

# Rangirajući pretraživač dokumenata - Infinity

Marko Budiselić

30.11.2015.

## 1 Pretprocesiranje

Bitan segment analize teksta je predprocesiranje. Predprocesiranje počinje s podjelom dokumenta na riječi (engl. *tokenization*). Implementiran je jednostavni tokenizator koji prvo zamijenjuje posebna znakove s prazninama, nakon toga odbacuje riječi duže od 25 znakova (25 je nastao subjektivnom procjenom) i na kraju se radi svođenje riječi na jednostavniji oblik (engl. *stemming*). Korišten je Snowball (Porter2) stemmer zato što je manje agresivan od Lancaster stemmer-a, a brži od Porter stemmer-a. Uvođenjem stemmer-a naraslo je vrijeme izvođenja pretprocesiranja i to za  $\sim 4$  puta (koristi se stemming Python modul). Trebalo bi koristiti neku brzu C/C++ implementaciju stemming postupka. Nije se koristio lematizator zato što je bitno očuvati riječi što je moguće više onakve kakve jesu. Primjerice, ne bi bilo dobro zamijeniti *is* s *be*, jer je korisnik možda baš htio dokumente gdje se pojavljuje *is*. Lematizator bi *is* riječ zamijenio s riječi *be* i tu bi se izbugila informacija (opet subjektivna procjena). Kada bi postojala potreba za lematizatorom tada bi se uz Python koristio spacy lematizator i to zato što je brži u odnosu na, primjerice, nltk lematizatora.

## 2 Algoritmi

### 2.1 Bag of words

Dvije su osnovne strukture podataka unutar algoritma. *Bag of words* unutar svakog dokumenta broji koliko se puta pojedina riječ pojavila unutar tog dokumenta. *Bag of documents* za svaku riječ broji koliko se puta pojavila u pojedinom dokumentu. Stvari su slične, ali omogućuju da se algoritam izvede u složenosti (*broj riječi upita \* broj dokumenata u kojima se pojavljuje riječ*). Pristigli upit najprije se pretprocesira, potom se za svaku riječ i za svaki dokument unutar kojeg se ta riječ pojavljuje računa normalizirana težina. Za svaku riječ gleda se koliko puta se ona pojavljuje unutar dokumenta i onda se ta vrijednost normalizira s veličinom dokumenta. Na taj način se postiže da kraći dokumenti u kojima broj pojavljivanje pojedine riječi čini veći postotak unutar dokumenta imaju snažniji utjecaj nego dokumenti u kojima se rijeđe pojavljuje dotična riječ. Na kraju se sortiraju dobiveni dokumenti prema izračunatim težinama. Dobre strane algoritma

su brzina, jednostavnost, normalizacija, a mana je što se ne uzima nikakva globalna informacija o zastupljenosti riječi unutar cijelog skupa dokumenata. Uvjet koji bi trebao biti zadovoljen je da ne postoji dokument unutar kojeg se nalaze sve riječi iz korpusa. U suprotnom će algoritam sporije raditi od očekivanog, ali pretpostavka je da se to neće dogoditi. Na danom skupu dokumenata, vrijeme izvođenja algoritma je reda veličine 1ms i ne bi trebalo postati drastično veće na skupu od, primjerice, reda veličine 1M dokumenata.

