

1. Projekt

Iskanje po zbirki dokumentov

Domen Gašperlin
Rok Grmek
Jakob Gaberc Artenjak
Anže Gregorc

Matemetično modeliranje, Fakulteta za računalništvo in informatiko

April, 2017

1 Opis problema

Namen projekta je izdelati iskalnik relevantnih dokumentov po ključnih besedah z metodo *latentnega semantičnega indeksiranja* (LSI), saj so metode, ki izberejo le dokumente, ki vsebujejo natanko iskane besede, precej nenatačne. Ljudje namreč uporabljamo veliko sopomenk, ki jih preproste metode ne povežejo. Metoda LSI zgradi model, ki združuje več besed v pojme in zato najde tudi dokumente, ki so relevantni, pa ne vsebujejo iskalne besede.

Izdelati je potrebno program, ki bo v dani zbirki za dane ključne besede poiskal najbolj relevantne dokumente.

2 Naloge

2.1 Iz zbirke dokumentov zgradite matriko **A** povezav med besedami in dokumenti.

Matrika **A**:

$$\mathbf{A} = \begin{array}{ccccc} & \text{doc1} & \text{doc2} & & \text{docD} \\ \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1D} \\ a_{21} & a_{22} & & \\ \vdots & & \ddots & \\ a_{B1} & & & a_{BD} \end{bmatrix} & \begin{array}{l} \text{beseda1} \\ \text{beseda2} \\ \vdots \\ \text{besedaB} \end{array} \end{array}$$

Vsak dokument ima v matriki svoj stolpec, vsaka beseda pa svojo vrstico. Element a_{ij} pa je frekvenca i -te besede v j -tem dokumentu.

Postopek gradnje:

```
number_of_docs = length(file_names); # shranimo stevilo vseh dokumentov
all_words = []; # inicializacija polja, ki bo vsebovala vse besede
num_of_words_in_docs = zeros(1, number_of_docs); # vektor, ki za vsak
    dokument hrani stevilo vseh besed
for i = 1:number_of_docs # sprehodimo se po vseh dokumentih
    # preberemo i-ti dokument
    doc = textread([path_to_docs, filesep, file_names{i}], '%s');
    # vse besede spremenimo na samo alfa numericne znake in v male crke
    for j = 1:length(doc)
        doc{j} = lower(doc{j}(isalnum(doc{j})));
    end
    # dodamo besede i-tega dokumenta v polje vseh besed
    all_words = [all_words; doc];
    # dodamo stevilo vseh besed v i-tem dokumentu
    num_of_words_in_docs(i) = length(doc);
end
```

Ko imamo zgrajeno polje vseh besed, moramo odstraniti podvojene besede in s tem ustvarimo polje, ki bo služilo kot stolpec v matriki A.

```
[unique_words, ~, numbers] = unique(all_words);
```

In sedaj imamo vse pripravljeno za gradnjo matrike A:

```
all_possible_numbers = (1:length(unique_words))'; # vektor od 1 do st
vseh unikatnih besed

# matrika A dimenzije (st. vseh unikatnih besed) x (st. vseh dokumentov)
A = zeros(length(unique_words), number_of_docs);
doc_end = 0;

# sprehodimo se po vseh dokumentih
for i = 1:number_of_docs
    # hranita stevilo, pri kateri se zacnejo in koncajo besede i-tega
    dokumenta v polju vseh besed
    doc_start = doc_end + 1;
    doc_end = doc_start + num_of_words_in_docs(i) - 1;

    # dodamo frekvence i-tega dokumenta v matriko A
    A(:, i) = histc(numbers(doc_start:doc_end, 1), all_possible_numbers);
end
```