

# Razrješavanje višeznačnosti riječi

Toni Benussi, Darko Jurić, Krešimir Mišura, Ante Trbojević

Sveučilište u Zagrebu, Fakultet elektrotehnike i računarstva

Unska 3, 10000 Zagreb, Hrvatska

{toni.benussi,darko.juric,kresimir.misura,ante.trbojevic}@fer.hr,

## Sažetak

Ovaj rad razmatra uporabu metoda strojnog učenja u razrješavanju višeznačnosti riječi (engl. *word sense disambiguation*). Na podacima

## 1. Uvod

Riječ, kao sastavnica svakog jezika, opisana je leksičkim izrazom i semantičkim značenjem. Jedan leksički izraz može imati više semantičkih značenja, pri čemu okolni kontekst sugerira semantičko značenje leksičkog izraza. Primjerice leksički izraz *jaguar* ima više semantičkih značenja (životinja iz porodice mačaka, britanski proizvođač automobila, britanska glazbena grupa, parfem), gdje se semantičko značenje riječi *jaguar* doznaje iz konteksta. Dakle svaka riječ u pisanom tekstu jedinstveno je označena leksičkim izrazom i semantičkim značenjem. Upravo iz razloga što odnos leksičkog izraza riječi i semantičkog značenja riječi nije injektivan, automatizirani postupci kao što je npr. pretraživanje teksta ili postavljanje upita nad bazom podataka i sl. vrlo često korisnicima daju neželjene rezultate. Danas lingvistička znanost i računarska znanost nastoje pronaći tehnike koje bi dovoljno dobro rješavale navedeni problem.

Tehnike strojnog učenja pružaju dovoljno dobru podlogu za realizaciju sustava koji bi bio sposoban razriješiti višeznačnosti riječi, tj. identificirati pravo značenje riječi ovisno o okolnom kontekstu. Stoga su u ovom radu istraženi postupci strojnog učenja za razrješavanje višeznačnosti riječi.

Slijedeći problem koji se pojavljuje je kako izlučiti značajke iz konteksta. U većini slučajeva za značajke se uzima  $n$  riječi koje se nalaze u neposrednoj blizini oko višeznačne riječi, gdje je  $n \in N$  neka fiksna konstanta. Pogrešno odabran  $n$  može pogoršati uspješnost klasifikacije, također intuitivno se može naslutiti kako nije svejedno koliko od tih  $n$  riječi se pojavljuje u tekstu s lijeve strane višeznačne riječi, a koliko pak s desne strane. Upravo ovaj rad nastoji ustanoviti te zavisnosti, te zajedno s priloženim rezultatima provedenih eksperimenata iznosi zaključke kako odabrati kvalitetne značajke.

## 2. Drugi odjeljak

Ovo je drugi odjeljak.

### 2.1. Prvi pododjeljak

Ovo je pododjeljak drugog odjeljka.

### 2.2. Drugi pododjeljak

Ovo je drugi pododjeljak prvog odjeljka. Na (pod)odjeljke se u tekstu možete referencirati ovako: “u odjeljku 2.1. pokazali smo da ...”.

### 2.2.1. Primjer pod-pododjeljka

Ovo je pod-pododjeljak. Izbjegavajte pod-pododjeljke po svaku cijenu.

### 2.3. Primjer pododjeljka s dugačkim naslovom koji prelazi u drugi redak

Još jedan primjer pododjeljka.

## 3. Veličina rada

Rad treba imati najmanje tri, a najviše šest stranica, uključivo popis literature.

## 4. Opis uzoraka

Označavanje vlastite baze podataka s višeznačnim riječima izuzetno je skup proces. Naime za relevantnu klasifikaciju potrebno je barem oko tisuću uzoraka za jednu višeznačnu riječ, pri tom označivač za svaki uzorak mora pročitati okolni kontekst i shvatiti semantičko značenje riječi te pridodati uzorku oznaku razreda, što predstavlja dugotrajan proces. Stoga je pribavljena besplatna baza podataka<sup>1</sup> s već označenim semantičkim značenjima za engleske riječi *interest* i *line*. Shodno tome u ovom radu bit će prikazani postupci razrješavanja višeznačnosti engleskih riječi *interest* i *line*. Tekstovi u bazi podataka su na engleskom jeziku i prikupljeni su iz *ACL/DCI Wall Street Journal* novina.

