

Programming Assignment 3

Implementing a data index and more

Jaka Kokošar, Danijel Maraž, and Toni Kocjan

Fakulteta za Računalništvo in Informatiko UL `dm9929@student.uni-lj.si`,
`jk0902@student.uni-lj.si`, `tk3152@student.uni-lj.si`

Povzetek The article covers the work done in the scope of the third programming assignment as part of the subject web information extraction and retrieval.

Keywords: Data Processing Indexing Retrieval

1 Introduction

After having collected web pages in the first assignment and thoroughly stripped them down to the data we are interested in in the second we were ready to continue in the final step which is constructing a data index and implementing the function of querying.

2 Data Processing

Glavna funkcija *preprocess* prejme kot argument rezultat funkcije *text*, ki sama prejme našo surovo html vsebino in tej odstrani nepotrebne tehnične html oznake. Preprocess nato:

- S funkcijo *remove_punctuation* odstrani ločila in več
- Z *nlTK.tokenize.word_tokenize* pretvori v vrsto besednih značk
- Odstrani se značke, ki niso alfabetične
- Vse velike začetnice se pretvorijo v male
- Odstrani se značke, ki so *stopword*
- S pomočjo pretvorbe v podatkovno strukturo množice se odstranijo duplikati

Nato ta vrne seznam ostalih značk.

3 Indexing

Funkcija *initiating_indexing* se požene in začne meriti čas gradnje indeksa. Na koncu izpiše na standardni izhod celoten porabljen čas. Glavno nalogo indeksiranja opravlja razred *BetterThanGoogle* nad katerim se pokliče funkcijo *create_index* in se mu kot argument poda relativno pot do datotek za sestavo indeksa.

3.1 BetterThanGoogle

3.2 Create_index

Funkcija kot argumenta prejme instanci razredov *Preprocessor* (ta služi za procesiranje besedila po opisu iz poglavja Data Processing) in *DBHandler* (ta služi za interakcijo z bazo). Nato pridobimo ime datoteke ter vsebino s pomočjo naše abstrakcije korpusov (*file_name* in *document*). Za tem iteriramo skozi vsak par ter spremenljivko *document* ustrezno obdelamo z razredom *preprocessor* (glej *_call_* od *preprocessor*). Za vsako značko, ki jo vrne *preprocessor*:

- Najdemo vse njene pojavitve v besedilu (*find_occurrences*)
- Pod pogojem, da smo našli vsaj eno pojavitev se izvede faza vnosa v indeks
- V indeks vnesemo ime značke, ime datoteke, število pojavitev značke, ter niz posameznih pojavitev ločen z vejico

Vredno je tudi omeniti, da program v log ves čas izpisuje koliko datotek je obdelal do sedaj in koliko mu jih še manjka.

4 Data Retrieval

Funkcija *initiating_search* prejme niz za katerega želimo iskati pojavitve v trenutnem indeksu. Ta ustvari nov objekt *SearchEngine* in mu poda instanco *DBHandler*. Potem se za dejansko iskanje na ustvarjenem objektu kliče funkcijo *perform_query*. Nato program izpiše na zaslon rezultate in čas porabljen za poizvedbo.

4.1 SearchEngine

Razred *SearchEngine* kot argument prejme *db_handler* (preko katerega se izvajajo vse interakcije z bazo).

Perform_query Funkcija kot argument prejme niz besed ločenih s presledki. Nato iz teh ustvari seznam besed, ter vsaki besedi velike črke zamenja z malimi. Za vsako se nato naredi poizvedba v bazi iz tabele *Posting* in s funkcijo *_find_occurrences_in_file* in več zankami ustvari terko *QueryResults(beseda, snipeți pojavitev besede)*. Na koncu funkcija vrne drugo terko s številom pojavitev na prvem mestu in seznamom terk *QueryResults* na drugem.

5 Rezultati poizvedb

5.1 predelovalne dejavnosti

Searching time: 22.426688194274902

Found 6367 results for "predelovalne dejavnosti"

