

# WIER 2021: Programming assignment 3, Report

---

Klemen Stanič, 63150267

Luka Kavčič, 63150139

## 1. Introduction

The goal of this assignment was to implement an algorithm that would first process the provided HTML pages and index them into a database. The built index is then used for querying against it in order to find the occurrences of the search term in all the pages. To get a better representation of the speedup an inverted index achieves compared to a more naive approach, we also implemented an algorithm that achieves the same results by sequentially searching files and compared the time complexity.

## 2. Implementation

### 2.A Data processing and indexing

Before we can index the data, the HTML pages must be cleaned of all the unnecessary words, that don't contain much information, such as stopwords. This process consists of:

- *HTML text extraction:* We use **beautifulsoup** to parse the webpages and extract only the textual information, disregarding all html tags, etc. At this point, we also remove all the *script* sections of the HTML, we clean up all trailing spaces, combine multiline paragraphs into a single line etc.
- *Word tokenization:* All the words are tokenized using the **nltk.tokenize** package. After the tokenization process, we transform the words to lower case. We remove all the stopwords using the provided stoplist.

The same process is used on search query terms and the *basic* version of this algorithm.

Indexing the cleaned data into the *sqlite* database:

- We open every file and clean it with the process we described above. We then iterate the words, and for each word (out of the preprocessed text), we count the number of occurrences of the word in the file, and find the indexes of these occurrences.
- We insert this data into IndexWord and Posting database tables.

### 2.B Data retrieval / querying

When a user enters string query, we first clean it using the same process we described above. Once cleaned, we iterate the query word by word and query the database for the word and append it to a list. This list, now containing all occurrences of all the words in the search query, is then sorted by the document name, in which it appeared in. Words that appear in the same document are aggregated along with their frequencies and indexes. The list is again sorted, this time using the word frequency. This list, combined with snippets near the word indexes, is then printed to the standard output.

### 2.C Naive approach

The algorithm opens each file sequentially and cleans the content of the file. We then iterate the list of (cleaned) words of the file and search for any words that match the words contained in search query. If the word matches the query search word, we save the words in the direct vicinity and increment the word's frequency. After the whole dataset is searched, we sort the results by frequency and print out the results.

### 3. Results

Results on the given search queries are presented in section 5. We limited the number of snippets per document to 2 and also truncated some results, to only show top hits. Full outputs, both for the basic and inverted index version are stored in folder `results`. TODO: comment results

#### 3.A Inverted index database, basic info:

Our database consists of two tables. First table, named `IndexWords`, contains all tokenized words found in the given documents, which includes 48 910 different words. During the importing phase we excluded stop words and did not save them to our database. The second table, named `Posting`, holds pairs of words and documents that contains them. We also save the frequency of appearance of the given word in the document as well as the place of appearance stored as indexes. Posting table holds 383 347 pairs.

|Word/Document frequencies|

Word	Document	Frequency
proizvodnja	../data/evem.gov.si/evem.gov.si.371.html	2266
gl	../data/evem.gov.si/evem.gov.si.371.html	1668
spada	../data/evem.gov.si/evem.gov.si.371.html	1338
dejavnosti	../data/evem.gov.si/evem.gov.si.371.html	1284
d.o.o	../data/podatki.gov.si/podatki.gov.si.340.html	967
xsd	../data/e-prostor.gov.si/e-prostor.gov.si.147.html	925
skupnost	../data/podatki.gov.si/podatki.gov.si.340.html	809
krajevna	../data/podatki.gov.si/podatki.gov.si.340.html	754
ministrstvo	../data/evem.gov.si/evem.gov.si.371.html	589

TODO : this TODO: speedup compare basic and sqlite

### 4. Conclusions

TODO: pretty good, but not really

#### (5.) Search query results

Search query *predelovalne dejavnosti*

Results for a query: "predelovalne dejavnosti"