Kao značajke uzorka uzimaju se riječi koje se nalaze unutar prozora konteksta promatranog uzorka. Prozor konteksta  $(l, r)$  (engl. *window of context*) sastoji se od  $l$  riječi koji se nalaze lijevo od višeznačne riječi (lijevi prozor konteksta) i  $r$  riječi koji se nalaze desno od višeznačne riječi (desni prozor konteksta), pri čemu se interpunkcijski i pravopisni znakovi ignoriraju. Sve riječi koje su uključene u prozor konteksta ostavljene su originalnoj leksičkoj formi, tj. nisu normalizirane, osim što su sva velika slova pretvorena u mala slova. Ova metoda reprezentiranja skupa značajki, opisana u (Pedersen, 2000), varijanta je metode „vreća riječi” (engl. *bag-of-words*) koja je prvi put opisana u (Gale et al., 1992), a razlikuje se po tome što razlikuje riječi koje se nalazi lijevo i desno od višeznačne riječi. Slika 1 opisuje izlučivanje značajki s definiranim prozorom konteksta  $(5, 3)$ .

Baza podataka sadrži 2368 instanci za riječ *interest* i 4148 instanci za riječ *line*, pri čemu svaka instanca sadrži nekoliko rečenica koje predstavljaju kontekst višeznačne riječi, iz kojih je zatim moguće izlučiti potrebne značajke.

<sup>1</sup><http://www.senseval.org/data.html>

Forward supplies can largely be stored in these same areas, and land forces are best held in reserve on our own soil. Drawing a line between military aid and military involvement may be harder, but it can be done if we keep the distinction clearly in mind.

Slika 1: Primjer uzorka višeznačne riječi *line* s prozorom konteksta (5, 3).

Također svaka instanca sadrži atribut koji definira semantičko značenje višeznačne riječi.

Nad bazom podataka napravljena je predobrada podataka tj. svi zapisi pretvoreni su u XML format u obliku kao što je prikazano na slici 2, gdje je `<tag key="division"/>` oznaka koja zamjenjuje višeznačnu riječ *line* i označuje semantičko značenje riječi u tom kontekstu.

```
<sentence> Forward supplies can largely
be stored in these same areas , and land
forces are best held in reserve on our own
soil. Drawing a <tag key="division"/>
between military aid and military
involvement may be harder , but it can be
done if we keep the distinction clearly in
mind. </sentence>
```

Slika 2: Primjer zapisa instance za višeznačnu riječ *line*

Riječ *interest* ima šest različitih semantičkih značenja (vidi tablicu 1), baš kao i riječ *line*, međutim za riječ *line* izdvojene su instance samo za tri semantička značenja (vidi tablicu 2). Naime distribucija semantičkih značenja u originalnoj bazi podataka za *interest* i *line* je vrlo neujednačen, stoga je stvoren umjetno uravnotežen skup instanci za *line* kako bi mogli istražiti koliko neke metode strojnog učenja kvalitetno klasificiraju na neuravnoteženom skupu podataka, a koliko na uravnoteženom skupu podataka.

Tablica 1: Distribucija semantičkih značenja za riječ *interest*

Semantičko značenje	Broj instanci
kamate	1252
udjel dionica u tvrtci	500
interes	361
prednost ili korist	178
pokazati zainteresiranost	66
prouzročiti zainteresiranost drugih	11

Kako je cilj ovog rada ustanoviti koji prozor konteksta ( $l, r$ ) odabrati ovisno o algoritmu strojnog učenja (razmatrani su k-NN, SVM, skup Bayeskovih klasifikatora) kako bi klasifikacija bila što uspješnija, stvoreno je 81 skupova uzoraka. Svaki skup uzoraka dobiven je drugačijim

