf{X}) = \mathbf{X}' \mathbf{X} - \mathbf{X}' \mathbf{v} \mathbf{v}' \mathbf{X} = \mathbf{X}' \mathbf{X} - \lambda \mathbf{u} \mathbf{u}'$$

$$\text{te je } \mathbf{X}^* = \mathbf{X} - \mathbf{v} \mathbf{v}' \mathbf{X} = \mathbf{X} - \lambda^{-1/2} \mathbf{v} \mathbf{u}' = \mathbf{X} - \mathbf{X} \mathbf{u} \mathbf{u}'.$$

Dakle, prvih $t = \text{rang}(\mathbf{X} \mathbf{X}') \leq \min(N, \ell)$ stupaca matrice \mathbf{U}_t može se odabratи kao

$$\mathbf{U}_t = \mathbf{X}' \mathbf{V}_t \Lambda_t^{-1/2}.$$

Proširenjem \mathbf{U}_t na \mathbf{U} i matrice $\Lambda_t^{-1/2} N \times \ell$ matricu čiji su dodatni elementi svi nula, dobili smo singular value decomposition(SVD) matrice \mathbf{X}'

$$\mathbf{X}' = \mathbf{U} \Sigma \mathbf{V}'$$

gdje je Σ $N \times \ell$ matrica sa svim elementima nula, osim vodećih na dijagonalni sa vrijednostima $\sigma_i = \lambda_i^{-1/2}$ koje zadovoljavaju $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_t > 0$ za $t = \text{rang}(\mathbf{X}) \leq \min(N, \ell)$, te \mathbf{U} i \mathbf{V} kvadratnim ortogonalnim matricama.

Prethodna definicija matrice \mathbf{U}_t implicira dualni reprezentaciju svojstvenog vektora \mathbf{u}_j matrice $\ell \mathbf{C}$ prikazanog pomoću odgovarajućeg svojstvenog vektora \mathbf{v}_j matrice \mathbf{K} pomnoženog sa $\lambda_j^{-1/2}$ tj.

$$\mathbf{u}_j = \lambda_j^{-1/2} \sum_{i=1}^l (v_j)_i \phi(x_i) = \sum_{i=1}^l \alpha_i^j \phi(x_i), \quad j = 1, \dots, t,$$

gdje su dualne varijable α_j za j -ti vektor \mathbf{u}_j dane sa $\alpha_j = \lambda_j^{-1/2} v_j$. Također možemo izračunati projekciju točke $\phi(\mathbf{x})$ na smjer \mathbf{u}_j u prostoru značajki

$$P_{U_j}(\phi(\mathbf{x})) = \mathbf{u}_j \cdot \phi(\mathbf{x}) = \left\langle \sum_{i=1}^l \alpha_i^j \phi(\mathbf{x}_i), \phi(\mathbf{x}) \right\rangle = \sum_{i=1}^l \alpha_i^j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle = \sum_{i=1}^l \alpha_i^j k(x_i, x).$$

Jezgrena metoda glavnih komponenata

Kernel PCA je primjena PCA algoritma u jezgrenom prostoru značajki te se tako iskorištava dualna reprezentacija. Iz prethodne jednadžbe znamo da je projekcija

$$P_{U_k}(\phi(\mathbf{x})) = (\mathbf{u}_j \cdot \phi(\mathbf{x}))_{j=1}^k = \left(\sum_{i=1}^l \alpha_i^j k(x_i, x) \right)_{j=1}^k \quad \text{sa} \quad \alpha_j = \lambda_j^{-\frac{1}{2}} v_j.$$

Koristeći jezgenu matricu \mathbf{K} te njenu dekompoziciju na $\mathbf{V} \Lambda_k \mathbf{V}^*$ jednostavno je prilagoditi PCA algoritam za prostor značajki (Shawe-Taylor et al., 2004.).

```

Ulaz    podaci S = { x1 x2 ... xl } ⊂ ℝ^n
Kij = k(xi, xj) jezgrena matrica
// centriranje jezgrene matrice, zamjenjuje centriranje
podataka
[V, Λ] = eig(K) dekompozicija jezgrene matrice
αj = λi^{-\frac{1}{2}} vj,    j = 1,...,k
x_i^* = \left( \sum_{i=1}^l \alpha_i^j k(x_i, x) \right)_{j=1}^k preslikavanje
Izlaz    S* = { x1^*.x2^*....xl^* } . dimenzija k

```

Algoritam jezgrene metode glavnih komponenata

2.2.2 Jezgrena kanonska korelacijska analiza⁶

Pretpostavimo da imamo 2 kopije korpusa, jednu napisanu na engleskom, a drugu na hrvatskom jeziku. Ovakav paralelni korpus nazivamo još i upareni set dokumenata. Zainteresirani smo za manje dimenzionalni prikaz dokumenata. Da imamo jedan jezik, mogli bi pomoću PCA dobiti dimenzije koje imaju najveće varijance. Time bi mogli riješiti problem sinonimije, jer ako se riječi često pojavljuju u dokumentima tj. vrlo su korelirane, nastoje se spojiti u istu dimenziju u novom prostoru. Za korpus sa dva prijevoda, možemo tražiti projekciju svakog prijevoda posebno tako da su projekcije maksimalno korelirane(Shawe-Taylor et al., 2004.).

Za upareni skupa podataka⁷, gdje je isti objekt \mathbf{x} prikazan na dva načina kao $\phi_a(\mathbf{x})$ i $\phi_b(\mathbf{x})$, $S = \{(\phi_a(\mathbf{x}_1), \phi_b(\mathbf{x}_1)), \dots, (\phi_a(\mathbf{x}_l), \phi_b(\mathbf{x}_l))\}$ nastojimo maksimizirati empiričku korelaciju između $\mathbf{x}_a = \mathbf{w}_a' \phi_a(\mathbf{x})$ i $\mathbf{x}_b = \mathbf{w}_b' \phi_b(\mathbf{x})$ projiciranih na \mathbf{w}_a i \mathbf{w}_b .

$$\begin{aligned} \text{Corr}(\mathbf{x}_a, \mathbf{x}_b) &= \frac{\text{E}[x_a, x_b]}{\sqrt{\text{E}[x_a, x_a] \text{E}[x_b, x_b]}} = \frac{\text{cov}(x_a, x_b)}{\sqrt{\text{var}(x_a, x_a) \text{var}(x_b, x_b)}} \\ &= \frac{\text{E}[\mathbf{w}_a' \phi(x_a) \phi(x_b)' \mathbf{w}_b]}{\sqrt{\text{E}[\mathbf{w}_a' \phi(x_a) \phi(x_a)' \mathbf{w}_a] \text{E}[\mathbf{w}_b' \phi(x_b) \phi(x_b)' \mathbf{w}_b]}} = \frac{\mathbf{w}_a' C_{ab} \mathbf{w}_b}{\sqrt{\mathbf{w}_a' C_{aa} \mathbf{w}_a \mathbf{w}_b' C_{bb} \mathbf{w}_b}} \end{aligned}$$

gdje smo na sljedeći način rastavili kovarijacijsku matricu

$$\begin{aligned} \mathbf{C} &= \frac{1}{l} \sum_{i=1}^l (\phi_a(\mathbf{x}_i), \phi_b(\mathbf{x}_i)) (\phi_a(\mathbf{x}_i), \phi_b(\mathbf{x}_i))' = \begin{pmatrix} \frac{1}{l} \sum_{i=1}^l \phi_a(x_i) \phi_a(x_i)' \frac{1}{l} \sum_{i=1}^l \phi_b(x_i) \phi_a(x_i)' \\ \frac{1}{l} \sum_{i=1}^l \phi_a(x_i) \phi_b(x_i)' \frac{1}{l} \sum_{i=1}^l \phi_b(x_i) \phi_b(x_i)' \end{pmatrix} \\ &= \begin{pmatrix} C_{aa} & C_{ba} \\ C_{ab} & C_{bb} \end{pmatrix} \end{aligned}$$

Ovaj optimizacijski problem sličan je optimizacijskom problemu PLS-a⁸. Vektori \mathbf{w}_a i \mathbf{w}_b određeni su opet samo svojim smjerovima jer njihovo množenje sa λ_a te λ_b daje

⁶ Eng. Kernel Canonical correlation analysis KCCA

⁷ Eng. paired or aligned dataset

$$\begin{aligned} \frac{\lambda_a \lambda_b w_a' C_{ab} w_b}{\sqrt{\lambda_a^2 w_a' C_{aa} w_a \lambda_b^2 w_b' C_{bb} w_b}} &= \frac{\lambda_a \lambda_b w_a' C_{ab} w_b}{\lambda_a \lambda_b \sqrt{w_a' C_{aa} w_a w_b' C_{bb} w_b}} \\ &= \frac{w_a' C_{ab} w_b}{\sqrt{w_a' C_{aa} w_a w_b' C_{bb} w_b}} \end{aligned}$$

te time možemo ograničiti doljnja dva izraza u nazivniku na vrijednost 1. Postavljamo sljedeći optimizacijski problem

$$\begin{aligned} \max_{w_a, w_b} \quad & w_a' C_{ab} w_b \\ \text{uz uvjet} \quad & w_a' C_{aa} w_a = 1 \quad \text{i} \quad w_b' C_{bb} w_b = 1 \end{aligned}$$

Ako primijenimo Lagrangeov multiplikator dobivamo

$$\max \quad w_a' C_{ab} w_b - \frac{\lambda_a}{2} (w_a' C_{aa} w_a - 1) - \frac{\lambda_b}{2} (w_b' C_{bb} w_b - 1)$$

Derivirajući po w_a i w_b

$$C_{ab} w_b - \lambda_a C_{aa} w_a = 0 \quad \text{i} \quad C_{ba} w_a - \lambda_b C_{bb} w_b = 0$$

oduzimajući prvu jednadžbu pomnoženu sa w_a' te drugu pomnoženu sa w_b' dobivamo

$$\lambda_a w_a' C_{aa} w_a - \lambda_b w_b' C_{bb} w_b = 0$$

Zaključujemo, ako uzmemo dva prethodna uvjeta da je $\lambda_a = \lambda_b$ te će se u dalnjem tekstu ova vrijednost obilježavati sa λ .

Problem možemo postaviti kao općeniti problem svojstvenih vrijednosti

$$\begin{pmatrix} 0 & C_{ab} \\ C_{ba} & 0 \end{pmatrix} \begin{pmatrix} w_a \\ w_b \end{pmatrix} = \lambda \begin{pmatrix} C_{aa} & 0 \\ 0 & C_{bb} \end{pmatrix} \begin{pmatrix} w_a \\ w_b \end{pmatrix}$$

Jezgrena kanonska korelacijska analiza

Naravno, problem želimo riješiti u dualnoj formulaciji. Nastaviti ćemo sa predstavljanjem w_a i w_b pomoću matrice X_a čiji su redci vektori $\phi_a(x_i)$, $i = 1, \dots, \ell$, te matrice X_b sa redcima $\phi_b(x_i)$

⁸ Eng. Partial last squares

$$\mathbf{w}_a = \mathbf{X}_a' \alpha_a \quad \text{and} \quad \mathbf{w}_b = \mathbf{X}_b' \alpha_b$$

te optimizacijski problem postaje

$$\max \alpha_a' \mathbf{X}_a \mathbf{X}_a' \mathbf{X}_b \mathbf{X}_b' \alpha_b$$

$$\text{uz uvjet } \alpha_a' \mathbf{X}_a \mathbf{X}_a' \mathbf{X}_a \mathbf{X}_a' \alpha_a = 1 \text{ te } \alpha_b' \mathbf{X}_b \mathbf{X}_b' \mathbf{X}_b \mathbf{X}_b' \alpha_b = 1$$

ili zapisano pomoću jezgrenih matrica

$$\max \alpha_a' \mathbf{K}_a \mathbf{K}_b \alpha_b$$

$$\text{uz uvjet } \alpha_a' \mathbf{K}_a^2 \alpha_a = 1 \text{ te } \alpha_b' \mathbf{K}_b^2 \alpha_b = 1$$

Pomakom u jezgreno definirani prostor značajki, dodatna fleksibilnost uvodi i opasnost od pretreniranosti. Time mislimo na pronalaženje lažnih korelacija koristeći vektore težine tako da su dvije projekcije potpuno poravnane. Na primjer, ako su podatci linearno nezavisni u oba prostora značajki možemo pronaći linearnu transformaciju koja preslikava ortogonalno svaki ulazni podatak u prostoru značajki. Sad je moguće pronaći i savršene korelacije između odabrane projekcije u jednom prostoru te proizvoljne projekcije u drugom. Maksimiziranje korelacije odgovara minimiziranju funkcije $\| \mathbf{w}_a' \phi_a(\mathbf{x}) - \mathbf{w}_b' \phi_b(\mathbf{x}) \|_2^2$. Korištenje jezgrene funkcije redovito dovodi do linearne nezavisnosti skupa za učenje te je potrebno uvesti kontrolu fleksibilnosti preslikavanja u prostore značajki. Regularizacijski parametri π_a i π_b glatko interpoliraju između maksimiziranja korelacije i maksimiziranja kovarijance. Koristeći Lagrangeovu metodu dolazimo do

$$\mathbf{K}_a \mathbf{K}_b \alpha_b - \lambda (1 - \pi_a) \mathbf{K}_a^2 \alpha_a - \lambda \pi_a \mathbf{K}_a \alpha_a = 0$$

$$\mathbf{K}_b \mathbf{K}_a \alpha_a - \lambda (1 - \pi_b) \mathbf{K}_b^2 \alpha_b - \lambda \pi_b \mathbf{K}_b \alpha_b = 0$$

Dakle formirali smo općeniti problem svojstvenih vrijednosti

$$\begin{pmatrix} 0 & K_a K_b \\ K_b K_a & 0 \end{pmatrix} \begin{pmatrix} \alpha_a \\ \alpha_b \end{pmatrix} = \lambda \begin{pmatrix} ((1 - \pi_a) K_a^2 + \pi_a K_a) 0 \\ ((1 - \pi_b) K_b^2 + \pi_b K_b) \alpha_b \end{pmatrix}$$

Standardan pristup rješavanju ovog problema je nepotpuna Cholesky dekompozicija na srednjoj matrici sa desne strane.

Regularizacijski parametar osim računskog problema određuje i funkcionalni prostor u kojem se traži rješenje. Sa većim vrijednostima regulacijskog parametra metoda je manje osjetljiva na ulazne podatke time i stabilnija (manje je vjerojatno da će pronaći lažne relacije među podacima). Sposobnost metode da uhvati koristan signal mjeri se

uspoređivanjem na stvarnim ulaznim i slučajnim podacima. Slučajni podaci dobivaju se slučajnim udruživanje parova podataka, na primjer (H_r , $\text{rand}(E)$) predstavlja hrvatsko-engleski korpus koji dobivamo premetanjem engleskih rečenica na stvarnom korpusu. Označimo dobivene korelacije sa $\text{KCCA}_\pi(H_r, E)$ sa regularizacijskim parametrom π , te J matricu sa svim jedinicama. Ako dolazi do pretreniranosti moguće je pronaći savršene korelacije te će vrijediti $\| J - \text{KCCA}_\pi(H_r, E) \| \approx 0$. Za male vrijednosti regularizacijskog parametra za bilo koje povezivanje dokumenata vrijednost funkcije mjere bit će blizu 0 te zaključujemo da to indicira pretreniranost. Kako se λ povećava funkcija mjere slučajno povezanih podataka penje se daleko od 0, dok ona pravilno uparenih dokumenata ostaje korelirana. Pravilni parametar π odabiremo tako da je $\| J - \text{KCCA}_\pi(H_r, \text{rand}(E)) \| > \| J - \text{KCCA}_\pi(H_r, E) \| > 0$ (Vinokourov et al., 2002.).

2.3 Jezgrene funkcije za klasifikaciju teksta

Najjednostavniji prikaz teksta u modelu vektorskog prostora⁹ je prikaz dokumenata kao „torbe riječi“¹⁰. Torba je skup u kojem su ponavljanja dopuštena, tako da se u obzir uzima pojava riječi i frekvencija pojavljivanja. Dokument je predstavljen riječima koje sadrži, ali se njihov redoslijed i konstrukcija rečenica ne uzimaju u obzir što dovodi do gubitka informacija. Riječ je bilo koji niz slova odvojen razmacima ili znakovima interpunkcije. Korišteni „izraz“ je sinonim za riječ. Torba, ili dokument, može se prikazati kao vektor izraza

$$\phi(d) = (\text{tf}(t_1, d), \text{tf}(t_2, d), \dots, \text{tf}(t_N, d)) \in \mathbb{R}^N,$$

u kojem svaki element $\text{tf}^{11}(t_i, d)$ označava broj pojavljivanja izraza t_i u dokumentu d . Obično $\phi(d)$ ima tisuće elemenata i često je taj broj veći od broja dokumenata u skupu za učenje. Za pojedini dokument ovakva reprezentacija je jako rijetka, ima samo nekolicinu nenul elementa, te se može koristiti i drukčiji, sažetiji zapisi.

Ako sa \mathbf{D} označimo dokument-izraz matricu čiji su redci dokument-vektori

$$\mathbf{D} = [\phi(d_1), \phi(d_2), \dots, \phi(d_l)]'$$

i analogna je prije korištenoj matrici \mathbf{X} , onda \mathbf{D}' predstavlja izraz-dokument matricu. Postoji direktna veza između matrica \mathbf{X} i \mathbf{D} gdje značajke postaju izrazi, a točke dokumenti. Također možemo definirati jezgrenu matricu

$$\mathbf{K} = \mathbf{DD}'$$

koja odgovara jezgri vektorskog prostora¹²

$$k(d_1, d_2) = \langle \phi(d_1), \phi(d_2) \rangle = \sum_{j=1}^N \text{tf}(t_j, d_1) \cdot \text{tf}(t_j, d_2)$$

Također je moguće koristiti nelinearna preslikavanja koristeći standardne jezgrene funkcije, npr. polinomialna jezgrena funkcija na normaliziranim podacima

⁹ Eng. Vector space model VSM

¹⁰ Eng. Bag-of-words

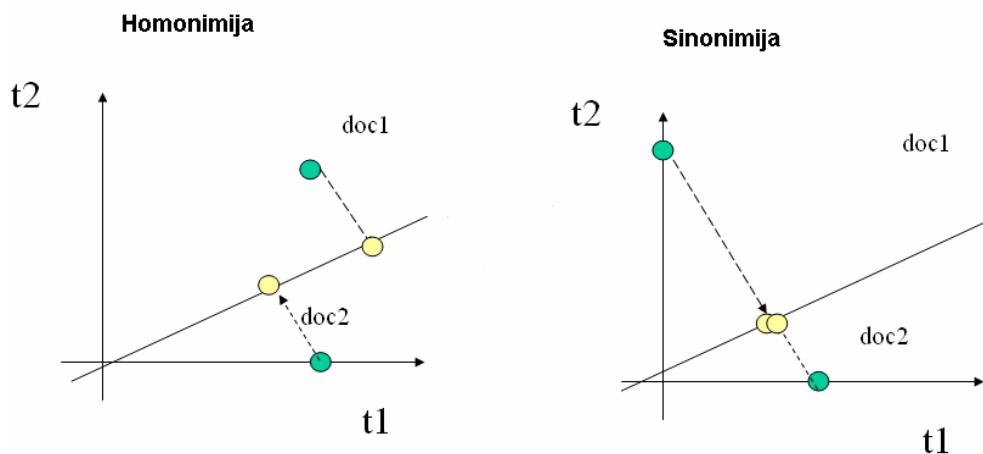
¹¹ Eng. Term frequency ili frekvencija pojavljivanja riječi u dokumentu

¹² Eng. Vector space kernel

$$k^*(d_1, d_2) = (k(d_1, d_2) + 1)^d = (\langle \phi(d_1), \phi(d_2) \rangle + 1)^d,$$

stupnja d koristi do n višerječnih izraza¹³ $0 \leq n \leq d$ za značajke. Iako se čini složenijom računski ne predstavlja veći problem od linearne. Također se može koristiti i bilo koja druga jezgrena funkcija.

Prikaz dokumenata u modelu vektorskog prostora ne uzima u obzir semantičke veze među izrazima, već izraze promatra kao izolirane pojave. Stoga su dokumenti koji govore o istoj temi, ali različitim izrazima preslikani u udaljena područja. Problemi sinonimije(različite riječi isto značenje) i homonimije(iste riječi različito značenje) mogu se riješiti dobro definiranim preslikavanjem koje udaljava dokumente koji sadrže homonime, a zbližava dokumente koji sadrže sinonime.



Slika 3. Homonimija i sinonimija

Neka rješenja semantičkih problema proizlaze iz same statistike riječi u dokumentu. Prvo je pridjeljivanje težina izrazima ili pridruženim koordinatama. Tu spada i uklanjanje stop riječi, riječi koje ne nose nikakvu informaciju o kategoriji te ih karakterizira velika frekvencija pojavljivanja. Definira se eksplisitna lista riječi koji će biti uklonjeni. Većinom su to prijedlozi, veznici i prilozi. Primjeri za hrvatski jezik su: *a, ah, ako, barem, bilo, četiri, ču, dokud, doskora, gdje, hoho, hopla, itekako, jesam, jutros, koga, manje, naprsto, nije, tebe, uzastopce, xxiii, zum, ...* Drugi problem je utjecaj dužine dokumenta. Što je duži dokument više riječi sadrži te norma dokument-vektora veća. Jezgrena funkcija koja odbacuje utjecaj norme dana je:

¹³ Eng. Monomials, multiword terms

$$\hat{k}(d_1, d_2) = \left\langle \frac{\phi(d_1)}{\|\phi(d_1)\|}, \frac{\phi(d_2)}{\|\phi(d_2)\|} \right\rangle = \frac{k(d_1, d_2)}{\sqrt{k(d_1, d_1)k(d_2, d_2)}}$$

Cilj nam je pronaći jezgenu semantičku funkciju koja uzima u obzir semantičke veze među riječima. Proučimo najjednostavnije linearno preslikavanje podataka $\tilde{\phi}(d) = \phi(d)\mathbf{S}$, gdje je \mathbf{S} matrica koja može biti kvadratna, dijagonalna, ili općenito $N \times k$ matrica. Jezgrena funkcija dobiva zatim novu formu

$$\tilde{k}(d_1, d_2) = \phi(d_1)\mathbf{S}\mathbf{S}'\phi(d_2)' = \tilde{\phi}(d_1)\tilde{\phi}(d_2)'.$$

Matrica \mathbf{S} zove se semantička matrica. Različiti izbori matrice vode k različitim varijantama modela vektorskog prostora¹⁴. Matrica \mathbf{S} može se definirati kao

$$\mathbf{S} = \mathbf{R}\mathbf{P}$$

gdje je \mathbf{R} dijagonalna matrica koja važe težine pojedinih riječi, i \mathbf{P} matrica bliskosti¹⁵ koja određuje semantičku blizinu različitih izraza.

Mjera *idf* važe težine pojedinih riječi kao funkciju inverzne frekvencije dokumenata. Ako prepostavimo da imamo l dokumenata, i $df^{16}(t)$ označava broj dokumenata koji sadrže izraz t , onda je mjera inverzne frekvencije dokumenata za izraz t dana sa

$$w(t) = \ln\left(\frac{l}{df(t)}\right)$$

Dakle, smanjuje se težine riječi koje se često pojavljuju te dolaze do izražaja diskriminantno bitne riječi. Za stop riječi, ako prepostavimo da se pojavljuju u skoro svim dokumentima vrijedi $df(t) = l$ te je $w(t) = 0$. Iako se na ovaj način implicitno možemo riješiti stop riječi ipak je poželjnije radi učinkovitosti definirati eksplisitnu listu izraza koji će biti uklonjeni. Matrica \mathbf{R} je dijagonalna sa elementima $R_{tt} = w(t)$. Jezgrena funkcija $\phi(d_1)\mathbf{R}\mathbf{R}'\phi(d_2)'$ koja koristi frekvenciju izraza kao i inverznu frekvenciju dokumenata često se naziva *tf-idf* reprezentacija.

¹⁴ Eng. Vector space model

¹⁵ Eng. Proximity matrix

¹⁶ Eng. Document frequency

Međutim *tf-idf* reprezentacija i dalje ne može prepoznati dva semantički povezana izraza. Semantička bliskost riječi uvodi se pomoću matrice \mathbf{P} od koje se zahtjeva da ima nedijagonalne elemente $P_{ij} > 0$ ako su izrazi i semantički povezani sa izrazom j . Ovakva jezgrena funkcija $\phi(d_1)\mathbf{P}\mathbf{P}'\phi(d_2)'$, odgovara reprezentaciji gušćih dokument-vektora $\phi(d)\mathbf{P}$ koji ima nenu null elemente za sve izraze semantički bliske izrazima koji se nalaze u dotičnom dokumentu.

Općeniti model vektorskog prostora¹⁷

Ova metoda nastoji uhvatiti izraz-izraz korelacije koristeći znanje o njihovom zajedničkom pojavljivanju u dokumentu. Dva izraza se smatraju semantički povezana ako se često pojavljuju u istim dokumentima. Prema ovome, dva dokumenta će biti slična i ako ne dijele zajedničke izraze nego izrazi koje sadrže često se zajedno pojavljuju u drugim dokumentima.

Dokument je predstavljen preslikavanjem:

$$\tilde{\phi}(d) = \phi(d)\mathbf{D}',$$

gdje je \mathbf{D} dokument – riječ matrica, što je jednako $\mathbf{P} = \mathbf{D}'$. Ovakvo preslikavanje nije odmah jasno, ali ako raspišemo jezgrenu funkciju

$$\tilde{k}(d_1, d_2) = \phi(d_1)\mathbf{D}'\mathbf{D}\phi(d_2)'$$

možemo primijetiti da matrica $\mathbf{D}'\mathbf{D}$ ima (i, j) ne nula element samo ako postoji dokument koji sadrži izraz i i izraz j ,

$$(\mathbf{D}'\mathbf{D})_{ij} = \sum_d \text{tf}(i, d) \text{tf}(j, d).$$

Ako je manje dokumenata od izraza ovakav postupak smanjuje dimenzionalnost preslikavanjem iz prostora veličine broja izraza u prostor veličine broja dokumenata. Iako zanimljiv, ovakav postupak je naivan u korištenju informacije zajedničkoj pojavljivanja.

¹⁷ GVSM eng. Generalised vector space model

Latentno semantičke jezgrena funkcije¹⁸

Latentno semantičko indeksiranje(LSI) slijedi isti princip kao GVSM, ali se razlikuju u tehnikama koje se koriste pri dobivanju informacija o zajedničkom pojavljivanju izraza. LSI koristi singular value decomposition (SVD).

Rastava matrice \mathbf{D}' je

$$\mathbf{D}' = \mathbf{U}\Sigma\mathbf{V}$$

gdje je Σ dijagonalna matrica, dok su \mathbf{U} i \mathbf{V} matrice čiji su stupci svojstveni vektori redom matrica $\mathbf{D}'\mathbf{D}$ i $\mathbf{D}\mathbf{D}'$. LSI preslikava dokumente u prostor razapet sa prvih k stupaca matrice \mathbf{U}

$$\mathbf{D} \rightarrow \phi(d)U_k.$$

Svojstvene vektore za skup dokumenta možemo zamisliti kao koncepte opisane linearnom kombinacijom izraza koji su odabrani na takav način da se dokumenti opisuju na najbolji način prikazani samo sa k takvih koncepata. Izrazi koji se često pojavljuju zajedno nastojat će se pojaviti u istom svojstvenom vektoru jer SVD ujedinjuje vrlo korelirane dimenzije da bi se dobio manji broj dimenzija sposobnih za rekonstrukciju cijelog vektora. Dakle, SVD koristi informacije o zajedničkom pojavljivanju izraza u sofisticiranim algoritmu.

Prema definiciji LSI-a vidimo da je jednak PCA algoritmu u prostoru značajki. Nova jezgrena funkcija postaje ona PCA-a algoritma

$$\tilde{k}(d_1, d_2) = \phi(d_1)\mathbf{U}_k\mathbf{U}_k' \phi(d_2)'$$

dok je matrica \mathbf{P} u ovom slučaju odabrana kao \mathbf{U}_k . Ako usporedimo sa GVSM te zapišemo njegovu jezgrenu funkciju u obliku:

$$\tilde{k}_{GVSM}(d_1, d_2) = \phi(d_1)\mathbf{D}'\mathbf{D}\phi(d_2)' = \phi(d_1)\mathbf{U}'\Sigma\mathbf{U}\phi(d_2)'$$

možemo vidjeti dvije prilagodbe. Prvo, provedena je redukcija dimenzionalnosti korištenjem samo k svojstvenih vektora te drugo, izbacivanjem matrice Σ projekcije podataka postaju ortonormalne.

Činjenica identičnosti LSI i PCA algoritma vodi do zaključka da se može također implementirati u dualnoj reprezentaciji ili

¹⁸ Eng. Latent semantic kernels

$$\phi(\mathbf{d})\mathbf{U}_k = \left(\lambda_j^{-\frac{1}{2}} \sum_{i=1}^l (\mathbf{v}_j)_i k(d_i, d) \right)_{j=1}^k,$$

gdje je k osnovna jezgrena funkcija koja može biti obična linearna jezgrena funkcija, ili složena jezgrena funkcija sa vaganjem težina izraza ili polinomialna jezgrena funkcija. Ovakva dualna reprezentacija naziva se latentno semantička jezgrena funkcija(LSK¹⁹) (Shawe-Taylor et al., 2004.).

Ako prikažemo matricu bliskosti

$$\mathbf{P} = \mathbf{U}_k \mathbf{U}_k'$$

vidimo da je ona kvadratna i određuje veze između različitih izraza. Kako se k povećava matrica \mathbf{P} teži jediničnoj matrici i jedinice samo na dijagonali pokazuje da su sve riječi semantički različite. Dakle, vrijednost k kontrolira količinu semantički bitnih informacija koje se uvode u reprezentaciju.

¹⁹ LSK eng. Latent semantic kernels

2.4 Istraživanje višejezičnih korpusa

Upareni korpsi²⁰ sastoje se od parova dokumenata koji su međusobni prijevodi na dva jezika. Oni su primjeri uparenih skupova za učenje korištenih pri CCA algoritmu. Prema prijevodu možemo se odnositi na dva načina: prvo, možemo ga smatrati kompleksnom oznakom za prvu verziju dokumenta, ili drugo: obje verzije možemo gledati kao dva pogleda na isti objekt.

Upareni korpus možemo koristiti pri učenju semantičkog preslikavanja koji će se koristiti samo na jednom prijevodu. Prijevodi se koristi kao složene oznake koje pronalaze semantičko preslikavanje potrebno za preslikavanje dokumenata u novi prostor bez semantički nebitni dijelova reprezentacije.

Latentno semantičke jezgrene funkcije mogu se također koristiti s ovakvim korpusima. Potrebno je spojiti dvije verzije dokumenata u jedan, višejezični tekst. Koristi se zatim PCA redukcija dimenzionalnost u k semantičkih dimenzija. Novi dokument bez prijevoda preslikava se u prostor svojstvenih vektora dobivenih spojenim dokumentima.

$$\mathbf{D} = \begin{pmatrix} D_x \\ D_y \end{pmatrix} = \mathbf{U}\Sigma\mathbf{V}^T$$

Latentno semantičke jezgrene funkcije imaju taj nedostatak da su ograničene spojem dvojezičnih dokumenata. Novi dokument bez prijevoda preslikava se u prostor svojstvenih vektora dobivenih spojem dokumenata tako da se i on proširi (dodaju se nule kao komponente drugog jezika). Zatim se novi dokument preslikava u prostor razapet sa k prvih svojstvenih vektora \mathbf{U}_k : $\mathbf{q}^* = \mathbf{U}_k^T \mathbf{q}$.