Results found in 0.02358245849609375s

Frequencies	Document	Snippet
1288	../data/evem.gov.si/evem.gov.si.371.html	... za infrastrukturo C PREDELOVALNE DEJAVNOSTI 10 Proizvodnja ... 32 Druge raznovrstne predelovalne dejavnosti 32.110 Kovanje ...
75	../data/evem.gov.si/evem.gov.si.377.html	... Defektolog v zdravstveni dejavnosti Dekan oziroma direktor ... Dietetik v zdravstveni dejavnosti Dimnikar Diplomirana medicinska ...
40	../data/podatki.gov.si/podatki.gov.si.340.html	... - NOSILEC DOPOLNILNE DEJAVNOSTI NA KMETIJI BREGAR ... šport CENTER INTERESNIH DEJAVNOSTI PTUJ CENTER JUDOVSKO ...
39	../data/evem.gov.si/evem.gov.si.452.html	... ... nerazvrščene (96.090) / Dejavnosti / eVEM Republika ...
31	../data/evem.gov.si/evem.gov.si.653.html	... Dovoljenje za opravljanje dejavnosti specializirane prodajalne z ... radijske ali televizijske dejavnosti Dovoljenje za izvajanje ...
29	../data/evem.gov.si/evem.gov.si.72.html	... od dohodka iz dejavnosti Republika Slovenija SPOT, ... od dohodka iz dejavnosti Davek od dohodka ...
29	../data/evem.gov.si/evem.gov.si.398.html	... usmerjene na opravljanje dejavnosti (npr.: pripravljala dela, ... za namene opravljanja dejavnosti ipd. V obdobju ...
23	../data/evem.gov.si/evem.gov.si.442.html	... ... telesa (96.040) / Dejavnosti / eVEM Republika ...
19	../data/evem.gov.si/evem.gov.si.460.html	... ... e-VEM eVEM>Dejavnosti>Drugje nerazvrščene predelovalne dejavnosti (32.990) Drugje ...
18	../data/evem.gov.si/evem.gov.si.28.html	... za opravljanje gospodarske dejavnosti. Lastnosti zasebnega zavoda ... niso za posamezne dejavnosti ali posamezne vrste ...
17	../data/evem.gov.si/evem.gov.si.42.html	... / Načini opravljanja dejavnosti Republika Slovenija SPOT, ... poslovni poti>Načini opravljanja dejavnosti Načini opravljanja dejavnosti ...
17	../data/evem.gov.si/evem.gov.si.265.html	... perutninskega (10.110) / Dejavnosti / eVEM Republika ... SKD šifra zajema dejavnosti in storitve, za ...
16	../data/evem.gov.si/evem.gov.si.574.html	... lesa (16.100) / Dejavnosti / eVEM Republika ... SKD šifra zajema dejavnosti in storitve, za ...
16	../data/evem.gov.si/evem.gov.si.564.html	... sokov (10.320) / Dejavnosti / eVEM Republika ... SKD šifra zajema dejavnosti in storitve, za ...
16	../data/evem.gov.si/evem.gov.si.522.html	... ... prometu (52.210) / Dejavnosti / eVEM Republika ...
16	../data/evem.gov.si/evem.gov.si.392.html	... redno poslovanje. 6. Dejavnosti podjetja Ob ustanovitvi gospodarske ... prihodkov) in druge dejavnosti, ki jih prav ...
16	../data/evem.gov.si/evem.gov.si.276.html	... veterinarskih (01.620) / Dejavnosti / eVEM Republika ... SKD šifra zajema dejavnosti in storitve, za ...
16	../data/evem.gov.si/evem.gov.si.272.html	...

```

gozdarstvo (02.400) / Dejavnosti / eVEM Republika ... SKD šifra zajema
dejavnosti in storitve, za ...
16      ../data/evem.gov.si/evem.gov.si.266.html      ...
izdelkov (10.130) / Dejavnosti / eVEM Republika ... SKD šifra zajema
dejavnosti in storitve, za ...
15      ../data/evem.gov.si/evem.gov.si.546.html      ...
sitarstvo (10.510) / Dejavnosti / eVEM Republika ... SKD šifra zajema
dejavnosti in storitve, za ...
15      ../data/evem.gov.si/evem.gov.si.267.html      ... maščob
(10.410) / Dejavnosti / eVEM Republika ... SKD šifra zajema dejavnosti in
storitve, za ...
15      ../data/evem.gov.si/evem.gov.si.140.html      ...
streliva (25.400) / Dejavnosti / eVEM Republika ... SKD šifra zajema
dejavnosti in storitve, za ...
...