Tablica 2: Distribucija semantičkih značenja za riječ *line*

Semantičko značenje	Broj instanci
tanak oblik, crta	373
umjetna podjela, granica	374
formacija ljudi ili stvari	349

izlučivanjem značajki iz skupa podataka, pri čemu se skupovi uzoraka razlikuju po veličini prozora konteksta ( $l, r$ ), odnosno po kombinaciji koliko je riječi izlučeno s lijeve strane višeznačne riječi (varijabla  $l$  u definiciji prozora konteksta), a koliko s desne strane višeznačne riječi (varijabla  $r$  u definiciji prozora konteksta). Veličina lijevog prozora konteksta i desnog prozora konteksta, tj. varijable  $l$  i  $r$  poprimaju vrijednosti iz skupa  $\{0, 1, 2, 3, 4, 5, 10, 25, 50\}$ . Ne postoji posebni razlog zašto su odabrane baš te vrijednosti, koje mogu poprimiti varijable  $l$  i  $r$ , već se slijedila preporuka iz rada (Pedersen, 2000), gdje su korištene iste vrijednosti. Nad svakim od tih skupova uzoraka provedeni su algoritmi: k-NN, SVM, naivni-Bayesov klasifikator, te su izračunate F1 mjere za svaku kombinaciju *skup uzoraka - algoritam*. Tablica 3 i tablica 4 prikazuju ukupan broj izlučenih značajki ovisno o kombinaciju ( $l, r$ ), tj. ovisno o skupu uzoraka. Primjećuje se kako skupovi uzoraka koji su dobiveni izlučivanjem s većim prozorom konteksta imaju više značajki, što je i očekivano. Naime veći prozor konteksta pohvatati će više različitih riječi pa će samim time i ukupan broj značajki biti veći.

## 5. Provođenje eksperimenata

Validacija i testiranje modela provedeno je unakrsnom provjerom. Poredak uzoraka u skupu uzoraka nasumično je ispremišan prije unakrsne provjere. Skup uzoraka zatim je podijeljen na pet podskupova, četiri podskupa služe za učenje modela, dok se posljednji podskup podijeli popola na još dva podskupa, od kojih jedan služi za validaciju, a drugi za testiranje. Prije podijele popola, podskup je nasumično ispremišan, s namjerom sprječavanja nepravilne distribucije uzoraka u skupu za validaciju ili skupu za testiranje. Nakon podijele petog podskupa na još dva skupa, vrši se validacija na dobivenom skupu za validaciju, dok se sa skupom za testiranje ne radi ništa. Nakon završene validacije slijedi slijedeća iteracija unakrsne provjere, tj. podskup za validaciju i testiranje ubacuje se u podskupove za učenje, a jedan od 4 prijašnja podskupa za učenje postaje skup za validaciju i testiranje. Takvih iterativnih koraka ima ukupno pet. Valja naglasiti kako se testiranje ne provodi odmah nakon validaciju u svakom iterativnom koraku unakrsne provjere, već se testiranja provodi nakon što je validacija provedena nad svakim od pet mogućih skupova za validaciju, nakon čega se izračunavaju optimalni parametri modela, te se tek nakon toga vrši testiranje.

### Program 1: Pseudokôd implementirane unakrsne provjere

```
1 unakrsnaProvjera(skup uzoraka, algoritam)
2 {
```

Tablica 3: Broj značajki za riječ *interest* ovisno o veličini prozora konteksta  $(l, r)$

$l \setminus r$	0	1	2	3	4	5	10	25	50
0	1	219	845	1489	1918	2269	3387	4357	4524
1	495	656	1165	1726	2118	2445	3509	4451	4617
2	1034	1151	1580	2071	2423	2719	3720	4624	4785
3	1515	1604	1980	2415	2716	2988	3915	4781	4940
4	1944	2017	2345	2733	2996	3245	4125	4953	5109
5	2301	2364	2662	3015	3261	3494	4332	5129	5282
10	3546	3587	3817	4075	4254	4442	5121	5813	5955
25	4906	4933	5111	5313	5455	5604	6168	6763	6885
50	5137	5162	5335	5530	5665	5810	6361	6938	7060

Tablica 4: Broj značajki za riječ *line* ovisno o veličini prozora konteksta  $(l, r)$

