

## 6. Full-text pretraga informacija

Sunday, 20 February 2022 16:09

### Information retrieval

Najpoznatiji sistem za pretragu informacija predstavljaju internet pretraživači (google)

Radi na osnovu poredjenja upita korisnika sa netekstuiranim dokumentima i vraćanje upita sortiranog po relevantnosti

Zbog velicine podataka imamo razne tehnike pretraživanja

Tehnike:

- Bazirana na pretrazi metapodataka dokumenta

  - Pretražuju se samo meta podaci ne sadržaj dokumenta

- Bazirana na pretrazi full-text indeksima (indeksima kreiranim na osnovu sadržaja)

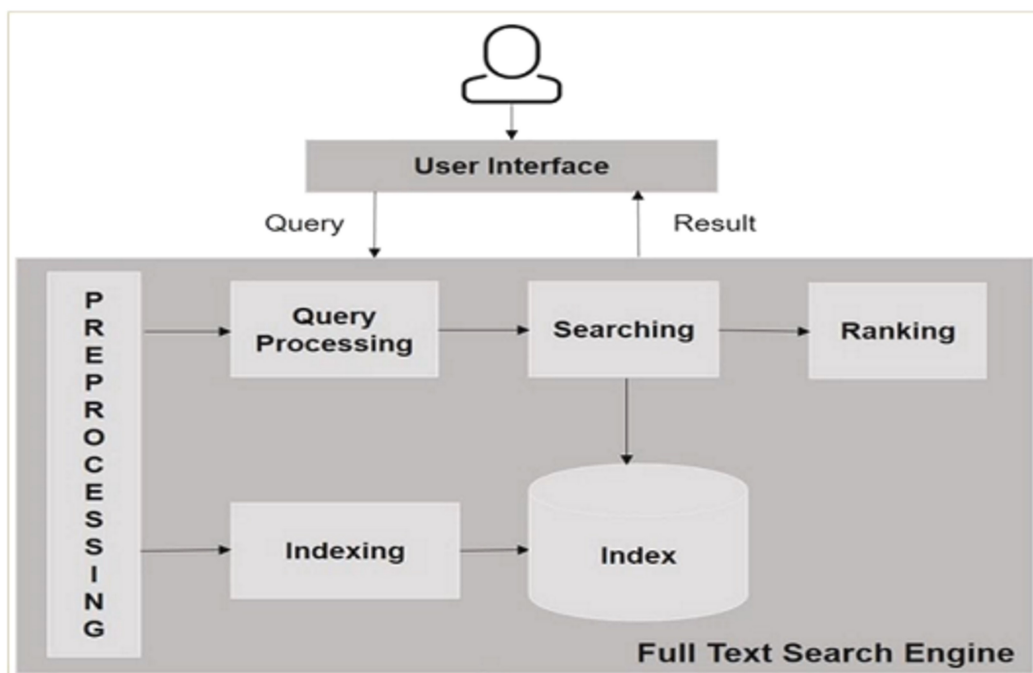
  - Pretražuje se kompletan sadržaj dokumenta

  - Bitni aspekti pretrage

    - Relevantnost -

    - Analiza - pretvara blok teksta u tokene (token je celina, npr rec)

- Mesavina ove dve



### Preprocesiranje

- Kroz ovo prolaze i dokumenti i korisnicka pretraga

- Ovo je samo analiza podataka kako bi se indeksirali

- A radi se i u korisnickoj pretrazi kako bi se sveli na iste mere i olaksali pretragu

  - Kao svodjenje svega na mala slova itd.