```

## Search query *trgovina*

Results for a query: "trgovina"

Results found in 0.011131763458251953s

Frequencies	Document	Snippet
364	../data/evem.gov.si/evem.gov.si.371.html	...
94	../data/evem.gov.si/evem.gov.si.651.html	... Druga
92	../data/evem.gov.si/evem.gov.si.21.html	... Moj e-
82	../data/podatki.gov.si/podatki.gov.si.340.html	... d.o.o.
13	../data/evem.gov.si/evem.gov.si.623.html	... ...
12	../data/evem.gov.si/evem.gov.si.630.html	... ...
12	../data/evem.gov.si/evem.gov.si.329.html	... ... in
10	../data/evem.gov.si/evem.gov.si.622.html	... ...
10	../data/evem.gov.si/evem.gov.si.327.html	... ...
10	../data/evem.gov.si/evem.gov.si.320.html	... ...
9	../data/evem.gov.si/evem.gov.si.620.html	... ...

```

9      ../data/evem.gov.si/evem.gov.si.343.html      ...  ...
prodajalnah s tekstilom Trgovina na drobno v ...
9      ../data/evem.gov.si/evem.gov.si.328.html      ...  ... in
plinastimi gorivi Trgovina na debelo s ...
8      ../data/evem.gov.si/evem.gov.si.450.html      ...  ...
popravila motornih koles; trgovina z njihovimi deli ...
8      ../data/evem.gov.si/evem.gov.si.334.html      ...  ...
semeni in krmo Trgovina na debelo z ...
8      ../data/evem.gov.si/evem.gov.si.323.html      ...  ...
debelo s tekstilom Trgovina na debelo s ...
8      ../data/evem.gov.si/evem.gov.si.316.html      ...  ... in
rudami (46.720) Trgovina na debelo s ...
7      ../data/evem.gov.si/evem.gov.si.617.html      ...  ...
preprogami in svetili Trgovina na debelo s ...
7      ../data/evem.gov.si/evem.gov.si.350.html      ...  ...
barvami in steklom Trgovina na drobno v ...
7      ../data/evem.gov.si/evem.gov.si.339.html      ...  ... z
rabljenim blagom Trgovina na drobno v ...
7      ../data/evem.gov.si/evem.gov.si.319.html      ...  ...
stroji, priključki, opremo Trgovina na debelo s ...
7      ../data/evem.gov.si/evem.gov.si.315.html      ...  ... in
materiali (46.460) Trgovina na debelo s ...
6      ../data/evem.gov.si/evem.gov.si.641.html      ...  ...
lastnimi motornimi gorivi Trgovina na drobno z ...
6      ../data/evem.gov.si/evem.gov.si.636.html      ...  ...
napravami in programi Trgovina na drobno v ...
6      ../data/evem.gov.si/evem.gov.si.632.html      ...  ... in
video napravami Trgovina na drobno v ...
6      ../data/evem.gov.si/evem.gov.si.619.html      ...  ... z
obdelovalnimi stroji Trgovina na debelo z ...
6      ../data/evem.gov.si/evem.gov.si.612.html      ...  ...
sadjem in zelenjavo Trgovina na debelo s ...
...

```

## Search query *social services*

Results for a query: "social services"

Results found in 0.00971531867980957s

Frequencies	Document	Snippet
5	../data/e-uprava.gov.si/e-uprava.gov.si.9.html	... culture
	Labour, retirement Social services, health, death ... employment	
	relationship etc.? Social services, health, death ...	
5	../data/e-uprava.gov.si/e-uprava.gov.si.45.html	... culture
	Labour, retirement Social services, health, death ... employment	
	relationship etc.? Social services, health, death ...	
1	../data/podatki.gov.si/podatki.gov.si.340.html	...
	recreation and spa services ltd. TERME MARIBOR, ...	

