#### Vektorski model

Dragan Ivanović dragan.ivanovic@uns.ac.rs

Katedra za informatiku, Fakultet tehničkih nauka, Novi Sad

2015.

#### Vektorski model

- Težinski faktori vezani za pojedine termove u odnosu na dokumente i upite su pozitivne ali ne celobrojne vrednosti
- I upit ima težinske faktore
- Ima rangiranja
- Ima parcijalnog poklapanja upita i dokumenta
- I upit i dokument se predstavljaju kao n-dimenzionalni vektor (n je broj termova u rečniku)
- Ugao koji zaklapaju vektori je obrnuto srazmeran relevantnosti dokumenta za postavljeni upit

• Do sada svi upiti su bili Bulovi

- Do sada svi upiti su bili Bulovi
  - dokumenti ili odgovaraju upitu, ili ne odgovaraju nema između

- Do sada svi upiti su bili Bulovi
  - dokumenti ili odgovaraju upitu, ili ne odgovaraju nema između
- Dobro za korisnike-eksperte sa preciznim razumevanjem svojih potreba i sadržaja kolekcije

- Do sada svi upiti su bili Bulovi
  - dokumenti ili odgovaraju upitu, ili ne odgovaraju nema između
- Dobro za korisnike-eksperte sa preciznim razumevanjem svojih potreba i sadržaja kolekcije
- Dobro za aplikacije one lako mogu da obrade hiljade rezultata

- Do sada svi upiti su bili Bulovi
  - dokumenti ili odgovaraju upitu, ili ne odgovaraju nema između
- Dobro za korisnike-eksperte sa preciznim razumevanjem svojih potreba i sadržaja kolekcije
- Dobro za aplikacije one lako mogu da obrade hiljade rezultata
- Nije dobro za većinu korisnika

- Do sada svi upiti su bili Bulovi
  - dokumenti ili odgovaraju upitu, ili ne odgovaraju nema između
- Dobro za korisnike-eksperte sa preciznim razumevanjem svojih potreba i sadržaja kolekcije
- Dobro za aplikacije one lako mogu da obrade hiljade rezultata
- Nije dobro za većinu korisnika
- Većina korisnika nije u stanju da piše Bulove upite (ili jeste, ali ih mrzi)

- Do sada svi upiti su bili Bulovi
  - dokumenti ili odgovaraju upitu, ili ne odgovaraju nema između
- Dobro za korisnike-eksperte sa preciznim razumevanjem svojih potreba i sadržaja kolekcije
- Dobro za aplikacije one lako mogu da obrade hiljade rezultata
- Nije dobro za većinu korisnika
- Većina korisnika nije u stanju da piše Bulove upite (ili jeste, ali ih mrzi)
- Većina korisnika ne želi da pregleda hiljade pogodaka

- Do sada svi upiti su bili Bulovi
  - dokumenti ili odgovaraju upitu, ili ne odgovaraju nema između
- Dobro za korisnike-eksperte sa preciznim razumevanjem svojih potreba i sadržaja kolekcije
- Dobro za aplikacije one lako mogu da obrade hiljade rezultata
- Nije dobro za većinu korisnika
- Većina korisnika nije u stanju da piše Bulove upite (ili jeste, ali ih mrzi)
- Većina korisnika ne želi da pregleda hiljade pogodaka
- Ovo posebno važi za pretragu na webu

 Bulovi upiti često rezultuju sa malo (=0) ili previše (1000+) pogodaka

- Bulovi upiti često rezultuju sa malo (=0) ili previše (1000+) pogodaka
- Upit 1: "standard user dlink  $650" \rightarrow 200,000$  hits

- Bulovi upiti često rezultuju sa malo (=0) ili previše (1000+) pogodaka
- Upit 1: "standard user dlink  $650" \rightarrow 200,000$  hits
- Upit 2: "standard user dlink 650 no card found": 0 hits

- Bulovi upiti često rezultuju sa malo (=0) ili previše (1000+) pogodaka
- Upit 1: "standard user dlink  $650" \rightarrow 200,000$  hits
- Upit 2: "standard user dlink 650 no card found": 0 hits
- Potrebna je veština da se napiše upit koji će vratiti razuman broj pogodaka

- Bulovi upiti često rezultuju sa malo (=0) ili previše (1000+) pogodaka
- Upit 1: "standard user dlink  $650" \rightarrow 200,000$  hits
- Upit 2: "standard user dlink 650 no card found": 0 hits
- Potrebna je veština da se napiše upit koji će vratiti razuman broj pogodaka
- Sa rangiranim skupom pogodaka nije važno koliko je velik rezultat

 Želimo da nađemo dokumente koji su najkorisniji za korisnika (sortiran rezultat)

- Želimo da nađemo dokumente koji su najkorisniji za korisnika (sortiran rezultat)
- Kako možemo da rangiramo dokumente u odnosu na upit?

- Želimo da nađemo dokumente koji su najkorisniji za korisnika (sortiran rezultat)
- Kako možemo da rangiramo dokumente u odnosu na upit?
- Dodelimo ocenu (score) recimo iz [0,1] svakom dokumentu

- Želimo da nađemo dokumente koji su najkorisniji za korisnika (sortiran rezultat)
- Kako možemo da rangiramo dokumente u odnosu na upit?
- Dodelimo ocenu (score) recimo iz [0, 1] svakom dokumentu
- Ocena je mera koliko se dokument i upit "poklapaju" (match)

• Treba nam način za dodelu ocene svakom paru upit/dokument

- Treba nam način za dodelu ocene svakom paru upit/dokument
- Počnimo od upita sa jednim termom

- Treba nam način za dodelu ocene svakom paru upit/dokument
- Počnimo od upita sa jednim termom
- Ako se term ne pojavljuje u dokumentu, ocena bi trebalo da bude 0

- Treba nam način za dodelu ocene svakom paru upit/dokument
- Počnimo od upita sa jednim termom
- Ako se term ne pojavljuje u dokumentu, ocena bi trebalo da bude 0
- Što češće se term pojavljuje u dokumentu, ocena je veća

