

Metody Obliczeniowe w Nauce i Technice

Laboratorium 7

Singular Value Decomposition

Wybrany przeze mnie zbiorem dokumentów są 2613 książek, artykułów udostępnionych w ramach projektu Gutenberg. Pobrane zostały przy pomocy biblioteki dostępnej biblioteki gutenberg. Kod do pobierania dokumentów znajduje się w pliku *download.py*.

Wyznaczony zbiór termów jest to iloczyn pierwszych 5000 słów (wcześniej zmodyfikowanych) z dokumentów oraz zbioru słów ze słownika *words.txt*. Operacje na słowach znajdują się w pliku *word_checking.py*. Natomiast tworzenie zbioru termów w pliku *do_matrix.py* a sam zbiór zapisany jest w pliku *union.txt*.

Częstość występowania słów została policzona również w pliku *do_matrix.py* a sama częstość występowania danych słów została zapisana w pliku *matrix.txt*, natomiast lista wszystkich dokumentów została zapisana w pliku *name.txt*.

W pliku *preapare.py* znajdują się wszystkie operacje na macierzy tj mnożenie przez IDF (inverse document frequency) oraz normalizacja macierzy. W tym pliku obliczane są również różne macierze SVD przy pomocy funkcji biblioteki Scipy dla rzadkich macierzy.

Obsługa zapytań znajduje się w pliku *finder.py* tam z wykorzystaniem policzonych wcześniej pliku odbywa się wyszukiwanie.

Wyniki z usuwaniem szumu oraz bez usuwania znacząco się różnią (tym bardziej im mniejsza jest wartość k), często fragment wzięty z tekstu nie był w pierwszej 10 najlepszych wyników dla sposobu bez usuwania szumu. Najlepsze wyniki zapytań otrzymywałem dla k rzędu 1400

Jeśli nie użyje IDF to znajduje się pewna mała pula dokumentów która w większości przypadków zajmuje 10 najwyższych miejsc. Natomiast jeśli użyje się IDF to problem ten nie występuje.

Program z usuwaniem szumu działa wolniej, jest to spowodowane ładowaniem większego pliku do pamięci oraz tym że macierz po SVD nie jest tak samo rzadka jak ta bez usuwania.