

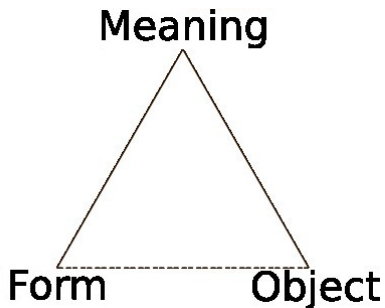
Przetwarzanie Języka Naturalnego

Lab 9 – Word Embedding

Aleksander Smywiński-Pohl

Wydział IEiT
Katedra Informatyki

16.05.2017





AGH

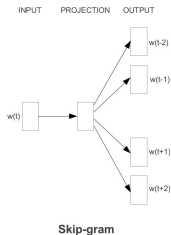
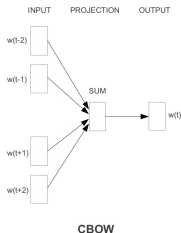
Jak opisać znaczenie słowa

- ⌘ definicja
- ⌘ konteksty (“a word is characterized by the company it keeps”)
- ⌘ relacje paradygmatyczne
- ⌘ relacje syntagmatyczne



- ✘ BOW, tf-idf itp.
- ✘ LSA (zliczanie)
- ✘ word embedding (przewidywanie)

- ✘ Continuous Bag of Words
- ✘ Skip-gram



Mikolov et al. "Efficient Estimation of Word Representations in Vector Space"



$$z = (z_1, z_2, \dots, z_k)$$

$$\sigma(z)_j = \frac{e^{z_j}}{\sum_{i=1}^k e^{z_i}}$$



AGH

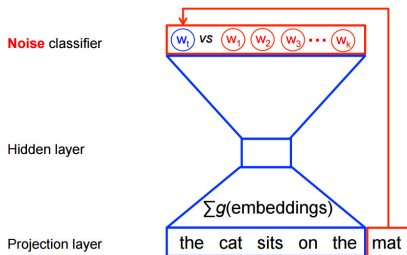
Algorytm

$$P(w_t|h) = \text{softmax}(\text{score}(w_t, h))$$

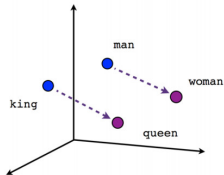


AGH

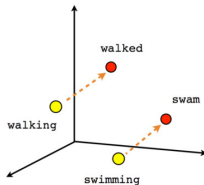
Algoritm c.d.



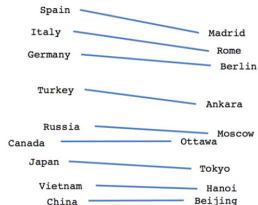
<https://www.tensorflow.org/tutorials/word2vec>



Male-Female



Verb tense



Country-Capital

<https://www.tensorflow.org/tutorials/word2vec>

- 1 wytrenować model word2vec na korpusie zawierającym minimum 100 mln słów w j. polskim (1 pkt)
- 2 wybrać trzy relacje morfosyntaktyczne dla różnych części mowy (np. dopełniacz–mianownik dla rzeczownika, czas teraźniejszy–czas przeszły dla czasownika) i przetestować ich zachowanie w wytrenowanym modelu (1 pkt)
- 3 wybrać trzy relacje semantyczne (np. państwo–stolica, część–całość) i przetestować ich zachowanie w wytrenowanym modelu (1 pkt)

- ✦ <https://www.tensorflow.org/tutorials/word2vec>
- ✦ Mikolov et al. “Efficient Estimation of Word Representations in Vector Space”
- ✦ Mikolov et al. “Distributed Representations of Words and Phrases and their Compositionality”