$l \setminus r$	0	1	2	3	4	5	10	25	50
0	1	224	663	1020	1299	1543	2433	3395	3647
1	235	422	825	1159	1427	1657	2531	3483	3731
2	554	691	1053	1363	1617	1836	2679	3609	3849
3	951	1058	1377	1657	1898	2109	2907	3803	4037
4	1326	1416	1703	1959	2186	2384	3143	4002	4231
5	1683	1766	2030	2267	2478	2659	3378	4198	4412
10	3185	3244	3447	3620	3783	3927	4494	5166	5346
25	5789	5828	5970	6094	6203	6296	6721	7205	7359
50	7046	7082	7205	7315	7408	7491	7854	8271	8406

```

3 lista_skupova = podijeli skup uzorak na
  pet jednakih dijelova
4 za svaki podskup i iz lista_skupova{
5   skup_za_ucenje = 4 podskupa iz
    lista_skupova među kojima nije i;
6   promijesaj i;
7   podijeli i na dva dijela
8   skup_za_validaciju = prva polovica
    podijeljenog skupa i
9
10  nauci_model(skup_za_ucenje, algoritam)
11  validacija(skup_za_validaciju, algoritam
    )
12 }
13
14 izracunaj parametre za algoritam
15
16 za svaki podskup i iz lista_skupova{
17   skup_za_ucenje = 4 podskupa iz
    lista_skupova među kojima nije i;
18   podijeli i na dva dijela
19   skup_za_testiranje = druga polovica
    podijeljenog skupa i
20
21   nauci_model(skup_za_ucenje, algoritam)
22   testiraj(skup_za_testiranje, algoritam)
23 }
24 }
```

Testiranje se provodi baš kao i validacija, tj. unakrsnom provjerom kroz pet koraka, ali je sada skup za validaciju

neaktivan, a skup za testiranje aktivan. Ovakav model učenja, validacije i testiranja odabran je s razlogom, naime želi se iskoristi što veći skup uzoraka i za validaciju i za testiranje. Naime ako bi se testiranje vršilo odmah nakon validacije u svakom iterativnom koraku unakrsne provjere, prilikom testiranja dobiveni rezultati ne bi bili potpuno objektivni, jer bi rezultati testiranja u tom koraku ovisili o trenutnom skupu za validaciju. TODO (srediti ovo obavezno)

## 6. Rezultati

### 6.1. Klasifikacija pomoću skupa naivnih Bayesovih klasifikatora

#### 6.1.1. Rezultati za riječ *interest*

Nakon provedene validacije odabire se devet najtočnijih klasifikatora iz svakog raspona kako bi se dobio ansambl Bayesovih klasifikatora. U tablici 5 prikazani su rezultati (F1 mjere) validacije, naivni Bayesovi klasifikatori koji su uključeni u ansambl su oni s prozorom konteksta: (1, 1), (3, 1), (10, 1), (1, 3), (4, 3), (10, 3), (2, 25), (5, 10), (10, 25).

Nakon provedenog testiranja F1 mjera ansambl naivnih Bayesovih klasifikatora iznosi

$$F1_{ansambl} = 0.7. \quad (1)$$

Valja primijetiti kako u ovom slučaju ansambl od devet naivnih klasifikatora ima manju točnost klasifikacije od svakog pojedinačnog naivnog Bayesovog klasifikatora od kojih je sastavljen, iz čega se može zaključiti da ansambl

Bayesovih klasifikatora ne mora nužno poboljšati točnost klasifikacije. Naime u ovom slučaju bilo bi bolje, umjesto ansambla, klasificirati npr. s naivnim Bayesovim klasifikatorom s definiranim prozorom konteksta (1, 1), kao što se može iščitati iz tablice 5 njegova F1 mjera iznosi

$$F1_{(1,1)} = 0.81. \quad (2)$$

Doduše treba uzeti u obzir da je mjera  $F1_{ansambl}$  izračunata na skupu za testiranje, dok je  $F1_{(1,1)}$  izračunata na skupu za validaciju, tj. te dvije mjere nisu baš usporedive jer nisu dobivene na temelju istih uzoraka.

Tablica 6: Matrica zabune ANBK-a za riječ *interest*

$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$	
71	0	0	0	28	108	$I_1$
1	0	0	0	0	6	$I_2$
1	0	0	0	4	37	$I_3$
3	0	0	0	7	92	$I_4$
0	0	0	0	213	58	$I_5$
0	0	0	0	0	556	$I_6$

Također ako se promotri tablica 6 lako se uočava kako niti jedan uzorak nije uspješno klasificiran u razrede  $I_2$ ,  $I_3$ ,  $I_4$ . Naime ti razredi imaju barem za red veličine manje uzoraka od najzastupljenijeg razreda  $I_6$  u skupu uzoraka (vidi tablicu 1). Ovakvu ne posve uspješnu, ali prihvatljivu klasifikaciju možemo opravdati neuravnoteženošću uzoraka između razreda u skupu uzoraka.