- Treba nam način za dodelu ocene svakom paru upit/dokument
- Počnimo od upita sa jednim termom
- Ako se term ne pojavljuje u dokumentu, ocena bi trebalo da bude 0
- Što češće se term pojavljuje u dokumentu, ocena je veća
- Razmotrićemo neke varijante kako ovo uraditi

• Uobičajena mera preklapanja dva skupa

- Uobičajena mera preklapanja dva skupa
- Neka su A i B skupovi (bar jedan je neprazan)

- Uobičajena mera preklapanja dva skupa
- Neka su A i B skupovi (bar jedan je neprazan)
- Jaccard-ov koeficijent:

$$\mathsf{jaccard}(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

- Uobičajena mera preklapanja dva skupa
- Neka su A i B skupovi (bar jedan je neprazan)
- Jaccard-ov koeficijent:

$$\mathsf{jaccard}(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

• jaccard(A, A) = 1

- Uobičajena mera preklapanja dva skupa
- Neka su A i B skupovi (bar jedan je neprazan)
- Jaccard-ov koeficijent:

$$\mathsf{jaccard}(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

- jaccard(A, A) = 1
- jaccard(A, B) = 0 if  $A \cap B = 0$

- Uobičajena mera preklapanja dva skupa
- Neka su A i B skupovi (bar jedan je neprazan)
- Jaccard-ov koeficijent:

$$\mathsf{jaccard}(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

- jaccard(A, A) = 1
- jaccard(A, B) = 0 if  $A \cap B = 0$
- A i B ne moraju imati isti broj elemenata.

- Uobičajena mera preklapanja dva skupa
- Neka su A i B skupovi (bar jedan je neprazan)
- Jaccard-ov koeficijent:

$$\mathsf{jaccard}(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

- jaccard(A, A) = 1
- jaccard(A, B) = 0 if  $A \cap B = 0$
- A i B ne moraju imati isti broj elemenata.
- Uvek se dodeljuje broj između 0 i 1.

#### Jaccard-ov koeficijent: primer

- Koja je ocena upit/dokument dobijena pomoću Jaccard-ovog koeficijenta za:
  - Upit: "ides of March"
  - Dokument: "Caesar died in March"

# Šta nije dobro kod Jaccard-ovog koeficijenta?

- Ne uzima u obzir frekvenciju terma (koliko puta se term pojavljuje)
- Retki termovi su informativniji od čestih; Jaccard ovo ne uzima u obzir
- Treba nam bolji način za normalizaciju dužine
- Kasnije ćemo koristiti  $|A \cap B|/\sqrt{|A \cup B|}$  (cosine) . . .
- ullet ....umesto  $|A\cap B|/|A\cup B|$  (Jaccard) za normalizaciju dužine

## Podsećanje: binarna matrica incidencije

	Ivanović	Milosavljević	-		
	D.	В.	S.	Μ.	
digitalan	1	1	1	1	
lucene	1	1	1	1	
dokument	1	1	1	1	
obrazovanje	0	0	1	0	
pretraga	1	1	1	1	
multimedijalan	0	1	1	1	
evaluacija	0	0	0	0	

. .

Svaki dokument je prikazan pomoću binarnog vektora  $\in \{0,1\}^{|V|}.$ 

### Podsećanje: binarna matrica incidencije

	Ivanović D.	Milosavljević B.	Gostojić S.	Zarić M.	
digitalan	1	1	1	1	
lucene	1	1	1	1	
dokument	1	1	1	1	
obrazovanje	0	0	1	0	
pretraga	1	1	1	1	
$\operatorname{multimedijalan}$	0	1	1	1	
evaluacija	0	0	0	0	

. .

Svaki dokument je prikazan pomoću binarnog vektora  $\in \{0,1\}^{|V|}$ .

#### Od sada ćemo koristiti brojačku matricu

	Ivanović	Milosavljević	Gostojić	Zarić	
	D.	B.	S.	Μ.	
digitalan	8	46	7	5	
lucene	11	68	3	2	
dokument	41	953	105	204	
obrazovanje	0	0	1	0	
pretraga	11	56	10	30	
multimedijalan	0	96	1	7	
evaluacija	0	0	0	0	

. .

Svaki dokument je prikazan pomoću vektora broja pojavljivanja  $\in \mathbb{N}^{|V|}$ .

### Od sada ćemo koristiti brojačku matricu

	Ivanović	Milosavljević	Gostojić	Zarić	
	D.	B.	S.	Μ.	
digitalan	8	46	7	5	
lucene	11	68	3	2	
dokument	41	953	105	204	
obrazovanje	0	0	1	0	
pretraga	11	56	10	30	
multimedijalan	0	96	1	7	
evaluacija	0	0	0	0	

. .

Svaki dokument je prikazan pomoću vektora broja pojavljivanja  $\in \mathbb{N}^{|V|}$ .

# Model "vreće sa re<u>čima"</u>

Ne uzimamo u obzir redosled reči u dokumentu

- Ne uzimamo u obzir redosled reči u dokumentu
- John is quicker than Mary i Mary is quicker than John su prikazani na isti način

- Ne uzimamo u obzir redosled reči u dokumentu
- John is quicker than Mary i Mary is quicker than John su prikazani na isti način
- Ovo se zove model "vreće sa rečima" (bag of words)

- Ne uzimamo u obzir redosled reči u dokumentu
- John is quicker than Mary i Mary is quicker than John su prikazani na isti način
- Ovo se zove model "vreće sa rečima" (bag of words)
- Korak nazad: pozicioni indeks razlikuje ova dva dokumenta

- Ne uzimamo u obzir redosled reči u dokumentu
- John is quicker than Mary i Mary is quicker than John su prikazani na isti način
- Ovo se zove model "vreće sa rečima" (bag of words)
- Korak nazad: pozicioni indeks razlikuje ova dva dokumenta
- Čuvanje informacije o poziciji ćemo ostaviti za drugi put

- Ne uzimamo u obzir redosled reči u dokumentu
- John is quicker than Mary i Mary is quicker than John su prikazani na isti način
- Ovo se zove model "vreće sa rečima" (bag of words)
- Korak nazad: pozicioni indeks razlikuje ova dva dokumenta
- Čuvanje informacije o poziciji ćemo ostaviti za drugi put
- Za sada koristimo model "vreće sa rečima"

• Frekvencija terma tf $_{t,d}$  terma t u dokumentu d definiše se kao broj pojavljivanja t u d.