Alternativna metoda je korištenje CCA. Ova metoda tretira dvije verzije dokumenta kao dva pogleda na isti semantički objekt te traži dvije različite projekcije u zajednički semantički prostor. Najprije se odabere k semantičkih dimenzija, $0 \leq k \leq \ell$, sa najvećom vrijednosti korelacije i zatim se novi dokument, ili upit, preslikaju u novi prostor $\mathbf{q}^* = \mathbf{A}^T \mathbf{Z}^T \mathbf{q}$,

²⁰ Eng. Pared or Aligned corpora

gdje je \mathbf{A} $\ell \times k$ matrica čiji su stupi rješenja KCCA algoritma za jedan od jezika tj. jezik koji se koristi u novom dokumentu. \mathbf{Z} je skup za učenje.

Također se KCCA može koristiti i u kategorizaciji teksta. Semantički vektori u jednom od jezika $\mathbf{W} = \mathbf{ZA}$ mogu se prenijeti i u druge aplikacije, npr. metodu potpornih vektora. Nova jezgrena funkcija postaje

$$\tilde{k}(d_1, d_2) = \phi(d_1)\mathbf{WW}'\phi(d_2)'$$

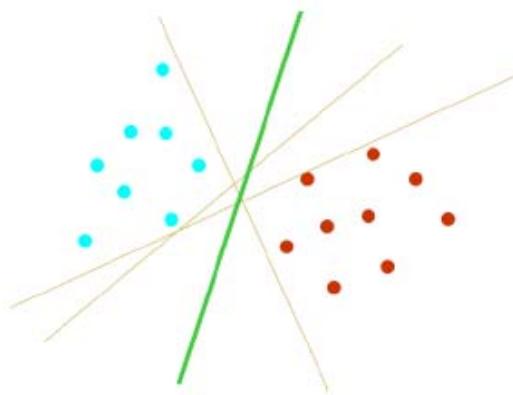
Matrica \mathbf{WW}' može se jedanput izračunati i spremiti za sljedeću upotrebu pri klasifikaciji podataka metodom potpornih vektora (Vinokourov et al., 2002.).

2.5 Metoda potpornih vektora

Nakon redukcija dimenzionalnosti zadatak nam postaje klasifikacija, tj. potrebno je odrediti kojem od dva razreda pripada testni uzorak. Skup za učenje možemo zapisati u obliku $\{ \mathbf{x}_i, \mathbf{y}_i \}, i = 1, \dots, n$ i $\mathbf{x}_i \in \mathbb{R}^d$, $y_i \in \{-1, +1\}$, gdje su \mathbf{x}_i ulazni vektori, a \mathbf{y}_i varijable odziva ili oznake.

Razmotrimo jednostavan primjer u kojem su podaci linearno razdvojivi tj. možemo nacrtati ravnu liniju $f(\mathbf{x}) = \mathbf{w}\mathbf{x} - b$ tako da se svi slučajevi sa $y_i = -1$ nalaze s jedne strane i za njih vrijedi $f(\mathbf{x}) < 0$, a slučajevi sa $y_i = +1$ nalaze se s druge strane i imaju vrijednost $f(\mathbf{x}) > 0$. Prema tome testne vektore možemo klasificirati prema njihovom predznaku funkcije.

Kako postoji beskonačno mnogo takvih hiperravnina koje rješavaju naš problem klasifikacije na skupu za učenje, postavlja se pitanje kako odabratи najbolju s obzirom na to da se sve razlikuju svojom klasifikacijom samo na testnim primjerima. Na primjer, možemo povući liniju vrlo blizu članova jednog razreda, recimo $y = -1$. Na dobivenim testnim primjerima nećemo imate velike pogreške sa slučajevima koji trebaju biti klasificirani kao $y = +1$, ali vrlo je velika vjerojatnost da ćemo pogriješiti sa slučajevima $y = -1$. Pametno bi bilo postaviti liniju čim dalje od $y = -1$ i $y = +1$ slučajeva iz skupa za učenje, tj. točno u njihovoј sredini.



Slika 4. Optimalna hiperravnina

Geometrijski, vektor \mathbf{w} usmjeren je ortogonalno na liniju definiranu sa $\mathbf{w}^T \mathbf{x} = b$. To se može protumačiti na sljedeći način. Za početak uzmimo $b = 0$. Sada je jasno da svi vektori, \mathbf{x} , kojima skalarni produkt sa \mathbf{w} nestaje zadovoljavaju ovu jednadžbu tj. svi vektori ortogonalni na \mathbf{w} zadovoljavaju ovu jednadžbu. Sada zamislimo da pomicemo hiperravninu iz ishodišta za neki vektor \mathbf{a} . Jednadžba ravnine postaje $(\mathbf{x} - \mathbf{a})^T \mathbf{w} = 0$ uz odmak $b = \mathbf{a}^T \mathbf{w}$ ili projekciju vektora \mathbf{a} na vektor \mathbf{w} . Bez gubitka na općenitosti možemo odabratи \mathbf{a} okomit

na hiperravninu, te u tom slučaju dužina $\|a\| = |b| / \|\mathbf{w}\|$ predstavlja najmanju udaljenost ishodišta i hiperravnine.

Definirajmo dvije nove hiperravnine paralelne sa ravninom kojom dijelimo grupe. One pretstavljaju ravnine koje prolaze najbližim primjerima za učenje sa obje strane. Nazivat ćemo ih potpornim hiperravninama²¹ jer podatkovni vektori koje sadrže podupiru ravninu.

Definirajmo i udaljenost tih hiperravnina od glavne, dijeleće, kao $d+$ i $d-$. Margina γ je definirana kao $d+ + d-$. Cilj nam je pronaći marginu koja je maksimalna, dok su podupiruće hiperravnine jednako udaljene od glavne. Za podupiruće hiperravnine možemo zapisati sljedeće jednadžbe:

$$\mathbf{w}^T \mathbf{x} = b + \delta$$

$$\mathbf{w}^T \mathbf{x} = b - \delta$$

Primjećujemo kako je problem pre-parametriziran te umjesto \mathbf{w} , b i δ možemo koristiti parametar α tako da i dalje jednadžba za \mathbf{x} bude zadovoljena. Da se riješimo ove dvosmislice potrebno je da $\delta = 1$ tj. postavljamo skalnu problema.

Zatim možemo i izračunati vrijednosti za $d+ = (|b+1| - |b|) / \|\mathbf{w}\| = 1 / \|\mathbf{w}\|$ (ovo je točno samo u slučaju ako $b \notin (-1,0)$ jer tada ishodište ne pada između hiperravnina. Ako je $b \notin (-1,0)$ potrebno je koristiti $d+ = (|b+1| + |b|) / \|\mathbf{w}\| = 1 / \|\mathbf{w}\|$). Dakle, margina je dva puta ta vrijednost $\gamma = 2 / \|\mathbf{w}\|$.

Sa gornjim definicijama možemo zapisati uvjete koje rješenje mora zadovoljiti,

$$\mathbf{w}^T \mathbf{x}_i - b \leq -1 \quad \forall y_i = -1$$

$$\mathbf{w}^T \mathbf{x}_i - b \geq +1 \quad \forall y_i = +1$$

ili u jednadžbi

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 \geq 0$$

Formulirajmo primarni problem SVM-a:

$$\begin{aligned} & \text{minimizirati} \quad \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{uz uvjet} \quad y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 \geq 0 \quad \forall i \end{aligned}$$

²¹ Eng. Support hiper-plane

ili riječima – nastojimo maksimizirati marginu, uz uvjet da svi slučajevi iz skupa za učenje padaju s vanjske strane potpornih hiperravnina. Slučajevi koji leže na potpornoj hiperravnini nazivaju se potporni vektori²² jer podupiru hiperravninu i dakle određuju rješenje problema.

Primarni problem može se riješiti kvadratnim programom. Ipak, ovako formulirani problem ne možemo riješiti jezgrenim načinom jer ne ovisi samo o skalarnom produktu podatkovnih vektora. Transformirat ćemo ga u dualni oblik tako da ga prvo zapišemo koristeći Lagrangiana,

$$\mathcal{L}(\mathbf{w}, \mathbf{b}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w}^\top \mathbf{x}_i - b) - 1]$$

Rješenje koje minimizira primarni problem uz dani uvjet dano je sa $\min_w \max \alpha \mathcal{L}(\mathbf{w}, \alpha)$ tj. dolazimo do problema sedla. Kad je originalna objektivna funkcija konveksna, (i samo tada), možemo zamijeniti minimizaciju i maksimizaciju. Koristeći ovo, dolazimo do novog uvjeta za vektor \mathbf{w} koji mora vrijediti u točci sedla. Derivirajući po \mathbf{w} i b dolazimo do

$$\mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0 \quad \Rightarrow \quad \mathbf{w}^* = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$\sum_i \alpha_i y_i = 0$$

Uvrštavanjem natrag u Lagrangian problem postaje dualni

$$\begin{aligned} \text{maksimizacija } \mathcal{L}D &= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ \text{uz uvjet } \sum_i \alpha_i y_i &= 0 \\ \alpha_i &\geq 0 \forall i \end{aligned}$$

Dualna formulacija je također kvadratični program, ali broj varijabli, α_i u ovom slučaju je jednak broju podatkovnih vektora, N . Ključna stvar je ipak da ovaj problem ovisi o \mathbf{x}_i samo preko skalarnog produkta $\mathbf{x}_i^\top \mathbf{x}_j$ te se uvodi supstitucija $\mathbf{x}_i^\top \mathbf{x}_j \rightarrow k(\mathbf{x}_i, \mathbf{x}_j)$.

Teorija dualnost jamči da će za konveksni problem, dualni problem biti konkavan, i nadalje da jedinstveno rješenje primarnog problema odgovara jedinstvenom rješenju dualnog problema. Vrijedi $\mathcal{L}P(\mathbf{w}^*) = \mathcal{L}D(\alpha^*)$.

²² Eng. Support vector

Nadalje moramo proučiti uvjete koji vrijede u točci sedla dakle moraju vrijediti i za rješenje. Zovu se KKT²³ uvjeti. Ovi uvjeti su općenito potrebni i dovoljni za problem konveksne optimizacije. Dobivaju se deriviranjem primarnog problema. Također potrebni su nam postavljeni uvjet nejednakost i ograničenje Lagrangeovih multiplikatora na nenegativnost. Te posljednje, važan "komplementarni uvjet" treba biti zadovoljen

$$\partial_w \mathcal{L}P = 0 \quad \rightarrow \quad \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0$$

$$\partial_b \mathcal{L}P = 0 \quad \rightarrow \quad \sum_i \alpha_i y_i = 0$$

$$y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 \geq 0 \quad \forall i$$

$$\alpha_i \geq 0 \forall i$$

$$\alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1] = 0$$

Posljednja jednadžba može biti iznenađujuća. Ukazuje na dva slučaja, ili vrijedi uvjet nejednakosti, ali nije zasićen: $y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 > 0$ te u tom slučaju vrijednosti za α_i moraju biti 0. Ako je uvjet nejednakost zasićen $y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 = 0$ vrijednost α_i mogu poprimiti bilo koju vrijednost $\alpha_i \geq 0$. Zasićen uvjet nejednakosti zovemo „aktivan“, dok ne zasićen nazivamo „neaktivan“. Svi slučajevi iz skupa za učenje sa $\alpha_i > 0$ predstavljaju aktivne uvjete na potpornoj hiperravnini te se zovu potporni vektori. Obično ih ima malo, i zovu ih još „rijetka“ rješenja (većina α nestaje).

Zainteresirani smo i za funkciju $f(\cdot)$ koju možemo koristiti za klasifikaciju:

$$f(\mathbf{x}) = \mathbf{w}^{*\top} \mathbf{x} - b^* = \sum_i \alpha_i y_i \mathbf{x}_i^T \mathbf{x} - b^*$$

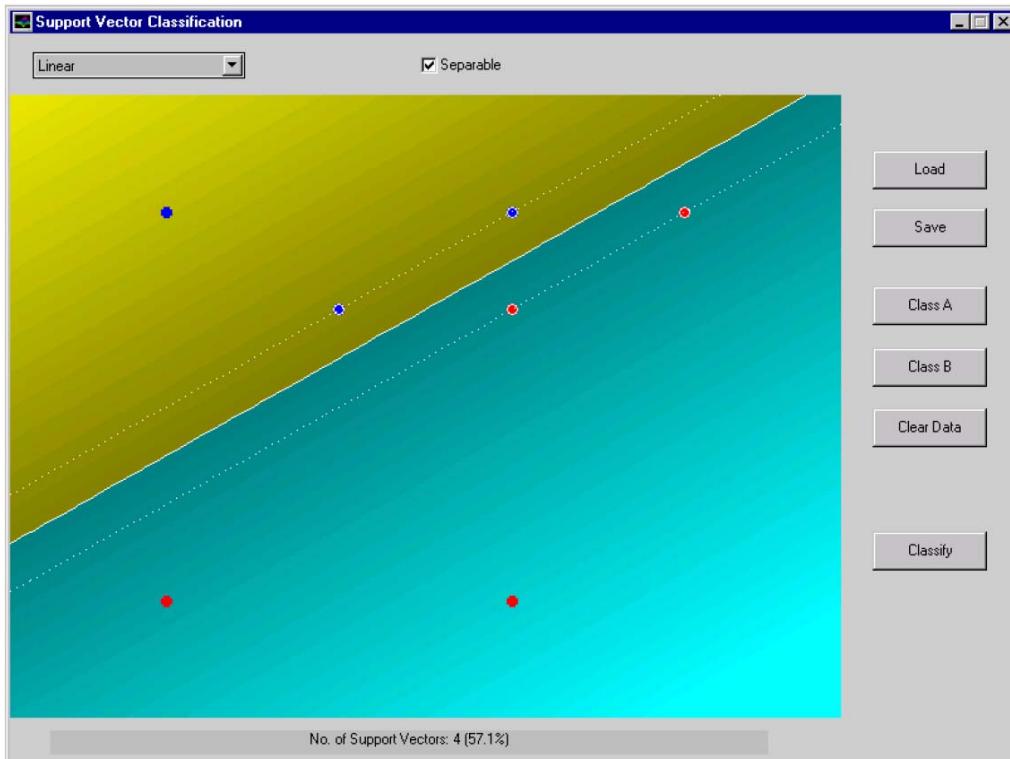
Primjenom KKT uvjeta dolazimo do rješenja za b^* koristeći „komplementarni uvjet“

$$b^* = \left(\sum_j \alpha_j y_j \mathbf{x}_j^T \mathbf{x} - y_i \right) \quad i \text{ je potporna vektor}$$

²³ Kratica stoji umjesto Karush-Kuhn-Tucker uvjeta

gdje smo koristili $y_i^2 = 1$. Koristeći jedan potporni vektor može se odrediti b , ali zbog numeričke stabilnosti bolje je koristiti njihov prosjek.

Najvažniji zaključaj je opet da se funkciju $f(\cdot)$ može izraziti samo preko skalarног produkta $\mathbf{x}_i^T \mathbf{x}_i$ koji se može zamijeniti sa jezgrenom matricom $k(\mathbf{x}_i, \mathbf{x}_j)$ za pomak u više dimenzionalni nelinearni prostor. Nadalje, kako je α rijetko, ne moramo koristiti puno elemenata jezgrene matrice za predviđanje klasifikacije novog ulaza \mathbf{x} .



Slika 5. Primjer optimalne hiperravnine sa SV-ima na marginama²⁴

Linearno neodvojivi slučajevi

Svi skupovi nisu linearno odvojivi te prema tome trebamo promijeniti gornje formulacije. Očito je da problem leži u uvjetu koji ne može uvijek biti ispunjen. Promijenimo taj uvjet koristeći „olabavljenu varijablu“ ξ_i

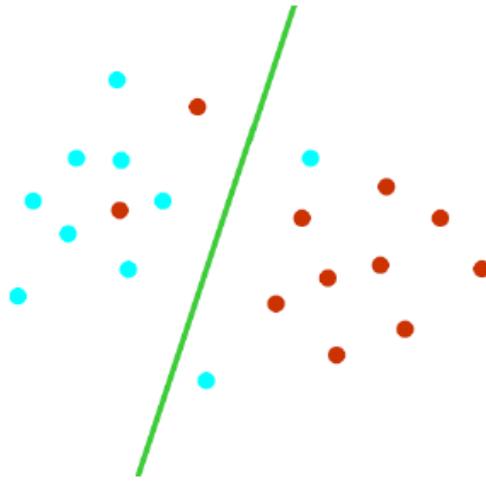
$$\mathbf{w}^T \mathbf{x}_i - b \leq -1 + \xi_i \quad \forall y_i = -1$$

²⁴ Slika je dobivena pomoću sučelja MATLAB SVM Toolbox, Steve Gunn

$$\mathbf{w}^T \mathbf{x}_i - b \geq +1 - \xi_i \quad \forall y_i = +1$$

$$\xi_i \geq 0 \forall i$$

Varijable ξ_i omogućuju kršenje uvjeta, ali potrebno je i njih kazniti – kako se ξ_i odabire sve većim, gornji uvjet postaje beskoristan. Funkcija kazne, u obliku $C(\sum_i \xi_i)^k$ vodi konveksnom optimizacijskom problemu za pozitivne cijelobrojne k . Za $k = 1, 2$ još je uvijek kvadratni program (QP). Odabrat ćemo $k = 1$. C kontrolira razmjenu između margine i kazne.



Slika 6. Općenita optimalna dijeleća hiperravnina

Da bude s krive strane potporne hiperravnine, podatkovnom slučaju potrebno je $\xi_i > 1$. Dakle, sumu $\sum_i \xi_i$ možemo interpretirati kao mjeru koliko su prekršaji teški, te je gornja granica za broj prekršaja.

Novi primarni problem postaje

$$\begin{aligned} \text{minimizirati} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i \\ \text{uz uvjet} \quad & y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 + \xi_i \geq 0 \quad \forall i \\ & \xi_i \geq 0 \forall i \end{aligned}$$

te vodi do Lagrangiana

$$\mathcal{L}(\mathbf{w}, b, \xi_i, \alpha_i, \mu_i) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 + \xi_i] - \sum_{i=1}^N \mu_i \xi_i$$

Iz čega izvodimo KKT uvjete

1. $\partial_w \mathcal{L}P = 0 \rightarrow \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i = 0$
2. $\partial_b \mathcal{L}P = 0 \rightarrow \sum_i \alpha_i y_i = 0$
3. $\partial_\xi \mathcal{L}P = 0 \rightarrow C - \alpha_i - \mu_i = 0$
4. uvjet 1 $y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 + \xi_i \geq 0 \quad \forall i$
5. uvjet 2 $\xi_i \geq 0 \forall i$
6. uvjet multiplikatora 1 $\alpha_i \geq 0 \forall i$
7. uvjet multiplikatora 2 $\mu_i \geq 0 \forall i$
8. „komplementarni uvjet“ 1 $\alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 + \xi_i] = 0$
9. „komplementarni uvjet“ 2 $\mu_i \xi_i = 0$

Iz ovih uvjeta dedukcijom dolazimo do sljedećega. Ako prepostavimo da je $\xi_i > 0$, tada je $\mu_i = 0$, dakle $\alpha_i = C(1)$ i $\xi_i = 1 - y_i (\mathbf{w}^T \mathbf{x}_i - b)$. Također kad je $\xi_i = 0$ slijedi $\mu_i > 0$ i $\alpha_i < C$. Nadalje za $\xi_i = 0$ vrijedi i $y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 = 0$, te $\alpha_i > 0$. U drugom slučaju, ako je $y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 > 0$ tada je $\alpha_i = 0$. Kao i prije, za podatke koji nisu potporni vektori i s prave strane ravnine vrijedi $\alpha_i = \xi_i = 0$ (svi uvjeti su neaktivni). Na potpornoj ravnini, još uvijek vrijedi $\xi_i = 0$, ali je sad $\alpha_i > 0$. Napokon, za podatke sa krive strane potporne hiperravnine α_i teži $\alpha_i = C$, a ξ_i uravnotežuje kršenje uvjeta tako da vrijedi $y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 + \xi_i = 0$.

Geometrijski možemo izračunati razmak između potporne ravnine i slučaja koji krši uvjet kao $\xi_i / \|\mathbf{w}\|$. Kako je ravnina definirana sa $y_i (\mathbf{w}^T \mathbf{x}_i - b) - 1 + \xi_i = 0$ paralelna sa potpornom hiperravninom te je udaljena od ishodišta za $|1 + y_i b - \xi_i| / \|\mathbf{w}\|$, a potporna ravnina je udaljena $|1 + y_i b| / \|\mathbf{w}\|$, dolazimo do gornje udaljenosti.

Naposlijetku, trebamo dualizirati problem zbog učinkovitog rješavanja pomoću jezgrenih metoda. Opet koristimo KKT uvjete da se riješim \mathbf{w} , b i ξ ,

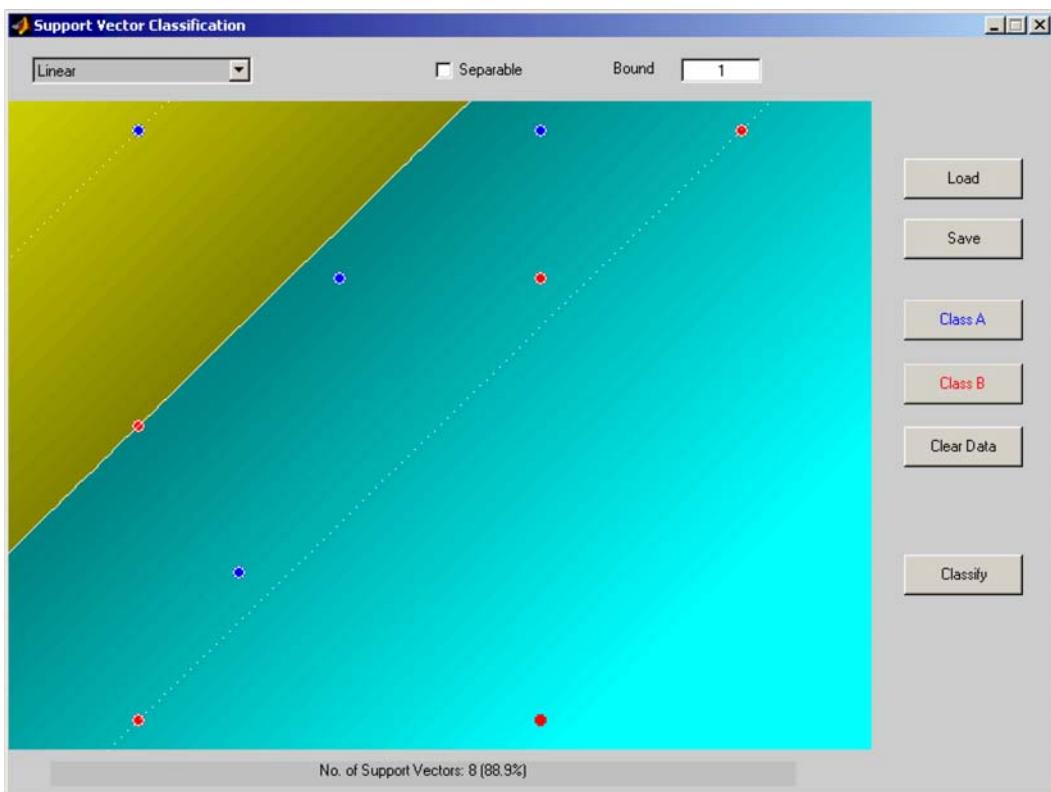
$$\text{maksimizacija } \mathcal{L}_D = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

$$\text{uz uvjet} \quad \sum_i \alpha_i y_i = 0$$

$$0 \leq \alpha_i \leq C \forall i$$

Iznenađujuće je da je to skoro isti kvadratni problem(QP) kao i prije, ali sa dodatnim uvjetom na multiplikatoru α_i koji je sad zatvoren u kutiju. Uvjet je izведен iz $\alpha_i = C - \mu_i$ te $\mu_i \geq 0$. Napomenimo još i da ovisi samo o skalarnom produktu vektora $\mathbf{x}_i^T \mathbf{x}_j$.

Parametar C uvodi dodatnu kontrolu nad klasifikatorom te ga je potrebno odrediti koristeći znanje o šumu u podacima. Uz $C = 1$ potporni vektori više ne trebaju ležati na margini te se orientacija hiperravnine i širina margine mijenjaju.

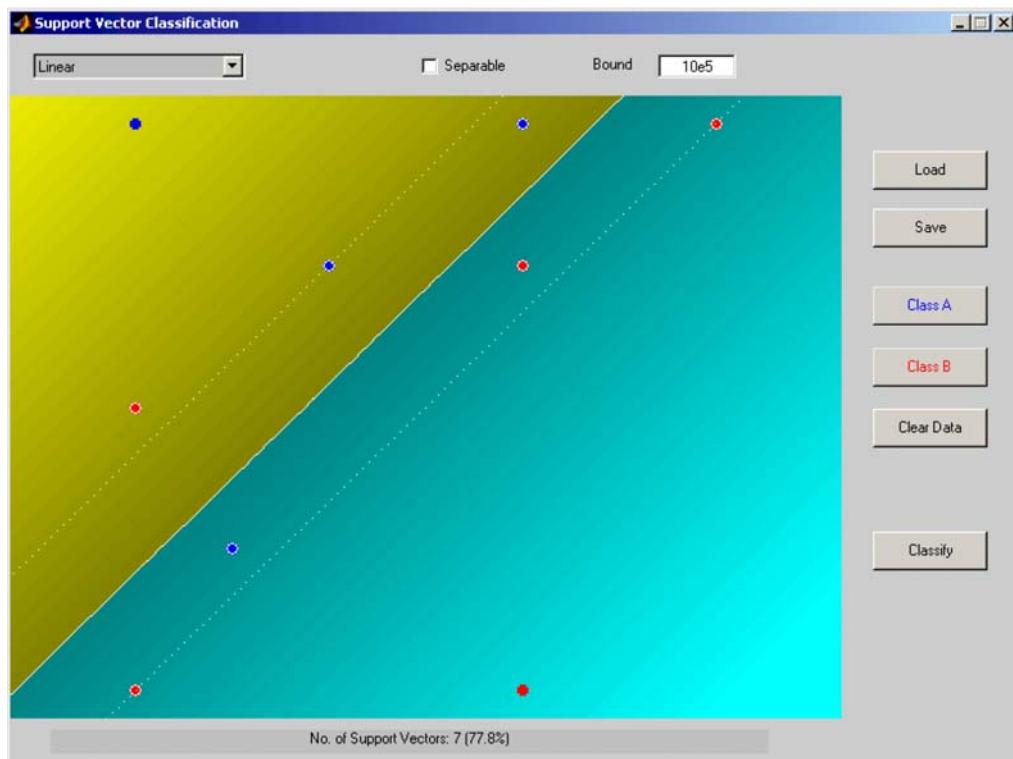


Slika 7. Primjer optimalne hiperravnina sa SV-ima na marginama ($C = 1$)

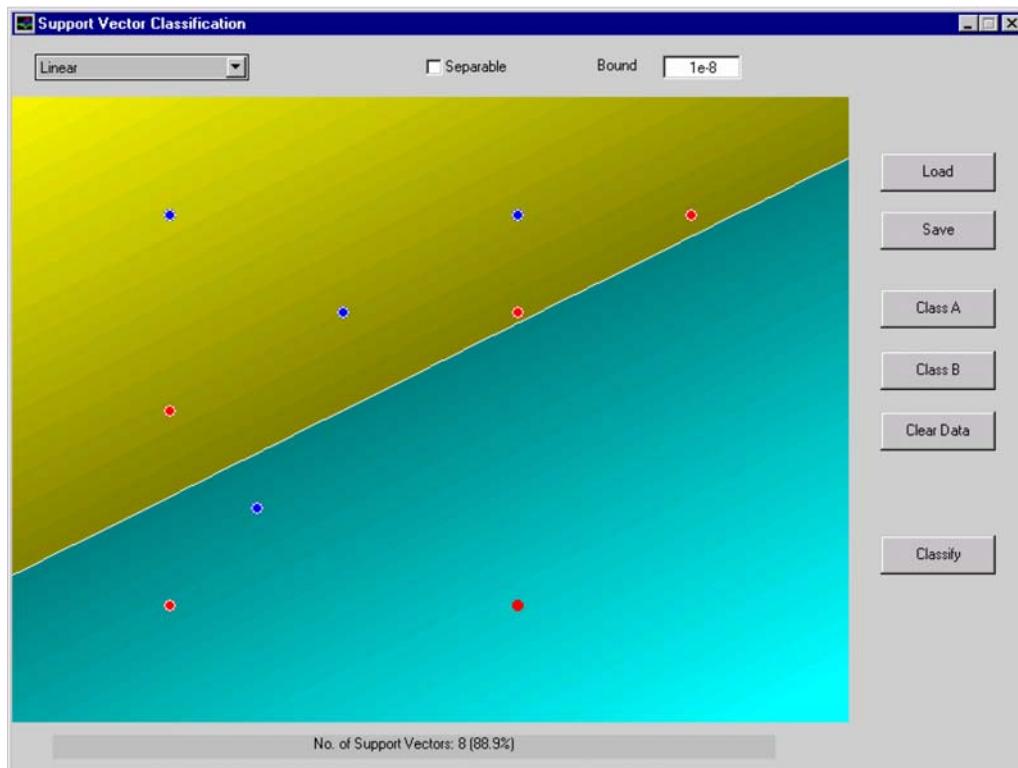
U ograničenje kad $\lim C \rightarrow \infty$ rješenje teži optimalnoj hiperravnini dobivenoj na linearno nerazdvojivim podacima.

U ograničenje kad $\lim C \rightarrow 0$ rješenje teži maksimiziranju margine, te će nakon nekog vremena svi Lagrangovi multiplikatori poprimiti vrijednost C. Klasifikator daje naglasak na minimizaciju pogrešno klasificiranih slučajeva, ali na taj način da povećava marginu, stvara marginu velike širine. Kao posljedica smanjenja C-a veličina margine se povećava.

Za određivanje parametra C koristi se postupak krosvalidacije.