Ovakav postupak za riječ *interest* proveden je i u radu (Pedersen, 2000), pri čemu je korištena identična metoda i identični način validacije i testiranja. U tom radu točnost<sup>2</sup> klasifikacije iznosi 0.88. Razlog, zbog čega je naš klasifikator neuspješniji, u odnosu na spomenuti rad, nije nam potpuno poznat. Pretpostavljamo da je razlog u različitom pristupu zaglađivanja apriornih vjerojatnosti s vjerojatnošću nula, tj. u različitim implementacijama Bayesovog klasifikatora (u ovom radu korištena je biblioteka sustava Weka).

### 6.1.2. Rezultati za riječ *line*

Rezultati validacije mogu se vidjeti u tablici 7. AMBK je sastavljen od Bayesovih klasifikatora s prozorom konteksta: (2, 2), (5, 2), (10, 1), (1, 4), (5, 4), (10, 4), (2, 10), (5, 10) i (10, 10). F1 mjera AMBK-a u ovom slučaju iznosi

$$F1_{ansambl} = 0.91, \quad (3)$$

što se može smatrati izuzetno uspješnom klasifikacijom. Također valja primijetiti kako u ovom slučaju ANBK ima veću točnost od bilo kojeg pojedinačnog NB klasifikatora, što se može pravdati uravnoteženošću skupa uzoraka po klasama.

Tablica 8: Matrica zabune ANBK-a za riječ *line*

$L_1$	$L_2$	$L_3$	
201	8	3	$L_1$
11	175	2	$L_2$
22	6	122	$L_3$



Slika 3: Ovo je opis slike. Puna rečenica koja završava točkom. Opis ide ispod slike. Opis treba biti kratak; detalje objasnite u tekstu.

## 6.2. Klasifikacija pomoću SVM-a

### 6.2.1. Rezultati za riječ *interest*

### 6.2.2. Rezultati za riječ *line*

## 7. Anonimizacija

Prije slanja rada na recenziju, rad treba (privremeno) anonimizirati. To uključuje četiri stvari:

1. Sakrijte imena autora navedena ispod naslova rada. Nemojte brisati imena autora jer ćete time poremetiti raspored teksta; samo zakomentirajte liniju *name*.
2. Uklonite bilo kakve tragove iz teksta na temelju kojih bi recenzent mogao naslutiti tko su autori teksta.
3. Uklonite zahvalu, ako je imate.
4. Provjerite da u generiranom PDF-dokumentu nisu uključeni metapodatci iz kojih bi bilo vidljivo tko je generirao PDF.

## 8. Ilustracije i tablice

### 8.1. Ilustracije

Ovo je primjer uključivanja ilustracije. Ilustracije u  $\LaTeX$ -kôd uključite *nakon* teksta koji se na njih poziva. Pustite  $\LaTeX$  da ilustraciju smjesti tamo gdje misli da je najbolje (to je najčešće pri vrhu stranice i najčešće tamo gdje je vi nikad ne biste smjestili). Na sliku se referencirate ovako: “na slici 3 prikazano je ...”. Koristite tildu (~) kako biste spriječili razdvajanje između riječi “slika” i broja slike.

### 8.2. Tablice

Postoje dvije vrste tablice: uska tablica koja stane unutar jednog stupca i široka tablica koja prelazi preko oba stupca.

<sup>2</sup>Kod klasifikacije u više od dvije klase, točnost je jednaka F1 mjeri.

Tablica 5: F1 mjere Bayesovih klasifikatora dobivene validacijom za riječ *interest*

$l \setminus r$	0	1	2	3	4	5	10	25	50
0	0,61	0,75	0,71	0,73	0,7	0,7	0,69	0,71	0,69
1	0,76	0,81	0,8	0,78	0,76	0,75	0,72	0,73	0,71
2	0,72	0,79	0,78	0,77	0,76	0,75	0,74	0,74	0,72
3	0,74	0,78	0,77	0,76	0,76	0,77	0,75	0,73	0,72
4	0,74	0,77	0,76	0,78	0,76	0,77	0,75	0,73	0,73
5	0,73	0,77	0,77	0,78	0,77	0,76	0,76	0,74	0,74
10	0,68	0,73	0,73	0,73	0,72	0,73	0,74	0,75	0,74
25	0,68	0,69	0,69	0,69	0,7	0,7	0,72	0,72	0,72
50	0,66	0,68	0,68	0,69	0,69	0,68	0,7	0,72	0,72