- Frekvencija terma t $_{t,d}$  terma t u dokumentu d definiše se kao broj pojavljivanja t u d.
- Hoćemo da koristimo tf kada računamo upit/dokument ocene.
   Kako?

- Frekvencija terma t $f_{t,d}$  terma t u dokumentu d definiše se kao broj pojavljivanja t u d.
- Hoćemo da koristimo tf kada računamo upit/dokument ocene.
   Kako?
- Sirova frekvencija terma nije ono što nam treba

- Frekvencija terma t $f_{t,d}$  terma t u dokumentu d definiše se kao broj pojavljivanja t u d.
- Hoćemo da koristimo tf kada računamo upit/dokument ocene.
   Kako?
- Sirova frekvencija terma nije ono što nam treba
- Dokument sa 10 pojava jednog terma je relevantniji od dokumenta sa jednom pojavom istog terma

- Frekvencija terma t $f_{t,d}$  terma t u dokumentu d definiše se kao broj pojavljivanja t u d.
- Hoćemo da koristimo tf kada računamo upit/dokument ocene.
   Kako?
- Sirova frekvencija terma nije ono što nam treba
- Dokument sa 10 pojava jednog terma je relevantniji od dokumenta sa jednom pojavom istog terma
- Ali nije 10 puta relevantniji

- Frekvencija terma t $f_{t,d}$  terma t u dokumentu d definiše se kao broj pojavljivanja t u d.
- Hoćemo da koristimo tf kada računamo upit/dokument ocene.
   Kako?
- Sirova frekvencija terma nije ono što nam treba
- Dokument sa 10 pojava jednog terma je relevantniji od dokumenta sa jednom pojavom istog terma
- Ali nije 10 puta relevantniji
- Relevantnost ne raste proporcionalno sa frekvencijom terma

• Logaritmska težina frekvencije terma t u d definiše se kao

$$\mathbf{w}_{t,d} = \left\{ egin{array}{ll} 1 + \log_{10} \mathrm{tf}_{t,d} & \mathrm{if} \ \mathrm{tf}_{t,d} > 0 \\ 0 & \mathrm{ina\check{c}e} \end{array} 
ight.$$

• Logaritmska težina frekvencije terma t u d definiše se kao

$$\mathbf{w}_{t,d} = \left\{ egin{array}{ll} 1 + \log_{10} \mathrm{tf}_{t,d} & \mathrm{if} \ \mathrm{tf}_{t,d} > 0 \\ 0 & \mathrm{ina\check{c}e} \end{array} 
ight.$$

ullet 0 o 0, 1 o 1, 2 o 1.3, 10 o 2, 1000 o 4, itd.

• Logaritmska težina frekvencije terma t u d definiše se kao

$$\mathbf{w}_{t,d} = \left\{ egin{array}{ll} 1 + \log_{10} \mathrm{tf}_{t,d} & \mathrm{if} \ \mathrm{tf}_{t,d} > 0 \\ 0 & \mathrm{ina\check{c}e} \end{array} 
ight.$$

- ullet 0  $\rightarrow$  0, 1  $\rightarrow$  1, 2  $\rightarrow$  1.3, 10  $\rightarrow$  2, 1000  $\rightarrow$  4, itd.
- Ocena za par upit/dokument: suma po termovima t za q i d: ocena  $=\sum_{t\in q\cap d}(1+\log \operatorname{tf}_{t,d})$

• Logaritmska težina frekvencije terma t u d definiše se kao

$$\mathbf{w}_{t,d} = \left\{ egin{array}{ll} 1 + \log_{10} \mathrm{tf}_{t,d} & \mathrm{if} \ \mathrm{tf}_{t,d} > 0 \\ 0 & \mathrm{ina\check{c}e} \end{array} 
ight.$$

- $0 \to 0$ ,  $1 \to 1$ ,  $2 \to 1.3$ ,  $10 \to 2$ ,  $1000 \to 4$ , itd.
- Ocena za par upit/dokument: suma po termovima t za q i d: ocena =  $\sum_{t \in q \cap d} (1 + \log \mathsf{tf}_{t,d})$
- Ocena je 0 akko nijedan term iz upita nije prisutan u dokumentu

• Retki termovi su infomativniji od čestih

# <u>Frekv</u>encija dokumenta

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)
  - Dokument koji sadrži ovaj term je verovatno relevantan

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)
  - Dokument koji sadrži ovaj term je verovatno relevantan
  - $\rightarrow$  Želimo veliku težinu za retke termove kao što je digitalizacija

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)
  - Dokument koji sadrži ovaj term je verovatno relevantan
  - → Želimo veliku težinu za retke termove kao što je digitalizacija
- Razmotrimo term koji je čest u kolekciji (npr. visoko, teško, analiza)

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)
  - Dokument koji sadrži ovaj term je verovatno relevantan
  - → Želimo veliku težinu za retke termove kao što je digitalizacija
- Razmotrimo term koji je čest u kolekciji (npr. visoko, teško, analiza)
  - Dokument koji sadrži ovaj term je verovatno relevantniji od onog koji ga ne sadrži, ali to nije siguran indikator relevatnosti

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)
  - Dokument koji sadrži ovaj term je verovatno relevantan
  - ullet ightarrow Želimo veliku težinu za retke termove kao što je  $\operatorname{digitalizacija}$
- Razmotrimo term koji je čest u kolekciji (npr. visoko, teško, analiza)
  - Dokument koji sadrži ovaj term je verovatno relevantniji od onog koji ga ne sadrži, ali to nije siguran indikator relevatnosti
  - Za česte termove želimo pozitivne težine za reči kao što su visoko, teško i analiza, ali manje težine nego za retke termove

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)
  - Dokument koji sadrži ovaj term je verovatno relevantan
  - → Želimo veliku težinu za retke termove kao što je digitalizacija
- Razmotrimo term koji je čest u kolekciji (npr. visoko, teško, analiza)
  - Dokument koji sadrži ovaj term je verovatno relevantniji od onog koji ga ne sadrži, ali to nije siguran indikator relevatnosti
  - Za česte termove želimo pozitivne težine za reči kao što su visoko, teško i analiza, ali manje težine nego za retke termove
- Koristićemo frekvenciju dokumenta da uzmemo to u obzir prilikom računanja ocene