1	../data/evem.gov.si/evem.gov.si.661.html	... Records and Related Services (AJ PES) and the ...
---	--	---

## Search query *Janez Janša*

Results for a query: "Janez Janša"

Results found in 0.013368606567382812s

Frequencies	Document	Snippet
-----		
6	../data/podatki.gov.si/podatki.gov.si.340.html	... HERMINA OSOJNIK DOLENC JANEZ - ZASEBNA AMBULANTA ... d.o.o., Brnik FERLEŽ JANEZ - NOTAR FERLIGOJ ...
2	../data/evem.gov.si/evem.gov.si.378.html	... 76 33 infoping@notarkaerjavecpong.si Janez Ferlež Lojzeta Fabjana ... 84 44 notar.valter.vindisping@gmailpong.com Janez Klemenc Prešernov trg ...
2	../data/evem.gov.si/evem.gov.si.362.html	... 76 33 infoping@notarkaerjavecpong.si Janez Ferlež Lojzeta Fabjana ... 84 44 notar.valter.vindisping@gmailpong.com Janez Klemenc Prešernov trg ...
1	../data/podatki.gov.si/podatki.gov.si.250.html	... je prof. dr. Janez Bogataj razdelil na ...
1	../data/e-prostor.gov.si/e-prostor.gov.si.11.html	... v zvezi z vsebino:Janez Košir, Geodetska uprava ...

## Search query *Študentski domovi*

Results for a query: "Študentski domovi"

Results found in 0.009506702423095703s

Frequencies	Document	Snippet
-----		
12	../data/podatki.gov.si/podatki.gov.si.340.html	... VIČ DIJAŠKI IN ŠTUDENTSKI DOM KOPER - ... CAPODISTRIA DIJAŠKI IN ŠTUDENTSKI DOM KRANJ DIJAŠKI ...
5	../data/podatki.gov.si/podatki.gov.si.401.html	... Navajanje vira: Zavodi, domovi in druge ustanove ...
5	../data/evem.gov.si/evem.gov.si.371.html	... da lahko ustanovi študentski dom domača ali ... namene (kot so domovi za starejše osebe, ...
4	../data/evem.gov.si/evem.gov.si.484.html	... Moj e-VEM eVEM>Dejavnosti>Počitniški domovi in letovišča (55.201) ...
3	../data/evem.gov.si/evem.gov.si.83.html	... preko pooblašene organizacije (študentski servisi, Zavod Republike ... ki jo priskrbi študentski servis. Napotnica je ...
3	../data/evem.gov.si/evem.gov.si.168.html	...

```

Moj e-VEM eVEM>Dejavnosti>Planinski domovi in mladinska prenočišča ...
3      ../data/evem.gov.si/evem.gov.si.165.html      ...
Moj e-VEM eVEM>Dejavnosti>Počitniški domovi Počitniški domovi Dejavnost ...
3      ../data/evem.gov.si/evem.gov.si.16.html      ...
Gostišče Prenosišče Počitniški domovi Počitniška stanovanja in ... gostom
(55.203) Planinski domovi in mladinska prenočišča ...
3      ../data/e-uprava.gov.si/e-uprava.gov.si.23.html      ... muzeji,
galerije, kulturni domovi... Zdravstvo Seznam institucij ... so bolnišnice,
zdravstveni domovi, lekarne... Sociala Seznam ...
2      ../data/podatki.gov.si/podatki.gov.si.398.html      ... in
šport Zavodi, domovi in druge ustanove ... z naslovom "Zavodi, domovi in
druge... Nadaljujte ...
2      ../data/podatki.gov.si/podatki.gov.si.397.html      ... in
šport Zavodi, domovi in druge ustanove ... z naslovom "Zavodi, domovi in
druge... Nadaljujte ...
2      ../data/podatki.gov.si/podatki.gov.si.296.html      ... in
šport Zavodi, domovi in druge ustanove ... z naslovom "Zavodi, domovi in
druge... Nadaljujte ...
2      ../data/podatki.gov.si/podatki.gov.si.133.html      ... in
šport Zavodi, domovi in druge ustanove ... z naslovom "Zavodi, domovi in
druge... Nadaljujte ...
2      ../data/podatki.gov.si/podatki.gov.si.129.html      ... in
šport Zavodi, domovi in druge ustanove ... z naslovom "Zavodi, domovi in
druge... Nadaljujte ...
2      ../data/podatki.gov.si/podatki.gov.si.125.html      ... in
šport Zavodi, domovi in druge ustanove ... z naslovom "Zavodi, domovi in
druge... Nadaljujte ...
2      ../data/podatki.gov.si/podatki.gov.si.119.html      ... in
šport Zavodi, domovi in druge ustanove ... z naslovom "Zavodi, domovi in
druge... Nadaljujte ...
1      ../data/evem.gov.si/evem.gov.si.73.html      ... preko
pooblaščne organizacije (študentski servisi, Zavod Republike ...
1      ../data/evem.gov.si/evem.gov.si.654.html      ... in
naselja Počitniški domovi Planšarija Osmica več ...
1      ../data/evem.gov.si/evem.gov.si.651.html      ... in
naselja Počitniški domovi Pogrebna in pokopališka ...
1      ../data/evem.gov.si/evem.gov.si.650.html      ... in
naselja Počitniški domovi Planšarija Osmica več ...
1      ../data/evem.gov.si/evem.gov.si.164.html      ... članico
univerze in študentski dom lahko ustanovijo ...
1      ../data/evem.gov.si/evem.gov.si.109.html      ... namene
(kot so domovi za starejše osebe, ...
1      ../data/e-uprava.gov.si/e-uprava.gov.si.56.html      ... primeru
telesne okvare. Domovi za starejše Namestitev ...
1      ../data/e-uprava.gov.si/e-uprava.gov.si.50.html      ... ob
poznejši upokojitvi, domovi za starejše, delna ...

