

Luka Dragar, Matej Vatovec, Aleksander Kovač

Inverted Index

Domača naloga pri predmetu Iskanje in ekstrakcija podatkov s spleta

MENTOR: prof. dr. Marko Bajec, doc. dr. Slavko Žitnik

Pri domači nalogi smo obdelali dano besedilo v večih dokumentih in dobljene podatke poindexirali. Najprej smo besedilo izluščili ter preobdelali, tako da smo ga tokenizirali, kasneje še normalizirali (izločili ločila, številke in izvedli lematizacijo ter besede shranili shranili v indeksno strukturo, s čimer smo poskrbeli za učinkovito strukturo iskanja in poizvedovanja po besedilu znotraj HMTL datotek.

• Indeksiranje podatkov

Program se sprehodi čez vsako HTML datoteko posebaj in naredi sledeče:

- izlušči samo besedilo (brez HTML značk...)
- tokenizira izluščeno besedilo
- kopira dobljeno tabelo tokenov in jo normalizira:
 - celotno besedilo spremeni v male črke
 - sprehodi se čez vsak token:
 - če je token ločilo ga zamenja z PUNCT
 - če je token številka (ali vsebuje številko) ga zamenja z NUM
 - če je token stopword ga zamenja z STOPWORD

tako dobimo dve tabeli. Primer:

```
['Galerija', 'MGML', ',', 'ki', 'se', 'nahaja', 'v' 'Ljubljani', ',', 'se', 'zapira', '.']  
['galerija', 'mgml', 'PUNCT', 'STOPWORD', 'STOPWORD', 'nahajati', 'STOPWORD', 'ljubljana',  
'PUNCT', 'STOPWORD', 'zampirati', 'PUNCT']
```

Prva tabela je tokenizirano besedilo, druga tabela pa normalizirano besedilo. Obe tabeli sta iste dožine zato da imajo vse besede iste indekse, kar olajša iskanje originalnega besedila pri poizvedbi. Pri indeksiranju se uporablja le druga tabela, torej normalizirano besedilo

Program se sprehodi čez vse tokene v normalizirani tabeli (druga tabela) in si za vsako besedo shrani število pojavitev in indekse posamezne pojavitve. Pri tem ignorira tokene PUNCT, STOPWORD in NUM.

Dobljene besede in njene pojavitve skupaj z imenom trenutne datoteke spravi v podatkovno bazo.

- **Iskanje podatkov brez obrnjenega indeksa**

Programu preko ukazne vrstice podamo parametre. Ti se tokenizirajo in normalizirajo po istem postopku, kot je opisan v prvem poglavju. Program se nato sprehodi čez vse HTML datoteke in iz vsake izlušči samo besedilo. Nato se dobljeno besedilo tokenizira in normalizira, kot je opisano v poglavju *Indeksiranje podatkov*. Sprehodimo se čez vsak token v normalizirani tabeli in pogledamo če se ta nahaja v podanih parametrih. Če se povečamo število njegovih pojavitev in shranimo njegov indeks (položaj tokena v tabeli). Ko imamo preštete vse pojavitve in shranjene indekse, lahko s pomočjo tokenizirane tabele in indeksov pridobimo originalni tekst ter ga izpišemo.

- **Iskanje podatkov z obrnjenim indeksom**

Vzpostavi se povezava z podatkovno bazo, kjer se nahaja indeks. Programu prav tako, kot v prejšnjem primeru preko ukazne vrstice podamo parametre. Ti se tokenizirajo in normalizirajo po istem postopku, kot je opisan v prvem poglavju. Nato se naredi poizvedba za vsako besedo v podanih parametrih:

```
SELECT documentName, sum(frequency), group_concat(indexes) FROM (" + subquery  
+ ") GROUP BY documentName ORDER BY sum(frequency) DESC
```

Sestavljena je iz dveh delov. Zunanji del prešteje, za vsako dobljeno stran, število ponovitev vseh besed in združi njihove indekse. V notranjem delu (subquery) pa se za vsako besedo zberejo pojavitve v bazi in se jih združi z ukazom UNION ALL. To naredi naslednja poizvedba:

```
SELECT documentName, frequency, indexes FROM Posting WHERE word = '"' + word  
+ '"'
```

Dobljene rezultate nato publikujemo v primeren izpis.

□ Analiza baze in rezultati

V naši SQLite podatkovni zbirki v tabeli IndexWord se nahaja 35029 zapisov, v tabeli »Posting« pa 352951 zapisov, kjer je se največkrat (upoštevajoč tudi besede, kjer jih izbrana knjižnica za lematizacijo ni uspela popraviti) pojavi beseda »proizvodnja« in sicer 2266-krat samo v dokument evem.gov.si.371.html, ki je hkrati tudi dokument z največ raznolikimi besedami (13301 zapisov), najmanj raznolikimi besedami (31) pa evem.gov.si.55.html. Besede v največ dokumentih so »pogoji« (1398), »uporabe« (1399), »domov« pa 1384. Iskani niz »predelovalne dejavnosti«, se pojavi največkrat – 1288 krat v dokumentu evem.gov.si.371.html. Iskani niz »social services« se pojavi najmanjkrat – 12 krat skupno v štirih dokumentih. Za preverbo rezultatov smo uporabili orodje [Posting | Table | inverted-index.db \(sqliteviewer.app\)](https://sqliteviewer.app/#/inverted-index.db/table/Posting/).