Slika 8. Primjer optimalne hiperravnine ($C = 10^5$)



Slika 9. Primjer optimalne hiperravnine ($C = 10^{-8}$)

2.5.1 Upotreba metode potpornih vektora

Klasifikator metodom potpornih vektora koristimo u sljedećim koracima:

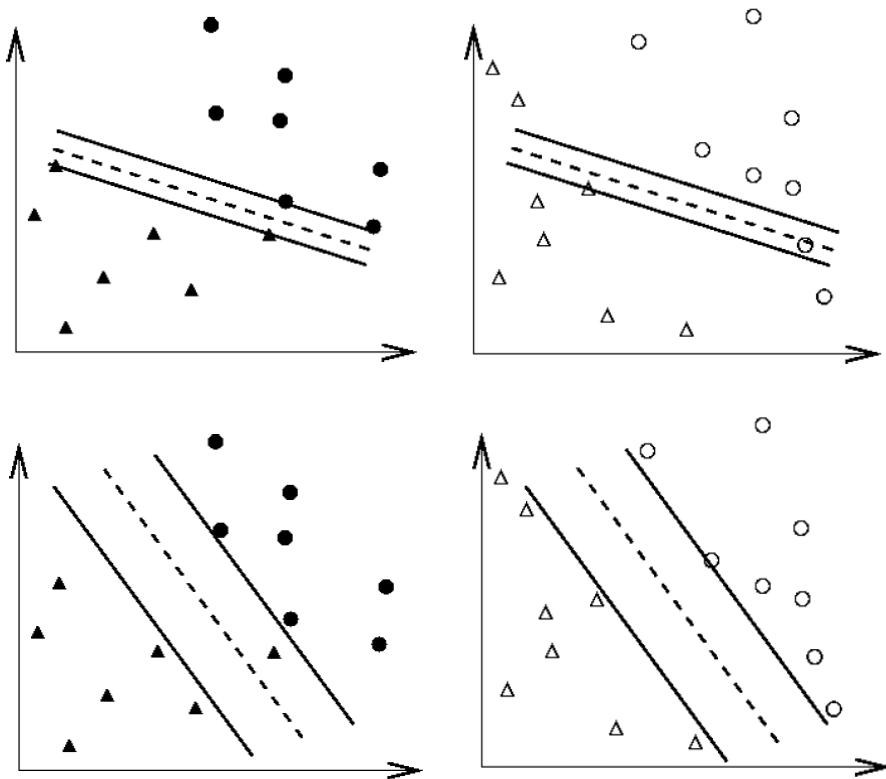
- Zapišemo podatke u formatu pogodnom za SVM program
- Provedemo jednostavno skaliranje podataka(ako je potrebno)
- Razmotrimo jezgrene funkcije, npr. RBF
- Koristeći krosvalidaciju dolazimo do najboljih parametara C i γ
- Koristeći najbolje parametre C i γ učimo na cijelom skupu za učenje²⁵
- Testiranje

Skaliranje podataka je vrlo važno, kao i skaliranje pri korištenju neuronskih mreža, tako i kod SVM-a. Najveća je prednost izbjegavanje dominacije atributa sa velikim numeričkim rasponom prema onima sa malim rasponom. Druga prednost je i izbjegavanje numeričkih teškoća pri izračunavanju. Kako jezgrene vrijednosti ovise o skalarnom produktu preslikanih vektora velike vrijednosti atributa moguće bi uzrokovati numeričke probleme. Preporuča se linearno skaliranje [-1, +1] ili [0, 1]. Također potrebno je skalirati i testni skup na isti način.

Pri odabiru jezgrene funkcije, RBF je prvi izbor. RBF nelinearno preslikava uzorke u višedimenzionalni prostor, tako da, za razliku od linearog jezgrenog modela, može riješiti i slučajeve kad su uzorci linearne neodvojivi. Nadalje, linearna jezgrena funkcija je samo specijalan slučaj RBF-a te se može dokazati da linearna jezgrena funkcija sa parametrom kazne C^* ima iste performanse kao RBF sa (C, γ) . I sigmoidalna jezgrena funkcija se ponaša kao RBF za određene parametre. Sljedeći razlog u odabiru stavlja RBF ispred polinomialne jezgrene funkcije jer ona koristi više hiperparametara koji utječu na složenost rješavanja. RBF se također lakše izračunava. Vrijednosti jezgrene funkcije su $0 \leq K_{ij} \leq 1$ za razliku od polinomialne funkcije gdje vrijednosti idu do beskonačnosti. Još moramo napomenuti da sigmoidalna jezgrena funkcija nije ispravna za neke parametre (Chang et al., 2005.).

²⁵ Najbolji parametri mogu biti neprilagođeni za cijeli skup za učenje, ali u praksi oni parametri dobiveni krosvalidacijom pogodni su za cijeli skup za učenje.

Pri korištenju RBF-a potrebno je podesiti dva parametra: C i γ . Potrebno je odabrati takve parametre koji čim točnije predviđaju oznake testnih podataka, ali to nužno ne znači da je potrebno dobiti veliku točnost na skupu za učenje, tj. točno predviđanje oznaka na skupu na kojem su one već poznate. Obično se skup za učenje dijeli na dva dijela u kojem jedan postaje nepoznat pri učenju klasifikatora. Točnost na ovakvom skupu preciznije govori o klasificiranju nepoznatih podataka. Ovakva procedura zove se krosvalidacija. Općenito se krosvalidacija izvršava na sljedeći način: podijeli se skup za učenje na v jednakih dijelova; sekvensijalno jedan podskup se testira koristeći klasifikator dobiven na preostalih $v - 1$ podskupu. Dakle, svaka instanca skupa za učenje predviđa se samo jedanput tako da je krosvalidacijska točnost postotak točno klasificiranih podataka. Također krosvalidacija može spriječiti pretreniranost. Na sljedećoj slici puni krugovi i trokuti predstavljaju skup za učenje, dok prazni skup za testiranje. Prvi klasifikator postiže loše rezultate na skupu za testiranja jer je pretreniran na skupu za učenje. Također ni točnost pri krosvalidaciji nije visoka. Drugi klasifikator bez pretreniranosti, daje bolje krosvalidacijske rezultate kao i testnu točnost (Chang et al., 2005.).



Slika 10. Pretrenirani klasifikator i bolji klasifikator

Tipični parovi C i γ isprobavaju se te se odabire onaj s najvećom krosvalidacijskom točnošću. Pokazalo se da je eksponencijalno rastuća sekvenca C i γ praktična metoda pri njihovom određivanju (npr. $C = 2^{-5}, 2^{-3}, \dots, 2^{15}$, $\gamma = 2^{-15}, 2^{-13}, \dots, 2^3$)

2.5.2 Biblioteka funkcija LIBSVM te OSU SVM

Prilikom implementacije metode potpornih vektora, koristištena je javno dostupna biblioteka funkcija LIBSVM²⁶ verzija 2.6 objavljena u travnju 2004., pisana u programskom jeziku C++ te Javi. Detalj o njoj mogu se pronaći na [<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>]. LIBSVM biblioteka implementira više tipova klasifikatora. Ova implementacija koristi tip nazvan C-SVC koji odgovara klasifikatoru sa slabom marginom. Korištena je RBF jezgrena funkcija.

Također pri klasifikaciji podataka reduciranih pomoću PCA-a i CCA-a korišteno je MATLAB sučelje na LIBSVM – OSU SVM, autori su Junshui Ma i Stanley Ahalt sa Ohio State University. Prema autorima, dovoljno je brz da može konkurirati svim SVM implementacijama u C-u, te nasleđuje LIBSVM – upravljanje memorijom što ga čini upotrebljivim i na ogromnim skupovima za učenje.

Pri dobivanju krosvalidacijskih parametara C i γ korišten je program grid.py napisan u programskom jeziku Python. Program grid.py nalazi se u na stranicama LIBSVM – a.

²⁶ Autori ove biblioteke su Chih-Chung Chang i Chih-Jen Lin, Department of Computer Science and Information Engineering, National Taiwan University.

2.5.3 Klasifikacija u više klase

Prethodno opisana metoda potpornih vektora klasificira podatke u dvije klase sa oznakama $y_i \in \{-1, +1\}$. Za višeklasno klasificiranje odabrana je metoda „jedan-prema-jedan“²⁷(Hsu et al.,2002.) u kojoj se stvara $k(k-1)/2$ klasifikator koji odvaja po dvije različite klase (Chang et al., 2005.).

U klasifikaciji se koristi strategija glasovanja, pri čemu svaki binarni klasifikator daje svoj glas. Svaki ulazni vektor klasificira se prema najvećem broju glasova. Za slučajevе u kojima dvije klase imaju jednak broj glasova, iako se ne smatra dobra tehnikom, odabire se klasa sa manjim indeksom.

²⁷ Eng. One-against-one

3 REZULTATI

3.1 Ulazni podaci

Kao ulazni podaci na raspolaganju su dvije baze, Vjesnik – baza članaka iz hrvatskih dnevnih novina Vjesnik u razdoblju od 2000. – 2003., te Croatia Weekly (brojevi 5-118) u razdoblju od 1998.-2000.

3.1.1 Croatia Weekly

Bazu članaka iz Hrvatsko-engleskoga paralelnoga korpusa za potrebe projekta prof. dr. Bojane Dalbelo-Bašić, Croatia Weekly, izradio je prof. Marko Tadić sa Zavoda za lingvistiku Filozofskoga fakulteta Sveučilišta u Zagrebu. Izvor su novine Croatia Weekly, brojevi 5-118, koji su izlazili u razdoblju od 1998.- 2000.

Baza sadrži 4 kategorije:

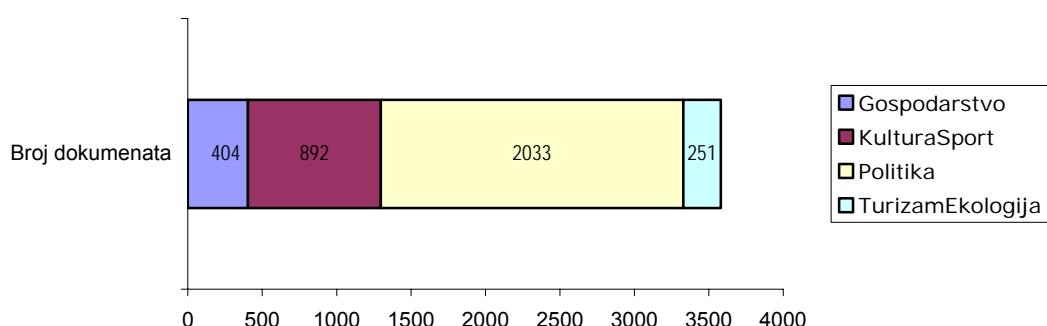
po - politika

go - gospodarstvo

te - turizam i ekologija

ks - kultura i sport

U korpusu se nalazi 3.580 članaka na hrvatskom jeziku i isto toliko članaka na engleskom jeziku. Članci na hrvatskom jeziku nisu normalizirani – iz njih su samo izbačene stop riječi. Jedan od problema ove baze je nejednako raspoređen broj dokumenata po kategorijama. Skup za učenje i skup za testiranje svaki sadrže 50% dokumenta svake kategorije. Za usporedbu, skup za učenje sadrži 1017 dokumenata politike, ali samo 128 dokumenata turizma i ekonomije.



Slika 11. Broj dokumenata po kategorijama za hrvatsko-engleski korpus.

3.1.2 Baza novinskih članaka Vjesnik

Baza Vjesnikovih članaka također je izradio prof. Marko Tadić sa Zavoda za lingvistiku Filozofskoga fakulteta Sveučilista u Zagrebu. Baza sadrži sve novinske članke (Vjesnik između 2000. i 2003. godine), preuzete iz Hrvatskog nacionalnog korpusa [<http://www.hnk.ffzg.hr/>].

Sažeta baza se sastoji od 10.000 dokumenata podijeljenih u kategorije:

ck - crna kronika

go - gospodarstvo

ku - kultura

sp - sport

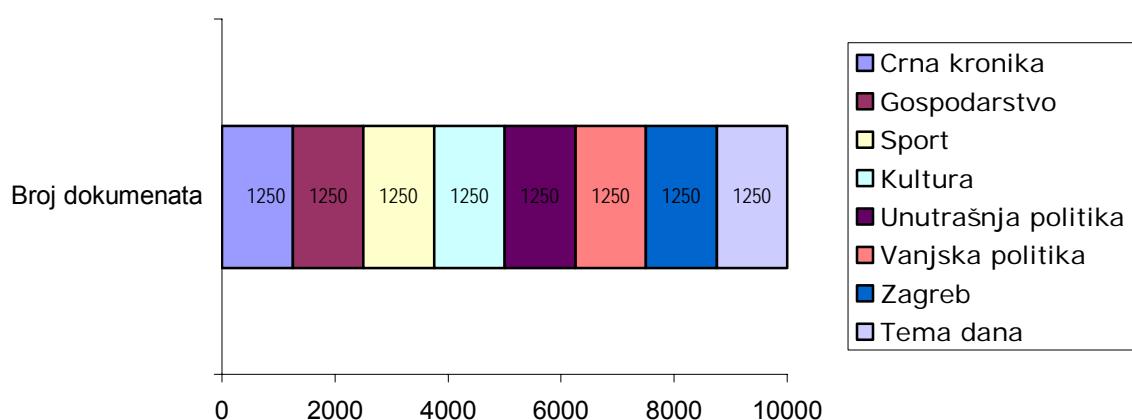
td - tema dana

un - unutarnja politika

vp - vanjska politika

zg - zagreb.

U svakoj kategoriji nalazi se 1.250 dokumenata.



Slika 12. Broj dokumenata po kategorijama za Vjesnikovu bazu članaka

3.1.3 Baza novinskih članaka Vjesnik AMNv1.0

Riječ je o bazi novinskih članaka Vjesnik na kojoj je provedena morfološka obrada postupkom Automatske morfološke normalizacije koju je razvio dipl. ing. Jan Šnajder, asistent na Zavodu za elektroniku, mikroelektroniku, računalne i inteligentne sustave Fakulteta elektrotehnike i računarstva.

Automatska morfološka normalizacija (AMN) je postupak kojim se pojavnice u tekstu svode na svoje morfološke norme. Osim što na taj način smanjujemo dimenzionalnost ulaznih podataka, povećavamo kvalitetu teksta koji želimo kategorizirati jer uklanjamo rasipanje značenja istog pojma na više leksički različitih oblika.

Postoje 3 vrste morfološke normalizacije:

1. *Flektivna*

Korištena je oznaka “-i” pri označavanju pripadnih baza (eng. Inflective). Eliminira efekt flektivne norme na način da sve oblike neke riječi svodi na jedan oblik, lemu ako je nju moguće automatski odrediti, inače odgovara kombinaciji više lema.

2. *Derivacijska*

Korištena je oznaka “-id” pri označavanju pripadnih baza. Nakon suođenja različitih oblika riječi na njihovu lemu, nalazi zajedničkog predstavnika više lema, prema derivacijskim pravilima morfologije. Derivacijska pravila morfologije opisuju tvorbu riječi (iz jedne leme u drugu) eventualno mijenjajući vrstu riječi. Smisao derivacijskih normi jest ostvarivanje još veće redukcije dimenzionalnosti ulaznih podataka od one koje se postiže primjenom infleksijskih pravila. Pored toga postiže se i veća koncentracija značenja nekog pojma. Npr. riječi “kompjuter” i “kompjuterski” imaju vrlo blisko značenje (vezane su uz isti pojam), zato je sigurno poželjno svesti ih na zajednički oblik prije samog postupka klasifikacije.

3. *Terminirajuća*

Korištena je oznaka “-idt” pri označavanju pripadnih baza. U slučajevima kad ne djeluju derivacijska pravila onda se reže riječ nakon 7-og slova.

Npr. neka je zadana riječ “kompjuterskog”. To je pridjev, i njegova lema je nominative jednine muškog roda: “kompjuterski”. To bi ujedno bila i infleksijska norma, ako je AMN postupkom ispravno pronađena. Osnova riječi “kompjuterski” je “kompjuter”, pa je derivacijska norma “kompjuter”.

Ovom se metodom broj različitih pojavnica u prosjeku smanjio na oko 45% početnog broja pojavnica.

Više informacije o ovoj metodi dostupno je na Internet adresi [<http://www.zemris.fer.hr/~jan/amn>].

3.2 Rezultati usporedbe PCA i CCA

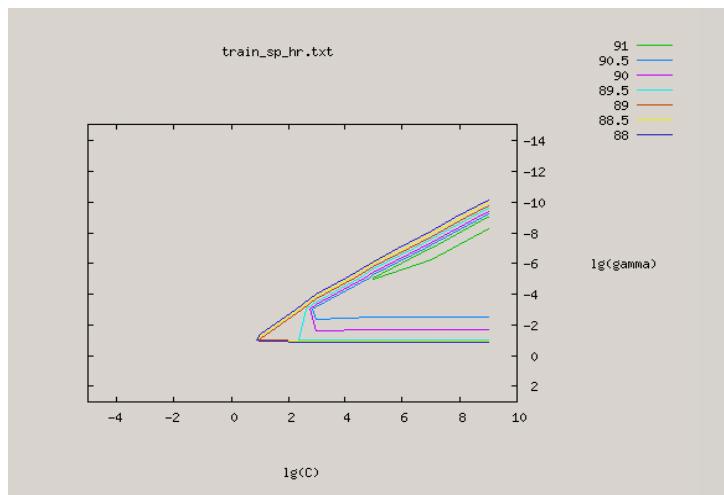
PCA i CCA algoritam implementirani su u MATLAB-u. Korištena je baza Croatia Weekly, hrvatsko-engleski paralelni korpus.

Bitno je spomenuti MATLAB-ovu funkciju `eigs` koja se povezuje na Fortranovu biblioteku ARPACK, i koristi procedure DSAUPD, DSEUPD, DNAUPD, DNEUPD, ZNAUPD, i ZNEUPD.

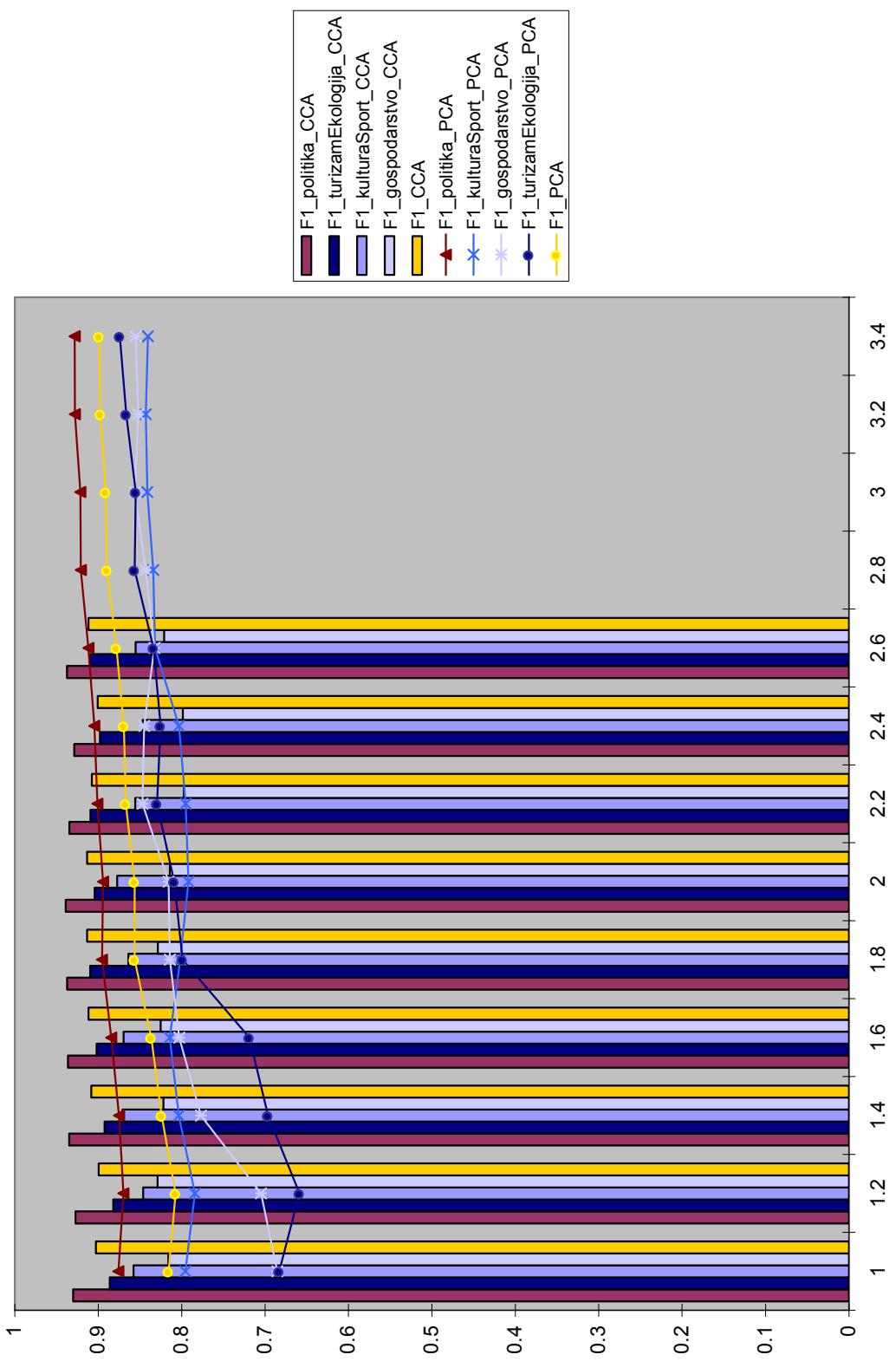
Eigs funkcija koristi Arnoldi metodu za izračunavanje svojstvenih vektora i rješavanje općeg problema svojstvenih vrijednosti rijetkih matrica – velikih matrica sa malim brojem nenu l elemenata. Njeno rješenje je aproksimacija međutim upotreba točne dekompozicije za izračunavanje svojstvenih vrijednosti iscrpljuje veliku memoriju računala.

Dok PCA algoritam nastoji pronaći smjerove maksimalne varijance među dokumentima na hrvatskom jeziku, CCA traži smjerove maksimalne kovarijance hrvatskih dokumenata sa engleskim te se skup za učenje i testiranje na hrvatskom jeziku preslikava na novo dobivene osi.

Parametri C i γ za metodu potpornih vektora dobiveni su pomoću programa `grid.py` napisanog u programskom jeziku Python, koji je sastavni dio LIBSVM-a. Program koristi tehniku krosvalidacije te skup za učenje dijeli na dva dijela: dok na jednom uči, drugi služi za testiranje i obrnuto. Na dobivenom konturnom plotu krosvalidacijske točnosti možemo očitati parametre koji nam odgovaraju. Najveća vrijednost krosvalidacijske točnosti 91.065% dobivena je za parametre $C = 128$, $\gamma = 0.0078125$ te su ti parametri korišteni u svim eksperimentima sa ovom bazom.



Slika 13. Konturni plot krosvalidacijske točnosti na skupu za učenje



Slika 14. Usporedba PCA i CCA redukcije dimenzionalnosti

(mjera $F1$ po kategorijama i $F1^{micro}$)

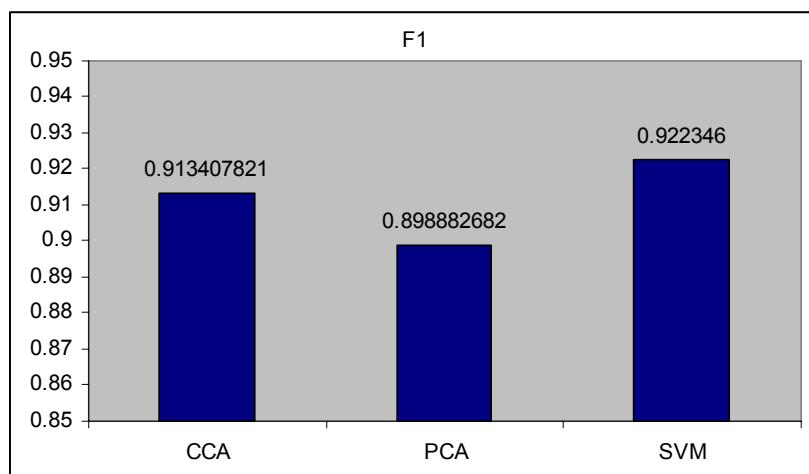
dimensionInost $k = 10^x$

Dimenzionalnost k je povećavana od $k = 10$ do $k = 500$ za CCA, te od $k = 10$ do $k = 1700$ za PCA. Napomenimo da je redukcija dimenzionalnosti ograničena odozgo i njen maksimum je $k = 1790$. To je očito ako se prisjetimo da se oba problema rješavaju pomoću jezgrene matrice dimenzionalnosti 1790×1790 ili *broj dokumenata x broj dokumenata*.

Regulacijski parametar CCA algoritma postavljen je na 1.5 uzastopnim ispitivanjem korelacijske matrice.

Cilj nam je procijeniti redukciju dimenzionalnosti općenito te također usporediti PCA i CCA metode. Kako je već navedeno PCA traži osi maksimalne varijance podataka dok CCA traži kovarijance između dokumenata na hrvatskom i engleskom jeziku. Prema rezultatima možemo zaključiti koliko uvođenje kompleksnih oznaka, tj. novog jezika pomaže u klasifikaciji. Prema (Fortuna, 2004.) novi jezik se može nadodati i umjetno, koristeći prijevodne programe kao npr. *Google Language Tools*²⁸.

Također možemo uspoređivati rezultate dobivene redukcijom dimenzionalnosti i bez redukcije dimenzionalnosti, primjenom metode potpornih vektora na dokumente. Prema rezultatima redukcija dimenzionalnosti daje neznatno gore rezultate nego bez redukcije.

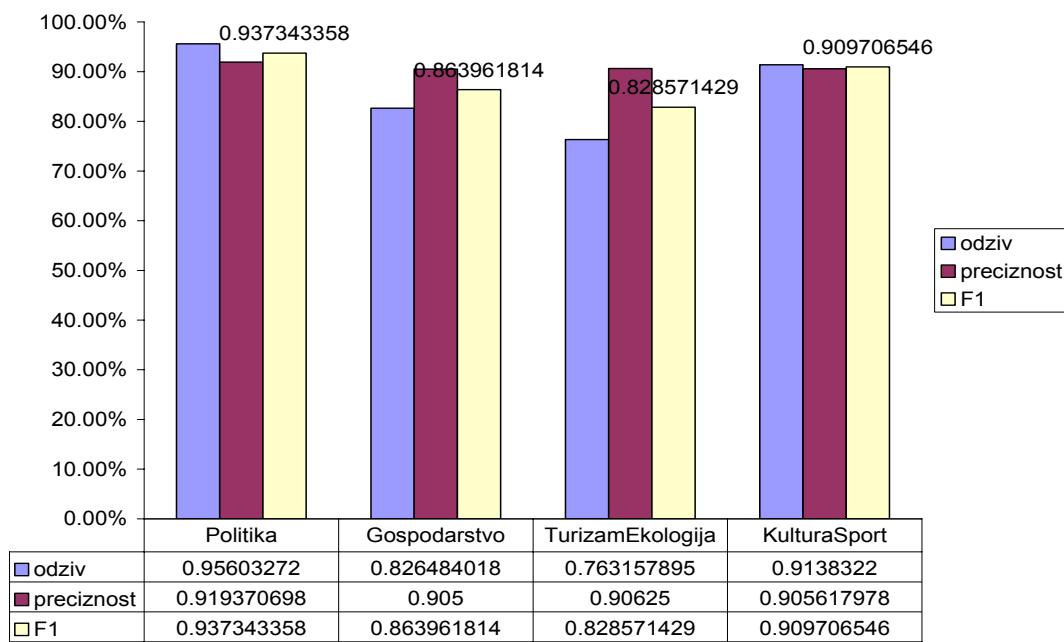


Slika 15. Usporedba mjere $F1^{micro}$ za redukciju dimenzionalnosti CCA ($k = 63$), PCA ($k = 1584$) te bez redukcije

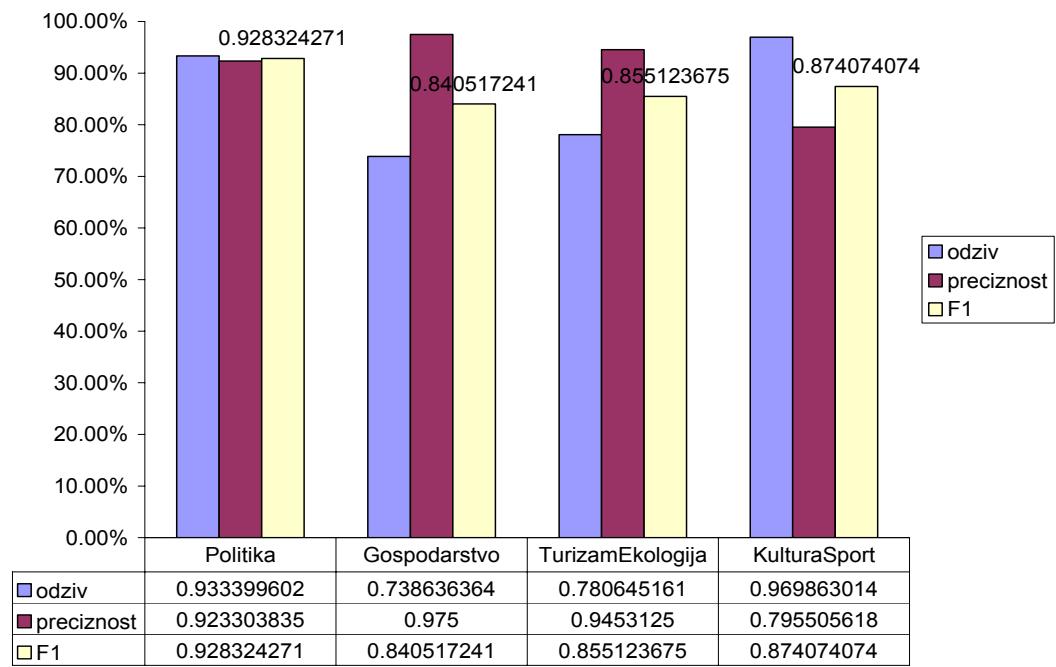
²⁸ http://www.google.com/language_tools

Povećanjem dimenzionalnosti prostora značajki rezultati mjere F1 za PCA i CCA postaju sve sličniji. Sa manjim vrijednostima k , PCA ima relativno lošije rezultate u usporedbi s CCA metodom.

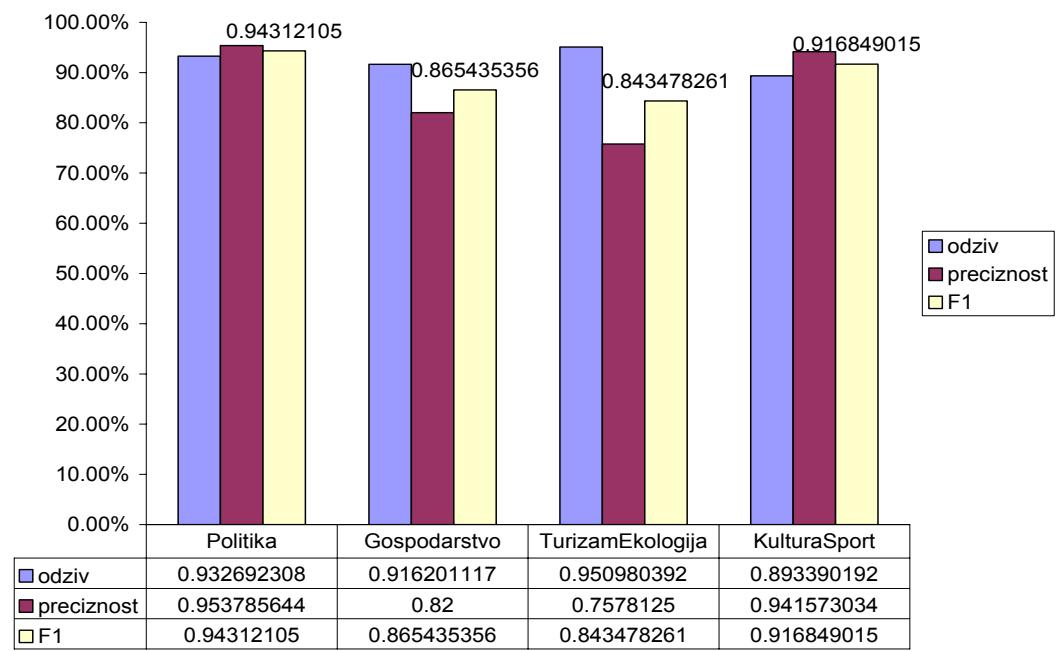
Na sljedećim slikama vidimo da u svim kategorijama mjera F1 CCA metode nadmašuje PCA. Najvećih varijacija među njima dolazi u kategoriji gospodarstvo. Kako je odziv za PCA mali, dok je preciznost velika možemo zaključiti da u toj kategoriji PCA propušta velik broj dokumenata koji se krivo klasificiraju. CCA metoda je obuhvatnija, ali također popušta na preciznosti – u tu kategoriju nadodaju se i dokumenti drugih kategorija.