- Posle se rankira i vraća korisniku rezultat

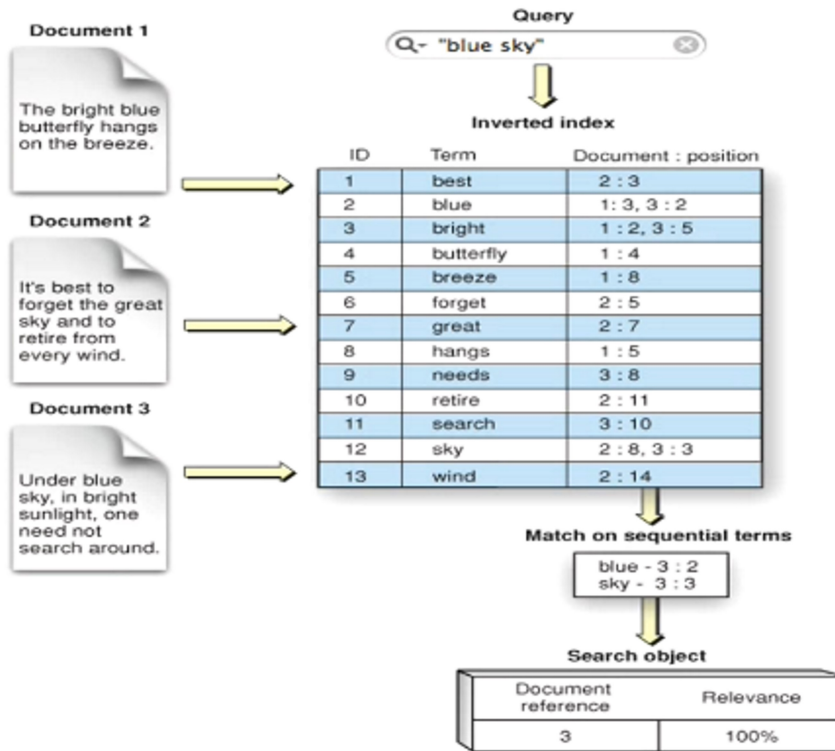
Dodatne transformacije mogu da se dese kao ako pogresno spelujemo rec ispravice se pa ce se potraziti

### Invertovani indeks

Struktura koja se koristi za pretragu full-text

- Vokabular - lista svih reci koje se pojavljuju u dokumentu

Pojavljivanja - lista koja sadrži info o svakoj reci vokabulara  
 Gde u dokumentu se pojavljuje rec  
 Koliko puta se pojavljuje rec  
 Dokument u kome se pojavljuje rec



Rec blue se pojavljuje u dokumentu 1 na poziciji 3

#### Pretraživanje

Rezultat upita je ono gde se pojavljuju obe reci u tekstu  
 100% je ovde al može da bude i manje kao za dokument 1 gde se blue javlja jednom tj 50% relevantnosti

Lucene indeks je najpopularnija biblioteka za invertovani indeks  
 Java jezik

#### Elasticsearch

Platforma za skladištenje i pretragu dokumenta bazirana na Lucene indeks-u  
 Pretraga sadržaja dokumenta  
 Implementirano Java, skladišti info kao JSON

#### Pojmovi:

Klaster - kolekcija jednog ili više servera zaduženi da zajedno čuvaju podatke i indeksuju ih i omoguće pretragu  
 Čvor - server na kome se čuvaju podaci  
 Indeks - kolekcija dokumenta koje imaju neke zajedničke karakteristike  
 Tip - Logičko particionisanje indeksa

Dokument - osnovna informativna jedinica (json)

Indeksiranje dokumenata - pribavljanje, obrada, analiza dokumenta i azuriranje indeksa

Shard - velike količine podataka se dele na male delove, nalaze se na raznim nodovima u klasteru

Svaki shard je nezavisan, ako otkaze jedan shard te informacije su izgubljene ali ostali čvorovi koji rade su nezavisni pa nastave da rade

Replike - svaki shard može da ima jednu ili više replika zbog otkaza

## TERMINOLOGY

MySQL	Elastic Search
Database	Index
Table	Type
Row	Document
Column	Field
Schema	Mapping
Index	Everything is indexed
SQL	Query DSL
SELECT * FROM table ...	GET http://...
UPDATE table SET ...	PUT http://...

Analizator ako ne navedemo dobijamo standardni  
Menja sva slova u mala  
Izbacuje nepotrebne reci kao clanovi, veznici itd.

Ima mnogo analizatora