- Retki termovi su infomativniji od čestih
- Razmotrimo term koji je redak u kolekciji (npr. digitalizacija)
  - Dokument koji sadrži ovaj term je verovatno relevantan
  - ightarrow Želimo veliku težinu za retke termove kao što je digitalizacija
- Razmotrimo term koji je čest u kolekciji (npr. visoko, teško, analiza)
  - Dokument koji sadrži ovaj term je verovatno relevantniji od onog koji ga ne sadrži, ali to nije siguran indikator relevatnosti
  - ullet ightarrow Za česte termove želimo pozitivne težine za reči kao što su visoko, teško i analiza, ali manje težine nego za retke termove
- Koristićemo frekvenciju dokumenta da uzmemo to u obzir prilikom računanja ocene
- Frekvencija dokumenta je broj dokumenata u kolekciji u kojima se pojavljuje dati term

•  $df_t$  je frekvencija dokumenta, odn. broj dokumenata u kojima se pojavljuje term t

- ullet df $_t$  je frekvencija dokumenta, odn. broj dokumenata u kojima se pojavljuje term t
- df je inverzna mera informativnosti terma

- ullet df $_t$  je frekvencija dokumenta, odn. broj dokumenata u kojima se pojavljuje term t
- df je inverzna mera informativnosti terma
- Definišemo idf težinu terma t kao:

$$\mathsf{idf}_t = \mathsf{log}_{10} \, \frac{\mathsf{N}}{\mathsf{df}_t}$$

- ullet df $_t$  je frekvencija dokumenta, odn. broj dokumenata u kojima se pojavljuje term t
- df je inverzna mera informativnosti terma
- Definišemo idf težinu terma t kao:

$$\mathsf{idf}_t = \mathsf{log}_{10} \, \frac{\mathsf{N}}{\mathsf{df}_t}$$

• idf je mera informativnosti terma

- ullet df $_t$  je frekvencija dokumenta, odn. broj dokumenata u kojima se pojavljuje term t
- df je inverzna mera informativnosti terma
- Definišemo idf težinu terma t kao:

$$\mathsf{idf}_t = \mathsf{log}_{10} \, \frac{\mathsf{N}}{\mathsf{df}_t}$$

- idf je mera informativnosti terma
- Koristićemo log  $N/df_t$  umesto  $N/df_t$  da "ublažimo" efekat idf-a

- ullet df $_t$  je frekvencija dokumenta, odn. broj dokumenata u kojima se pojavljuje term t
- df je inverzna mera informativnosti terma
- Definišemo idf težinu terma t kao:

$$\mathsf{idf}_t = \mathsf{log}_{10} \, \frac{\mathsf{N}}{\mathsf{df}_t}$$

- idf je mera informativnosti terma
- Koristićemo log  $N/df_t$  umesto  $N/df_t$  da "ublažimo" efekat idf-a
- Koristimo logaritmovanje i za frekvenciju terma i za frekvenciju dokumenta

#### Primeri za idf

Izračunati idf $_t$  koristeći formulu: idf $_t = \log_{10} \frac{1,000,000}{\mathrm{df}_t}$ 

term	df <sub>t</sub>	idf <sub>t</sub>
XMIRS	1	6
digitalizacija	100	4
nedelja	1000	3
analiza	10.000	2
ispod	100.000	1
i	1.000.000	0

# Uticaj idf-a na rangiranje

• idf utiče na rangiranje samo ako upit ima bar dva terma

# Uticaj idf-a na rangiranje

- idf utiče na rangiranje samo ako upit ima bar dva terma
- Na primer, u upitu "digitalizacija dokumenata", idf težina povećava relativnu težinu za digitalizacija i smanjuje relativnu težinu za dokumenata

# Uticaj idf-a na rangiranje

- idf utiče na rangiranje samo ako upit ima bar dva terma
- Na primer, u upitu "digitalizacija dokumenata", idf težina povećava relativnu težinu za digitalizacija i smanjuje relativnu težinu za dokumenata
- idf nema uticaja na rangiranje rezultata upita sa jednim termom

Reč	cf	df
osiguranje	10440	3997
pokušati	10422	8760

• Frekvencija kolekcije terma t je broj pojavljivanja t u kolekciji

Reč	cf	df
osiguranje	10440	3997
pokušati	10422	8760

- Frekvencija kolekcije terma t je broj pojavljivanja t u kolekciji
- Koja reč je bolji term za upit (i trebalo bi da dobije veću težinu?

Reč	cf	df
osiguranje	10440	3997
pokušati	10422	8760

- Frekvencija kolekcije terma t je broj pojavljivanja t u kolekciji
- Koja reč je bolji term za upit (i trebalo bi da dobije veću težinu?
- Ovaj primer sugeriše da je df bolji za težine nego cf:

Reč	cf	df
osiguranje	10440	3997
pokušati	10422	8760

- Frekvencija kolekcije terma t je broj pojavljivanja t u kolekciji
- Koja reč je bolji term za upit (i trebalo bi da dobije veću težinu?
- Ovaj primer sugeriše da je df bolji za težine nego cf:
- Želimo da manji broj dokumenata koji sadrži osiguranje dobije veći značaj u odnosu na gomilu dokumenata koji sadrže pokušati pri upitu koji sadrži ova dva terma

# Frekvencije terma, kolekcije i dokumenta

Veličina	Simbol	Definicija
frekv. terma	$tf_{t,d}$	broj pojavljivanja t u d
frekv. dokumenta	$df_t$	broj dokumenata u kolekciji koji sadrže <i>t</i>
frekv. kolekcije	cf <sub>t</sub>	ukupan broj pojavljivanja <i>t</i> u kolekciji

# Frekvencije terma, kolekcije i dokumenta

Veličina	Simbol	Definicija
frekv. terma	$tf_{t,d}$	broj pojavljivanja t u d
frekv. dokumenta	$df_t$	broj dokumenata u kolekciji koji sadrže <i>t</i>
frekv. kolekcije	cf <sub>t</sub>	ukupan broj pojavljivanja <i>t</i> u kolekciji

• Veza između df i cf?