Slika 16. Odziv, preciznost i F1 za CCA ($k = 63$)



Slika 17. Odziv, preciznost i F1 za PCA ($k = 1584$)



Slika 18. Odziv, preciznost i F1 za SVM

Vjerojatnosti

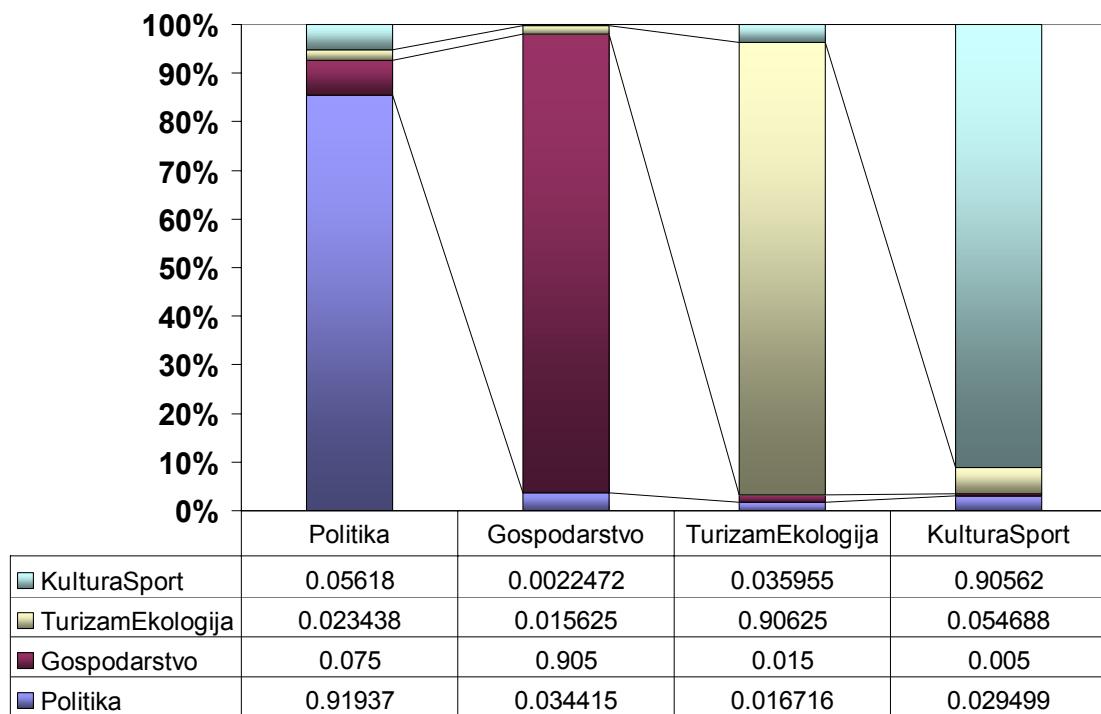
Metoda potpornih vektora kao izlaz predviđa oznake razreda, ali ne daje nikakve vjerojatnosne informacije. Korištene biblioteke funkcija metode potpornih vektora, LIBSVM, proširene su ovakvima proračunima koji su spremljeni u Confusion Matrix. Svaki element matrice daje vjerojatnost

$$\text{ConfMatrix}[i, j] = P(X \in j | X \in i)$$

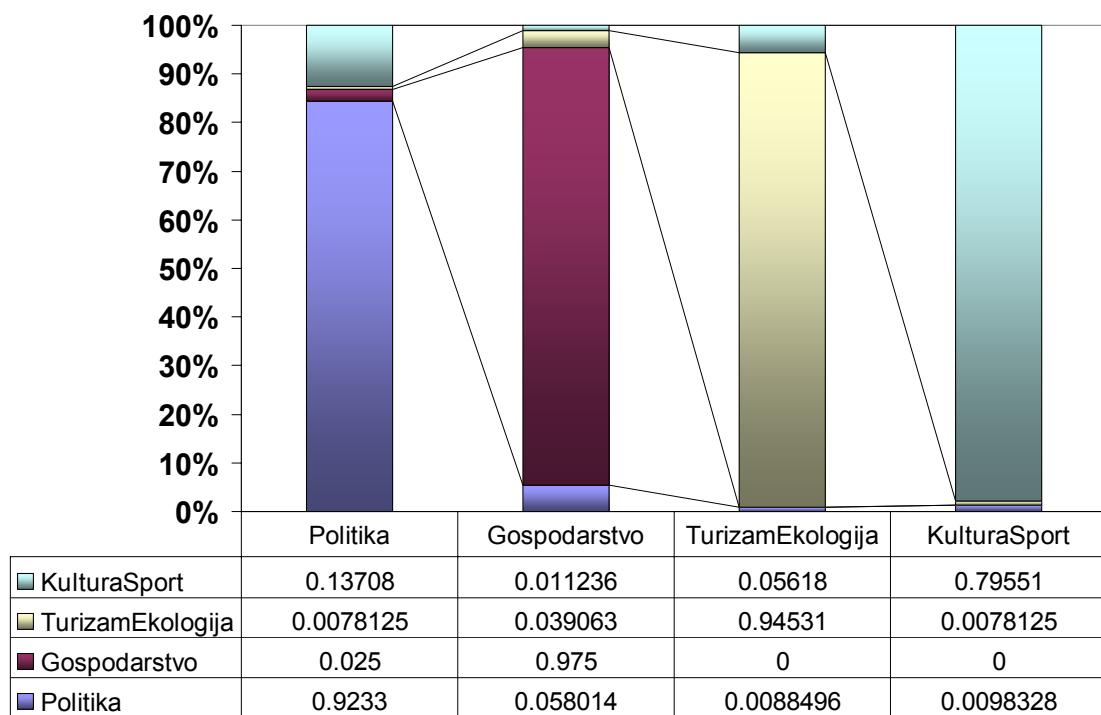
ili vjerojatnost klasificiranja dokumenta kategorije i u kategoriju j .

Više o vjerojatnostima i njihovom izračunavanju može se pronaći u (Chang et al. , 2005.)

Zaključke koje dobijemo vjerojatnostima možemo uspoređivati sa odzivom i preciznošću. Tako za PCA kategorije gospodarstvo i turizam i ekologija imaju skoro jednako nizak odziv i visoku preciznost. U Confusion matrici to znači redak popunjen većinom postocima blizu 0 i jednim postotkom malo manjim od 1 koji odgovara dotičnoj kategoriji – ili malo će se drugih kategorija klasificirati u dotičnu – velika preciznost. Također stupac označava odziv, a raznoliko popunjen za ove dvije kategorije daje mali odziv. Veliki odziv za PCA ima kultura i sport – većina dokumenta te kategorije predviđanjem će ostati u toj kategoriji. Mala preciznost te kategorije objašnjavamo velikim brojem dokumenata politike za koje se predviđa da pripadaju kulturi i sportu.



Slika 19. Confusion Matrix za $k = 63$ CCA



Slika 20. Confusion Matrix za $k = 1584$ PCA

4 KONCEPTNI VEKTORI

U ovom poglavlju proučiti ćemo razdiobu veliko dimenzionalnih i rijetkih tekstualnih skupova podataka kao npr. Vjesnikove baze u različite konceptualne kategorije.

4.1 Konceptni vektori i objektivna funkcija

Dokument vektore $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ koji popunjavaju *rječ x dokument* matricu, ako prepostavimo da su normalizirani, možemo promatrati kao točke jedinične kugle u \mathbb{R}^m . Nadalje, sve su komponente dokument vektora pozitivne i nalaze se u $\mathbb{R}^{m \geq 0}$. Za ovakve vektore skalarni produkt je prirodna mjera sličnosti. Za dva vektora \mathbf{x} i \mathbf{y} u $\mathbb{R}^{m \geq 0}$ vrijedi

$$\mathbf{x}^\top \mathbf{y} = \|\mathbf{x}\| \|\mathbf{y}\| \cos \theta(\mathbf{x}, \mathbf{y}) = \cos \theta(\mathbf{x}, \mathbf{y})$$

gdje je $\theta(\mathbf{x}, \mathbf{y})$ kut između dva vektora i $0 \leq \theta(\mathbf{x}, \mathbf{y}) \leq \pi/2$.

Prodot $\mathbf{x}^\top \mathbf{y}$ naziva se «kosinusna sličnost». Svojstvo dokument vektora da su rijetki omogućava lako i učinkovito računanje konceptne sličnosti.

Za dokument vektore $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ prepostavimo razdiobu dokumenata u k različitim grupama $\pi_1, \pi_2, \dots, \pi_k$ tako da

$$\bigcup_{j=1}^k \pi_j = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \quad \text{i vrijedi } \pi_j \cup \pi_l = \emptyset \text{ ako } j \neq l.$$

Za svaki $1 \leq j \leq k$ računa se srednji vektor ili centroid grupe dokument vektora sadržanih u grupi π_j formulom

$$\mathbf{m}_j = \frac{1}{n_j} \sum_{x \in \pi_j} \mathbf{x}$$

gdje je n_j broj dokument vektora u grupi π_j . Kako centroidi ne moraju imati jediničnu normu, njihov smjer se zadržava u konceptnim vektorima

$$\mathbf{c}_j = \frac{\mathbf{m}_j}{\|\mathbf{m}_j\|}$$

Za svaki jedinični vektor \mathbf{z} u \mathbb{R}^m iz Cauchy – Schwarz nejednakost vrijedi

$$\sum_{x \in \pi_j} \mathbf{x}^\top \mathbf{z} \leq \sum_{x \in \pi_j} \mathbf{x}^\top \mathbf{c}_j$$

Prema tome konceptne vektore možemo promatrati kao one vektore koji su najbliži u kosinusnoj sličnosti prema svim dokumentima u grupi π_j .

Produkt $\sum_{x \in \pi_j} \mathbf{x}^T \mathbf{c}_j$ možemo koristiti kao mjeru koherencije ili kvalitete grupe. Ako su svi dokumenti u grupi identični, tada će prosječna koherencija grupe biti najveća tj. imat će vrijednost 1. S druge strane, ako vektori u grupi jako variraju, prosječna koherencija će biti blizu 0. Svako particioniranje u grupe $\{\pi_j\}_{j=1}^k$ mjerimo pomoću *objektivne funkcije*²⁹ (Dhillon et al., 2001.)

$$Q(\{\pi_j\}_{j=1}^k) = \sum_{j=1}^k \sum_{x \in \pi_j} \mathbf{x}^T \mathbf{c}_j$$

4.2 Spherical k-means algoritam

Sljedeći korak je pronaći particioniranje dokument vektora $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ u k različitim grupa $\pi_1^*, \pi_2^*, \dots, \pi_k^*$ tako da se maksimizira objektivna funkcija tj. tražimo rješenje sljedećeg optimizacijskog problema

$$\{\pi_j\}_{j=1}^k = \arg \max_{\{\pi_j\}_{j=1}^k} Q(\{\pi_j\}_{j=1}^k)$$

Kako je problem NP težak za pronalaženje rješenja koristi se heuristički iterativni algoritam - spherical k-means algoritam:

1. Počinje se sa proizvoljnom razdiobom dokument vektora $\{\pi_j^{(0)}\}_{j=1}^k$. Indeks iteracije t jednak je 0, a $\{\mathbf{c}_j^{(0)}\}_{j=1}^k$ su konceptni vektori dobiveni trenutnim particioniranjem.
2. Za svaki dokument vektor \mathbf{x}_i , i $1 \leq i \leq n$ potrebno je pronaći konceptni vektor njemu najbliži i izračunati novo particioniranje $\{\pi_j^{(t+1)}\}_{j=1}^k$

$$\pi_j^{(t+1)} = \{ \mathbf{x} \in \{\mathbf{x}_i\}_{i=1}^n : \mathbf{x}^T \mathbf{c}_j^{(t)} > \mathbf{x}^T \mathbf{c}_l^{(t)}, 1 \leq l \leq n, l \neq j \}, \quad 1 \leq j \leq k$$

Riječima, $\pi_j^{(t+1)}$ je skup svih dokument vektora koji su najbliži konceptnom vektoru $\mathbf{c}_j^{(t)}$. Ovako definirane grupe nazivaju se *Voronoi* ili *Dirichlet*-ove grupe.

²⁹ Eng. objective function

3. Izračunavaju se novi konceptni vektori prema novom grupiranju $\{ \mathbf{c}_j^{(t+1)} \}_{j=1}^k$
4. Ako je ispunjen uvjet zaustavljanja tada je $\pi_j^{(t+1)} = \mathbf{c}_j^{(t+1)}$ i $\mathbf{c}_j^{(t+1)} = \mathbf{c}_j^{(t+1)}$ za $1 \leq j \leq k$.
Inače t se poveća za 1 i vraćamo se na korak 2.

Jedan od primjera kriterija zaustavljanja je $|Q(\{ \pi_j^{(t)} \}_{j=1}^k) - Q(\{ \pi_j^{(t+1)} \}_{j=1}^k)| < \varepsilon$ tj. algoritam se zaustavlja ako je promjena objektivne funkcije nakon jedne iteracije manja od praga.

Nadalje vrijednost objektivne funkcije ne smanjuje se iteracijama

$$Q(\{ \pi_j^{(t)} \}_{j=1}^k) \leq Q(\{ \pi_j^{(t+1)} \}_{j=1}^k), \text{ za } \forall t \geq 0$$

Dokaz:

$$\begin{aligned} Q(\{ \pi_j^{(t)} \}_{j=1}^k) &= \sum_{j=1}^k \sum_{x \in \pi_j^{(t)}} \mathbf{x}^T \mathbf{c}_j^{(t)} = \sum_{j=1}^k \sum_{l=1}^k \sum_{x \in \pi_j^{(t)} \cap \pi_l^{(t+1)}} \mathbf{x}^T \mathbf{c}_j^{(t)} \\ &\leq \sum_{j=1}^k \sum_{l=1}^k \sum_{x \in \pi_j^{(t)} \cap \pi_l^{(t+1)}} \mathbf{x}^T \mathbf{c}_l^{(t)} \\ &= \sum_{l=1}^k \sum_{j=1}^k \sum_{x \in \pi_j^{(t)} \cap \pi_l^{(t+1)}} \mathbf{x}^T \mathbf{c}_l^{(t)} \\ &= \sum_{l=1}^k \sum_{x \in \pi_l^{(t+1)}} \mathbf{x}^T \mathbf{c}_l^{(t)} \leq \sum_{l=1}^k \sum_{x \in \pi_l^{(t+1)}} \mathbf{x}^T \mathbf{c}_l^{(t+1)} = Q(\{ \pi_j^{(t+1)} \}_{j=1}^k) \end{aligned}$$

Također postoji $\lim_{t \rightarrow \infty} Q(\{ \pi_j^{(t)} \}_{j=1}^k)$ i vrijedi

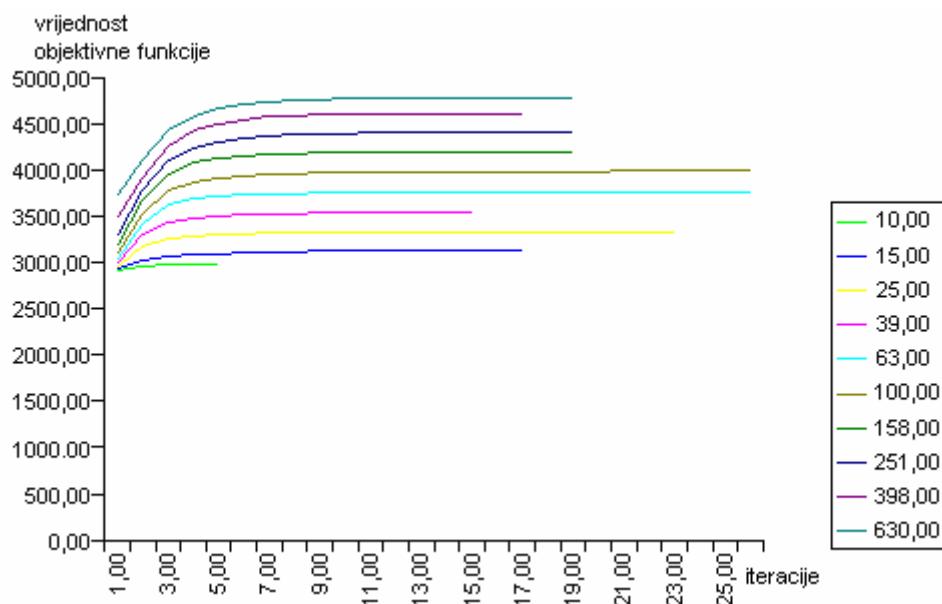
$$Q(\{ \pi_j^{(t)} \}_{j=1}^k) = \sum_{j=1}^k n_j^{(t)} \|\mathbf{m}_j^{(t)}\| \leq \sum_{j=1}^k n_j^{(t)} = n$$

Objektivna funkcija je rastuća i ograničena odozgo konstantom tj. iteracijama konvergira i granica postoji. Ove dokaze može se pronaći u (Dhillon et al., 2005.)

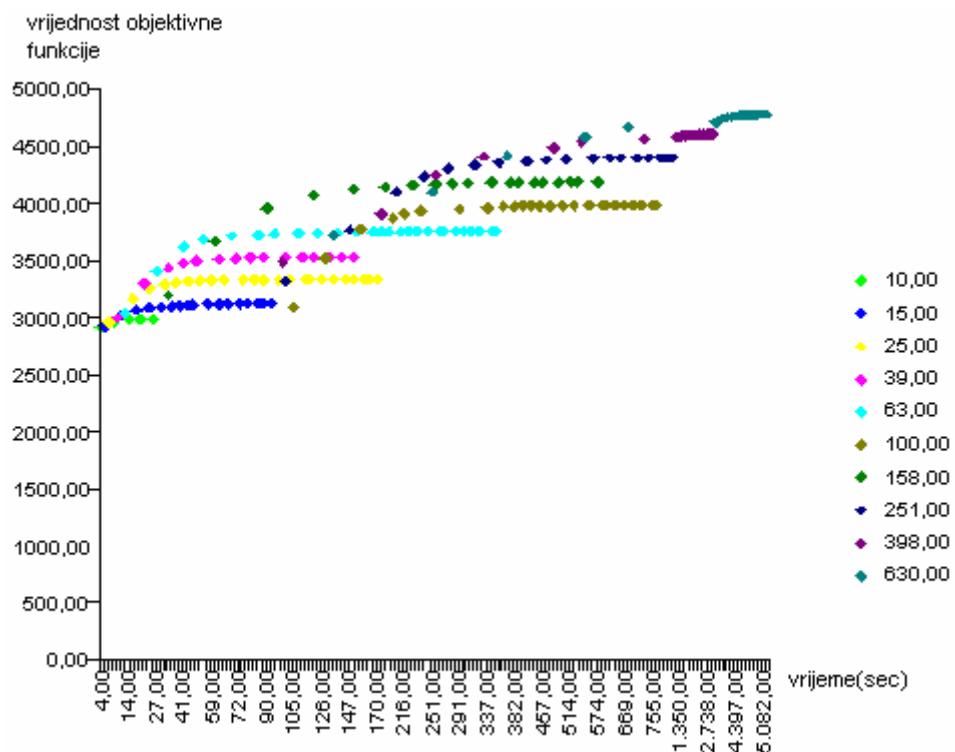
Spherical k-means algoritam kao i ostali gradijentni algoritmi može zaglaviti u lokalnom maksimumu. Ipak u eksperimentima daje dobre rezultate.

4.3 Eksperimentalni rezultati

4.3.1 Objektivna funkcija je rastuća



Slika 21. Promjena objektivne funkcije po iteracijama za različiti k (broj grupa)



Slika 22. Promjena objektivne funkcije u vremenu(sekunde) za različiti k (broj grupa)

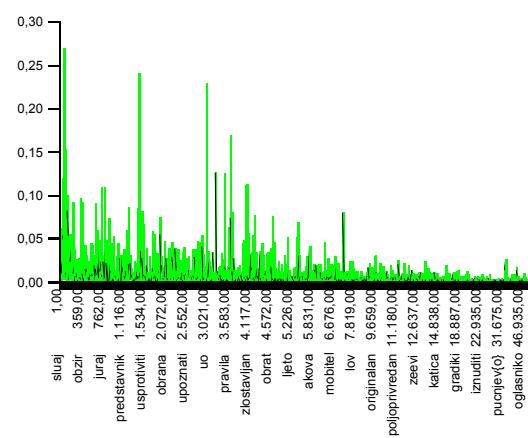
4.3.2 Konceptni vektori su lokalni, rijetki i teže ortonormalnosti

U prethodnom tekstu $\{\pi_j\}_{j=1}^k$ označava zadnje partitioniranje dokumenata u k grupa. Za svaku grupu $1 \leq j \leq k$ definiramo skup riječi W_j , riječ $1 \leq w \leq m$ nalazi se u W_j ako je težina te riječi u konceptnom vektoru \mathbf{c}_j veća nego u ostalim konceptnim vektorima ili

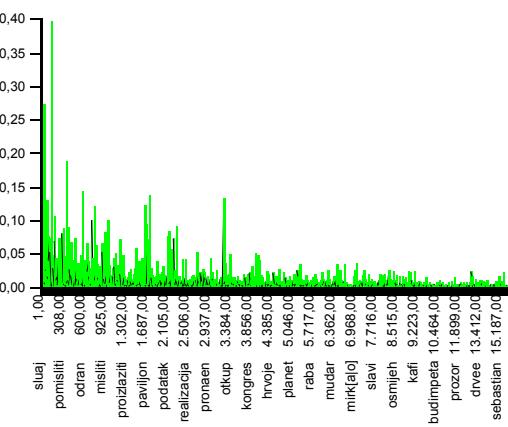
$$W_j = \{w : 1 \leq w \leq m, \mathbf{c}_{wj} > \mathbf{c}_{wl}, 1 \leq l \leq k, l \neq j\}.$$

Skup riječi može se upotrijebiti za opisati grupe dokument vektora. U eksperimentu koji slijedi Vjesnikova baza razdijeljena je u 8 grupa, onoliko koliko ima kategorija različitih vrsta članaka. Svaki skup riječi lokalizira se oko jedne vrste članaka, sporta, kulture, unutrašnje politike, crne kronike, gospodarstva, teme dana, vanjske politike, Zagreba. U tablici 1. prikazani su konceptni vektori za svih 8 kategorija i prvih 15 riječi kako se pojavljuju u konceptnom vektoru. Različitim bojama obilježene su skupine riječi (crvenom je obilježena riječ pripadajuće skupine riječi). Primjećujemo da je većina težina konceptnog vektora koncentrirana ili lokalizirana u pridruženom skupu riječi. Sada možemo reći da su konceptni vektori lokalizirani.

Riječ	Vrijednost	Sljedeća vrij
godina	0,27	0,36
sud	0,24	0,24
policija	0,23	0,23
kaznen	0,17	0,17
rei	0,15	0,19
zagreb	0,15	0,26
djelo	0,13	0,13
upanijski	0,12	0,12
svjedok	0,12	0,12
sat	0,12	0,12
optuen	0,11	0,11
ubojsztvo	0,11	0,11
dan	0,11	0,13
automobil	0,11	0,11
kuna	0,11	0,19



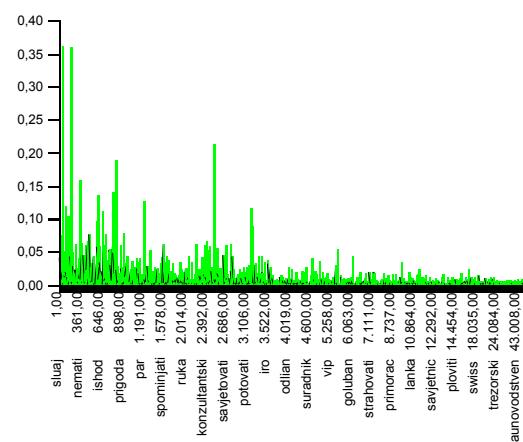
Crna kronika



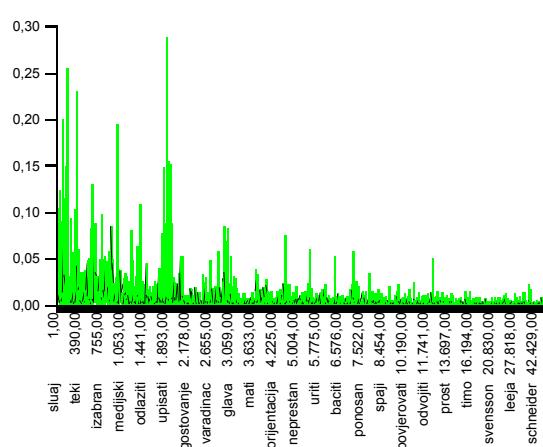
Kultura

Riječ	Vrijednost	Sljedeća vrij
hrvatski	0,40	0,48
godina	0,27	0,36
moi	0,19	0,28
zagreb	0,16	0,26
nov	0,14	0,16
film	0,14	0,14
knjiga	0,13	0,13
imati	0,13	0,20
izložba	0,12	0,12
velik	0,12	0,12
djelo	0,11	0,13
kultura	0,11	0,11
rad	0,10	0,14
ovjek	0,10	0,10
festival	0,09	0,09

Riječ	Vrijednost	Sljedeća vrij
godina	0,36	0,36
hrvatski	0,36	0,48
banka	0,21	0,21
kuna	0,19	0,19
moi	0,16	0,28
trite	0,14	0,14
nov	0,14	0,16
cijena	0,13	0,13
imati	0,12	0,20
dionica	0,12	0,12
tvrtka	0,11	0,11
zagreb	0,11	0,26
trebatи	0,10	0,16
poslovan	0,10	0,10
rast	0,09	0,09



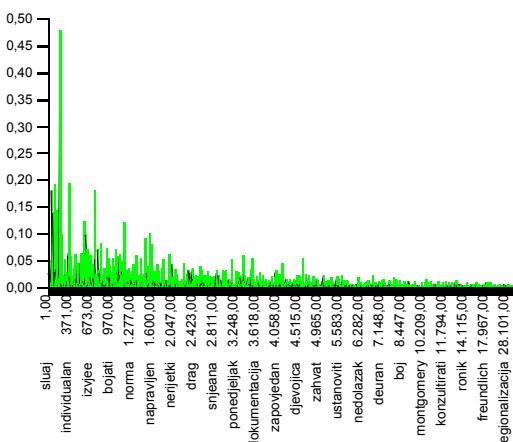
Gospodarstvo



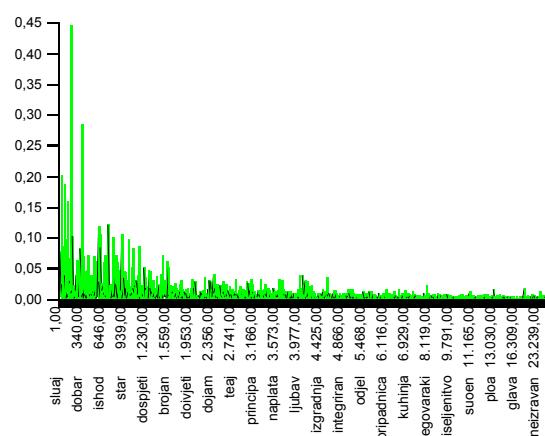
Sport

Riječ	Vrijednost	Sljedeća vrij
igra	0,15	0,15
igra	0,20	0,20
utakmica	0,29	0,29
hrvatski	0,26	0,48
moi	0,23	0,28
imati	0,20	0,20
zagreb	0,17	0,26
igrati	0,16	0,16
klub	0,15	0,15
momad	0,15	0,15
pobjeda	0,13	0,13
godina	0,12	0,36
trener	0,12	0,12
trebatи	0,11	0,16
minuta	0,11	0,11
dan	0,10	0,13

Riječ	Vrijednost	Sljedeća vrij
hrvatski	0,48	0,48
moi	0,19	0,28
vlada	0,19	0,19
predsjednik	0,18	0,19
godina	0,18	0,36
trebati	0,14	0,16
rei	0,14	0,19
zagreb	0,14	0,26
imati	0,13	0,20
kazati	0,12	0,12
nov	0,12	0,16
zakon	0,12	0,12
ministar	0,10	0,10
prav	0,10	0,10
stranka	0,10	0,11

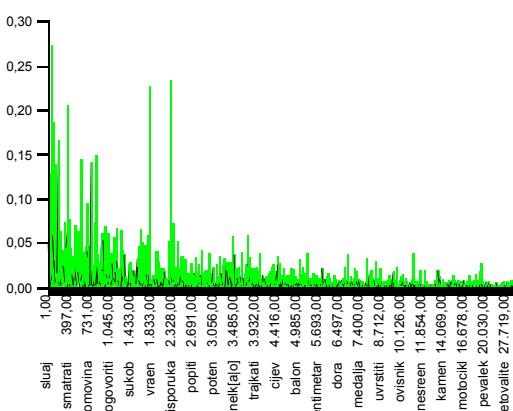


Tema dana



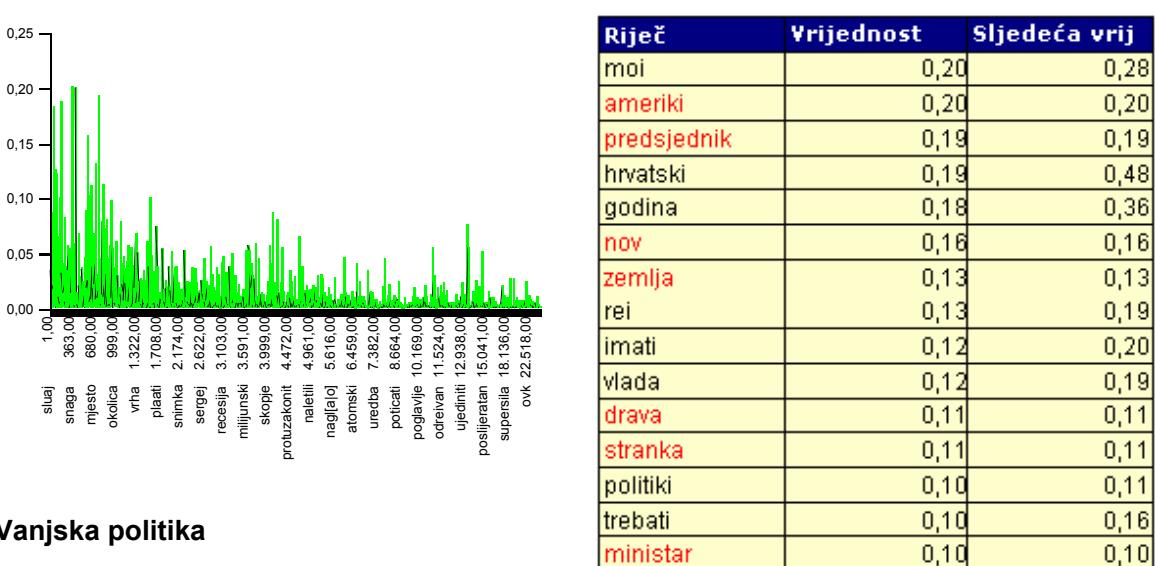
Unutrašnja politika

Riječ	Vrijednost	Sljedeća vrij
hrvatski	0,44	0,48
moi	0,28	0,28
godina	0,20	0,36
imati	0,19	0,20
trebati	0,16	0,16
vlada	0,15	0,19
predsjednik	0,12	0,19
nov	0,12	0,16
politiki	0,11	0,11
moz	0,10	0,10
stranka	0,10	0,11
prav	0,10	0,10
drava	0,10	0,11
morati	0,10	0,10
ovjek	0,10	0,10



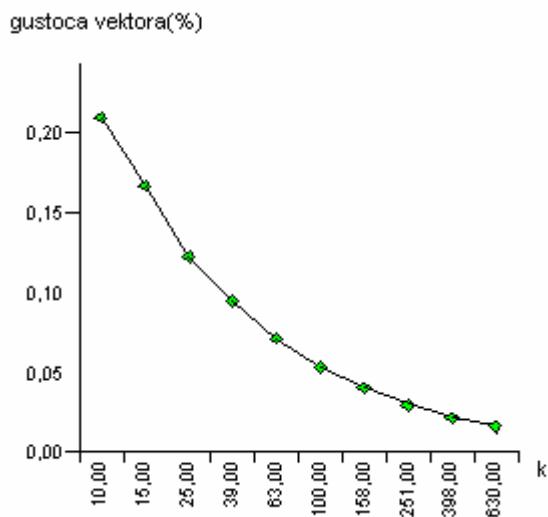
Zagreb

Riječ	Vrijednost	Sljedeća vrij
godina	0,27	0,36
zagreb	0,26	0,26
grad	0,23	0,23
gradski	0,23	0,23
moi	0,21	0,28
rei	0,19	0,19
hrvatski	0,17	0,48
zagrebaki	0,16	0,16
kuna	0,15	0,19
nov	0,14	0,16
rad	0,14	0,14
trebati	0,14	0,16
imati	0,13	0,20
dan	0,13	0,13
sat	0,10	0,12



Tablica 1. Osam konceptnih vektora koji odgovaraju grupiranju Vjesnikove baze u 8 grupa.