Tablica 7: F1 mjere Bayesovih klasifikatora dobivene validacijom za riječ *line*

$l \setminus r$	0	1	2	3	4	5	10	25	50
0	0,32	0,55	0,63	0,67	0,66	0,63	0,66	0,66	0,65
1	0,67	0,68	0,73	0,76	0,8	0,79	0,77	0,73	0,72
2	0,72	0,73	0,79	0,78	0,79	0,79	0,82	0,77	0,76
3	0,77	0,74	0,78	0,75	0,78	0,8	0,83	0,79	0,81
4	0,78	0,77	0,81	0,8	0,79	0,8	0,83	0,83	0,85
5	0,8	0,81	0,83	0,83	0,82	0,83	0,86	0,84	0,85
10	0,81	0,86	0,83	0,86	0,87	0,87	0,9	0,9	0,9
25	0,76	0,77	0,81	0,81	0,82	0,84	0,85	0,83	0,83
50	0,72	0,76	0,76	0,76	0,79	0,81	0,82	0,82	0,81

Tablica 9: Ovo je opis tablice. Opis ide iznad tablice.

Zaglavlje1	Zaglavlje2
Jedan	Tekst u prvom retku
Dva	Tekstu u drugom retku
Tri	Tekst u trećem retku
	Tekst u četvrtom retku

### 8.2.1. Uske tablice

Primjer uske tablice je tablica 9. Nipošto nemojte koristiti okomite crte u tablici. Te crte nemaju smisla i loše izgledaju.

### 8.3. Široke tablice

Tablica 10 je primjer široke tablice koja ide preko oba stupca. Slično se mogu napraviti i slike koje idu preko oba stupca.

## 9. Matematičke formule

Matematičke formule koje se pojavljuju unutar rečenice pišite unutar tzv. *inline* matematičke okoline:  $2 + 3$ ,  $\sqrt{16}$ ,  $h(x) = \mathbf{1}(\theta_1 x_1 + \theta_0 > 0)$ . Veće formule pišite u tzv. *displayed* matematičkoj okolini:

$$b_k^{(i)} = \begin{cases} 1 & \text{ako } k = \operatorname{argmin}_j \|\mathbf{x}^{(i)} - \mu_j\| \\ 0 & \text{inače} \end{cases}$$

Matematičke izraze na koje se kasnije pozivate pišite unutar okoline *equation*:

$$J = \sum_{i=1}^N \sum_{k=1}^K b_k^{(i)} \|\mathbf{x}^{(i)} - \mu_k\|^2 \quad (4)$$

Sada se možete pozvati na (4). Ako se odlomak nastavlja nakon formule

$$f(x) = x^2 + \varepsilon \quad (5)$$

kao ovaj ovdje, onda obavezno na početku nastavka odlomka koristite naredbu *noindent* kako biste spriječili uvlačenje retka na početku odlomka.

Ako se matematička formula prostire kroz više redaka, koristite naredbu *align*.

$$\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$$

$$\operatorname{Var}(X + Y) = \operatorname{Var}(X) + \operatorname{Var}(Y)$$

Matematičke nazive koji se sastoje od više slova trebete pisati unutar narebe *mathit*, u suprotnom  $\LaTeX$  između slova umeće razmak kao da je riječ o množenju. Usporedite *Consistent(h, D)* i *Consistent(h, D)*.

Ako vam treba neki matematički simbol, znate kako izgleda, ali ne znate kako se zove odgovarajuća naredba u  $\LaTeX$ -u, iskušajte *Detexify*.<sup>3</sup> (Usput rečeno, riječ je o klasifikatoru k-NN pisanome u Haskellu.)

<sup>3</sup><http://detexify.kirelabs.org/>

Tablica 10: Opis široke tablice.