# Frekvencije terma, kolekcije i dokumenta

Veličina	Simbol	Definicija
frekv. terma	$tf_{t,d}$	broj pojavljivanja t u d
frekv. dokumenta	$df_t$	broj dokumenata u kolekciji koji sadrže <i>t</i>
frekv. kolekcije	cf <sub>t</sub>	ukupan broj pojavljivanja <i>t</i> u kolekciji

Veza između tf i cf?

• tf-idf težina terma je proizvod njegove tf težine i njegove idf težine

 tf-idf težina terma je proizvod njegove tf težine i njegove idf težine

•

$$w_{t,d} = (1 + \log \mathsf{tf}_{t,d}) \cdot \log \frac{N}{\mathsf{df}_t}$$

 tf-idf težina terma je proizvod njegove tf težine i njegove idf težine

•

$$w_{t,d} = (1 + \log \mathsf{tf}_{t,d}) \cdot \log rac{\mathsf{N}}{\mathsf{df}_t}$$

Najpoznatija težina u IR

 tf-idf težina terma je proizvod njegove tf težine i njegove idf težine

•

$$w_{t,d} = (1 + \log \mathsf{tf}_{t,d}) \cdot \log rac{\mathsf{N}}{\mathsf{df}_t}$$

- Najpoznatija težina u IR
- Poznata i kao: tf.idf, tf x idf

• Dodeli tf-idf težinu svakom termu t za svaki dokument d:  $w_{t,d} = (1 + \log \operatorname{tf}_{t,d}) \cdot \log \frac{N}{\operatorname{df}_t}$ 

- Dodeli tf-idf težinu svakom termu t za svaki dokument d:  $w_{t,d} = (1 + \log \operatorname{tf}_{t,d}) \cdot \log \frac{N}{\operatorname{df}_t}$
- N je ukupan broj dokumenata

- Dodeli tf-idf težinu svakom termu t za svaki dokument d:  $w_{t,d} = (1 + \log \operatorname{tf}_{t,d}) \cdot \log \frac{N}{\operatorname{df}_{\star}}$
- N je ukupan broj dokumenata
- Raste sa brojem pojavljivanja u dokumentu

- Dodeli tf-idf težinu svakom termu t za svaki dokument d:  $w_{t,d} = (1 + \log \operatorname{tf}_{t,d}) \cdot \log \frac{N}{\operatorname{df}_{\star}}$
- N je ukupan broj dokumenata
- Raste sa brojem pojavljivanja u dokumentu
- Raste sa retkošću terma u kolekciji

## Binarna ightarrow brojačka ightarrow težinska matrica

	Ivanović	Milosavljević	Gostojić	Zarić	
	D.	B.	S.	Μ.	
digitalan	0.25	1.51	0.12	0.05	
lucene	0.75	3.2	0.28	0.18	
dokument	0.12	2.3	0.51	0.83	
obrazovanje	0	0	0.5	0	
pretraga	0.25	1.8	0.24	0.26	
multimedijalan	0	8.25	0.12	0.23	
evaluacija	0	0	0	0	

. .

Svaki dokument je predstavljen vektorom realnih vrednosti tf-idf težina  $\in \mathbb{R}^{|V|}$ 

# Binarna ightarrow brojačka ightarrow težinska matrica

	Ivanović	Milosavljević	Gostojić	Zarić	
	D.	B.	S.	М.	
digitalan	0.25	1.51	0.12	0.05	
lucene	0.75	3.2	0.28	0.18	
dokument	0.12	2.3	0.51	0.83	
obrazovanje	0	0	0.5	0	
pretraga	0.25	1.8	0.24	0.26	
multimedijalan	0	8.25	0.12	0.23	
evaluacija	0	0	0	0	

. . .

Svaki dokument je predstavljen vektorom realnih vrednosti tf-idf težina  $\in \mathbb{R}^{|V|}$ 

ullet Svaki dokument se reprezentuje kao vektor realnih brojeva sa tf-idf težinama  $\in \mathbb{R}^{|V|}$ 

- ullet Svaki dokument se reprezentuje kao vektor realnih brojeva sa tf-idf težinama  $\in \mathbb{R}^{|V|}$
- ullet Tako imamo |V|-dimenzionalni vektorski prostor

- Svaki dokument se reprezentuje kao vektor realnih brojeva sa tf-idf težinama  $\in \mathbb{R}^{|V|}$
- $\bullet$  Tako imamo |V|-dimenzionalni vektorski prostor
- Termovi su ose prostora

- Svaki dokument se reprezentuje kao vektor realnih brojeva sa tf-idf težinama  $\in \mathbb{R}^{|V|}$
- Tako imamo |V|-dimenzionalni vektorski prostor
- Termovi su ose prostora
- Dokumenti su tačke ili vektori u ovom prostoru

- Svaki dokument se reprezentuje kao vektor realnih brojeva sa tf-idf težinama  $\in \mathbb{R}^{|V|}$
- $\bullet$  Tako imamo |V|-dimenzionalni vektorski prostor
- Termovi su ose prostora
- Dokumenti su tačke ili vektori u ovom prostoru
- Visoka dimenzionalnost: desetak miliona dimenzija kada se primeni na web pretraživač

- Svaki dokument se reprezentuje kao vektor realnih brojeva sa tf-idf težinama  $\in \mathbb{R}^{|V|}$
- $\bullet$  Tako imamo |V|-dimenzionalni vektorski prostor
- Termovi su ose prostora
- Dokumenti su tačke ili vektori u ovom prostoru
- Visoka dimenzionalnost: desetak miliona dimenzija kada se primeni na web pretraživač
- Vrlo retki vektori većina vrednosti je 0

• Ideja 1: upite, kao i dokumente, predstaviti kao vektore u vektorskom prostoru

- Ideja 1: upite, kao i dokumente, predstaviti kao vektore u vektorskom prostoru
- Ideja 2: rangirati dokumente prema njihovoj blizini sa upitom

- Ideja 1: upite, kao i dokumente, predstaviti kao vektore u vektorskom prostoru
- Ideja 2: rangirati dokumente prema njihovoj blizini sa upitom
- blizina = sličnost

