

Miha Kranjc, Marko Adžaga, Matic  
Hrastelj

## Domača naloga 1 - Poročilo

Iskanje in ekstrakcija podatkov s spleta,  
Domača naloga 1

MENTOR: prof. dr. Marko Bajec, Timotej  
Knez

### Posebnosti, odločitve in opis implementacije

Implementacija pajka je ločena po komponentah, ki opravljajo različne naloge. Glavni del je razred *Crawler*, ki nadzoruje celoten proces ekstrakcije. Na začetku izvajanja ustvari določeno število niti oziroma delavcev. Vsaka nit izvaja neskončno zanko, ki iz podatkovne baze prevzame spletno stran, iz nje izlušči podatke ter rezultate zapiše nazaj v podatkovno bazo. Dostopi do baze so zavarovani s pomočjo ključavnice. Samo luščenje podatkov se izvaja v ločenem razredu *Extractor*, s pomočjo knjižnic *Selenium* in *BeautifulSoup*. Med obiskovanjem spletnih strani vedno upoštevamo potencialen časovni zamik pri nalaganju ter vse možne izjeme, ki se lahko zgodijo pri poizvedovanju. Rezultati se naknadno vpišejo nazaj v podatkovno bazo. Tokrat nam ni potrebno skrbeti za hkraten dostop, saj ne more priti do konflikta. V podatkovno bazo v tabelo *page* smo dodali tudi stolpec, ki vsebuje hash vrednost HTML strani, za lažje primerjanje med stranmi.

### Parametri pajka in težave

Implementiran pajek ima zgolj tri parametre, ki jih lahko določi uporabnik. V datoteki *main.py* lahko nastavimo semenska spletna mesta oziroma izhodišča (trenutno so podana štiri spletna mesta iz navodil naloge), število niti z argumentom *worker\_count*, s katerimi pajek prečesava splet ter minimalni pretečeni čas med klici na spletno mesto *default\_access\_period*, ki je privzeto nastavljen na 5 sekund. Opcijsko se lahko nastavijo še parametri za povezavo na podatkovno bazo v samem razredu pajka *Crawler*. Parametrov, ki pa jih uporabnik ne more nastaviti in se nastavijo ob vzpostavitvi pajka, pa je kar precej. Od ključavnice, do slovarja dostopnih časov vseh spletnih mest in samih ekstraktorjev.

Med samo implementacijo in testiranjem smo imeli tudi kar nekaj težav. Nekatere omembe vredne težave so bile problem s certifikatom SSL, kjer smo morali dovoliti, da ni potreben, nekaj težav smo imeli tudi s ključavnico, ter tudi s počasnim dodajanjem zapisov v bazo, kjer smo morali uporabiti poseben ukaz za hitrejša vnašanja.

## Statistika

### Spletna mesta iz seznama semen:

- Število spletnih mest: **4**
- Število spletnih strani (HTML, BINARY, FRONTIER, DUPLICATE): **29.888**
- Število spletnih strani (HTML): **23.303**
- Število duplikatov: **768**
- Število binarnih elementov glede na tip:
  - o PDF: **25.475**
  - o DOC: **3.157**
  - o DOCX: **8.803**
  - o PPT: **11**
  - o PPTX: **198**
- Število slik: **52.831**
- Povprečno število slik na spletno stran: **2,4**

### Celotna podatkovna baza crawldb:

- Število spletnih mest: **318**
- Število spletnih strani (HTML, BINARY, FRONTIER, DUPLICATE): **133.301**
- Število spletnih strani (HTML): **52.735**
- Število duplikatov: **6.674**
- Število binarnih elementov glede na tip:
  - o PDF: **55.343**
  - o DOC: **7.565**
  - o DOCX: **14.499**
  - o PPT: **78**
  - o PPTX: **259**
- Število slik: **488.941**
- Povprečno število slik na spletno stran: **11,2**

## Vizualizacija povezav

Zaradi prevelikega števila skupnih url-jev smo se odločili vizualizirati samo url-je z spletne domene <https://www.e-prostor.gov.si>. Vizualizacijo smo dodatno omejili z poddomenami: dostopi, inspire, podrocja, novica, en, skupno.