Zaglavlje1	Zaglavlje2	Zaglavlje 3
A	Neki vrlo dugačak tekst koji je širi od jednog stupca	128
B	Neki vrlo dugačak tekst koji je širi od jednog stupca	3123
C	Neki vrlo dugačak tekst koji je širi od jednog stupca	–32

## 10. Pseudokôd i programski kôd

Za pseudokôd je nabolje koristiti paket *Algorithm2e*. Alternativno, a posebice za programski kôd, možete koristiti pakete *listings* ili *fancyvrb*. Navođenje programskog kôda izbjegnite, osim ako za to nemate valjan razlog.

## 11. Navodi literature

Navodi literature pišu se u zagradi s prezimenom autora prvog autora te godinom izdanja (Chomsky, 1973). Više navoda se pišu jedan iza drugoga, u zajedničkoj zagradi i međusobno odvojeni točka-zarezom (Chomsky, 1973; Chave, 1964; Feigl, 1958). Navodi se najčešće pišu na kraju rečenice, a svakako prije interpunkcijskog znaka za kraj rečenice.

Ako rad ima više od dva autora, piše se ime samo prvog autora, nakon čega slijedi kratica *et al.*, koja znači *et alia* tj. i drugi (?). Ako su samo dva autora, pišu se prezimena oba autora (?).

Ako je ime autora ugrađeno u rečenicu, onda se ime autora piše izvan zagrada, a u zgrade ide samo godina izdanja. Npr. Chomsky (1973) je predložio da...". Razlika je dakle u tome referencirate li se na sâm rad ili na autora koji je napisao taj rad.

Popis literature dan je u abecednom poretку na kraju članka. Oblik navoda ovisi o vrsti bibliografske jedinice: konferencijski radovi (Chave, 1964), knjige (Butcher, 1981), članci u časopisu (?), doktorske disertacije (Croft, 1978) i poglavlja knjige (Feigl, 1958).

Sve ovo dobivate automatski ako koristite BibTeX. U datoteku `su2010.bbl` upišite navode literature, a zatim se na njih pozivajte putem njihovih simboličkih oznaka.

## 12. Tuđice

Umjesto tuđica (a to će biti uglavnom engleske riječi) trebale pisati hrvatski prijevod, a tuđice trebale navesti u kurzivu u zagradi. Na primjer: stablo odlučivanja (engl. *decision tree*). Nemojte pisati velika početna slova engleskog naziva, osim ako se ne radi o vlastitom imenu. Ako želite odmah uvesti i kraticu, dodajte i nju u zagradu. Na primjer: umjetna neuronska mreža (engl. *artificial neural network*, ANN). Iznimno, ako za neki engleski naziv nema odgovarajućeg prijevoda, možete koristiti engleski naziv, ali onda obavezno uvijek u kurzivu.

Nikada, baš nikada, engleske nazive ne koristite u naslovima odjeljaka.

## 13. Dodatne tipografske napomene

Slijede neke dodatne tipografske napomene.

### 13.1. Naglašavanje

Nipošto nemojte koristiti **masna slova**. Ako baš želite nešto naglasiti, koristite *kurziv*.

### 13.2. Zgrade

Ispred zgrade obavezno se piše bjelina (to nije poziv funkcije), a također i iza zgrade, osim, naravno, ako iza zgrade slijedi interpunkcijski znak.

### 13.3. Navodnici

Navodnici se pišu "ovako" (različiti su lijevi i desni navodnici), a ne "ovako" niti 'ovako' niti nešto treće. Doduše, navodnici se rijetko koriste u znanstvenom tekstu.

### 13.4. Trotočje

Trotočje (tzv. elipsa) "... " nema što tražiti u znanstvenom tekstu. Koristite kratice "itd." i "i sl.". U okviru matematičkog teksta, naravno, tri točke imaju drugo značenje (skraćeni prikaz niza) i koristite ih ako vam trebaju.

### 13.5. Razmaci nakon kratice

Nakon kratice koja završava točkom, a koja nije kraj rečenice, stavite tildu (~) umjesto razmaka. U suprotnom će L<sup>A</sup>T<sub>E</sub>X misliti da se radi o kraju rečenice i staviti će nešto veći razmak. Pogledajte razliku između

"tzv. elipsa" i

"tzv. elipsa"

(mala je, ali postoji). Tilda ujedno sprječava rastavljanje rečenice na tom mjestu.

### 13.6. Brojke

Brojke manje od 10 raspišite riječima, osim ako se radi vrijednostima parametara, rezultatima ili sl. Dakle, umjesto "u 3 eksperimenta", pišite "u tri eksperimenta".