U tablici 1. također možemo primijetiti da su konceptni vektori rijetki, tj. imaju mali broj nenegativnih elemenata. Od 100 000 elemenata, popunjeno je 20 000 odnosno 20%. Kako broj grupe raste tako je manji broj dokument vektora koji se nalaze unutar pojedine grupe i možemo očekivati sve veću rijetkost konceptnih vektora. Dokument vektori su oko 99% rijetki.

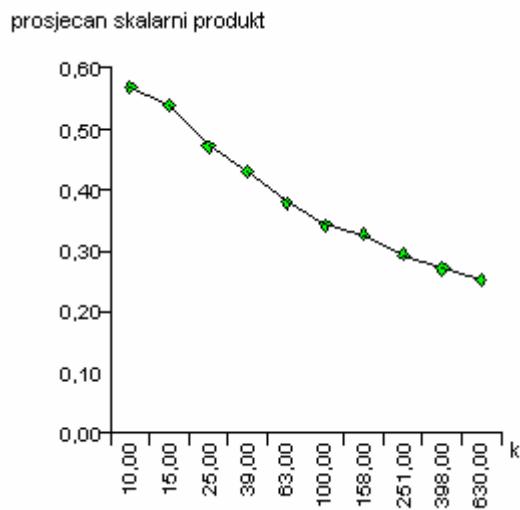


Slika 23. Povećanjem broja konceptnih vektora oni postaju sve rijetki

Za svaki $k \geq 2$ možemo izračunati prosječan skalarni produkt konceptnih vektora $\{ \mathbf{c}_j^+ \}_{j=1}^k$ kao

$$\frac{2}{k(k-1)} \sum_{j=1}^k \sum_{l=j+1}^k (\mathbf{c}_j^+)^T \mathbf{c}_l^+$$

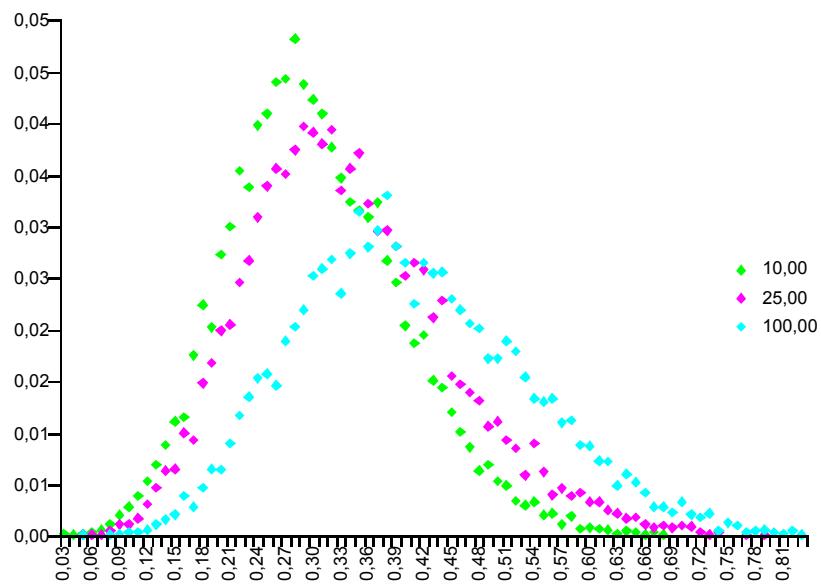
Prosječan skalarni produkt poprima vrijednosti između [0,1] gdje 0 odgovara ortonormalnim vektorima. Kako se broj k grupa povećava tako prosječan skalarni produkt teži prema 0 ili konceptni vektori teže ortonormalnosti.



Slika 24. Konceptni vektori teže ortonormalnosti

4.3.3 Struktura unutar i između grupa

Kako smo prije vidjeli svaki dokument vektor je najbliži upravo konceptnom vektoru grupe u kojoj se on nalazi. Strukturu unutar grupe možemo promatrati pomoću funkcije vjerojatnost kosinusnih sličnosti unutar grupe ili skalarnog produkta $\mathbf{x}^T \mathbf{c}_j$ dokument-vektora i pripadajućeg konceptnog vektora.

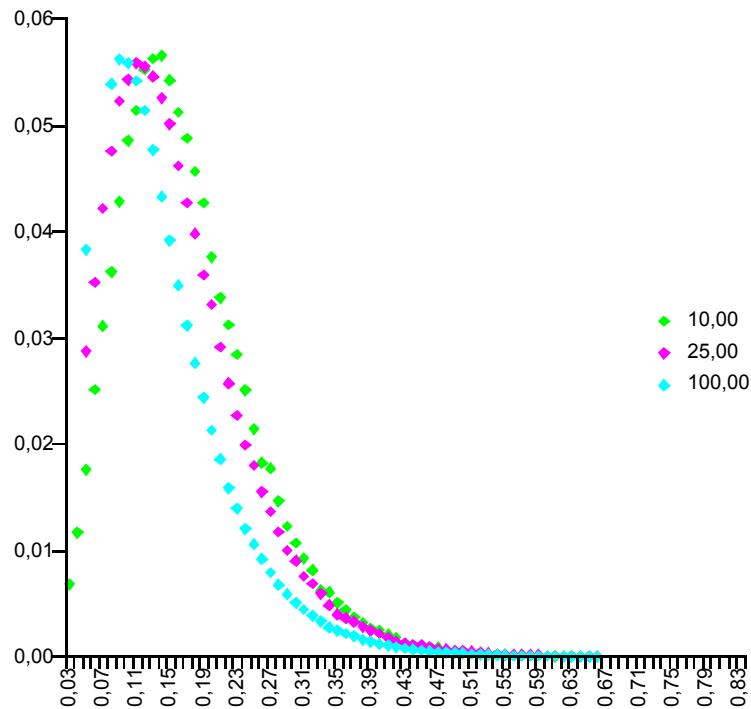


Slika 25. Funkcija vjerojatnosti $\mathbf{x}^T \mathbf{c}_j$ unutar grupe

Za k je 10 vidimo da se većina vrijednosti nalazi u intervalu [0.1, 0.5] iz čega bi mogli zaključiti da ne postoje vektori bliski svojim konceptnim vektorima. Dakle konceptne vektore ne trebamo zamišljati okružene točkama dokument - vektora, zapravo postoji veliki prazni prostor između dokument - vektora i pridruženih konceptnih vektora. Povećanjem broja konceptnih vektora k sve je veća kompaktnost unutar grupe, kao što je bilo za očekivati, dok se očekivanje funkcije vjerojatnosti lagano odmiče prema većim vrijednostima. Dakle, kako se k povećava i prazan prostor oko konceptnih vektora se smanjuje.

S druge strane, iako dokument vektori nisu bliski s pripadajućim konceptnim vektorom, još su udaljeniji od ostalih konceptnih vektora. To svojstvo nam zapravo omogućava smisleno grupiranje dokumenata. Strukturu između grupa, ili skalarni produkt $\mathbf{x}^T \mathbf{c}_j$ dokument vektora i ostalih konceptnog vektora možete vidjeti na slici 5. Struktura

između grupa se progresivno se udaljava od unutra grupne strukture i teži prema 0. To navodi na očitu sličnost sa fraktalima (Mandelbrot, 1988).



Slika 26. *Funkcija vjerojatnosti $\mathbf{x}^T \mathbf{c}_j$ između grupa*

5 UPORABA RAZLIČITIH METRIKA

Samo srce datamining i information-retrivial zadataka je funkcija udaljenosti kojom mjerimo sličnost među podacima. Možemo reći da većina algoritama kritično ovisi o upotrebi dobre metrike na ulaznom prostoru. Na primjer, K-means, algoritam najbližih susjeda i jezgreni algoritmi kao SVM, trebaju dobru metriku koja reflektira bitne veze među

Minkowsky metrika

$$D(x,y) = \left(\sum_{i=1}^m |x_i - y_i|^r \right)^{\frac{1}{r}}$$

Euklidska metrika

$$D(x,y) = \sqrt{\sum_{i=1}^m (x_i - y_i)^2}$$

Manhattan/city block metrika

$$D(x,y) = \sum_{i=1}^m |x_i - y_i|$$

Mahalanobis metrika

$$D(x,y) = [\det V]^{\frac{1}{m}} (x - y)^T V^{-1} (x - y)$$

Kvadratna metrika

$$D(x,y) = (x - y)^T Q (x - y)$$

podatcima. Ovaj problem najviše se očituje u nenadziranim algoritmima kao grupiranju i povezan je sa još jednim problemom na tom području: Ako su tri algoritma korištena za grupiranje kolekcije dokumenata, jedna grupa prema autorstvu, druga prema temi, a treća prema stilu pisanja, tko može reći koje je «pravo» rješenje? Još gore, ako imamo samo algoritam koji grupira po temi, a mi želimo po stilu pisanja, relativno je malo mehanizama kojim možemo pravilno grupirati podatke. Prepostavimo da korisnik ukaže na određene točke u ulaznom prostoru (\mathbb{R}^n) kao «slične». Možemo li naučiti metriku udaljenosti prema ovim vezama, metriku koja dodjeljuje male udaljenosti sličnim

parovima? U prethodnom primjeru to bi značilo da iz grupe dokumenata pisanih istim stilom, možemo odrediti bitne značajke u određivanju stila.

5.1 Euklidska mjera udaljenosti

Podaci u našem slučaju su vektori u \mathbb{R}^m prostoru sa elementima koji predstavljaju broj pojavljivanja i -te značajke. Ako su značajke riječi, kao u ovom slučaju, to zovemo «bag of words», ali općenito značajke ne moraju predstavljati riječi te ćemo u dalnjem tekstu

upotrebljavati općenitiji izraz – «značajke». Za mjerjenje udaljenosti među vektorima najčešće se upotrebljava Euklidska udaljenost ili dot produkt³⁰

$$\begin{aligned} d(\mathbf{u}, \mathbf{v})^2 &= (\mathbf{u} - \mathbf{v})^\top (\mathbf{u} - \mathbf{v}) \\ &= \sum_{j=1}^n (\mathbf{u}_j - \mathbf{v}_j)^2 \end{aligned}$$

koja se najčešće koristi s heurističkim mjerjenjem težina značajki³¹ kao npr. izračunavanje *tfidf* u preprocesiranju. Postoje dva bitna problema sa ovom udaljenošću:

1. Zanemarene su korelacije među riječima
2. Mjerjenje značajki je neizbjegljivo proizvoljno.

Prvi problem je posebno važan kod jezika, jer lingvističke značajke (rijeci) imaju snažne međusobne korelacije. Ovakve korelacije ne mogu se razmatrati običnim dot produkтом. Ovakav problem moguće je riješiti specifičnim jezgrenim funkcijama.

Drugi problem je jednostavniji. Dok *tfidf* često daje dobre rezultate u praksi postoji više opcija, osobito u *tf* kao logaritmi i kvadratni korijeni, ali ne postoji princip po kojem se one odabiru. Nadalje, ne postoji teoretska osnova koja daje optimalnost u pogledu funkcija udaljenosti.

Gornji problem korelacija među značajkama i težina značajki može se gledati kao problem definiranja prikladne matrice u prostoru značajki, temeljene na distribuciji podataka.

5.2 Mahalanobisova udaljenost

Prepostavimo da imamo skup točka $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^n \subseteq \mathbb{R}^m$ i postoje informacije o sličnosti S : $(\mathbf{x}_i, \mathbf{x}_j) \in S$ ako su \mathbf{x}_i i \mathbf{x}_j slični,

Pitanje je kako naučiti funkciju udaljenosti $d(\mathbf{x}, \mathbf{y})$ tako da se poštuje sličnost, tj. da su slične točke blizu?

³⁰ Kada su vektori normalizirani, $|\mathbf{u}| = |\mathbf{v}| = 1$, Euklidska udaljenost je $(\mathbf{u} - \mathbf{v})^\top (\mathbf{u} - \mathbf{v}) = |\mathbf{u}|^2 + |\mathbf{v}|^2 - 2\mathbf{u} \cdot \mathbf{v} \propto -$

$\mathbf{u} \cdot \mathbf{v} = -\cos(\mathbf{u}, \mathbf{v})$; (Manning and Schutze, 1999)

³¹ feature weighting

Proučimo funkciju udaljenosti u obliku:

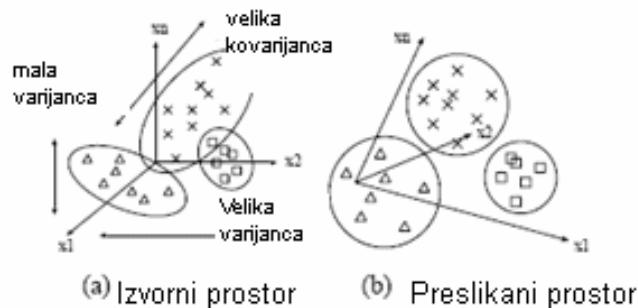
$$d(\mathbf{x}, \mathbf{y}) = d_A(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_A = \sqrt{(\mathbf{x} - \mathbf{y})^T A (\mathbf{x} - \mathbf{y})}$$

Da bi gornja funkcija bila metrika – i zadovoljavala nenegativnost i nejednakost trokuta – potrebno je da A bude pozitivna, semi-definitna. Postavljajući $\mathbf{A} = \mathbf{I}$ daje Euklidsku udaljenost; ako ograničenim A na dijagonalnu matricu, učimo metriku u kojoj različite osi imaju različite «težine»; općenito A parametrizira obitelj Mahalanobisovih udaljenosti na \mathbb{R}^n ³².

Ako zapišemo gornju jednadžbu u drugom obliku (matrica A simetrična)

$$d_A(\mathbf{x}, \mathbf{y})^2 = (\mathbf{A}^{1/2}(\mathbf{u} - \mathbf{v}))^T (\mathbf{A}^{1/2}(\mathbf{u} - \mathbf{v}))$$

dolazimo do zaključka da je učenje takve metrike jednako preslikavanju podataka, koje svaki vektor \mathbf{x} preslika u $\mathbf{A}^{-1/2}\mathbf{x}$ i zatim se primjenjuje Euklidova metrika u $\mathbf{A}^{-1/2}$ -preslikanom prostoru. Možemo zamisliti da su podaci u grupama distribuirani kako je prikazano na slici 1(a), u elipsoidnom obliku koji ima velike (ko)varijance za neke, a male (ko)varijance za druge dimenzije. Nadalje, grupa nije ni poravnana sa osima koordinatnog sustava. Kad pronađemo matricu A koja minimalizira šumove u grupi, koja smanjuje velike varijance, a povećava male stvarajući od elipsoidnog oblika kružni u novom $\mathbf{A}^{-1/2}$ -preslikanom prostoru (slika 1(b)), možemo očekivati da smo dobili novu, pouzdaniju metriku.



Slika 27. Geometrija prostora značajki

³² Ako se koristi funkcija preslikavanja ϕ , a time funkcija udaljenosti postaje $\sqrt{(\phi(x) - \phi(y))^T A (\phi(x) - \phi(y))}$, može se naučiti i nelinearna metrika

5.3 Učenje funkcije udaljenosti

Učenje funkcije udaljenosti možemo grubo podijeliti na učenje metrike i jezgreno učenje.

5.3.1 Učenje metrike

Učenje metrike pokušava pronaći optimalnu linearu transformaciju grupe podatkovnih vektora koja bolje karakterizira sličnost među vektorima. Transformacija je sama linearna, ali podatkovni vektori mogu se prije preslikati koristeći nelinearnu funkciju $\phi(x)$. Ova transformacija jednaka je dodjeljivanju težina značajkama vektora pa se učenje metrike naziva i «feature weighting». Za skup podatkovnih vektora $X = \{x_i\}_{i=1}^n$ u \mathbb{R}^m , učenje metrike nastoji naučiti funkciju udaljenosti $d_A(x_i, x_j)$ između vektora x_i i x_j . Matematički funkciju udaljenosti možemo izraziti kao

$$d_A(x, y) = \sqrt{(\phi(x) - \phi(y))^T A (\phi(x) - \phi(y))}$$

gdje je A pozitivna (semi-)definitna matrica koja zadovoljava svojstva metrike – nenegativnost i nejednakost trokuta. Općenito, A je parametar Mahalanobisovih udaljenosti u \mathbb{R}^m . Izbor funkcije preslikavanja ϕ i matrice skaliranja A je upravo to što razlikuje algoritme za učenje metrike udaljenosti.

Mnogi algoritmi temelje se na korištenju kontekstualnih informacija ili «side informations». Kontekstualne informacije mogu biti informacije dobivene od korisnika o sličnim karakteristikama podskupa podataka. Na temelju takvih informacija, (Bar-hillel et al., 2003.) koriste Relevant Component Analysis (RCA) u izračunavanju Mahalanobisove metrike. Autori koriste relaciju ekvivalentnosti za kontekstualnu informaciju. Računaju

$$\mathbf{C} = \frac{1}{P} \sum_{j=1}^{|G|} \sum_{i=1}^{|S_j|} (x_{ji} - m_j)(x_{ji} - m_j)^T$$

gdje je m_j centroid j -te grupe, dok $|G|$ i $|S_j|$ predstavljaju broj grupa i broj uzoraka u j -toj grupi. Matrica $\mathbf{W} = \mathbf{C}^{-1/2}$ koristi se za transformaciju, a inverz matrice \mathbf{C} kao Mahalanobisova matrica.

(Aggarwal, 2003.) razmatra funkcije udaljenosti osjetljive na pojedine karakteristike podataka. Korišteni modeli su parametrizirani Minkowski model:

$$D(x_i, x_j, \lambda) = \left(\sum_{r=1}^m \lambda_r |x_{ir} - x_{jr}|^p \right)^{1/p}$$

i parametrizirani kosinusni model

$$\cos(\mathbf{x}_i, \mathbf{x}_j, \lambda) = \sum_{r=1}^m \frac{\lambda_r x_{ir} x_{jr}}{\sqrt{\sum_{r=1}^m \lambda_r^2 x_{ir}^2 \sum_{r=1}^m \lambda_r^2 x_{jr}^2}}$$

Oba modela nastoje minimizirati pogrešku s obzirom na λ_r . Parametrizirani Minkowski model može se promatrati kao dodjeljivanje težina značajkama na ulaznom prostoru. Slično, i parametrizirani kosinusni model možemo promatrati kao skalarni produkt na ulaznom prostoru.

(Xing et al., 2003) problem u (Bar-hillel et al., 2003.) promatraju kao problem konveksne optimizacije i razrađuju tehniku za dobivanje dijagonalne i pune Mahalanobisove matrice. Za razliku od RCA ovdje se koriste iterativne metode osjetljive na postavke parametara i također računski zahtjevne.

Jednostavan način definiranja kriterija za željenu metriku je već spomenuti zahtjev da parovi točaka $(\mathbf{x}_i, \mathbf{x}_j) \in S$ budu blizu tj. želimo minimizirati $\sum_{(xi, xj) \in S} \|\mathbf{x}_i - \mathbf{x}_j\|_A^2$. Trivijalno

rješenje problema je za $\mathbf{A} = 0$, dodajući uvjet $\sum_{(xi, xj) \in D} \|\mathbf{x}_i - \mathbf{x}_j\|_A^2 \geq 1$ osiguravamo da S ne preslika sve točke u jednu. Ovdje D označava skup parova točaka za koje se ili zna da «nisu slične» ako je ta informacija eksplicitno dostupna ili da se uzimaju svi parovi točaka koji se ne nalaze u S . Dolazimo do optimizacijskog problema

$$\begin{aligned} \min_A \quad & \sum_{(xi, xj) \in S} \|\mathbf{x}_i - \mathbf{x}_j\|_A^2 \\ \text{st.} \quad & \sum_{(xi, xj) \in D} \|\mathbf{x}_i - \mathbf{x}_j\|_A^2 \geq 1 \\ & A \geq 0 \end{aligned}$$

Izbor konstante 1 na lijevoj strani nejednadžbe je proizvoljan, ali i nebitan; ako konstantu zamijenimo bilo kojom pozitivnom konstantom to rezultira samo mijenjanjem matrice \mathbf{A} u $c^2 \mathbf{A}$.

S jedne strane (Xing et al., 2003.) gleda na ovaj optimizacijski problem kao konveksni. U slučaju da želimo izračunati dijagonalnu matricu $\mathbf{A} = \text{diag}(A_{11}, A_{22}, \dots, A_{mm})$ dolazimo do učinkovitog algoritma koji koristi Newton-Raphsonovu metodu. Definirajmo

$$g(A) = g(A_{11}, \dots, A_{nn}) = \sum_{(xi, xj) \in S} \|\mathbf{x}_i - \mathbf{x}_j\|_A^2 - \log(\sum_{(xi, xj) \in D} \|\mathbf{x}_i - \mathbf{x}_j\|_A^2)$$

i slijedi da minimaliziranje funkcije g (uz $A \succeq 0$ ³³) je ekvivalentno rješavanju gornjeg problema. Dakle, možemo koristiti Newton-Raphsonovu metodu za optimiziranje funkcije g .

Ako tražimo punu matrica \mathbf{A} , Newton-Rapsonova metoda postaje presložena (složenost je $O(n^6)$). Polazimo od algoritma gradijentnog spusta i iterativnih projekcija

```

Iterate
  Iterate
     $A := \arg \min_{A'} \{ \|A' - A\|_F : A' \in C_1 \}$ 
     $A := \arg \min_{A'} \{ \|A' - A\|_F : A' \in C_2 \}$ 
  until  $A$  converges
   $A := A + \alpha(\nabla_A g(A))_{\perp \nabla_A f}$ 
until convergence

```

gdje je $\| \cdot \|_F$ već ranije spominjana Frobeniusova norma matrice $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^p \sum_{j=1}^q |a_{ij}|}$.

Postavljamo sljedeći problem

$$\max_{\mathbf{A}} \quad g(\mathbf{A}) = \sum_{(x_i, x_j) \in D} \| \mathbf{x}_i - \mathbf{x}_j \|_{\mathbf{A}}^2$$

$$\text{st.} \quad f(\mathbf{A}) = \sum_{(x_i, x_j) \in S} \| \mathbf{x}_i - \mathbf{x}_j \|_{\mathbf{A}}^2 \leq 1 \\ \mathbf{A} \geq 0$$

Koristi se gradijentni silazak za rješavanje funkcije g iza kojeg slijede iterativne projekcije koje osiguravaju da vrijede oba dva postavljena uvjeta tj. nakon gradijentnog

koraka $\mathbf{A} = \mathbf{A} + \alpha \nabla_{\mathbf{A}} g(\mathbf{A})$, projicira se matrica \mathbf{A} u $C_1 = \{\mathbf{A} : \sum_{(x_i, x_j) \in S} \| \mathbf{x}_i - \mathbf{x}_j \|_{\mathbf{A}}^2 \leq 1\}$ i C_2

$= \{\mathbf{A} : \mathbf{A} \geq 0\}$ kako je prikazano u algoritmu. Razlog ovakvog pristupa je računski nezahtjevno projiciranje matrice \mathbf{A} u C_1 i C_2 . Prvi uvjet $\arg \min_{\mathbf{A}'} \{ \| \mathbf{A}' - \mathbf{A} \|_F : \mathbf{A}' \in C_1 \}$ lako je riješiti skupom rijetkih linearnih jednadžbi i složenosti je $O(n^2)$. Sljedeći korak projekcije najprije rastavlja matricu \mathbf{A} na $\mathbf{A} = \mathbf{X}^T \Lambda \mathbf{X}$, gdje je $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ dijagonalna matrica svojstvenih vrijednosti matrica \mathbf{A} koju zamjenjujemo sa $\Lambda' = \text{diag}(\max\{0, \lambda_1\}, \max\{0, \lambda_2\}, \dots, \max\{0, \lambda_n\})$, a matricu \mathbf{A} sa $\mathbf{A} = \mathbf{X}^T \Lambda' \mathbf{X}$.

(Mochihashi,2004.) kako sam i kaže, dijeli isti cilj i nasljeđuje djelo (Xing et al.,2003.). Za razliku od parova sličnih točaka i minimaliziranja udaljenosti među njima, (Mochihashi,2004.) traži matricu \mathbf{M} koja minimalizira udaljenost između pojedinih

³³ Da ovaj uvjet bude ispunjen, potrebno je modificirati Newtonov pomak sa $H^{-1} \nabla g$ sa $\alpha H^{-1} \nabla g$, gdje je α parametar koraka koji daje najviši spust uz uvjet da je $A_{ij} \geq 0$

podatkovnih vektora i pridruženom centroidu grupe kojoj pripada. Matematički napisan ovaj kvadratni problem uz \mathbf{s}_i koji označava i -ti podatkovni vektor i \mathbf{c}_i centroid i -te grupe od njih N, je

$$\begin{aligned} \mathbf{M} &= \arg \min_{\mathbf{M}} \sum_{i=1}^N \sum_{sj \in X_i} d_M(\mathbf{s}_j, \mathbf{c}_i)^2 \\ &= \arg \min_{\mathbf{M}} \sum_{i=1}^N \sum_{sj \in X_i} (\mathbf{s}_j - \mathbf{c}_i)^T \mathbf{M} (\mathbf{s}_j - \mathbf{c}_i) \end{aligned}$$

uz uvjet da je $|\mathbf{M}| = 1$. Uvjet, kako je prije navedeno, isključuje mogućnost da je $\mathbf{M} = 0$, a 1 je proizvoljna konstanta. Promjenom te konstante sa c dobivamo novo rješenje $c^2 \mathbf{M}$.

Matrica koja rješava problem minimizacije definiran prije je

$$\mathbf{M} = |\mathbf{A}|^{1/n} \mathbf{A}^{-1}$$

gdje je $\mathbf{A} = [a_{kl}]$

$$a_{kl} = \sum_{i=1}^N \sum_{sj \in X_i} (\mathbf{s}_{jl} - \mathbf{c}_{il}) (\mathbf{s}_{jk} - \mathbf{c}_{ik})$$

Kad je \mathbf{A} singularna, možemo umjesto \mathbf{A}^{-1} koristiti njen Moore-Penrose pseudoinverz \mathbf{A}^+ ³⁴. Općenito se \mathbf{A} sastoji od jezičnih značajki, vrlo je rijetka i često singularna. Dakle, \mathbf{A}^+ je često potrebna u računanju.

U problemima grupiranja gdje su značajke lingvističke ovakav način dobivanja funkcije udaljenosti je računski nemoguć. U matrica M imala bi *broj riječi x broj riječi* elemenata. Ipak, u kombinaciji sa metodama redukcije dimenzionalnosti, npr. SVD-om, ovakva metrika bila bi upotrebljiva.

Općenito, učenje metrike nastoji naučiti dobru funkciju udaljenosti pomoću težina značajki na ulaznom prostoru. Iako se značajke može preslikati pomoću nelinearne funkcije, takva transformacija je eksplicitna i zbog računske složenosti ovakav pristup je nepraktičan. Jezgrena metoda učenja funkcije udaljenosti uspješno zaobilazi problem računske složenosti.

³⁴ Moore-Penrose pseudoinverz \mathbf{A}^+ jedinstvena je matrica koja ima svojstvo da vrijedi $x = \mathbf{A}^+ y$ kao rješenje najmanjih kvadrata za $\mathbf{Ax} = y$ i ako \mathbf{A} je singularna.

5.3.2 Jezgreno učenje

Jezgrene metode implicitno preslikavaju skup podatkovnih vektora $X = \{x\}_{i=1}^n$ na ulaznom prostoru I u drugi, više dimenzionalni (moguće i beskonačni) projicirani prostor P pomoću funkcije preslikavanja ϕ . Jezgrena funkcija k definirana je kao skalarni produkt između dva vektora u projiciranom prostoru P , kao

$$k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle.$$

Jezgrene metode koriste skalarni produkt (K) kao mjeru sličnosti. Potrebno je odabrat i pravilni K koji mora biti pozitivan, (semi-) definitivan i simetričan. Koriste se polinomne jezgrene funkcije, Gaussov RBF, Laplace RBF. U P , udaljenost se može izračunati pomoću jezgrenog trika

$$d(x_i, x_j) = \sqrt{K(x_i, x_i) + K(x_j, x_j) - 2K(x_i, x_j)}$$

gdje ispravnost funkcije udaljenosti slijedi iz ispravnosti jezgrene matrice K . Zbog velike uloge jezgrene funkcije, loš odabir vodi lošim performansama ove, jezgrene metode. Umjesto da se koristi predefinirani K mnogi se istražuje i učenje K na skupu podataka za učenje.

(Cristianini et al., 2002) predložili su «kernel alignment» za mjerjenje sličnosti između dvije jezgrene funkcije. Geometrijski, «alignment» je definiran kao kosinusni kut između dvije jezgrene matrice, nakon transformiranja te dvije matrice u vektore. Vezano uz to, predložili su i idealiziranu jezgenu matricu (K^*). Prepostavimo da je $y(x_i) \in \{-1, 1\}$ je oznaka klase za x_i . K^* je definiran kao

$$K^*(x_i, x_j) = 1, \text{ ako } y(x_i) = y(x_j)$$

$$0, \text{ ako } y(x_i) \neq y(x_j)$$

a pretstavlja ciljnu jezgenu matricu kojoj početna treba težiti. Učenje jezgrene matrice proučavali su (Kwok et al., 2003) i (Zhang, 2003.) te (Wu et al., 2004). Daljnje ulaženje u ovu problematiku izvan je teme ovoga rada.