Found 1570 results in

"data/evem.gov.si/evem.gov.si.371.html":

```
0: '... - \n \n * **C** PREDELOVALNE DEJAVNOSTI\n\n
... '
1: '... ruge raznovrstne predelovalne dejavnosti\n\n ... '
2: '... gje ščnerazvrene predelovalne dejavnosti\n\n ... '
3: '... v industrijskem predelovalnem procesu, gl ... '
4: '... e iz čpodroja C (Predelovalne dejavnosti)_ ... '
5: '... e ustrezne šifre dejavnosti /storitve in ... '
6: '... h za opravljanje dejavnosti.\n\nV iskalnik ... '
7: '... ih je 645 od 645 dejavnosti\n\nIzpisanih j ... '
8: '... zpisanih je od dejavnosti\n\n * **A** K ... '
9: '... h za opravljanje dejavnosti:\n\n Pridelav ... '
10: '... h za opravljanje dejavnosti:\n\n Pridelav ... '
11: '... h za opravljanje dejavnosti:\n\n Pridelava ... '
12: '... h za opravljanje dejavnosti:\n\n ... '
13: '... h za opravljanje dejavnosti:\n\n ... '
14: '... h za opravljanje dejavnosti:\n\n ... '
15: '... h za opravljanje dejavnosti:\n\n ... '
16: '... h za opravljanje dejavnosti:\n\n ... '
17: '... **\n\n#### · Lista dejavnosti, ki se čobia ... '
18: '... obrtni čnain\n\nZa dejavnosti, ki so ščuvr ... '
19: '... ščvrene na Listo dejavnosti, ki se čobia ... '
20: '... registraciji te dejavnosti poslovni sub ... '
```

5.2 trgovina

Searching time: 5.710252046585083

Found 1158 results for "trgovina"

Found 368 results in

"data/evem.gov.si/evem.gov.si.371.html":

```

0: '... .110_\n          * _trgovina na debelo s ...'
1: '... .890_\n          * _trgovina na debelo z ...'
2: '... .890_\n          * _trgovina na debelo s ...'
3: '... .380_\n          * _trgovina na drobno s ...'
4: '... alij\n          * Trgovina na debelo z ...'
5: '... jami\n          * Trgovina na drobno z ...'
6: '... .500_\n          * _trgovina na debelo s ...'
7: '... .460_\n          * _trgovina na drobno s ...'
8: '... a čščienje tal v trgovinah in industri ...'
9: '... .320_\n          * _trgovina (odkup in pr ...'
10: '... .220_\n          * _trgovina na debelo z ...'
11: '... o \n \n * **G** TRGOVINA; ŽVZDREVANJE ...'
12: '... IL\n\n          * **45** Trgovina z motornimi ...'
13: '...          * **45.110** Trgovina z avtomobili ...'
14: '... da:**\n\n          * trgovina na debelo in ...'
15: '... :_**\n\n          * _trgovina na debelo al ...'
16: '...          * **45.190** Trgovina z drugimi mo ...'
17: '... da:**\n\n          * trgovina na debelo al ...'
18: '... :_**\n\n          * _trgovina na debelo al ...'
19: '...          * **45.310** Trgovina na debelo z ...'
20: '... da:**\n\n          * trgovina na debelo s ...'

```

5.3 social services

Searching time: 1.4749388694763184

Found 12 results for "social services"

Found 5 results in

"data/e-uprava.gov.si/e-uprava.gov.si.45.html":

```

0: '... retirement\n * Social services , heal ...'
1: '... hip etc.?\n\n### Social services , heal ...'
2: '... tain financial social assistance? Ho ...'

```

```
3: '... ent\n * Social services , health , death ...'
4: '... .?\n\n### Social services , health , death ...'
```

Found 5 results in

```
"data/e-uprava.gov.si/e-uprava.gov.si.9.html":
0: '... retirement\n * Social services , heal ...'
1: '... hip etc.?\n\n### Social services , heal ...'
2: '... tain financial social assistance? Ho ...'
3: '... ent\n * Social services , health , death ...'
4: '... .?\n\n### Social services , health , death ...'
```

Found 1 results in

```
"data/evem.gov.si/evem.gov.si.661.html":
0: '... ords and Related Services (AJ PES) and ...'
```

Found 1 results in

```
"data/podatki.gov.si/podatki.gov.si.340.html":
0: '... creation and spa services ltd.\n\nTERME ...'
```

5.4 MJU

Searching time: 0.6603918075561523

Found 28 results for "MJU"

Found 5 results in

```
"data/podatki.gov.si/podatki.gov.si.295.html":
0: '... 017/18** , ki ga MJU organizira s ...'
1: '... \n\n\n\n\n\nhtml\n\n##### MJU\n\n**Podrobnos ...'
```

Found 2 results in

"data/evem.gov.si/evem.gov.si.68.html":

- 0: '... za javno upravo (MJU), ki na podl ...'
- 1: '... racija posreduje\nMJU, ki uredi za ...'

Found 2 results in

"data/podatki.gov.si/podatki.gov.si.105.html":

- 0: '... ronska špota: gp.mju@gov.si\n * ...'
- 1: '... 1)\n\nhtml\n\n##### MJU\n\n**Podrobnos ...'

5.5 •