### 13.7. Matematički simboli

Sve matematičke simbole, makar bili jednoslovni, pišite unutar matematičke okoline: *N* umjesto N, *k* umjesto k, itd. Nikada ne započinjite rečenicu matematičkim simbolom. Umjesto toga rečenicu započnite odgovarajućom imenicom. Npr. umjesto "*N* je broj koji..." napišite "Vrijednost *N* je broj koji..."

### 13.8. Crte i spojnice

Između dijelova polusloženica koristite pišete jednu crticu odnosno spojnicu: aritmetičko-logička operacija, EM-algoritam, VC-dimenzija, web-stranica. Između dijelova umetnute rečenice, piše se crta, koja se u L<sup>A</sup>T<sub>E</sub>X-u dobiva

pisanjem dviju crtica: “Prvi eksperiment je pokazao – premda ne posve uvjerljivo – da točnost klasifikatora ne ovisi samo o ...”.

### 13.9. Fusnote

Ovo je primjer fusnote.<sup>4</sup> Oznaku fusnote stavljajte nakon interpunkcijskog znaka, ako se odnosi na rečenicu ili na zadnju riječ, inače je stavite nakon riječi na koju se odnosi. Fusnote pokušajte izbjegavati, osim za poveznice (v. dolje).

### 13.10. Poveznice

Poveznice na web-stranice pišite neproporcionalnim fontom: `ktlab.fer.hr`. Svakako pokušajte sve poveznice smjestiti u fusnote;<sup>5</sup> tekst u kojem se pojavljuju poveznice izgleda loše zato jer su poveznice često predugačke.

### 13.11. Natuknice

Natuknice međusobno odvojite točka-zarezom:

1. Prva natuknica,
2. Druga natuknica,
3. Zadnju natuknicu završite točkom.

Ako neka od natuknica u sebi već sadrži zarez, onda sve natuknice završite točka zarezom, osim zadnje:

- Prva natuknica, koja u sebi ima zarez;
- Druga natuknica;
- Zadnja natuknica.

Natuknice koristite umjereno (zbog prostora).

### 13.12. Siročad

Po svaku cijenu nastojte izbjeći odlomke kod kojih u zadnjem retku visi samo jedna riječ, kao što je slučaj s ovim odlomkom. To lako možete riješiti preoblikovanjem nekih rečenica, no taj posao ostavite za kraj (može i nakon recenzije).

Isto vrijedi i za posljednje retke odlomaka ili (još gore) odsječaka koji usamljeno prelaze na iduću stranicu.

## 14. Zaključak

Zaključak je posljednji numerirani odsječak rada. Zaključak bi trebao biti veličine do najviše pola visine stupca. Zaključak je dobro razdijeliti u više odlomaka (npr. dva ili tri).

## Zahvale

Po želji možete prije popisa literature uključiti ovaj odjeljak i zahvaliti onima koji su vam na bilo koji način pomogli u izradi rada, a nisu autori rada. Ako nemate kome zahvaliti, obrišite ovaj odjeljak.

## Literatura

- Judith Butcher. 1981. *Copy-editing*. Cambridge University Press, 2nd edition.
- K. E. Chave. 1964. Skeletal durability and preservation. In J. Imbrie and N. Newel, editors, *Approaches to paleoecology*, pages 377–87, New York. Wiley.
- N. Chomsky. 1973. Conditions on transformations. In S. R. Anderson and P. Kiparsky, editors, *A festschrift for Morris Halle*, New York. Holt, Rinehart & Winston.
- W. B. Croft. 1978. *Organizing and searching large files of document descriptions*. Ph.D. thesis, Cambridge University.
- F. Feigl, 1958. *Spot Tests in Organic Analysis*, chapter 6. Publisher publisher, 5th edition.
- W. Gale, K. Church, and D. Yarowsky. 1992. A method for disambiguating word senses in a large corpus. *Computers and the Humanities*.
- T. Pedersen. 2000. A simple approach to building ensembles of naive bayesian classifiers for word sense disambiguation. *Proceeding Proceedings of the 1st North American chapter of the Association for Computational Linguistics conference*.

---

<sup>4</sup>Ovo je primjer teksta fusnote.

<sup>5</sup>`ktlab.fer.hr`