5.4 Jezgreni k-means algoritam

Kao i kod spherical k-mean algoritma potrebno je definirati objektivnu funkciju koju čemo u ovom slučaju minimalizirati. Za ulazni skup dokument-vektora $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$, centroida grupa $\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_k$, i funkciju f koja svakom dokument-vektoru pridružuje jednu grupu nastojimo minimalizirati

$$\sum_{j=1}^k \|\phi(\mathbf{x}_i) - \mathbf{c}_{f(x_i)}\|^2.$$

Početni centroidi su slučajno odabrani, da bi se gornja funkcija minimalizirala mijenjajući funkciju f kao i centroide grupa. Postupak je iterativan. Najprije se točke pomicu u grupu čijem su centroidu najbliže. Jasno je da će to minimalizirati gornju funkciju. Zatim se u sljedećem koraku mijenja centroid grupa sa centrom mase svih točaka koje se nalaze unutar te grupe. Također se može dokazati da i ovaj korak minimalizira gornju funkciju. Nakon konačnog broja iteracija algoritam će konvergirati.

Definirajmo matricu \mathbf{A} dimenzija $\ell \times N$ kao

$$A_{ik} = 1, \text{ ako je } \mathbf{x}_i \text{ nalazi u grupi } k \\ 0, \text{ inače}$$

U svakom retku matrice \mathbf{A} mora postojati točno jedan element, dok suma svakog retka daje broj dokument-vektora u jednoj grupi. Koordinate centroida možemo izračunati

$$\mathbf{X}' \mathbf{A} \mathbf{D}$$

Gdje je \mathbf{X} matrica dokument-vektora u redcima, a \mathbf{D} je $N \times N$ dijagonalna matrica sa elementima koji su inverz stupčaste sume matrice \mathbf{A} (broj elemenata dodijeljenih pojedinoj grupi). Udaljenost vektora $\phi(\mathbf{x})$ od centroida dana je sa

$$\begin{aligned} \|\phi(\mathbf{x}) - \mathbf{c}_k\|^2 &= \|\phi(\mathbf{x})\|^2 - 2 \langle \phi(\mathbf{x}), \mathbf{c}_k \rangle + \|\mathbf{c}_k\|^2 \\ &= k(\mathbf{x}, \mathbf{x}) - 2 (\mathbf{k}' \mathbf{AD})_k + (\mathbf{DA}' \mathbf{X} \mathbf{X}' \mathbf{AD})_{kk} \end{aligned}$$

gdje je \mathbf{k} vektor skalarnih produkata između vektora $\phi(\mathbf{x})$ i vektora za učenje.

Dakle grupa kojoj se vektor $\phi(\mathbf{x})$ dodaje definirana je sa

$$\operatorname{argmin}_{1 \leq k \leq N} \|\phi(\mathbf{x}) - \mathbf{c}_k\|^2 = \operatorname{argmin}_{1 \leq k \leq N} (\mathbf{DA}' \mathbf{K} \mathbf{AD})_{kk} - 2 (\mathbf{k}' \mathbf{AD})_k$$

gdje je \mathbf{K} jezgrena matrica.

U svakoj se iteraciji na ovaj način računa grupa kojoj je dodijeljen pojedini vektor. Kao i u većini jezgrenih algoritama, tako i u ovome, algoritam dobiva informacije o ulaznim podacima i o prostoru u koji se preslikavaju preko jezgrevne matrice \mathbf{K} . Možemo reći da jezgrena matrica ima centralnu ulogu u jezgrenim algoritmima i dakle njen odabir jako je važan (Shawe-Taylor et al., 2004.).

5.5 Odabir jezgrene funkcije

Jezgrena matrica sadrži sve bitne informacije o položaju ulaznih podataka u prostoru značajki. Razmotrit ćemo dva načina pogrešnog odabira jezgrene matrice. Ako je jezgrena funkcija preopćenita i ne daje važnost određenim sličnostima među ulaznim podacima. Jezgrena funkcija daje iste težine svakom ulaznom paru, time nedijagonalni elementi jezgrene matrice postaju mali, a dijagonalni teže 1. Jezgrena matrica svedena je na koncept identiteta. Dakle, došli smo do pretreniranosti jer možemo točno klasificirati skup za učenje, dok jezgrena matrica nema svojstvo generalizacije na novim podacima. S druge strane, ako je matrica potpuno uniformna, tada je svaki ulaz sličan svakom drugom ulazu. Ovaj problem preslikava svaki ulaz u isti vektor u prostoru značajki i vodi podtreniranosti podataka. Geometrijski, prvi problem odgovara preslikavanju ulaza u ortogonalne vektore u prostoru značajki, dok su u drugom problemu sve točke preslikane u istu sliku.

Idealan odabir jezgrene matrice temeljio bi se na našem a priori znanju problema i pojednostavnio bi učenje na odabir funkcije uzorka u prostoru značajki. Nažalost, nije uvijek moguće pravilno odabrati jednu jezgrenu matricu a priori, već moramo odabrati obitelj jezgrenih matrica koje reflektiraju naše a priori očekivanje i ostavljamo si na izbor odabir jezgrene matrice koju ćemo koristiti. Dakle, sustav koji uči mora sad riješiti dva problema: odabrati jezgrenu matricu iz obitelji jezgrenih matrica i poslije ili konkurentno odabrati funkciju uzorka u prostoru značajki.

Mnogi pristupi se mogu koristiti u rješavanju dvojnog problema učenja. Najjednostavniji koriste malo podataka iz skupa za učenje, i često ne uzimaju o označenosti učenja s učiteljem. Razrađenje metode koje koriste označenosti trebaju mjeru «dobrote» u rješavanju jezgrenog dijela problema. Dakle, uvodi se pojam sličnosti među jezgrenim matricama i odabire se jezgrena matrica koja je najbliža idealnoj jezgrenoj matrici

$$K^*(x_i, x_j) = 1, \text{ ako } y(x_i) = y(x_j)$$

$$0, \text{ ako } y(x_i) \neq y(x_j).$$

Mjera sličnosti među jezgrama, ili u slučaju idealne jezgrene matrice između jezgrene matrice i cilja, mora zadovoljavati osnovna svojstva: mora biti simetrična, maksimizirati se kad su argumenti jednaki i minimalizirati se kad se primjeni na dvije nezavisne jezgrene matrice. Nadalje, u slučaju idealne jezgrene matrice ograničit ćemo se samo na skup za učenje jer se idealna jezgrena matrica može izračunati samo na tom

skupu. Dakle, potrebno je dokazati da se procjena može dobiti koristeći samo manji skup podataka.

(Christianini et al.,2001.) uvode pojam jezgrenog alignmenta, mjere sličnosti između dvije jezgrene funkcije ili jezgrene funkcije i ciljne funkcije.

Alignment jezgrene funkcije k_1 i jezgrene funkcije k_2 na skupu $S = \{x_1, x_2, \dots, x_n\}$ definiran je kao

$$A(S, k_1, k_2) = \frac{\langle K_1, K_2 \rangle_F}{\sqrt{\langle K_1, K_1 \rangle_F \langle K_2, K_2 \rangle_F}}$$

gdje su K_i jezgrene matrice dobivene primjenom funkcije k_i na skupu S . Alignment možemo promatrati i kao kosinus kuta između dva bi-dimenzionalna vektora K_1 i K_2 . Ako je K_2 idealni kernel definiran sa $K_2 = yy'$, gdje je y vektor oznaka $\{-1, 1\}$ na skupu S , tada

$$A(S, K, yy') = \frac{\langle K, yy' \rangle_F}{\sqrt{\langle K, K \rangle_F \langle yy', yy' \rangle_F}} = \frac{\langle K, yy' \rangle_F}{m\sqrt{\langle K, K \rangle_F}} \text{ jer vrijedi} \\ \langle yy', yy' \rangle_F = m^2$$

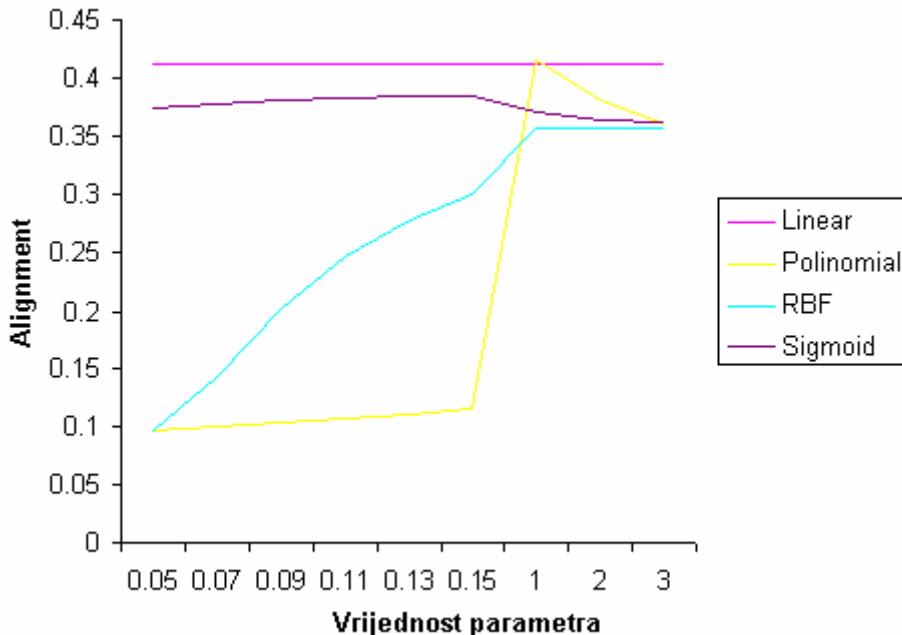
Dakle, mjeru sličnosti između dvije jezgre transformirali smo u mjeru dobrote pojedine jezgrene funkcije.

(Christianini et al. ,2001.) pokazali su da ako je očekivanje alignmenta visoko, onda postoji funkcija koja dobro generalizira. Dakle, možemo optimirati alignment na skupu za učenje i očekivati da će na testnom skupu jezgrena funkcija također davati visoki alignment i imati dobre performanse.

Najprije ćemo odabrati obitelj(linearna, RBF,...) jezgrenih funkcija. Zatim je potrebno odabrati parametre tih funkcija koji će davati najbolje rezultate na Vjesnikovoj bazi.

Linearna	$F(x_i, x_j) = (x_i \cdot x_j)$
RBF	$F(x_i, x_j) = \exp(-\gamma \ x_i - x_j\ ^2)$
Sigmoidalna	$F(x_i, x_j) = \tanh(a^* x_i \cdot x_j)$
Polinomialna	$F(x_i, x_j) = (x_i \cdot x_j + 1)^d$

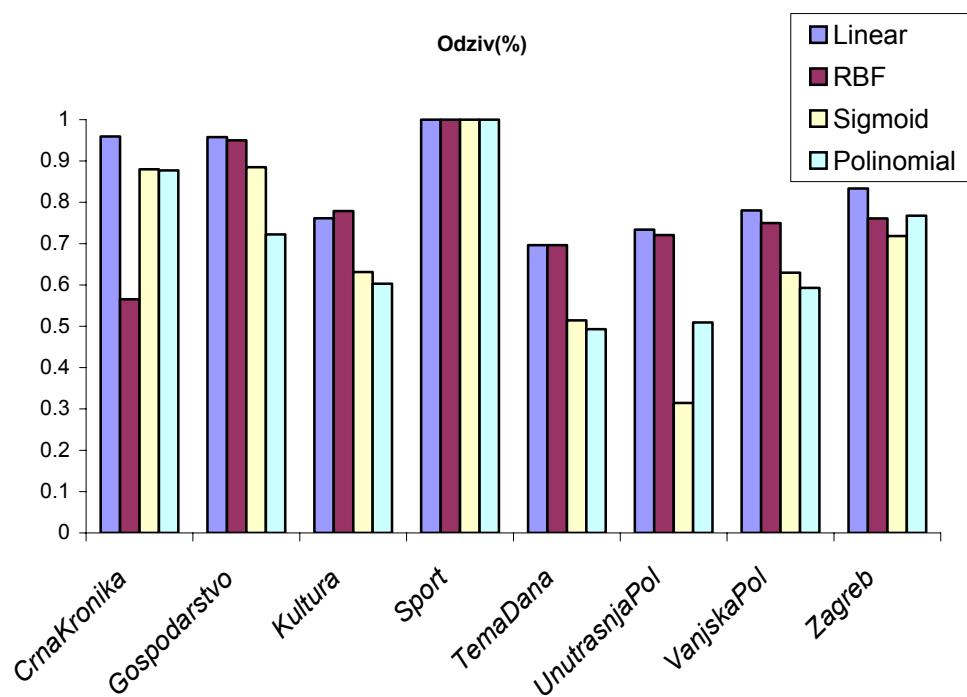
Rezultati



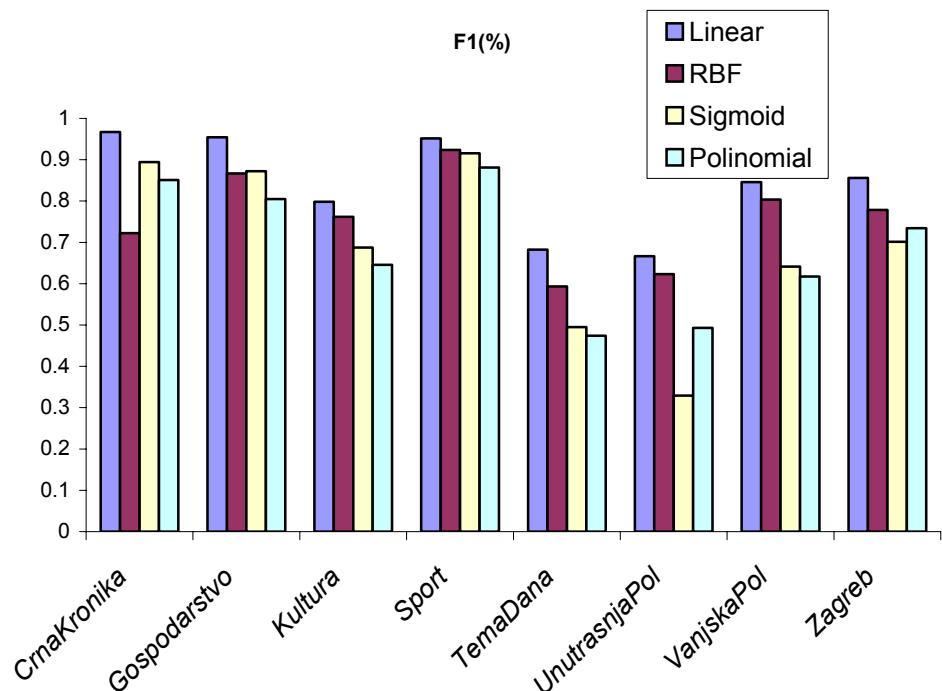
Slika 28. Alignment za različite parametre σ i različite jezgrene funkcije

Prema rezultatima vidimo da linearna funkcija (koja zapravo odgovara kosinusu kuta između vektora) ima najveći alignment sa našim idealnim kernelom dobivenim na skupu za učenje. Ta funkcija nije ovisna o parametrima. Linearnoj funkciji najviše se približava polinomialna za parametar 1, ili $(\mathbf{x}_i \cdot \mathbf{x}_j + 1)^1$ za koju možemo reći da je skoro identična linearnej. Ako primijenimo jezgrevi k-means na skupu za učenje, možemo vidjeti da točnost grupiranja na tom skupu sa različitim jezgrenim matricama potvrđuje prepostavke postavljene sa alignmentom.

Isto tako uspješnost jezgrene funkcije kao metrike možemo ocijeniti i na taj način da gledamo $\phi(\mathbf{x}) - \mathbf{c}_k$, ili udaljenost vektora u prostoru značajki od njegovog konceptnog vektora na skupu za učenje. Odziv označava koliki je broj vektora koji pripada dotičnoj kategoriji pravilno klasificiran, ili u našem slučaju najbliži centroidu upravo svoje grupe. Prema rezultatima za sport(kategorija 4) možemo vidjeti da su svi dokumenti u toj kategoriji u svim metrikama i promjenama jezgrene matrice zapravo najbliži uvijek svom centroidu. Druge kategorije mijenjaju se izborom jezgrene funkcije, ali većinom vidimo da linearna jezgrena funkcija daje najbolje rezultate. Ako pogledamo mjeru F1, tj. i preciznosti i odzivu damo istu važnost, vidimo da linearna funkcija nadmašuje sve druge.



Slika 29. Postotak vektora koji su najbliži svom centroidu upotreboom različitih metrika



Slika 30. Mjera F1, ocjena odziva i preciznosti vektora i njihove udaljenosti od pripadnih centroida

6 KONCEPTNO INDEKSIRANJE

Konceptno indeksiranje također koristi grupe dokument vektora za preslikavanje u niže dimenzionalni prostor. Postoje nenađizirana i nadzirana tehnika redukcije dimenzionalnosti.

Konceptno indeksiranje na sljedeći nenađizirani način računa redukciju dimenzionalnosti. Kao i kod konceptne dekompozicije, ako je k željeni broj dimenzija, koristi se k konceptnih vektora kao osi novog k -dimenzionalnog prostora. Uz n dokument vektora $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ u \mathbb{R}^m možemo definirati *rječ x dokument* matricu kao

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n],$$

i konceptnu matricu $m \times k$ tako da su stupci matrice konceptni vektori

$$\mathbf{C} = [\mathbf{c}_1 \ \mathbf{c}_2 \ \dots \ \mathbf{c}_k].$$

Svaki konceptni vektor predstavlja osi novog prostora, i k -dimenzionalna reprezentacija dokument vektora dobiva se projiciranjem u ovaj prostor. Ta projekcija može se matrično zapisati, ako je \mathbf{d}_i i -ti dokument vektor tada $\mathbf{d}\mathbf{C}$ predstavlja k -dimenzionalnu reprezentaciju tog vektora, a $\mathbf{X}^T\mathbf{C}$ predstavlja preslikani cijeli skup dokumenata. Slično, upit \mathbf{q} se preslikava se $\mathbf{q}\mathbf{C}$. Sličnost između dva dokumenta u reduciranom prostoru računa se kosinusnom sličnosti između reduciranih vektora.

U slučaju nadzirane redukcije dimenzionalnosti, konceptno indeksiranje koristi već postojeće klase dokumenta i unutra njih traži grupe sličnih dokumenta. U najjednostavnijem slučaju svaka grupa pripada jednoj klasi iz skupa dokumenata. U ovom slučaju rang niže dimenzionalnog prostora bit će jednak broju klasa. Općenito tražimo niže dimenzionalan prostor ranga k koji je veći od broja klasa I . Najprije se grupira u I grupa tako da se stvori grupa za svaku klasu dokumenata, a onda se one dalje particioniraju dok ne dobijemo k grupa. Bitno je primjetiti da će svaka grupa dokumenata sadržavati dokumente samo jedne klase. Nakon što su grupe određene konceptno indeksiranje ne razlikuje se od opisanog u prethodnom odjeljku – kao u slučaju nenađizirane redukcije dimenzionalnosti.

Za razumijevanje konceptnog indeksiranja mora se kao prvo shvatiti bit konceptnih vektora i drugo značenje reducirane dimenzionalne reprezentacije svakog dokumenata. Kako smo u prethodnom poglavlju opisali konceptni vektor sumarizira kontekst dokument vektora u njegovoj grupi. Također svaka grupa nastoji obuhvatiti čim sličnije vektore, tj. dokument vektori koji se nalaze u različitim grupama manje su slični nego oni u istoj grupi. Za k konceptnih vektora i dokument \mathbf{d} i -ta koordinata reduciranog dimenzionalnog prostora je kosinusna sličnost između dokumenta \mathbf{d} i i -tog konceptnog vektora. Dimenzije

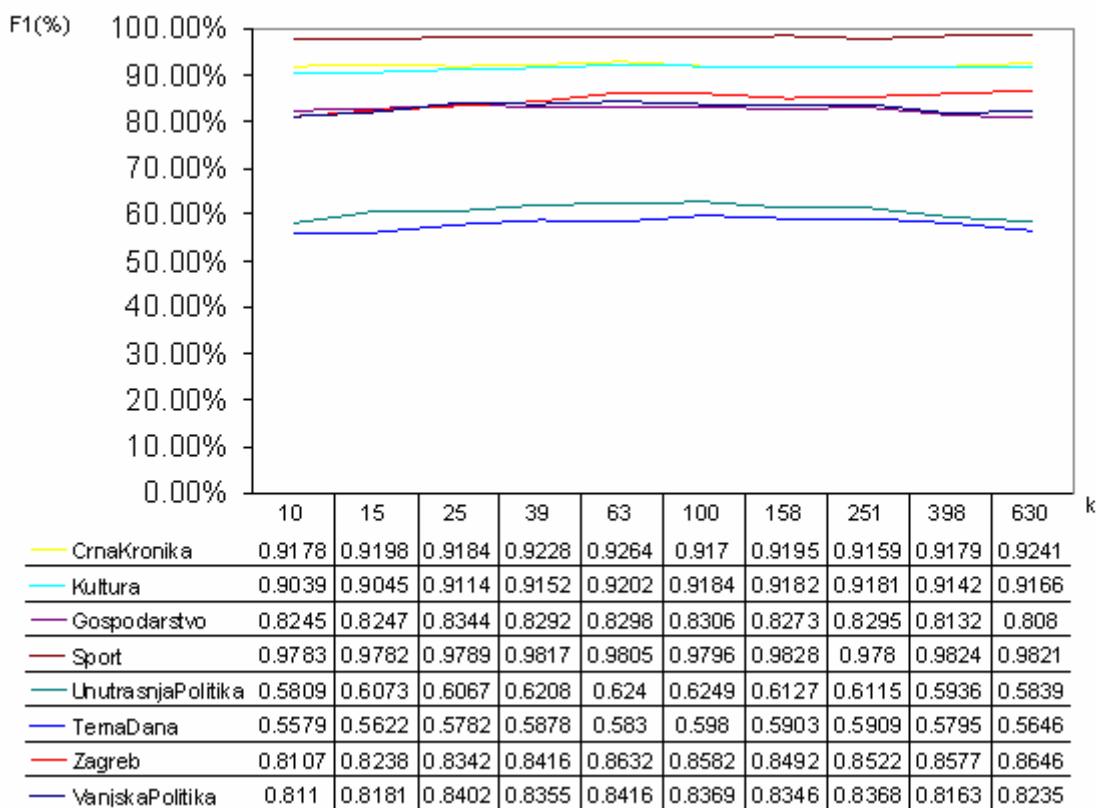
reduciranog prostora odgovaraju stupnju do kojeg dokument vektor odgovara konceptu obuhvaćenom konceptnim vektorom. Ova interpretacija preslikavanja dokument-vektora u niže dimenzionalni prostor razlog je za nazivanje ovakve redukcije dimenzionalnosti *konceptno indeksiranje* (Karypis et al., 2000).

6.1 Eksperimentalni rezultati

6.1.1 Nadzirana i nenadzirana redukcija dimenzionalnosti

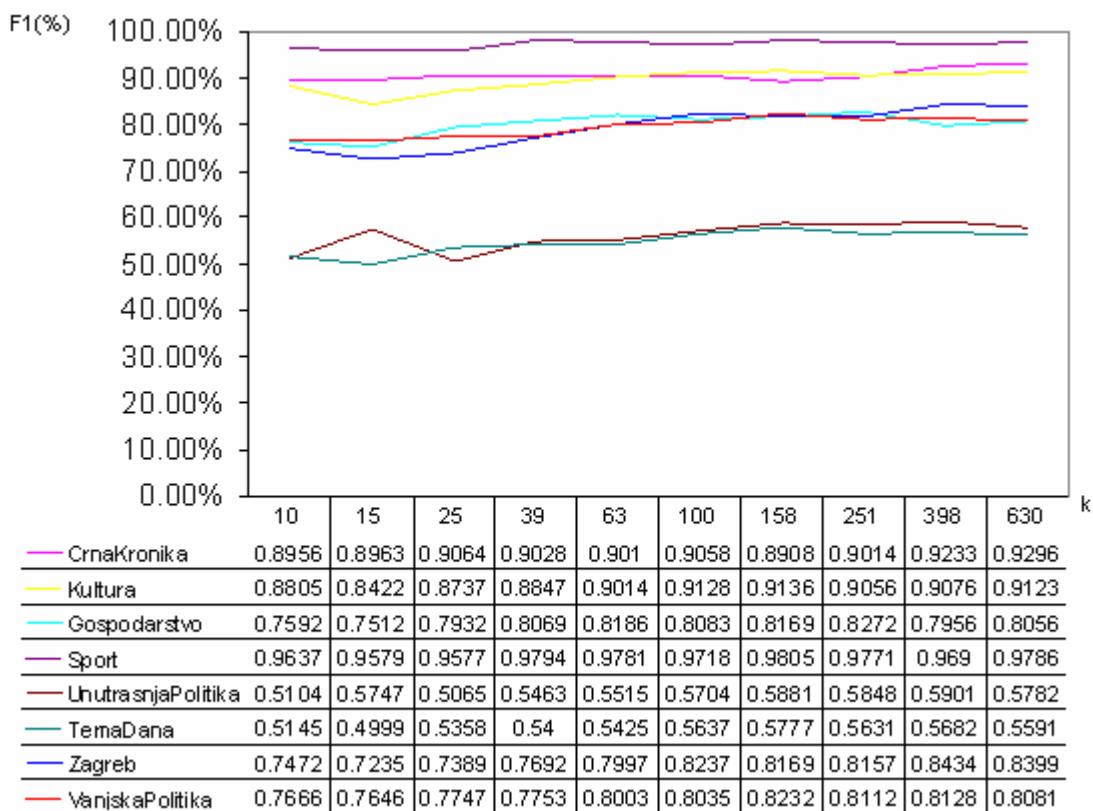
Koristeći konceptno indeksiranje kao redukciju dimenzionalnosti, ulaz u SVM smanjen je sa matrice veličine $m \times n$ ($100\ 000 \times 10\ 000$) na matricu $k \times n$ ($k \times 10\ 000$). Zajedno sa smanjenjem ulaznih podataka možemo očekivati lošije rezultate koji će se postupno popravljati povećanjem k . Parametri SVM-a postavljeni su na $C = 1000$, $\gamma = 0.01$ postupkom krosvalidacije.

Prvo su testirani nenadzirana i nadzirana redukcija dimenzionalnosti.



Slika 31. Mjera $F1$ po kategorijama sa k konceptnih vektora

Nadzirana redukcija dimenzionalnosti

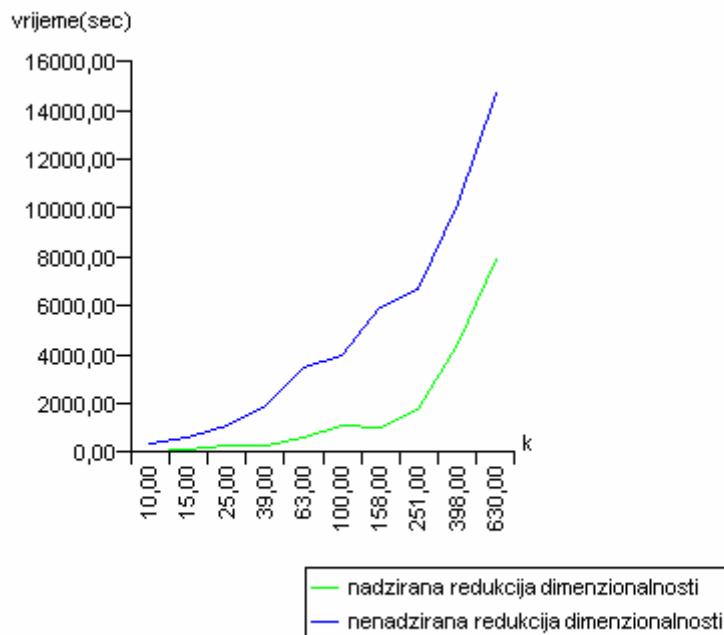


Slika 32. Mjera F1 po kategorijama sa k konceptnih vektora

Nenadzirana redukcija dimenzionalnosti

Uspoređivanjem nadzirane i nenadzirane redukcije dimenzionalnosti, možemo vidjeti da zapravo dobivamo klasifikaciju iste točnosti. Mjera F1 za nenadziranu redukciju dimenzionalnosti je manja ili jednaka odgovarajućoj mjeri F1 nadzirane redukcije dimenzionalnosti po kategorijama. Jedina bitna razlika je vremenska. Kako kod nenadzirane redukcije trebamo uspoređivati svaki dokument vektor sa svim konceptnim vektorima kroz iteracije, a kod nadzirane samo sa konceptnim vektorima njegove kategorije, možemo zaključiti i dokazati da će nenadzirana redukcija dimenzionalnosti biti znatno sporija.

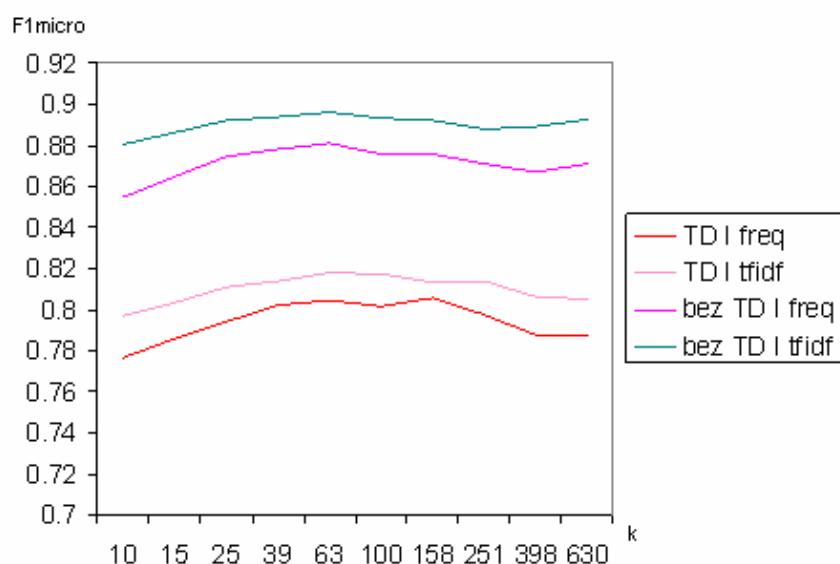
U sljedećim testiranjima korištena je samo nadzirana redukcija dimenzionalnosti te su se mijenjale matrice *rječ x dokument* dobivene različitom morfološkom normalizacijom. Tri vrste matrica „-i“, „-id“, „-idt“ označavaju 3 vrste morfološke normalizacije kojom su dobivene: flektivna, derivacijska te terminirajuća.



