

Rangirajući pretraživač dokumenata - Infinity

Marko Budiselić

30.11.2015.

1 Pretprocesiranje

2 Algoritmi

2.1 Bag of words

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 određeni dokument.

4 Upute za pokretanje

`https://infinity.buda.link`

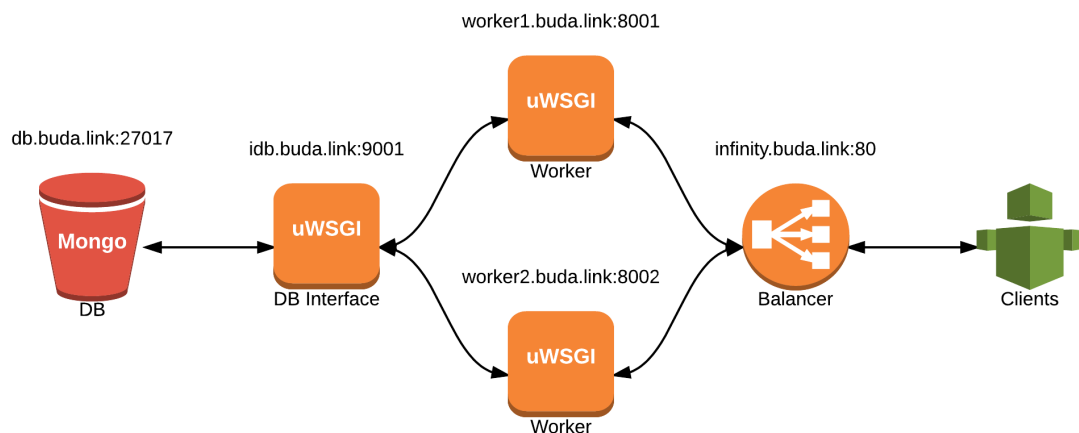
TODO: source setup.py

TODO: objasniti konzolu (history)

5 Produkcijaska okolina

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

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



Slika 1: Produkcijaska okolina

6 Web sučelje

TODO: Angular + materializecss

TODO: nekakva slika sučelja

P.S.

infinity > gugol