The screenshot displays the SQLite Viewer Web App interface in a web browser. The URL is <https://sqliteviewer.app/#/inverted-index.db/table/Posting/>. The page title is "SQLite Viewer Web App". Below the title, it states: "SQLite Viewer Web is a free, web-based SQLite Explorer, inspired by DB Browser for SQLite and Airtable. Use this web-based SQLite Tool to quickly and easily inspect .sqlite files. Your data stays private: Everything is done client-side and never leaves your browser." There is a logo for the application on the right.

The main interface shows a table named "Posting" from the database "inverted-index.db". The table has two columns: "word" and "frequency". The "word" column contains the value "trgovina" for all 16 visible rows. The "frequency" column contains various document names, all starting with "./sites/evem.gov.si/". The table is sorted by frequency in descending order. The first row has a frequency of 215, and the last row has a frequency of 1. The interface includes a search bar at the top right, a "Reset Filters" button, and a "Records: 215*" indicator. The bottom of the interface shows a pagination bar indicating "Page 1 / 3".

word	frequency
trgovina	./sites/evem.gov.si/evem.gov.si.371.html
trgovina	./sites/evem.gov.si/evem.gov.si.651.html
trgovina	./sites/evem.gov.si/evem.gov.si.21.html
trgovina	./sites/podatki.gov.si/podatki.gov.si.340.html
trgovina	./sites/evem.gov.si/evem.gov.si.623.html
trgovina	./sites/evem.gov.si/evem.gov.si.329.html
trgovina	./sites/evem.gov.si/evem.gov.si.630.html
trgovina	./sites/evem.gov.si/evem.gov.si.320.html
trgovina	./sites/evem.gov.si/evem.gov.si.327.html
trgovina	./sites/evem.gov.si/evem.gov.si.622.html
trgovina	./sites/evem.gov.si/evem.gov.si.328.html
trgovina	./sites/evem.gov.si/evem.gov.si.343.html
trgovina	./sites/evem.gov.si/evem.gov.si.620.html
trgovina	./sites/evem.gov.si/evem.gov.si.316.html
trgovina	./sites/evem.gov.si/evem.gov.si.323.html
trgovina	./sites/evem.gov.si/evem.gov.si.334.html

□ Hitrosti poizvedb

Preizkusili smo dve poizvedbi za največkrat pojavljeno besedo “proizvodnja”. S pomočjo inverznega indeks (podatkovne baze), je poizvedba trajala **22 milisekund**, brez indeksa pa je trajala **83 sekund**.

Primer poizvedbe za »avtor franci«.

```
['avtor', 'franci']
SELECT documentName, sum(frequency), group_concat(indexes) FROM (SELECT documentName, frequency, indexes FROM Posting WHERE word = 'avtor' UNION ALL SELECT documentName, frequency, indexes FROM Posting WHERE word = 'franci') GROUP BY documentName ORDER BY sum(frequency) DESC
Search took: 0.00404953956040039 seconds
Frequencies Document                               Snippet
-----
0 podatki.gov.si/podatki.gov.si.10.html             širši javnosti. Avtor: Ministrstvo ... na 5 min Avtor: DARS ... v Vipavski dolini Avtor: DARS ... o gostoti prometa Avtor: DARS ... na cestah. Avtor: DARS ... javnem sektoru. Avtor: Adam ... leto 2015. Avtor: Virostatiq.com ... slovenskim sopredjem. Avtor: Virostatiq.com ... 2013 naprej. Avtor: M3U
4 even.gov.si/even.gov.si.12.html                   skleneta naročnik in avtor. Naročnik ... pravna oseba. Avtor je fizična ... avtorsko delo. Avtor je lahko ... ali je avtor že sicer
3 podatki.gov.si/podatki.gov.si.10.html             širši javnosti. Avtor: Ministrstvo ... na 5 min Avtor: DARS ... v Vipavski dolini Avtor: DARS
3 even.gov.si/even.gov.si.85.html                   ki jih plača avtor. Od ... če je avtor zavarovanec po ... dohodka, če avtor ni zavarovanec
2 even.gov.si/even.gov.si.73.html                   je, da avtor in podjetnik ... morata plačati tudi avtor in podjetnik
2 e-uprava.gov.si/e-uprava.gov.si.9.html             Bled, avtor: Franci
2 e-uprava.gov.si/e-uprava.gov.si.45.html            Bled, avtor: Franci
1 e-uprava.gov.si/e-uprava.gov.si.6.html             park Topla, avtor: Tomo
1 e-uprava.gov.si/e-uprava.gov.si.5.html             Bled, avtor: Mirko
1 e-uprava.gov.si/e-uprava.gov.si.43.html            Bled, avtor: Mirko
1 e-uprava.gov.si/e-uprava.gov.si.42.html            park Topla, avtor: Tomo
1 e-prostor.gov.si/e-prostor.gov.si.40.html          Žal pa avtor ni imel
```