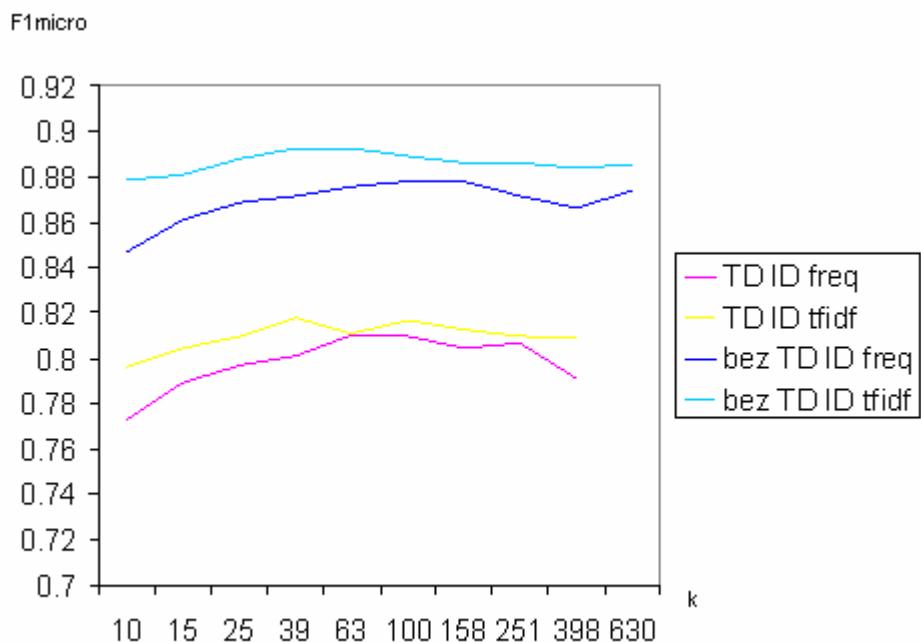
Slika 33. Vremenska ovisnost(sekunde) algoritma redukcijom dimenzionalnosti u k grupa

Usporedba nadzirane i nenadzirane redukcije dimenzionalnosti

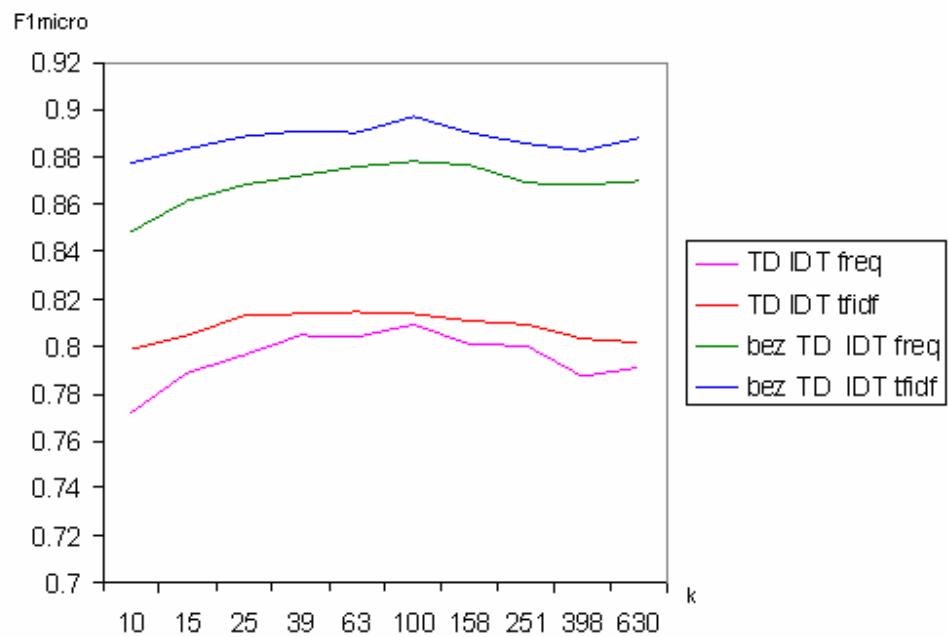
6.1.2 Utjecaj morfološke normalizacije na hrvatski jezik



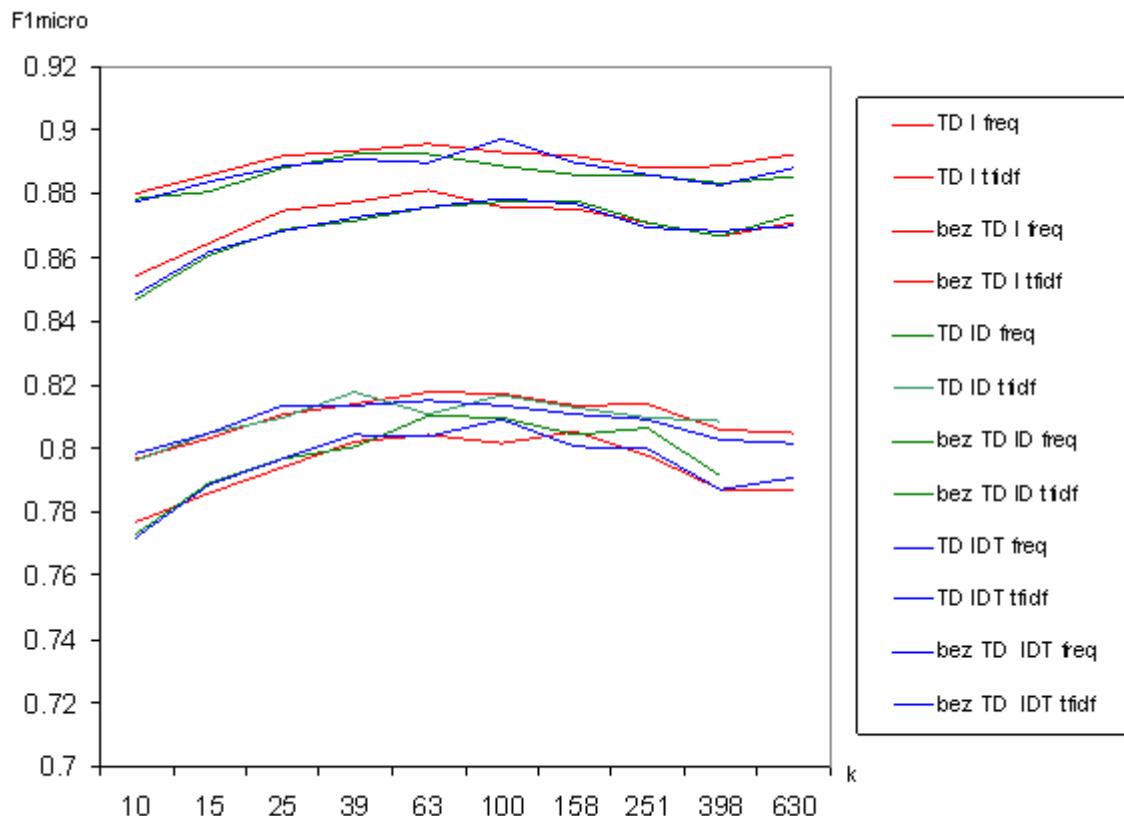
Slika 34. Mjera $F1^{micro}$ za baze sa i bez teme dana(TD, bez TD), korištenjem frekvencija ili mjere tfidf za flektivnu morfološku normalizaciju



Slika 35. Mjera $F1^{micro}$ za baze sa i bez teme dana(TD, bez TD), korištenjem frekvencija ili mjere tfidf za derivacijsku morfološku normalizaciju



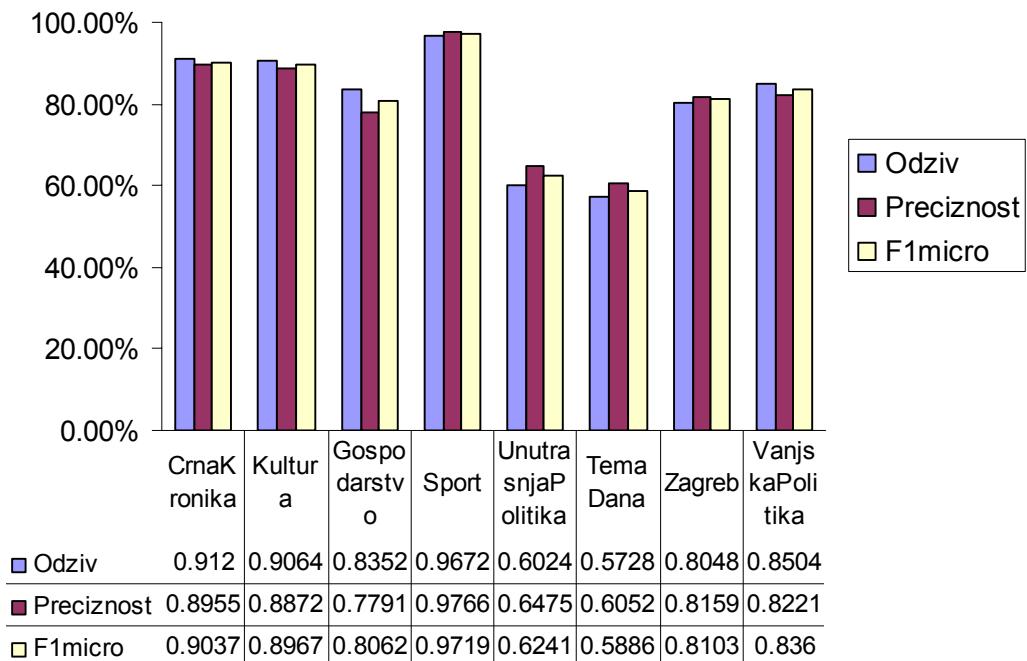
Slika 36. Mjera $F1^{micro}$ za baze sa i bez teme dana(TD, bez TD), korištenjem frekvencija ili mjere tfidf za terminirajuću morfološku normalizaciju



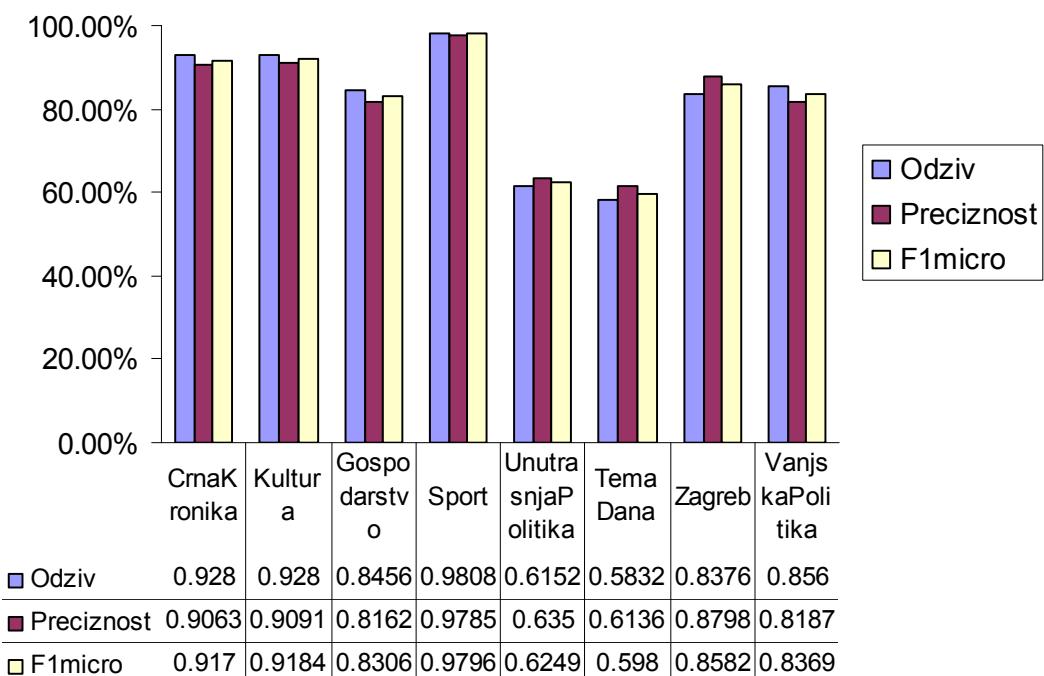
Slika 37. Mjera $F1^{micro}$ za baze sa i bez teme dana(TD, bez TD), sve 3 morfološke normalizacije(I, ID, IDT) te korištenjem frekvencija ili mjere tfidf po broju korištenih konceptnih vektora

Sve morfološke obrade daju približno iste rezultate, iako ih možemo sortirati po neznatnoj uspješnosti: flektivna, terminirajuća, derivacijska. Velike su promjene u rezultatima sa i bez kategorije tema dana, te u oba slučaja tfidf pokazuje se kao bolji odabir.

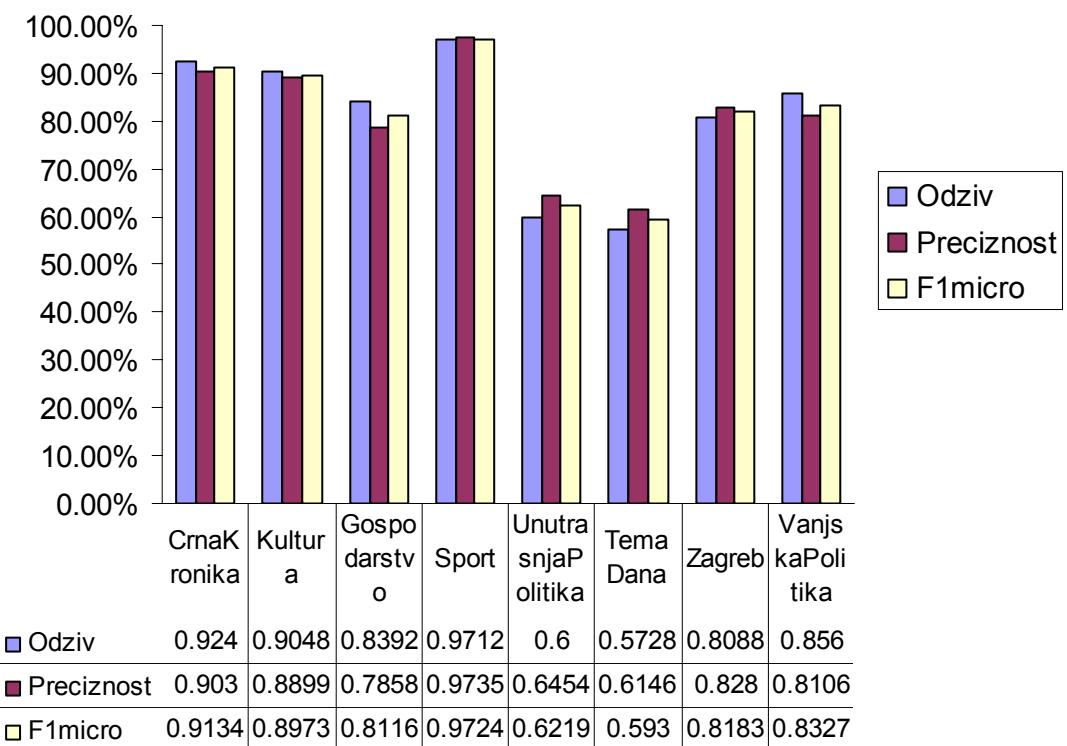
6.1.3 Mjere uspješnosti za različite morfološke normalizacije, sa i bez kategorije tema dana (broj konceptnih vektora: 100)



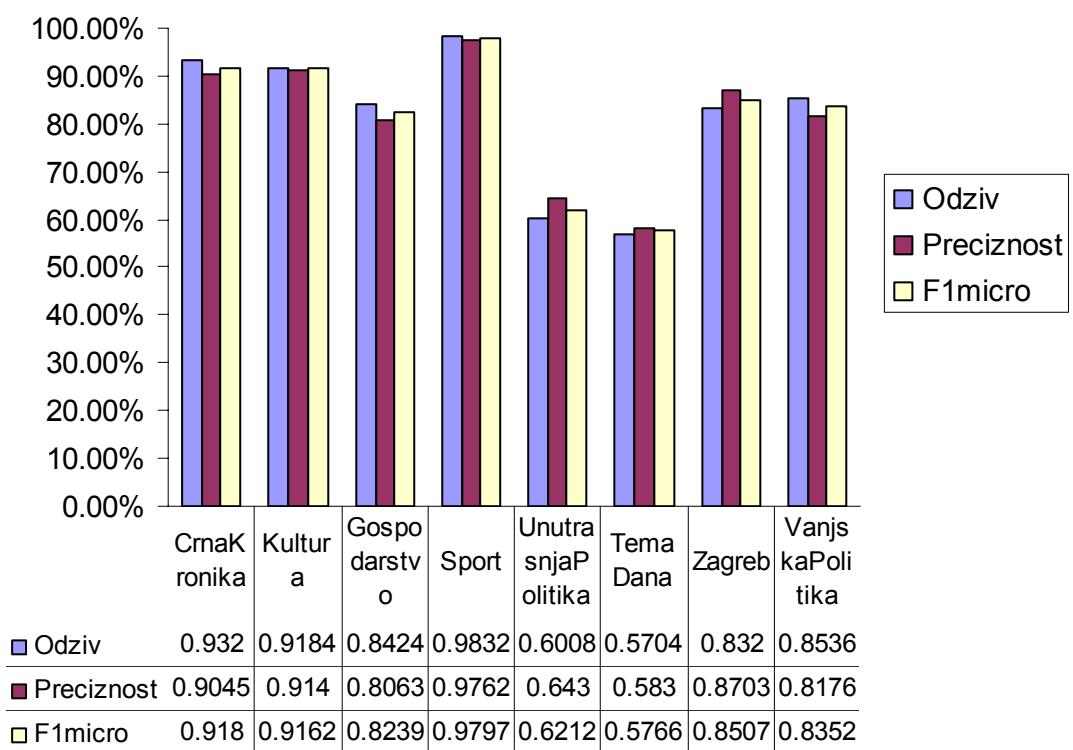
Slika 38. Mjera $F1^{micro}$ za flektivnu morfološku normalizaciju, frekvencija



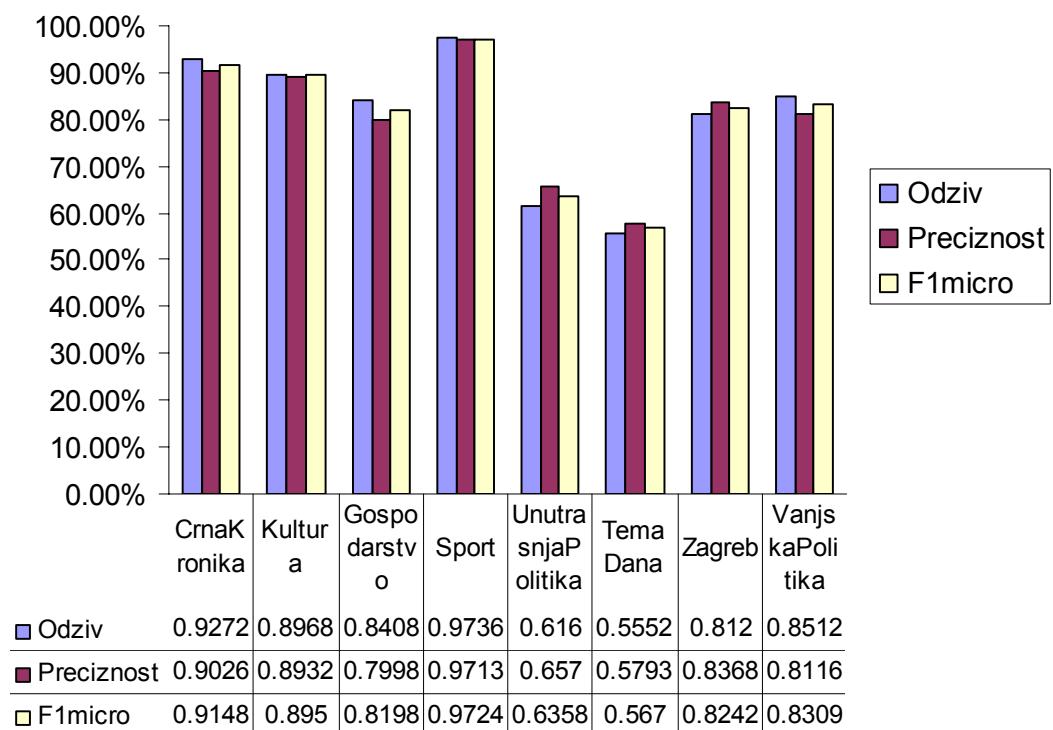
Slika 39. Mjera $F1^{micro}$ za flektivnu morfološku normalizaciju, tfidf



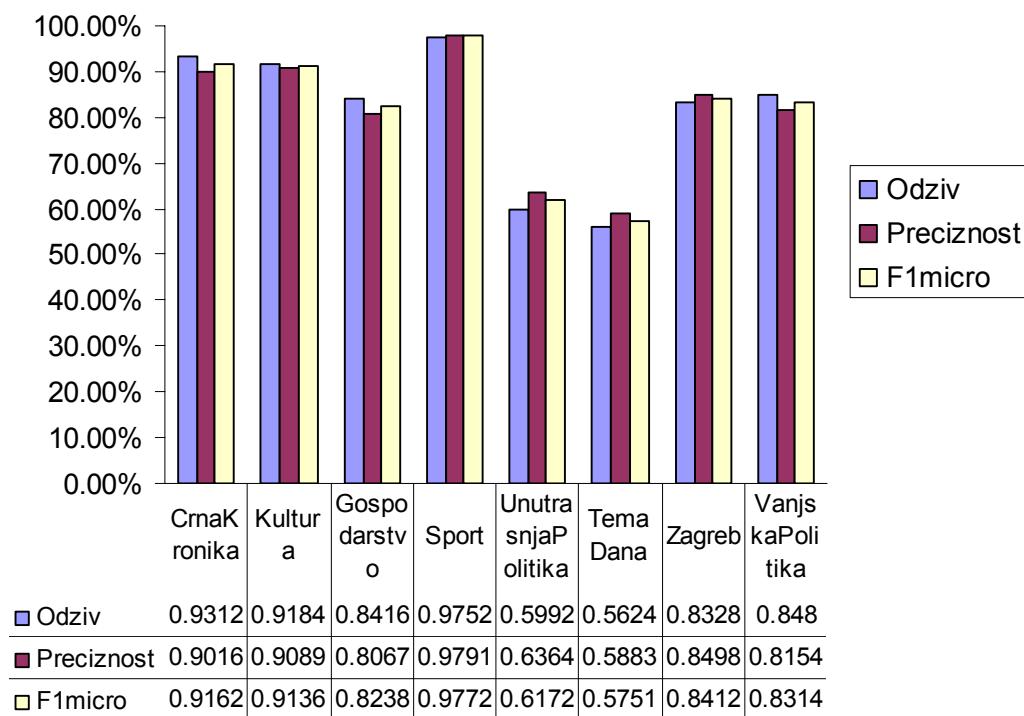
Slika 40. Mjera $F1^{micro}$ za derivacijsku morfološku normalizaciju, frekvencija



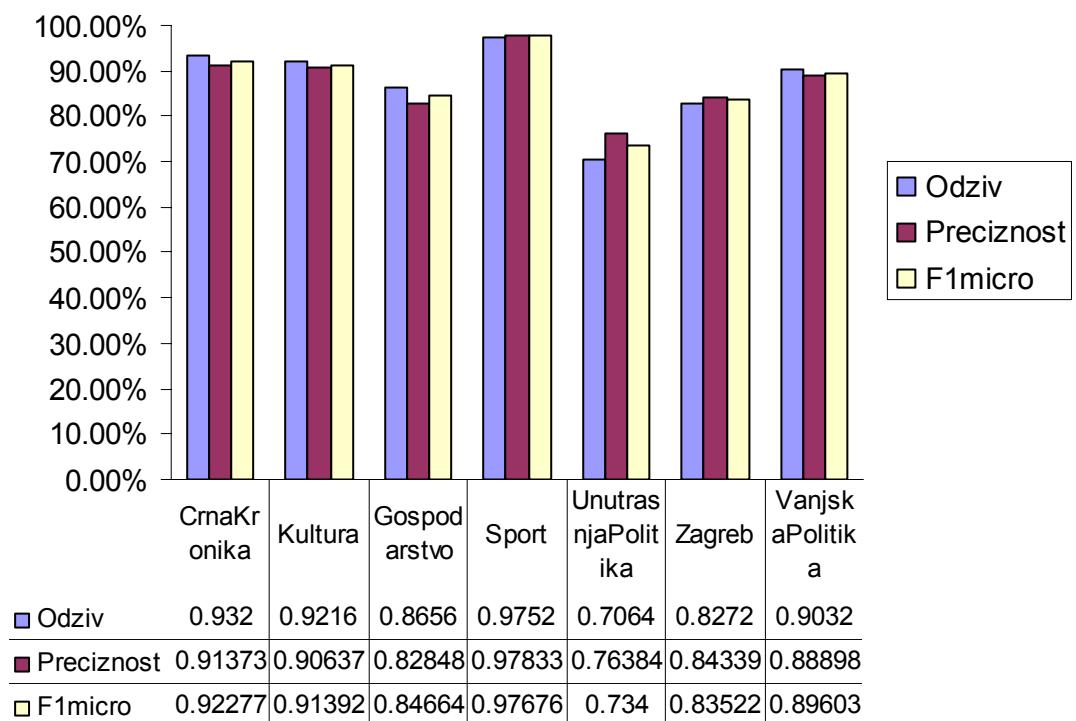
Slika 41. Mjera $F1^{micro}$ za derivacijsku morfološku normalizaciju, tfidf



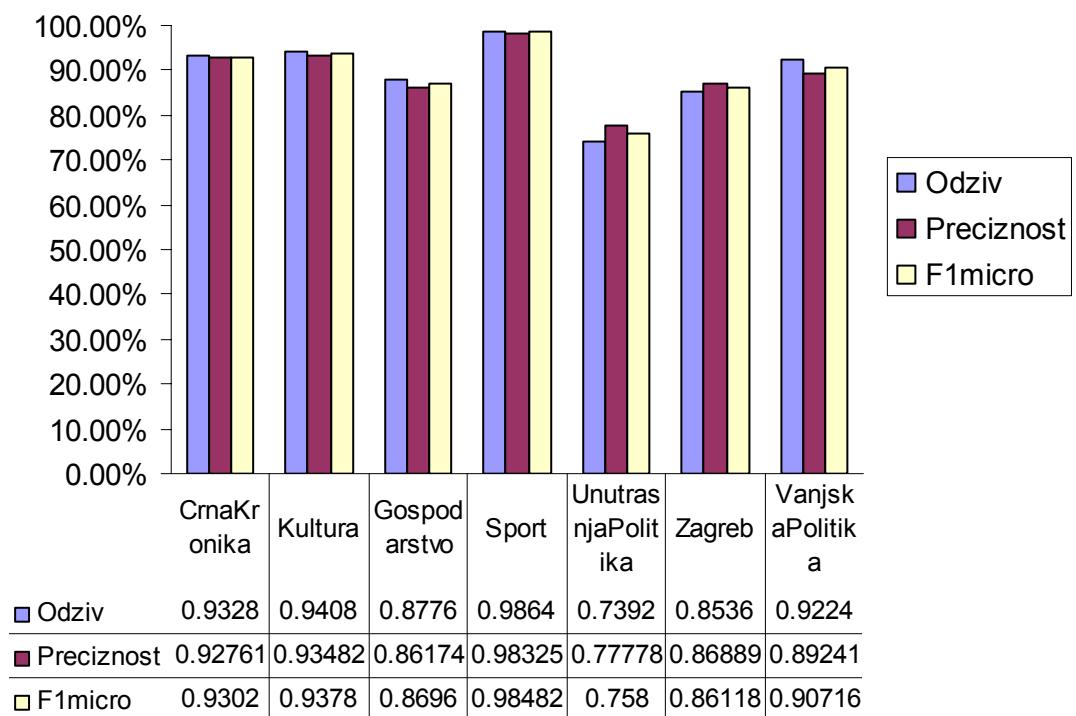
Slika 42. Mjera $F1^{micro}$ za terminirajuću morfološku normalizaciju, frekvencija



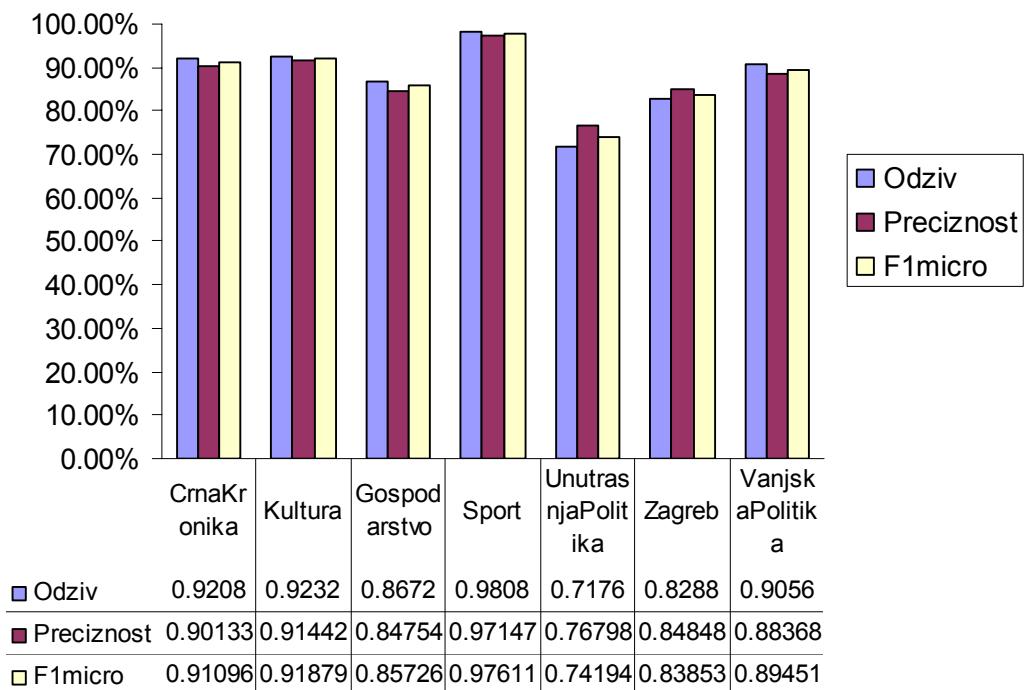
Slika 43. Mjera $F1^{micro}$ za terminirajuću morfološku normalizaciju, tfidf



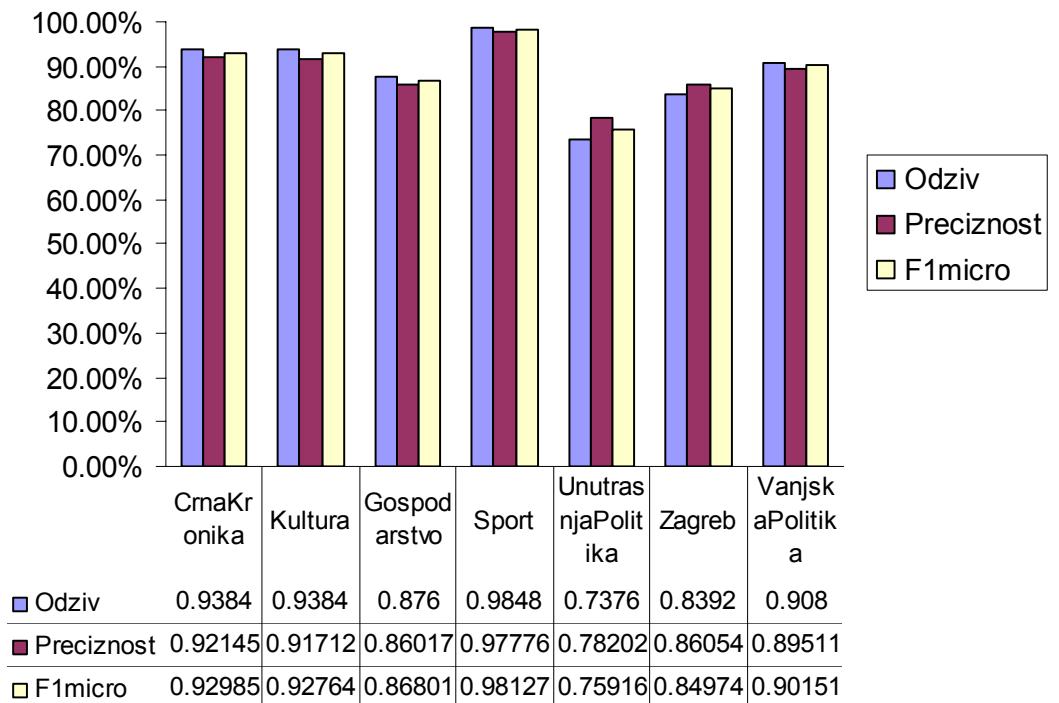
Slika 44. Mjera $F1^{micro}$ za flektivnu morfološku normalizaciju, frekvencija, bez teme dana



