

5. zadanie zamerané na Elastic:

Odovzdanie do 12.12. 23:59. Zadania 1-8 je dokopy za 7,5 boda tak isto ako aj 9 – 11.

Dokopy teda za 15 bodov. Odovzdávate dokument s popisom a queries separátne každú ako json čo posielate a json čo dostanete ako odpoveď.

1. Rozbehajte si 3 inštancie Elasticsearch-u
 2. Vytvorte index pre Tweety, ktorý bude mať "optimálny" počet shardov a replík pre 3 nody (aby tam bola distribúcia dotazov vo vyhľadávaní, aj distribúcia uložených dát)
 3. Vytvorte mapping pre normalizované dáta z Postgresu - Tweet musí obsahovať údaje rovnaké ako máte už uložené v PostgreSQL. Dbajte na to, aby ste vytvorili polia v správnom dátovom type (polia ktoré má zmysel analyzovať analyzujte správne, tie ktoré nemá, aby neboli zbytočne analyzované (keyword analyzer)) tak aby index nebol zbytočne veľký. Mapovanie musí byť striktné.
 4. Pre index tweets vytvorte 3 vlastné analyzéry (v settings) nasledovne:
 1. Analyzér "englando". Tento analyzér bude obsahovať nasledovné:
 - 1.2. filtre: english_possessive_stemmer, lowercase, english_stop, english_stemmer,
 - 1.3. char_filter: html_strip
 - 1.4. tokenizer: štandardný
 - ukážku nájdete na stránke elastic.co pre anglický analyzér
 2. Analyzér custom_ngram:
 - 2.2. Filtre: lowercase, asciifolding, filter_ngrams (definujte si ho sami na rozmedzie 1-10)
 - 2.3. char_filter: html_strip
 - 2.4. tokenizer: štandardný
 3. Analyzér custom_shingles:
 - 3.2. Filtre: lowercase, asciifolding, filter_shingles (definujte si ho sami a dajte token_separator: "")
 - 3.3. char_filter: html_strip
 - 3.4. tokenizer: štandardný
- Do mapovania pridajte:
1. každý anglický text (rátajme že každý tweet a description u autora je primárne v angličtine) nech je analyzovaný novým analyzérom "englando"
 2. Priradte analyzery
 - a. author.name nech má aj mapovania pre custom_ngram, a custom_shingles,
 - b. author.screen_name nech má aj custom_ngram,
 - c. author.description nech má aj custom_shingles. Toto platí aj pre mentions, ak tam tie záznamy máte.
 3. Hashtagy indexujte ako lowercase
5. Vytvorte bulk import pre vaše normalizované Tweety.
 6. Importujete dáta do Elasticsearchu prvych 5000 tweetov
 7. Experimentujte s nódami, a zistite koľko nódov musí bežať (a ktoré) aby vám Elasticsearch vedel pridávať dokumenty, mazať dokumenty, prezerať dokumenty a vyhľadávať nad nimi? Dá sa nastaviť Elastic tak, aby mu stačil jeden nód?

8. Upravujte počet retweetov pre vami vybraný tweet pomocou vášho jednoduchého scriptu (v rámci Elasticsearchu) a sledujte ako sa mení `_seq_no` a `_primary_term` pri tom ako zabíjate a spúšťate nódy.
9. Zrušte repliky a importujete všetky tweety
10. Vyhľadajte vo vašich tweetoch spojenie "gates s0ros vaccine micr0chip". V query použite `function_score`, kde jednotlivé medzikroky sú nasledovné:

Query:

1. Must - vyhľadajte vo viacerých poliach (konkrétne: `author.name` (pomocou `shingle`), `content` (cez analyzovaný anglický text), `author.description` (pomocou `shingles`), `author.screen_name` (pomocou `ngram`)) spojenie "gates s0ros vaccine micr0chip", zapojte podporu pre preklepy, operátor je OR.
 - 2.1 tieto polia vo vyhľadávaní boost-nite nasledovne - `author.name` * 6, `content` * 8, `author.description` * 6, `author.screen_name` * 10.
3. Filter - vyfiltrujte len tie, ktoré majú `author.statuses_count` > 1000 a tie, ktoré majú hashtag „qanon“
4. Should – boost-nite 10 krat tie, ktoré obsahujú v `mentions.name` (tento objekt je typu `nested`) cez `ngram` string "real".
5. Nastavte podmienené váhy cez `functions` nasledovne:
 - 5.1. `retweet_count`, ktorý je väčší rovný ako 100 a menší rovný ako 500 na 6,
 - 5.2. `author.followers_count` väčší ako 100 na 3

Zobrazte agregácie pre výsledky na konci. Vytvorte bucket hashtags podľa hashtagov a spočítajte hodnoty výskytov (na webe by to mohli byť facety).

11. Konšpiračné teórie podľa Elasticu. Pracujte zo všetkými tweetami, ktoré máte. Následne pre všetky týždne zistite pomocou vnorených agregácií, koľko `retweet_count` sumárne majú tweety ktoré majú hashtagy z prvého zadania. Teda na základe hashtagov znova rozdeľte tweety do konšpiračných teórií ale pomocou agregácií.