```

## Search query *Vozniško dovoljenje*

Results for a query: "vozniško dovoljenje"



Results found in 0.012945175170898438s

Frequencies	Document	Snippet
188	../data/evem.gov.si/evem.gov.si.371.html	... mora uporabnik pridobiti dovoljenje pri pristojnem organu. ... je potrebno pridobiti dovoljenje, ki ga izda ...
107	../data/evem.gov.si/evem.gov.si.653.html	... medicine - licenca Dovoljenje za opravljanje dejavnosti ... zdravili na drobno Dovoljenje Ministrstva za zdravje ...
16	../data/evem.gov.si/evem.gov.si.398.html	... iz knjige sklepov, dovoljenje ATPV za spremembo ... črk. OZS, Obrtno dovoljenje Kdo si mora ...
10	../data/evem.gov.si/evem.gov.si.84.html	... imeti veljavno enotno dovoljenje za prebivanje in ... Srbije pa tudi delovno dovoljenje po bilateralnem sporazumu, ...
7	../data/evem.gov.si/evem.gov.si.43.html	... in pridobiti ustrezno dovoljenje, preden začne opravljati dejavnost, ... mora imeti uporabno dovoljenje Poslovni prostor lahko ...
7	../data/evem.gov.si/evem.gov.si.312.html	... distributerji, ki imajo dovoljenje za opravljanje prometa ... o registraciji ali dovoljenje za nujne primere ...
7	../data/e-uprava.gov.si/e-uprava.gov.si.56.html	... enoti. Če se ... Vozniško dovoljenje Z vozniškim ... Če se ... Vozniško dovoljenje Z vozniškim dovoljenjem ...
6	../data/evem.gov.si/evem.gov.si.599.html	... in imajo veljavno dovoljenje organa, pristojnega za ... nadaljnjem besedilu: veletrgovci). Dovoljenje za promet z ...
6	../data/evem.gov.si/evem.gov.si.539.html	... dejavnosti potrebujete vstopno dovoljenje in ustrezen kader. ... je potrebno pridobiti dovoljenje za gospodarski ribolov. ...
6	../data/evem.gov.si/evem.gov.si.373.html	... enote Prošnja za dovoljenje za prebivanje na ... X Prošnja za dovoljenje za prebivanje na ...
5	../data/evem.gov.si/evem.gov.si.582.html	... o registraciji ali dovoljenje za nujne primere ... nujne primere ali dovoljenje za vzporedno trgovanje ...
5	../data/evem.gov.si/evem.gov.si.157.html	... ne potrebuje začetno dovoljenje, ki preverja izpolnjevanje ... imeti mora ustrezno uporabno dovoljenje, upoštevati se morajo minimalni ...
5	../data/e-uprava.gov.si/e-uprava.gov.si.33.html	... izkaznica, potni list, vozniško dovoljenje). Tip dokumenta ... listina Osebna izkaznica Vozniško dovoljenje Obmejna prepustnica ...
4	../data/evem.gov.si/evem.gov.si.610.html	... za to izdano dovoljenje. Različne vrste dovoljenj ... in registracija Priglasitev Nacionalno dovoljenje (avtorizacija) Dovoljenje po poenostavljenem ...
4	../data/evem.gov.si/evem.gov.si.584.html	... o registraciji ali dovoljenje za nujne primere ... nujne primere ali dovoljenje za vzporedno trgovanje ...
4	../data/evem.gov.si/evem.gov.si.572.html	... mora vlogo za dovoljenje za obnovo vinograda ... po potrebi pridobiti dovoljenje za nujne primere ...



