

VILNIAUS UNIVERSITETAS
MATEMATIKOS IR INFORMATIKOS FAKULTETAS
INFORMATIKOS KATEDRA

Baigiamasis bakalauro darbas

Dokumentų klasterizacija
(Document clustering)

Atliko: 4 kurso 1 grupės studentas

Dominykas Ablingis (parašas)

Darbo vadovas:

lekt. Rimantas Kybartas (parašas)

Recenzentas:

doc. dr. Vardenis Pavardenis (parašas)

Vilnius
2018

Turinys

Sąvokų apibrėžimai	2
Įvadas	3
1. Pagrindinė tiriamoji dalis	4
2. Algoritmų testavimas/kokybės vertinimas	5
3. Rezultatų vizualizacija	6
4. Susijusios, aktuali problemos	7
4.1. keyword extraction.....	7
4.2. cluster labeling.....	7
4.3. document classification	7
4.4. dimensionality reduction	7
Išvados	8
Conclusions	9

Sąvokų apibrėžimai

- cluster
- clustering
- document
- corpus
- term – terminas dar vadinamas (tokens, words, terms or attributes)
- token
- feature
- online
- offline
- supervised
- unsupervised
- matmenų redukcija (angl. *dimensionality reduction*)

Įvadas

Šiais laikais kai kiekvienas žmogus turintis prieigą prie interneto gali dalintis informacija, daugybė knygų yra skaitmenizuojamos kiekvieną dieną ir mokslo institucijos dalinasi savo moksline informacija, pasiekamos informacijos kiekis didėja su kiekviena diena ir pagrindinė problema nebėra informacijos trūkumas, o atradimas ko reikia. Tam spresti buvo ir yra kuriama įvairūs mechanizmai: paieškos varikliai...Šiame darbe autorius nagrinės viena iš šios problemos sprendimo metodų, klasterizacijos algoritmus ir jų panaudojimą tekstiniams dokumentams.

Šiandieninėje visuomenėje prieinamos informacijos kiekis didėja su kiekviena diena ir pagrindinė problema yra ne informacijos trūkumas, o jos gausa ir galimybė surasti reikiamą informaciją. Paieškos programos gali pateikti didelius kiekius tekstinių dokumentų, bet reikiamo dokumento paieška užima daug laiko ir be to, ne visada gauta informacija atitinka paiešką. Tekstinių dokumentų paieškos procesui palengvinti ir pagreitinti gali būti taikomi klasterizavimo metodai, kurie yra pakankamai gerai žinomi ir jau seniai naudojami duomenų klasterizavimui. Klasterizavimo metu dokumentai pagal savo turinį suskirstomi į klasterius vadovaujantis tam tikrais kriterijais, pvz., pagal temą, pagal dokumento naujumą ir pan.

1. Pagrindinė tiriamoji dalis

Darbo tikslas – mokslinės literatūros apie dokumentų apdorojimą ir klasterizavimą apžvalga ir analizė

Darbo uždaviniai:

1. Panaudojimo sritys
2. Dokumentų klasterizavimui naudojamų įrankių ir metodų apžvalga.
3. Dokumentų klasterizavimo proceso žingsnių aprašymas
4. Kylantys iššūkiai
5. Susijusios problemos

Kursiniame darbe buvo siekta susipažinti su moksline literatūra, išanalizuoti egzistuojančius dokumentų klasterizavimo metodus ir įrankius, kurie bus išbandyti taikomojo pobūdžio kursiniame darbe apie lietuviškų dokumentų klasterizavimą.

Šiame darbe taip pat bus paminėtos kitos su dokumentu klasterizacija susijusios problemos ir ši informacija bus išdėstyta eilės tvarka kaip būtų sprendžiamos užduotys.

1. Duomenų išgavimas iš skirtingų tekstinių dokumentų formatų
2. Duomenų apdorojimas
3. Algoritmai
4. Algoritmų testavimas
5. Rezultatų vizualizacija

2. Algoritmų testavimas/kokybės vertinimas

viena iš fundamentalių neprižiūrime mokymosi problemų yra modelių testavimas. skirtingai nei „prižiurime mokymasi“ kur svarbu atdidėti dali duomenų su kuriais nėra mokomasi o tik testuojama sugeneruoti modeliai, „neprižiurime mokymasi“ mes to negalim atlikti...bet egzistuoja keltatas metodu kaip galima patikrinti sudarytus klasterius.

- pirma galima sugeneruotus klasterius leisti tikrinti **žmonėms** . tam yra keli būdai. galima tiesiog duoti sugeneruotus klasterius ir bandyti nuspresti ar jie tinkami. taipat galima parainkti du atsiktinius dokumentus ir spėti ar jie turētu būti viename ar atskiruose klasteriuose, ir tada palyginti su kompiuterio sugeneruotu rezultatu.
- taip yra metodų kaip atlikti testavimą automatiškai. vienas jų paimti 2(ar daugiau) dokumentus iš skirtingų klasterių ir apkeisti juos vietomis, tada patikrinti klasteriu //stipruma. tai atlike dokybe kartų galime spręsti kaip sėkmingai sekėsi klasterizacija, jeigu apmainant jų kokybė nukrito tai reikškia, kad dokumentai sėkmingai suklaserizuoti, bet jeigu nesikeite tai reiške, kad klasteriai mažai vienas nuo kito skiresi ir klasreizacija nesekminga.

3. Rezultatų vizualizacija

dažnai norėdami geriau pažinti (ar testuoti) klasterizacijos rezultatus mes galime juo vizualizuoti. vizualizacijos gali būti įvairios ir dažnai priklauso nuo algoritmo rūšies, bet dažniausiai naudojama **point cloud** vizualizacijos. jose matosi pagal kokius parametrus (ašis) buvo klasterizuojama ir kaip atrodo sudaryti klasteriai, taip pat galima pridėti papildomų indikatorių, priklausančių nuo algoritmo. pavyzdžiui klasteriams sudarytiems k-means metodu galima nubraižyti atitinkamus „centrinius taškus“.

4. Susijusios, aktuolios problemos

4.1. keyword extraction

4.2. cluster labeling

4.3. document classification

4.4. dimensionality reduction

dimensionality reduction methods can be considered a subtype of soft clustering; for documents, these include latent semantic indexing (truncated singular value decomposition on term histograms) and topic models.

Išvados

Conclusions

Šiame skyriuje pateikiamos išvados (reziumė) anglų kalba.