- Ideja 1: upite, kao i dokumente, predstaviti kao vektore u vektorskom prostoru
- Ideja 2: rangirati dokumente prema njihovoj blizini sa upitom
- blizina = sličnost
- blizina ≈ negativno rastojanje

- Ideja 1: upite, kao i dokumente, predstaviti kao vektore u vektorskom prostoru
- Ideja 2: rangirati dokumente prema njihovoj blizini sa upitom
- blizina = sličnost
- blizina ≈ negativno rastojanje
- Podsećanje: ovo radimo da prevaziđemo problem "ili jesi ili nisi" Bulovog modela

• Prvi pokušaj: rastojanje između dve tačke

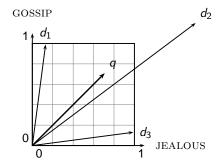
- Prvi pokušaj: rastojanje između dve tačke
- (= rastojanje između krajnjih tačaka dvaju vektora)

- Prvi pokušaj: rastojanje između dve tačke
- (= rastojanje između krajnjih tačaka dvaju vektora)
- Euklidsko rastojanje?

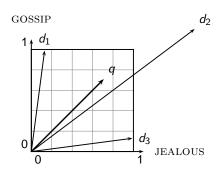
- Prvi pokušaj: rastojanje između dve tačke
- (= rastojanje između krajnjih tačaka dvaju vektora)
- Euklidsko rastojanje?
- Euklidsko rastojanje je loša ideja . . .

- Prvi pokušaj: rastojanje između dve tačke
- (= rastojanje između krajnjih tačaka dvaju vektora)
- Euklidsko rastojanje?
- Euklidsko rastojanje je loša ideja . . .
- ...jer je veliko za vektore različitih dužina.

# Zašto je rastojanje loša ideja



# Zašto je rastojanje loša ideja



Euklidsko  $\vec{q}$  i  $\vec{d_2}$  je veliko iako je distribucija termova u upitu q i dokumentu  $d_2$  vrlo slična

# Ugao umesto rastojanja

• Rangiraćemo dokumente prema uglu koji zaklapaju sa upitom

- Rangiraćemo dokumente prema uglu koji zaklapaju sa upitom
- Eksperiment: uzmimo dokument d i dodajmo ga još jednom na njegov kraj; nazovimo to d'

- Rangiraćemo dokumente prema uglu koji zaklapaju sa upitom
- Eksperiment: uzmimo dokument d i dodajmo ga još jednom na njegov kraj; nazovimo to d'
- "Semantički" d i d' imaju isti sadržaj

- Rangiraćemo dokumente prema uglu koji zaklapaju sa upitom
- Eksperiment: uzmimo dokument d i dodajmo ga još jednom na njegov kraj; nazovimo to d'
- "Semantički" d i d' imaju isti sadržaj
- Ugao između dokumenata je 0, što odgovara maksimalnoj sličnosti

- Rangiraćemo dokumente prema uglu koji zaklapaju sa upitom
- Eksperiment: uzmimo dokument d i dodajmo ga još jednom na njegov kraj; nazovimo to d'
- "Semantički" d i d' imaju isti sadržaj
- Ugao između dokumenata je 0, što odgovara maksimalnoj sličnosti
- Euklidsko rastojanje između d i d' je veliko (u svakom slučaju > 0)

Sledeće dve stvari su ekvivalentne

- Sledeće dve stvari su ekvivalentne
  - Rangiraj dokumente prema uglu između upita i dokumenta u rastućem redosledu

- Sledeće dve stvari su ekvivalentne
  - Rangiraj dokumente prema uglu između upita i dokumenta u rastućem redosledu
  - Rangiraj dokumente prema cos(query,document) u opadajućem redosledu

- Sledeće dve stvari su ekvivalentne
  - Rangiraj dokumente prema uglu između upita i dokumenta u rastućem redosledu
  - Rangiraj dokumente prema cos(query,document) u opadajućem redosledu
- Kosinus je monotono opadajuća funkcija ugla u intervalu [0°, 180°]

• Kako da izračunamo kosinus?

- Kako da izračunamo kosinus?
- Vektor se može normalizovati deljenjem svake komponente njegovom dužinom – ovde koristimo L<sub>2</sub> normu:

$$||x||_2 = \sqrt{\sum_i x_i^2}$$

- Kako da izračunamo kosinus?
- Vektor se može normalizovati deljenjem svake komponente njegovom dužinom ovde koristimo  $L_2$  normu:  $||x||_2 = \sqrt{\sum_i x_i^2}$
- Ovo premešta sve vektore u jediničnu sferu . . .

- Kako da izračunamo kosinus?
- Vektor se može normalizovati deljenjem svake komponente njegovom dužinom ovde koristimo  $L_2$  normu:  $||x||_2 = \sqrt{\sum_i x_i^2}$
- Ovo premešta sve vektore u jediničnu sferu . . .
- ...jer je nakon normalizacije:  $||x||_2 = \sqrt{\sum_i x_i^2} = 1.0$

- Kako da izračunamo kosinus?
- Vektor se može normalizovati deljenjem svake komponente njegovom dužinom ovde koristimo  $L_2$  normu:  $||x||_2 = \sqrt{\sum_i x_i^2}$
- Ovo premešta sve vektore u jediničnu sferu . . .
- ullet ...jer je nakon normalizacije:  $||x||_2 = \sqrt{\sum_i x_i^2} = 1.0$
- Kao rezultat, i kratki i dugački dokumenti imaju težine istog reda veličine

- Kako da izračunamo kosinus?
- Vektor se može normalizovati deljenjem svake komponente njegovom dužinom ovde koristimo  $L_2$  normu:  $||x||_2 = \sqrt{\sum_i x_i^2}$
- Ovo premešta sve vektore u jediničnu sferu . . .
- ullet ...jer je nakon normalizacije:  $||x||_2 = \sqrt{\sum_i x_i^2} = 1.0$
- Kao rezultat, i kratki i dugački dokumenti imaju težine istog reda veličine
- Efekat na dva dokumenta d i d' (d dodat na samog sebe) sa prethodnog primera: imaju identične vektore nakon normalizacije

$$\cos(\vec{q}, \vec{d}) = \sin(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}||\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