## 2.2 Vector space

## 2.3 Binary independence

# 3 Implementacija

Svaki algoritam tipa je *IRAlgorithm*. Taj tip posjeduje četiri metode. *config* putem koje se prima nekakva konfiguracija algoritma, ako takva za dani algoritam postoji. *preprocess\_all* kroz koju se radi predprocesiranje svih dostupnih dokumenata. *preprocess\_one* putem koje se radi predprocesiranje samo jednog novog dokumenta. I, *run* koja vraća konkretne rezultate. Kako bi se izbjeglo pokretanje predprocesiranja svaki put kada se treba pokrenuti neki algoritam uveden je razred *AlgorithmBox* koji posjeduje instance svih dostupnih algoritama u varijabli *available\_algorithms*. U varijabli *prepared\_algorithms* nalaze se samo one instance algoritama za koje je već proveden proces predprocesiranja. Varijabla *prepared\_algorithms* se prilikom svakog postavljanja novih dokumenata kroz *setter files* postavi na prazni riječnik. Na taj način se postiže da se kod idućeg dohvaćanja nekog algoritma ponovno forsira proces predprocesiranja. Ukoliko se dodao samo jedan dokument onda se samo pozove metoda *preprocess\_one* nad svim algoritmima koji su unutar riječnika *prepared\_algorithms*. Mana ovog pristupa je što svaki algoritam za sebe drži kopiju svih dokumenata i ponavljanje predprocesiranja za svaki algoritam posebno. Prednost je, svojstvo da svaki algoritam ima izolirani skup podataka s kojim može činiti što je već potrebno kako bi algoritam valjano radio. Odlučio sam se za taj pristup zato što je memorija manji problem od nekakve jednostavnosti korištenja i važnosti da je svaki algoritam za sebe izoliran. Jedna od mana trenutne implementacija je i što se dokumenti pamte unutar riječnika. Optimalnije bi bilo pamtit i dokumente unutar liste, te imati potporne strukture podataka kao što su riječnici u kojima bi se pamtilo primjerice na kojoj poziciji u listi je određen dokument.

# 4 Upute za pokretanje

<https://infinity.buda.link>

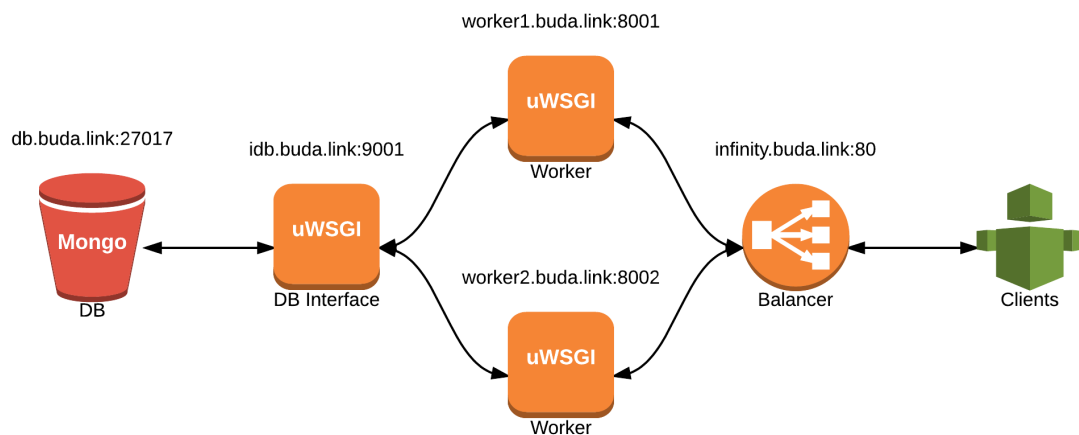
TODO: source setup.py

TODO: objasniti konzolu (history)

## 5 Producerska okolina

Svi upiti klijenata dolaze na Nginx load balancer koji ima izrazito veliku propusnost. Nakon toga load balancer prosledjuje upit na worker instance koje svaka za sebe imaju cijeli dataset i znaju vratiti odgovarajuci rezultat. Dataset worker instance preuzimaju od db interface instance, a ne direktno iz baze. Trenutno je u deploymentu samo jedna instanca sucelja prema bazi, ali u produkciji tu moze biti opet load balancer i vise instanci interface-a prema bazi podataka. Baza podataka je MongoDB, takodjer samo jedna instanca, no u praksi tu moze doci mongo replica set ili mongo shard cluster.

Kao sto se vidi na slici 1, sve instance imaju simbolicka imena (FQDN). To je takodjer izrazito bitno jer se time postize transparentnost pristupa i migracijska transparentnost. U konkretnoj implementaciji sva imena su definirana u `/etc/hosts` file-u na deploy stroju, no u produkciji ce imena biti definirana na redundantnim DNS serverima.



Slika 1: Producerska okolina

## 6 Web sučelje

TODO: Angular + materializecss

TODO: nekakva slika sučelja

P.S.

infinity > gugol