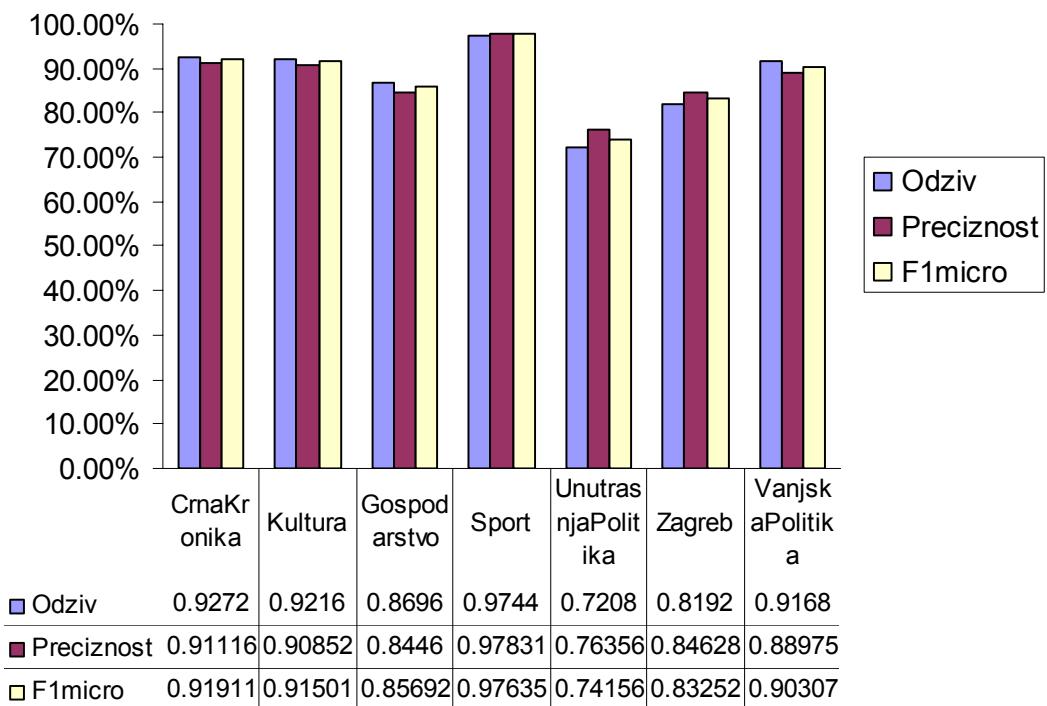
Slika 45. Mjera $F1^{micro}$ za flektivnu morfološku normalizaciju, tfidf, bez teme dana



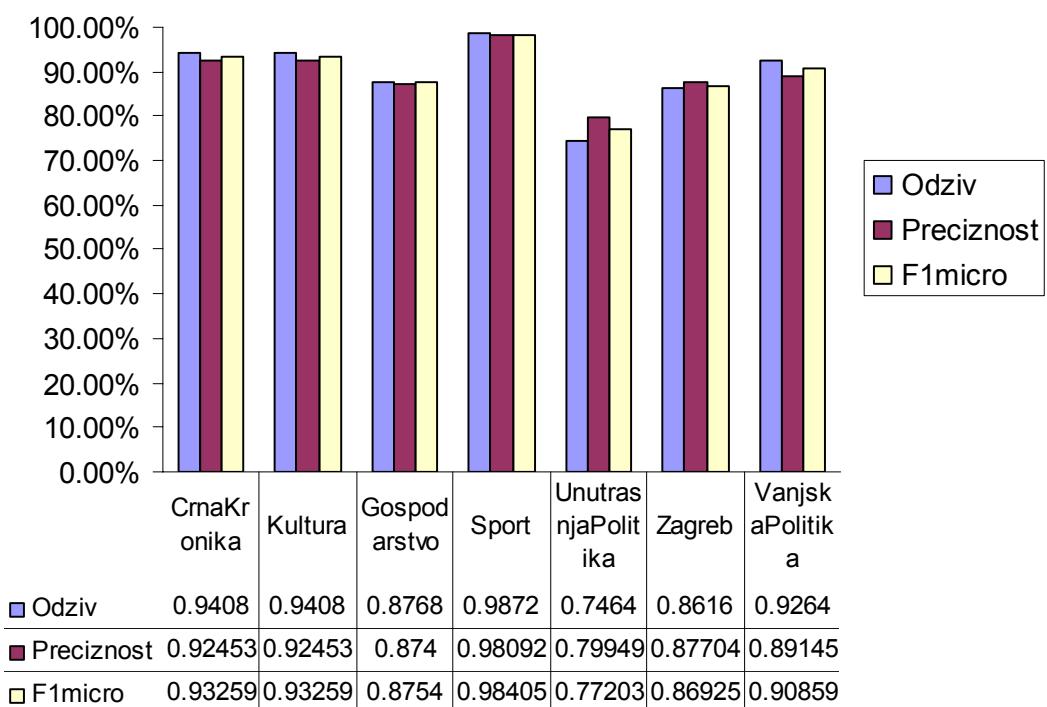
Slika 46. Mjera $F1^{micro}$ za derivacijsku morfološku normalizaciju, frekvencija, bez teme dana



Slika 47. Mjera $F1^{micro}$ za derivacijsku morfološku normalizaciju, tfidf, bez teme dana



Slika 48. Mjera $F1^{micro}$ za terminirajuću morfološku normalizaciju, frekvencija, bez teme dana



Slika 49. Mjera $F1^{micro}$ za terminirajuću morfološku normalizaciju, tfidf, bez teme dana

7 ZAKLJUČAK

Nakon završetka ovog rada, dok pišem zadnju stranicu, glavna misao mi je kako i na čemu je moguće doraditi metode koje sam obuhvatila i proučila. Istovremeno se nadam da će ovaj diplomski biti dobar temelj za istraživanja koja slijede.

Donedavno područje dubinske analize teksta na hrvatskom jeziku nije ni postojalo. Buđenje i nagli razvoj duguje text mining timu sa Fakulteta elektrotehnike i računarstva u Zagrebu, a ponajviše sjajnoj mentorici i voditeljici prof. dr. sc. Bojani Dalbelo-Bašić. Stalno pojavljivanje novih baza svjedoči o potrebama ovog područja, a njihova sve veća veličina tjera na daljnju optimizaciju, istraživanje i usavršavanje ovog područja. Kolegama koje interesira ovo područje ostavljam prvu stepenicu i još mnoga otvorena pitanja koja se mogu poboljšavati. Smatram da je zbog bogatstva hrvatskog jezika napredak na ovom području moguć s razvojem lingvistike, sintakse i semantike i njihovom sve boljom računalnom obradom.

8 LITERATURA

M. Malenica. *Primjena jezgrenih metoda u kategorizaciji teksta*, diplomski rad, FER, Zagreb, rujan 2004.

Marko Tadić. *Hrvatski nacionalni korpus na Internetu*, 2005.
[<http://www.hnk.ffzg.hr/>]

I.S.Dhillon, D.S. Modha. *Concept decomposition for large sparse text data using clustering*, Machine Learning, 42:1, 2001, pp. 143-175

Jing Gao, and Jun Zhang. *Text Retrieval Using Sparsified Concept Decomposition Matrix*, Technical Report 412-04, Department of Computer Science, University of Kentucky, Lexington, KY, 2004. [<http://cs.engr.uky.edu/~jzhang/pub/MINING/gao3.ps.gz>]

N. Cristianini, J. Shawe-Taylor. *An Introduction to Support Vector Machines (and other kernel based learning methods)*, Cambridge University Press, Cambridge, 2000.

T. Joachims. *Learning to classify text using support vector machines*, Kluwer, Massachusetts, 2001.

M.A. Hearst. *Trends and Controversies: Support Vector Machines*, IEEE Intelligent Systems, 13(4),, 1998 [www.computer.org/intelligent/ex1998/pdf/x4018.pdf]

G. Karypis, E.-H. Han. *Concept indexing: A fast dimensionality reduction algorithm with applications to document retrieval and categorization*, Technical Report TR-00-0016, University of Minnesota, 2000.
[<http://citeseer.ist.psu.edu/karypis00concept.html>]

Golub, G. H.,C. F. Van Loan. *Matrix computations*. Baltimore, MD, USA: The Johns Hopkins University Press, 1996.

Mandelbrot, B. B. *Fractal Geometry of Nature*. W. H. Freeman & Company, 1988.

E. Xing, A. Ng, M. Jordan, and S. Russell. *Distance metric learning, with application to clustering with side-information*. Advances in Neural Information Processing Systems 15, 2003. [<http://books.nips.cc/papers/files/nips15/AA03.pdf>]

Gang Wu, Navneet Panda, Edward Y. Chang. *Formulating Distance Functions via the Kernel Trick*, ACM International Conference on Knowledge Discovery and Data Mining (KDD), Chicago, August 2005 (future plan)

A. Bar-hillel, T. Hertz, N. Shental, and D. Weinshall. *Learning distance functions using equivalence relations*. In Proceedings of the Twentieth International Conference on Machine Learning, August 2003.[
<http://www.hpl.hp.com/conferences/icml2003/papers/88.pdf>]

C. C. Aggarwal. *Towards systematic design of distance functions for data mining applications*. The Ninth ACM SIGKDD International Conference on Knowledge Discovery in Data and Data Mining, 2003.[<http://citeseer.ist.psu.edu/717882.html>]

Nello Cristianini, John Shawe-Taylor, Andre Elisseeff, Jaz Kandola. *On Kernel-target alignment*, Neural Information Processing Systems 2001. [www.support-vector.net/papers/nips01_alignment.ps]

Daichi Mochihashi, Genichiro Kikui, Kenji Kita. *Learning Nonstructural Distance Metric by Minimum Cluster Distortions*, Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing (EMNLP-2004), pp.341-348, Barcelona, July 2004. [acl.ldc.upenn.edu/acl2004/emnlp/pdf/Mochihashi.pdf]

Rong Zhang and Alexander I. Rudnicki, *A Large Scale Clustering Scheme for Kernel K-Means*, icpr, vol. 04, no. 4, p. 40289, 16 2002.

Z. Zhang. *Learning metrics via discriminant kernels and multidimensional scaling : Towards expected Euclidean representation*. In Proceedings of the Twentieth International Conference on Machine Learning, August 2003.

J. T. Kwok and I. W. Tsang. *Learning with idealized kernels*. In Proceedings of the Twentieth International Conference on Machine Learning, pages 400–407, August 2003.

Alexei Vinokourov, John Shawe-Taylor, Nello Cristianini. *Inferring a Semantic Representation of Text via Cross-Language Correlation Analysis*. In Advances of Neural Information Processing Systems 15, 2002.[www.ecs.soton.ac.uk/~av/alexei_nips02.pdf]

Chih-Chung Chang and Chih-Jen Lin. *LIBSVM: a Library for Support Vector Machines*, 2005. [<http://www.csie.ntu.edu.tw/~cjlin/papers/libsvm.pdf>]

Blaž Fortuna. *Kernel Canonical Correlation Analysis With Applications*. SIKDD 2004 at multiconference IS 2004, 12-15 Oct 2004, Ljubljana, Slovenia 2004. [eprints.pascal-network.org/archive/00000736/01/BlazFortuna-Kcca.pdf]

J.Shawe-Taylor, N.Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.

C.-W. Hsu and C.-J. Lin. *A comparison of methods for multi-class support vector machines*. IEEE Transactions on Neural Networks, 13(2):415–425, 2002.
[\[http://citeseer.ist.psu.edu/hsu01comparison.html\]](http://citeseer.ist.psu.edu/hsu01comparison.html)

Lehoucq, R.B., D.C. Sorensen, and C. Yang, *ARPACK Users' Guide: Solution of Large-Scale Eigenvalue Problems with Implicitly Restarted Arnoldi Methods*, SIAM Publications, Philadelphia, 1998.
[\[http://www.caam.rice.edu/software/ARPACK/UG/ug.html\]](http://www.caam.rice.edu/software/ARPACK/UG/ug.html)

Steve Gunn. *Support Vector Machiens for Classification and Regression*. ISIS Technical Report ISIS-1-98, Image Speech & Intelligent Systems Research Group, University of Southampton, May. 1998

<http://www.kernel-methods.net/index.html> Dodatak knjige J.Shawe-Taylor, N.Cristianini. *Kernel Methods for Pattern Analysis*. Cambridge University Press, 2004.

<http://www.zemris.fer.hr/projects/textmining> Stranice text mining projekta na Fakultetu elektrotehnike i računarstva u Zagrebu, voditelj projekta: Prof. dr. sc. Bojana Dalbelo-Bašić

9 PRIMJER RADA METODE GLAVNIH KOMPONENTA I KANONSKE KORELACIJSKE ANALIZE

Smisao ovog poglavlja je približiti navedene metode čitatelju. Kako se ne bi udaljavali od podataka korištenih pri analizi odabrani su neki od naslova članaka iz kategorije politika i kategorije gospodarstvo baze čanaka Croatia Weekly. Slijede naslovi tih članaka zajedno sa kategorijom i prijevodom na engleski jezik:

Skup za učenje:

1 go	Svjetska financijska kriza mimoći će Hrvatsku	World financial crisis passing Croatia
2 go	Zagrebačka banka dobila veliki kredit ebrd-a	Zagreb bank secures loan of ebdr
3 go	Svjetska banka odobrila 101 dolara kredita hrvatske željeznice	World bank lends millions for croatia railroad company
4 go	Svjetska banka potpora hrvatskoj 200 dolara	World bank million of support for Croatia
5 go	Bankarska kriza mora biti riješena sljedeće godine	Financial crisis must be resolved next year
6 go	Hrvatsko informatičko tržište vrijedi 300 dolara	Croatian information technology market worth million usd
7 go	Kredit ebrd-a 168 dem za sanaciju zagrebačkog odlagališta otpada	Ebrd grants 168 million dm loan for Zagreb landfill rehabilitation
8	Najmanje 330 dem od prodaje pbz-a	330 million dm for state coffers from pbz

go	proračunu do kraja godine	sale
9 po	Sporazum otvaranju graničnih prijelaza mora se poštivati	Agreement opening border crossing observed
10 po	Bošnjaci odbili potpisati parafirani sporazum	Bosniaks refuse to sign agreement
11 po	Preko graničnih prijelaza u Istri ušlo dosta turista	Border crossing istria tourists
12 po	Uskrsnja najezda stranih turista na Jadran	Eastern invasion of foreign tourists
13 po	Bošnjaci pristaju na dogovor o humanom preseljenju	Bosniaks agree to humane relocation
14 po	Dan potpisivanja daytonskog sporazuma državni praznik	Day of signing dayton agreement is state holiday
15 po	Pardew zagovara reviziju daytonskog sporazuma	Pardew jolted revision of dayton agreement
16 po	Mjesec dana za dogovor o izbornom zakonu	Month day for agreement about election

Skup za testiranje:

1 go	Zagrebačka banka prodaje tvornicu duhana zagreb	Zagreb bank sale factory of tabaco Zagreb
2 go	Svjetska banka strukturne promjene preduvjet gospodarskog rasta	World bank structural prerequisite economic growth
3 po	U srpnju sporazum o imovinsko-pravnim odnosima	Agreement about property rights relations in july
4 po	Oess-ov sporazum sa srj spominje icty	Osce agreement with Yugoslavia fails to mention icty

Crvenom bojom označene su riječi koje se pojavljuju više od jedanput. Za izgradnju dokument-riječ matrice potrebno je parsirati i lematizirati riječi koje se nalaze u skupu za učenje. Odbacuju se stop riječi, brojevi i interpunkcijski znakovi. Prolaskom kroz sve primjere iz skupa za učenje generira se lista svih izraza koji se pojavljuju te im se dodjeljuje broj koji će ih dalje zamjenjivati kao indeks u dokument-riječ matrici. Ovakav postupak potreban je za dokumente na hrvatskom i posebno na engleskom jeziku. Broj riječi u engleskom i hrvatskom neće biti isti jer prijevod nije dostavan.

- | | |
|----------------|---------------------|
| 1. svjetska | 1. <i>world</i> |
| 2. financijska | 2. <i>financial</i> |
| 3. kriza | 3. <i>crisis</i> |
| 4. mimoći | 4. <i>passing</i> |
| 5. hrvatska | 5. <i>croatia</i> |
| 6. zagrebačka | 6. <i>zagreb</i> |
| 7. banka | 7. <i>bank</i> |
| 8. dobiti | 8. <i>secures</i> |
| 9. veliki | 9. <i>loan</i> |
| 10. kredit | 10. <i>ebdr</i> |
| 11. erdb-a | 11. ... |
| 12. ... | |

Na temelju rječnika tvore se vektori za svaki dokument. Dimenzionalnost vektora je broj riječi koje su se pojavile u skupu za učenje, a svaki element vektora predstavlja frekvenciju pojavljivanja riječi u razmatranom dokumentu. Kako se u naslovima na hrvatskom jeziku pojavljuje 65 riječi, 65 je indeks zadnje riječi("zakon") i dimenzionalnost svih vektora. Slaganjem dokument-vektora u matricu dobijemo dokument-riječ matricu.

	svjetska	financijsk	kriza	mimoći	hrvatska	Zagrebač		zakon
1	1	1	1	1	1	0		0
2	0	0	0	0	0	1		0
3	1	0	0	0	0	0		0
4	1	0	0	0	0	0		0
5	0	0	1	0	0	0		0
6	0	0	0	0	0	0		0
7	0	0	0	0	0	1		0
8	0	0	0	0	0	0		0
9	0	0	0	0	0	0		0
10	0	0	0	0	0	0		0
11	0	0	0	0	0	0		0
12	0	0	0	0	0	0		0

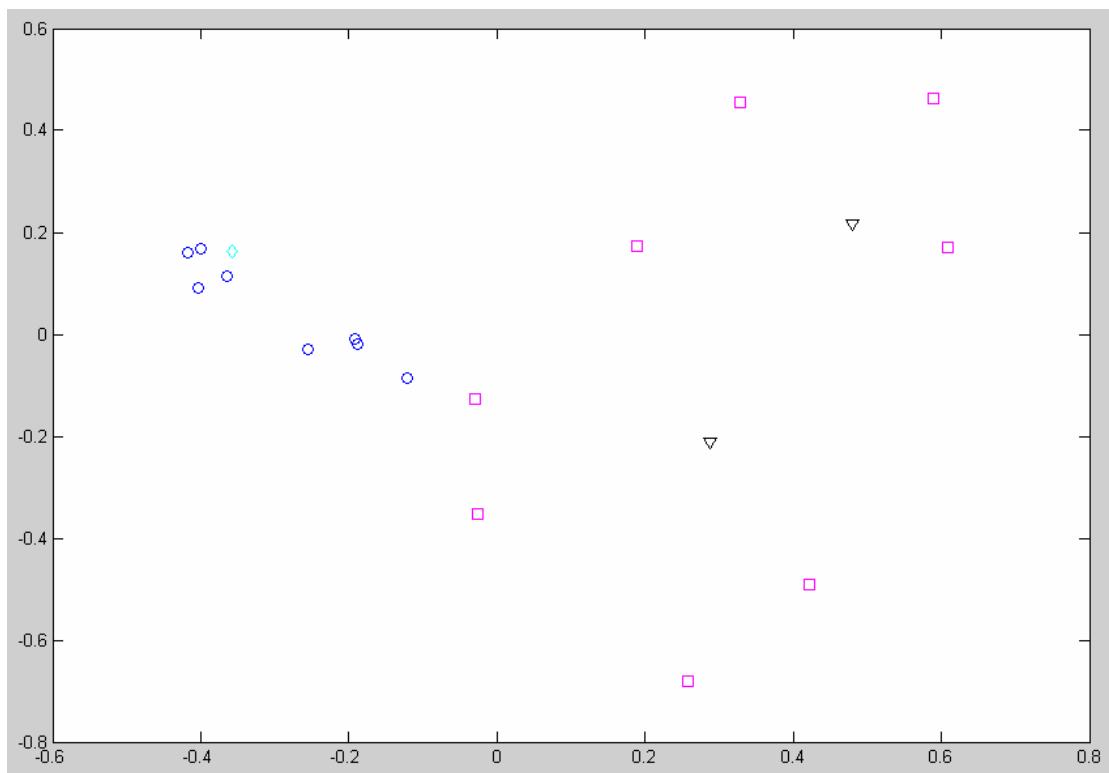
13	0	0	0	0	0	0		0
14	0	0	0	0	0	0		0
15	0	0	0	0	0	0		0
16	0	0	0	0	0	0		1

Dimenzionalnost matrice dokument-riječ za odabране naslove Croatia Weekly

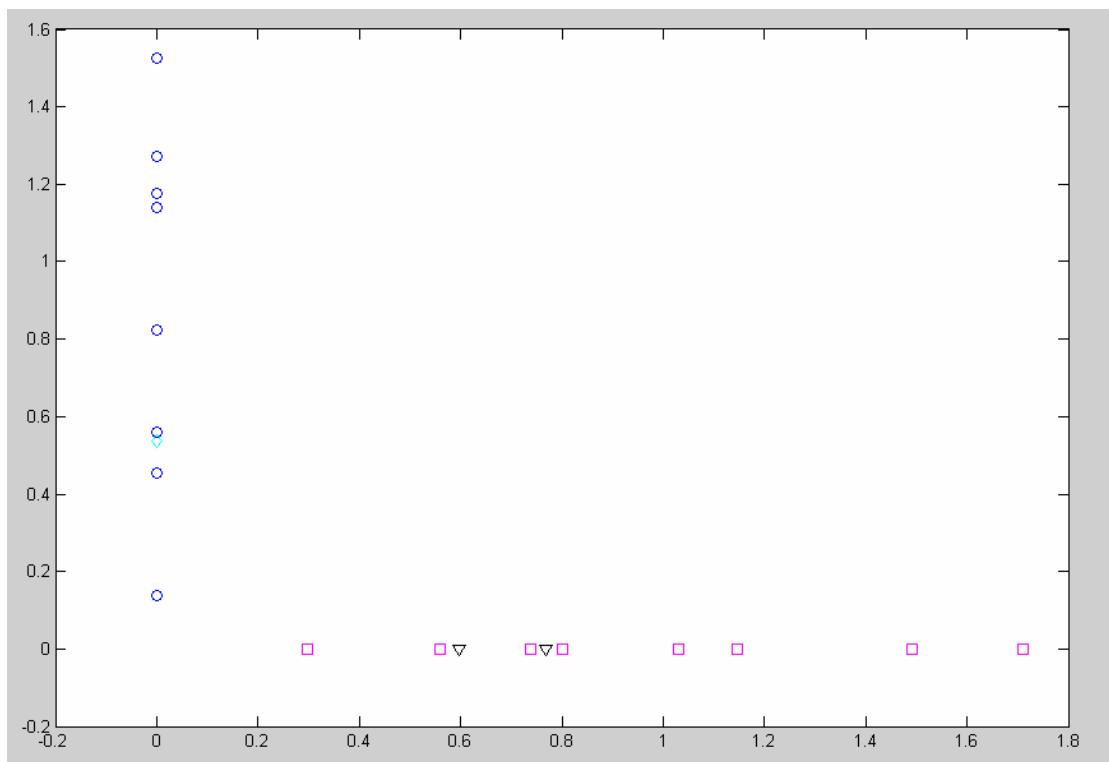
	Hrvatski	Engleski
Skup za učenje	16 x 65	16 x 58
Skup za treniranje	4 x 65	4 x 58

Ovakve matrice ulaz su u metodu glavnih komponenti i kanonsku korelacijsku analizu, metode za redukciju dimenzionalnosti. Radi boljeg prikaza i razumljivosti preslikavat ćemo dokument-vektore u dvodimenzionalni prostor. Potrebno je pronaći dva vektora u višedimenzionalnom prostoru riječi na koje će se dokument-vektori preslikati. Svaka metoda ima svoj način pronalaska takvih vektora. Vidjet ćemo koliko će vektori pronađeni metodom glavnih komponenata zadovoljiti naše potrebe. Kako ta metoda nema pristupa oznakama dokumenata(kategorijama) preostaje joj da potraži smjerove maksimalne varijance podataka i nada se da smjerovi male varijance predstavljaju samo šum. Kanonska korelacijska analiza traži smjerove maksimalne korelacije među prikazima dokumenta na dva jezika. Engleski jezik nam zapravo služi kao kompleksna oznaka kad već nemam pristupa stvarnim oznakama kategorije. Neke riječi kao npr. dogovor i sporazum u engleskom jeziku se preslikavaju u jednu riječ – agreement, te ne treba spominjati pozitivne promjene koje će to uvesti na ovako malom skupu i povezati dokumente koji nemaju nikakvih sličnosti na hrvatskom jeziku. Ako pažljivije pogledate naslove i njihove prijevode vidjet ćete da engleski nije doslovan prijevod hrvatskog i da se neke riječi dodatno pojavljuju u engleskom i povezuju inače odvojene dokumente.

Na slici su prikazani preslikani dokumenti u novom dvodimenzionalnom prostoru određenom vektorima dobivenim metodom glavnih komponenata i kanonskom korelacijskom analizom. Rozom i plavom prikazani su redom dokumenti kategorije gospodarstvo te kategorije politika iz skupa za učenje. Crno je skup za testiranje kategorije gospodarstvo, a svjetlo plavo skup za testiranje kategorije politika.



Slika 50. Dvodimenzionalni preslikani dokument-vektori dobiveni metodom glavnih komponenata



Slika 51. Dvodimenzionalni preslikani dokument-vektori dobiveni kanonskom korelacijskom analizom

Skup za učenje i treniranje redukcijom dimenzionalnosti kanonskom korelacijskom analizom postaje jasno i savršeno odvojiv. Za $x = 0$ možemo klasificirati dokument u klasu politika sa 100% točnošću, te također za $y = 0$ možemo tvrditi da je kategorija dokumenta gospodarstvo. Vektori pronađeni kanonskom korelacijskom analizom jasno su pridijeljeni svaki svojoj kategoriji, te su okomiti na sve dokument-vektore suprotne kategorije. Za metodu glavnih komponenata to nije slučaj. Dokument-vektori su linearna kombinacija vektora pronađenih ovom metodom. Samo sa jednim, prvim vektorom, možemo jasno odrediti klasifikaciju ako postavimo granicu za $x \geq 0$. Negativne vrijednosti komponenata vektora dobivene su zbog centriranja ulaznih podataka.

Ako se sjećamo definicije matrice bliskosti i jezgrene funkcije,

$$\tilde{k}(d_1, d_2) = \phi(d_1)U_k U_k' \phi(d_2)'$$

$$\tilde{k}(d_1, d_2) = \phi(d_1)WW' \phi(d_2)'$$

vidimo da su matrice WW' i $U_k U_k'$ rječ x rječ matrice koja određuju veze između različitih riječi. Osim što želimo povezati dokumente u kojima su se pojavile iste riječi, povezujemo i riječi koje se često pojavljuju u sličnim dokumentima što možemo vidjeti u matricama WW' za kanonsku korelacijsku analizu i $U_k U_k'$ za metodu glavnih komponenata. Matrica WW' savršeno odvaja riječi različitih kategorija malim vrijednostima i povezuje riječi iste kategorije. Dvije riječi koje se obje nalaze u istoj kategoriji(npr. svjetska-banka) za kanonsku korelacijsku analizu imaju vrijednosti 0.15557, 0.037886, 0.046873, 0.037886, 0.037886,... dok za usporedbu one koje se nalaze u različitim kategorijama(npr. sporazum-banka) imaju vrijednosti -1.6328e-017, 7.8959e-018, -1.1772e-018... U metodi glavnih komponenata razlika između kategorija nije tako jasno određena.

	Svjetska	financijska	kriza	mimoći
Svjetska	0.15557	0.037886	0.046873	0.037886
Financijska	0.037886	0.0092262	0.011415	0.0092262
Kriza	0.046873	0.011415	0.014122	0.011415
Mimoći	0.037886	0.0092262	0.011415	0.0092262
Hrvatska	0.037886	0.0092262	0.011415	0.0092262
Zagrebačka	0.063146	0.015378	0.019025	0.015378
Banka	0.14712	0.035827	0.044325	0.035827
Dobila	0.029429	0.0071666	0.0088666	0.0071666
Veliki	0.029429	0.0071666	0.0088666	0.0071666
Kredit	0.11953	0.029108	0.036012	0.029108
Ebrd-a	0.063146	0.015378	0.019025	0.015378
...
Mjesec	-1.6328e-017	7.8959e-018	-1.1772e-018	7.8959e-018
Izborni	-1.6328e-017	7.8959e-018	-1.1772e-018	7.8959e-018
Zakon	-1.6328e-017	7.8959e-018	-1.1772e-018	7.8959e-018

Matrica bliskosti za prvih par riječi (kanonska korelacijska analiza)

	Svjetska	financijska	kriza	mimoći
Svjetska	0.86791	0.14411	0.081873	0.14411
Financijska	0.14411	0.025203	0.012693	0.025203
Kriza	0.081873	0.012693	0.0083587	0.012693
Mimoći	0.14411	0.025203	0.012693	0.025203
Hrvatska	0.14411	0.025203	0.012693	0.025203
Zagrebačka	-0.21098	-0.066344	0.0021886	-0.06634
Banka	0.68845	0.098315	0.076226	0.098315
Dobila	-0.03536	-0.020587	0.0070463	-0.02058
Veliki	-0.03536	-0.020587	0.0070463	-0.02058
Kredit	0.095499	-0.019134	0.033694	-0.01913
Ebrd-a	-0.21098	-0.066344	0.0021886	-0.06634
...
Mjesec	-0.07797	-0.010901	-0.008798	-0.01090
Izborni	-0.07797	-0.010901	-0.008798	-0.01090
Zakon	-0.07797	-0.010901	-0.008798	-0.01090

Matrica bliskosti za prvih par riječi (metoda glavnih komponenata)

Različite boje odvajaju riječi po dokumentima, tako se prvi 5 riječi označenih crvenom bojom nalazi u prvom dokumentu, sljedećih 6 u drugom dokumentu itd. Zelenom bojom obilježene su posljednje 3 riječi koje pripadaju zadnjem dokumentu i bitno je napomenuti da pripadaju kategoriji politika za razliku od prva dva dokumenta koja govore o gospodarstvu. Metoda glavnih komponenata nije uspjela povezati sve riječi koje definiraju kategorije. Tako su samo skup riječi svjetska-banka jako povezani jer se pojavljuju u drugim dokumentima zajedno.

Za kanonsku korelačijsku analizu vrijedi: kako su spoj riječi „svjetska banka“ pojavljuje u mnogo dokumenta, a „zagrebačka banka“ je spomenuta u jednom logično je povezivanje svjetska – zagrebačka. Svjetska - dobila ima manju vrijednosti od svjetska - zagrebačka, svjetska - kredit ima veću jer se eksplicitno pojavljuje u jednom od dokumenata.

Kanonska korelačijska analiza je uspjela točno odvojiti riječi, tj. riječi jednog i drugog svojstvenog vektora se uopće ne miješaju i točno definiraju svoje kategorije jer su svi dokumenti jedne kategorije bili međusobno povezani, na hrvatskom ili ako ne na engleskom prijevodu. Metoda glavnih komponenata je s druge strane kao komponente vektora izmiješala riječi jedne i druge kategorije, tako da postoji veza svjetska - zakon jača od svjetska - veliki iako se svjetska i veliki pojavljuju u kategoriji gospodarstvo, a zakon u politici.

Metoda glavnih komponenata naravno ne gleda oznake pojedinih dokumenta nego traži osi najveće varijance te u ovom primjeru možemo zaključiti kako je kanonska korelačijska analiza kompleksnim oznakama(engleskim jezikom) uspjela odrediti bolji potprostor.