 $\bullet$   $q_i$  je tf-idf težina terma i u upitu

$$\cos(\vec{q}, \vec{d}) = \sin(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}||\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

- q<sub>i</sub> je tf-idf težina terma i u upitu
- d<sub>i</sub> je tf-idf težina terma i u dokumentu

$$\cos(\vec{q}, \vec{d}) = \sin(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}||\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

- q<sub>i</sub> je tf-idf težina terma i u upitu
- d<sub>i</sub> je tf-idf težina terma i u dokumentu
- $|\vec{q}|$  i  $|\vec{d}|$  su dužine  $\vec{q}$  i  $\vec{d}$ .

$$\cos(\vec{q}, \vec{d}) = \sin(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{|\vec{q}||\vec{d}|} = \frac{\sum_{i=1}^{|V|} q_i d_i}{\sqrt{\sum_{i=1}^{|V|} q_i^2} \sqrt{\sum_{i=1}^{|V|} d_i^2}}$$

- q<sub>i</sub> je tf-idf težina terma i u upitu
- d<sub>i</sub> je tf-idf težina terma i u dokumentu
- $|\vec{q}|$  i  $|\vec{d}|$  su dužine  $\vec{q}$  i  $\vec{d}$ .
- Ovo je kosinusna sličnost  $\vec{q}$  i  $\vec{d}$  ......ili, ekvivalentno, kosinus ugla između  $\vec{q}$  i  $\vec{d}$ .

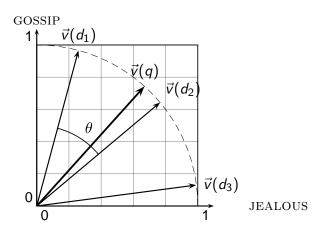
#### Kosinus za normalizovane vektore

Za normalizovane vektore kosinus je jednak skalarnom proizvodu

#### Kosinus za normalizovane vektore

- Za normalizovane vektore kosinus je jednak skalarnom proizvodu
- $\cos(\vec{q}, \vec{d}) = \vec{q} \cdot \vec{d} = \sum_i q_i \cdot d_i$  (ako su  $\vec{q}$  i  $\vec{d}$  normalizovani)

## Ilustracija kosinusne sličnosti



Pretpostavimo da imamo tri disertacije u kolekciji:

ID: Ivanović D.

MB: Milosavljević B.

GS: Gostojić S.

Pretpostavimo da imamo tri disertacije u kolekciji:

ID: Ivanović D.

MB: Milosavljević B.

GS: Gostojić S.

frekv. terma (broj)

term	ID	MB	GS
indeksiranje	5	58	0
dokument	41	953	105
obrazovanje	0	0	1
multimedijalan	0	96	1

frekv. terma (broj)

term	ID	MB	GS
indeksiranje	5	58	0
dokument	41	953	105
obrazovanje	0	0	1
$\operatorname{multimedijalan}$	0	96	1

frekv. terma (broj)

log frekv.

term	ID	MB	GS
indeksiranje	5	58	0
dokument	41	953	105
obrazovanje	0	0	1
multimedijalan	0	96	1

<b>.</b>	ID	MD	CC
term	ID	MB	GS
indeksiranje	1,7	2,76	0
dokument	2,61	3,98	3,02
obrazovanje	0	0	1
multimedijalan	0	2,98	1

frekv. dok

term	idf
indeksiranje	0,176
$\operatorname{dokument}$	0
obrazovanje	0,477
$\operatorname{multimedijalan}$	0,176

frekv. dok

term	idf
indeksiranje	0,176
$\operatorname{dokument}$	0
obrazovanje	0,477
$\operatorname{multimedijalan}$	0,176

 $tf_idf$ 

term	ID	MB	GS
indeksiranje	0,29	0,49	0
dokument	0	0	0
obrazovanje	0	0	0,48
multimedijalan	0	0,52	0,18

 $tf_idf$ 

term	ID	MB	GS
indeksiranje	0,29	0,49	0
dokument	0	0	0
obrazovanje	0	0	0,48
multimedijalan	0	0,52	0,18

 $tf_idf$ 

term	ID	MB	GS
indeksiranje	0,29	0,49	0
$\operatorname{dokument}$	0	0	0
obrazovanje	0	0	0,48
$\operatorname{multimedijalan}$	0	0,52	0,18

*tf<sub>i</sub>df* & normalizacija

term	ID	MB	GS
indeksiranje	1	0,69	0
dokument	0	0	0
obrazovanje	0	0	0,94
multimedijalan	0	0,73	0,35

tfidf

term	ID	MB	GS
indeksiranje	0,29	0,49	0
$\operatorname{dokument}$	0	0	0
obrazovanje	0	0	0,48
$\operatorname{multimedijalan}$	0	0,52	0,18

*tf<sub>i</sub>df* & normalizacija

			J	
	term	ID	MB	GS
ĺ	indeksiranje	1	0,69	0
	dokument	0	0	0
	obrazovanje	0	0	0,94
	multimedijalan	0	0,73	0,35

• Razmotrimo kako bismo odgovorili na upit: ... multimedijalnih ... indeksiranje ... multimedijalnog , pri čemu ... predstavljaju delove upita koji će nakon pretprocesiranja upita biti izbačeni.