4           ../data/evem.gov.si/evem.gov.si.467.html           ...  
dejavnosti se ne potrebuje začetno dovoljenje, ki preverja izpolnjevanje  
... imeti mora ustrezno uporabno dovoljenje, upoštevati se morajo minimalni  
prostorski ...

4           ../data/evem.gov.si/evem.gov.si.464.html           ...  
dejavnosti se ne potrebuje začetno dovoljenje, ki preverja izpolnjevanje  
... imeti mora ustrezno uporabno dovoljenje, upoštevati se morajo minimalni  
...

4           ../data/evem.gov.si/evem.gov.si.359.html           ... se  
nanaša posredništvo, dovoljenje za promet z ... ali začasno izredno  
dovoljenje za promet z ...

4           ../data/evem.gov.si/evem.gov.si.158.html           ...  
dejavnosti se ne potrebuje začetno dovoljenje, ki preverja izpolnjevanje  
... imeti mora ustrezno uporabno dovoljenje, upoštevati se morajo minimalni  
prostorski ...

4           ../data/evem.gov.si/evem.gov.si.152.html           ...  
dejavnosti se ne potrebuje začetno dovoljenje, ki preverja izpolnjevanje  
... imeti mora ustrezno uporabno dovoljenje, upoštevati se morajo ...

4           ../data/evem.gov.si/evem.gov.si.140.html           ...  
strelišča, ko dobi dovoljenje Ministra, pristojnega za ... je potrebno  
posebno dovoljenje, ki ga izda ...

3           ../data/podatki.gov.si/podatki.gov.si.423.html       ...  
prodajaln, ki imajo dovoljenje za izdajo zdravil ... prodajaln, ki imajo  
dovoljenje za izdajo zdravil ...

3           ../data/evem.gov.si/evem.gov.si.68.html           ...  
položaja točke SPOT Dovoljenje za pridobitev položaja ... kršitev. MGRT  
odvzame dovoljenje za opravljanje postopkov ...

3           ../data/evem.gov.si/evem.gov.si.645.html           ... zdravil  
prek medmrežja Dovoljenje za izdajo zdravil ... promet z zdravili  
Dovoljenje za opravljanje dejavnosti ...

3           ../data/evem.gov.si/evem.gov.si.642.html           ... pogoje  
in imeti dovoljenje za opravljanje dejavnosti ... in dokazila Stalno  
dovoljenje za opravljanje dejavnosti ...

3           ../data/evem.gov.si/evem.gov.si.596.html           ... osebe,  
ki imajo dovoljenje za izvajanje sevalne ... sevalne dejavnosti in  
dovoljenje za uporabo vira ...

...