

# Przetwarzanie Języka Naturalnego

## Lab 8 – LSA

Aleksander Smywiński-Pohl

Wydział IEiT  
Katedra Informatyki

11.04.2017

- ✘ każda składowa wektora to częstotliwość (lub waga) odpowiadającego jej słowa w danym tekście
- ✘ „bag of words”: nie uwzględniamy kolejności występowania wyrazów
- ✘ Wady:
  - wyrazy występujące tylko w jednym tekście niosą z sobą dużo informacji (vide prawo Zipfa), ale nie mają dużego wpływu na podobieństwo (vide np. metryka cosinusowa)
  - wyrazy występujące często nie niosą ze sobą informacji, a mają mocny wpływ na pozorne podobieństwo
  - zgodność tekstów sprowadza się do zgodności wyrazów

- ✘ usuwamy *hapax legomena*
- ✘ usuwamy wyrazy, które występują w więcej niż 70% tekstów
- ✘ w macierzy term-document wartość w danej komórce zawiera wagę danego wyrazu w danym tekście

- ✘ czasem nazywane Latent Semantic Indexing (LSI)
- ✘ metoda analizy podobieństwa między dokumentami i wyrazami oparta na tworzeniu zbioru *pojęć (concepts)*
- ✘ założenie: słowa bliskoznaczne pojawiają się w podobnych fragmentach tekstu
- ✘ rozkład macierzy term-document przy pomocy dekompozycji głównych składowych (ang. *singular value decomposition - SVD*)
- ✘ zmniejszona zostaje liczba wierszy (wyrazów) przy zachowaniu podobieństwa między kolumnami (dokumentami)



$A$  – macierz term-document o wymiarach  $n \times m$

$$A = U\Sigma V^T$$

$U$  – macierz pojęć o wymiarach  $n \times l$  (wektory własne)

$\Sigma$  – przekątniowa macierz wartości własnych o wymiarach  $l \times l$

$V$  – macierz dokumentów o wymiarach  $m \times l$  (wektory własne)

Wymiary nowego układu współrzędnych wyznaczonego przez wektory własne to *pojęcia* lub *tematy* (*concepts, topics*).



Wybierając  $k$  największych wartości własnych dokonujemy redukcji wymiarów:

$$A' = U' \Sigma' V'^T$$

Wymiary macierzy  $U'$ ,  $\Sigma'$ ,  $V'$  to teraz kolejno:  $n \times k$ ,  $k \times k$ ,  $m \times k$

Możemy teraz porównywać wyrazy i dokumenty w przestrzeni o mniejszej liczbie wymiarów.

Zalety:

- ✕ oszczędność reprezentacji ( $k$  jest często rzędu setek)
- ✕ zwiększona skuteczność (usuwany jest szum)



Wady:

- ✘ pojęcia, jako wektory, często nie mają zrozumiałej dla człowieka postaci (niejasne komponenty powstałe przy redukcji wymiarów, ujemne wartości wag, etc.), na przykład:
  - $[(auto), (motor), (kwiat)] \rightarrow [(1.34 * auto + 0.28 * motor), (kwiat)]$
  - $[(auto), (butelka), (kwiat)] \rightarrow [(1.34 * auto + 0.28 * butelka), (kwiat)]$
- ✘ ponieważ waga dla każdego słowa to pewien punkt w przestrzeni, LSA jest nieczułe na polisemię
- ✘ konsekwencje wynikające z użycia modelu „bag of words”

- 1 zbudować model LSA np. przy użyciu biblioteki gensim (<http://radimrehurek.com/gensim/tutorial.html>). Proszę pamiętać o sprowadzeniu wyrazów do formy podstawowej (1 pkt.)
- 2 napisać program, znajdujący najbardziej zbliżone notatki do notatki wzorcowej (1 pkt.)
- 3 porównać wyniki wyszukiwania notatek przy użyci LSA z modelami z zadania 6 (1 pkt. – należy wykorzystać dane z zadania 6)

Materialy: [http:](http://home.agh.edu.pl/~zkaleta/pjn/lab6.tar.gz)

[//home.agh.edu.pl/~zkaleta/pjn/lab6.tar.gz](http://home.agh.edu.pl/~zkaleta/pjn/lab6.tar.gz)