*tf<sub>i</sub>df* & normalizacija

term	ID	MB	GS	upit
indeksiranje	1	0,69	0	0,62
dokument	0	0	0	0
obrazovanje	0	0	0,94	0
multimedijalan	0	0,73	0,35	0,79

*tf<sub>i</sub>df* & normalizacija

term	ID	MB	GS	upit
indeksiranje	1	0,69	0	0,62
dokument	0	0	0	0
obrazovanje	0	0	0,94	0
multimedijalan	0	0,73	0,35	0,79

•  $cos(ID,upit) = 1 * 0,62 + 0 * 0 + 0 * 0 + 0 * 0,79 \approx 0,62$ 

*tf<sub>i</sub>df* & normalizacija

		,		
term	ID	MB	GS	upit
indeksiranje	1	0,69	0	0,62
dokument	0	0	0	0
obrazovanje	0	0	0,94	0
multimedijalan	0	0,73	0,35	0,79

- $cos(ID,upit) = 1 * 0.62 + 0 * 0 + 0 * 0 + 0 * 0.79 \approx 0.62$
- $cos(MB,upit) = 0,69 * 0,62 + 0 * 0 + 0 * 0 + 0,73 * 0,79 \approx 1$

*tf<sub>i</sub>df* & normalizacija

		-		
term	ID	MB	GS	upit
indeksiranje	1	0,69	0	0,62
dokument	0	0	0	0
obrazovanje	0	0	0,94	0
multimedijalan	0	0,73	0,35	0,79

- $cos(ID,upit) = 1 * 0.62 + 0 * 0 + 0 * 0 + 0 * 0.79 \approx 0.62$
- $cos(MB,upit) = 0,69*0,62+0*0+0*0+0,73*0,79 \approx 1$
- $cos(GS,upit) = 0 * 0,62 + 0 * 0 + 0,94 * 0 + 0,35 * 0,79 \approx 0,28$

*tf¡df* & normalizacija

		,		
term	ID	MB	GS	upit
indeksiranje	1	0,69	0	0,62
$\parallel$ dokument	0	0	0	0
obrazovanje	0	0	0,94	0
$\parallel$ multimedijalan	0	0,73	0,35	0,79

- $cos(ID,upit) = 1 * 0.62 + 0 * 0 + 0 * 0 + 0 * 0.79 \approx 0.62$
- $cos(MB, upit) = 0.69 * 0.62 + 0 * 0 + 0 * 0 + 0.73 * 0.79 \approx 1$
- $cos(GS,upit) = 0 * 0,62 + 0 * 0 + 0,94 * 0 + 0,35 * 0,79 \approx 0,28$
- Koja disertacija najbolje odgovara upitu?

### Izračunavanje kosinusne ocene

```
CosineScore(q)
     float Scores[N] = 0
     float Length[N]
    for each query term t
     do calculate w_{t,q} and fetch postings list for t
 5
         for each pair (d, \mathsf{tf}_{t,d}) in postings list
         do Scores[d] += w_{t,d} \times w_{t,a}
 6
     Read the array Length
    for each d
     do Scores[d] = Scores[d]/Length[d]
     return Top K components of Scores[]
10
```

## Komponente tf-idf težina

Frekv. terma		Frekv. dokumenta		Normalizacija	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{max_t(tf_{t,d})}$	p (prob idf)	$\max\{0,\log \frac{N-\mathrm{df}_t}{\mathrm{df}_t}\}$	u (pivoted unique)	1/u
b (boolean)	$\begin{cases} 1 & \text{if } \operatorname{tf}_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			b (byte size)	$1/\mathit{CharLength}^{lpha}, \ lpha < 1$
L (log ave)	$\frac{1 + \log(\operatorname{tf}_{t,d})}{1 + \log(\operatorname{ave}_{t \in d}(\operatorname{tf}_{t,d}))}$				

# Komponente tf-idf težina

Frekv. terma		Frekv. dokumenta		Normalizacija	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{max_t(tf_{t,d})}$	p (prob idf)	$\max\{0,\log \tfrac{N-\mathrm{df}_t}{\mathrm{df}_t}\}$	u (pivoted unique)	1/u
b (boolean)	$egin{cases} 1 &  ext{if } \operatorname{tf}_{t,d} > 0 \ 0 &  ext{oth erwise} \end{cases}$			b (byte size)	$1/\mathit{CharLength}^{lpha}$ , $lpha < 1$
L (log ave)	$\frac{1 + \log(\operatorname{tf}_{t,d})}{1 + \log(\operatorname{ave}_{t \in d}(\operatorname{tf}_{t,d}))}$				

Najbolja poznata kombinacija komponenti težine

## Komponente tf-idf težina

Frekv. terma		Frekv. dokumenta		Normaliz acija	
n (natural)	tf <sub>t,d</sub>	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \ldots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{max_t(tf_{t,d})}$	p (prob idf)	$\max\{0,\log \frac{N-\mathrm{df}_t}{\mathrm{df}_t}\}$	u (pivoted unique)	1/u
b (boolean)	$egin{cases} 1 &  ext{if } \operatorname{tf}_{t,d} > 0 \ 0 &  ext{otherwise} \end{cases}$			b (byte size)	$1/\mathit{CharLength}^{lpha}, \ lpha < 1$
L (log ave)	$\frac{1 + \log(\operatorname{tf}_{t,d})}{1 + \log(\operatorname{ave}_{t \in d}(\operatorname{tf}_{t,d}))}$				

Default: bez težina

• Mogu da se koriste različite težine za upite i dokumente

- Mogu da se koriste različite težine za upite i dokumente
- Notacija: qqq.ddd

- Mogu da se koriste različite težine za upite i dokumente
- Notacija: qqq.ddd
- Primer: Itn Inc

- Mogu da se koriste različite težine za upite i dokumente
- Notacija: qqq.ddd
- Primer: Itn.Inc
- upit: logaritamski tf, idf, bez normalizacije

- Mogu da se koriste različite težine za upite i dokumente
- Notacija: qqq.ddd
- Primer: Itn.Inc
- upit: logaritamski tf, idf, bez normalizacije
- dokument: logaritamski tf, bez idf, kosinusna normalizacija

Predstavi svaki dokument kao tf-idf vektor

- Predstavi svaki dokument kao tf-idf vektor
- Predstavi upit kao ponderisani tf-idf vektor

- Predstavi svaki dokument kao tf-idf vektor
- Predstavi upit kao ponderisani tf-idf vektor
- Izračunaj kosinusnu sličnost između upita i svakog dokumenta

- Predstavi svaki dokument kao tf-idf vektor
- Predstavi upit kao ponderisani tf-idf vektor
- Izračunaj kosinusnu sličnost između upita i svakog dokumenta
- Rangiraj dokumente prema sličnosti

- Predstavi svaki dokument kao tf-idf vektor
- Predstavi upit kao ponderisani tf-idf vektor
- Izračunaj kosinusnu sličnost između upita i svakog dokumenta
- Rangiraj dokumente prema sličnosti
- ullet Prikaži najboljih K (npr. K=10) dokumenata korisniku