

0. Ембеддинг - информативное векторное представление объектов

Задачи NLP

- Sentiment analysis
- Spam filtering
- Fake news detection
- Topic prediction
- #hashtag prediction

text classification problems

1) bag of words - каждое предл/мечет - это 5 букв. написана в эмбеддинге

Препроцессинг данных

- токенизация: разбиение на части (на слова/слои/буквы)
- нормализация:
 - стемминг - отбрасыва до корня слова (примерно) (примеры: портмер, лапкастер, портмер 2)
 - лемматизация - отбрасыва до 1-ой формы слова для каждой части речи.
(примеры можно посмотреть в nltk)
- удаление цифр, стоп-слов и т.д.

TF-IDF

tf - term frequency - частота слова в документе

idf - inverse document frequency - логарифм отношения кол-ва док-тов к числу док-тов, где встречается это слово

Word representations via matrix factorization

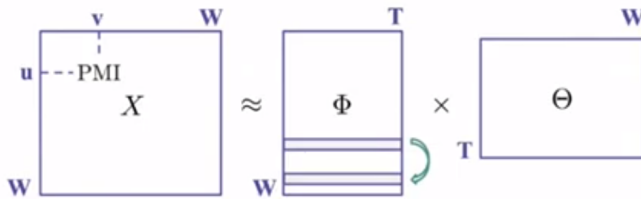
- **Input:** PMI, word cooccurrences, etc.
- **Method:** dimensionality reduction (SVD)
- **Output:** word similarities

: Pointwise Mutual Information (PMI)

$$PMI = \log \frac{p(u, v)}{p(u)p(v)} = \log \frac{n_{uv}n}{n_u n_v}$$

Much better solution: **Positive PMI (pPMI)**

$$pPMI = \max(0, PMI)$$



37

Word 2 vec

Embeddings: word2vec

Source Text

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

The quick brown fox jumps over the lazy dog. →

Training Samples

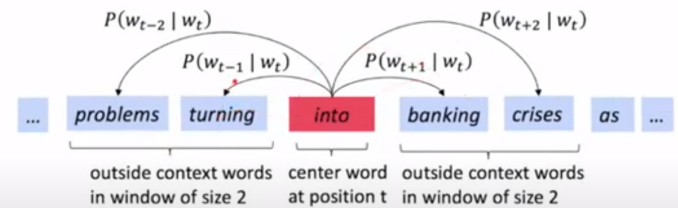
(the, quick)
(the, brown)

(quick, the)
(quick, brown)
(quick, fox)

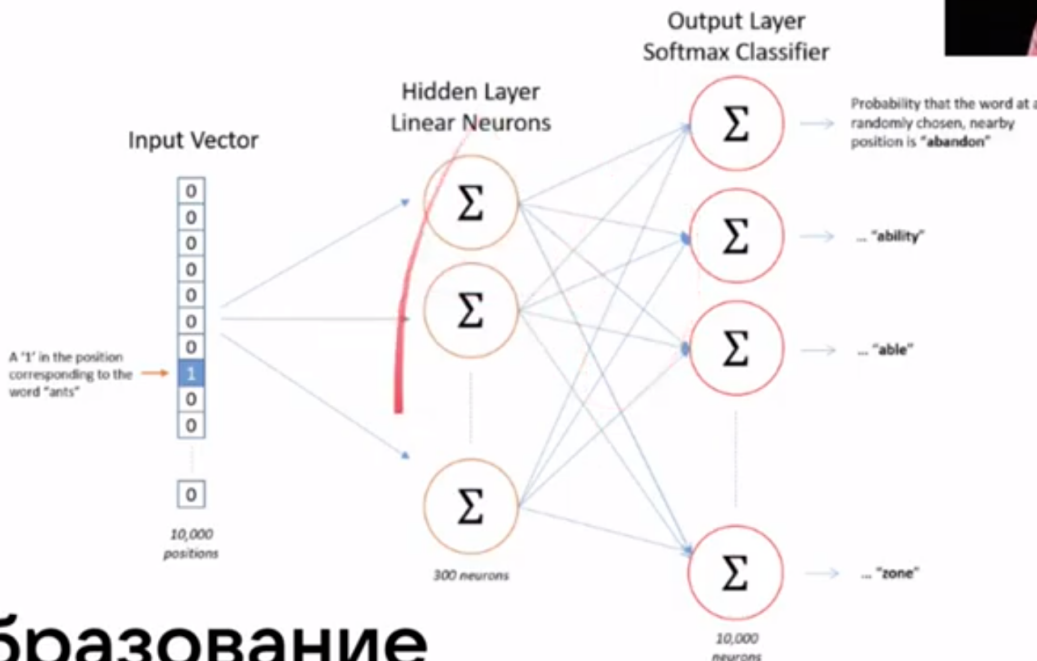
(brown, the)
(brown, quick)
(brown, fox)
(brown, jumps)

(fox, quick)
(fox, brown)
(fox, jumps)
(fox, over)

Embeddings: word2vec



Embedding



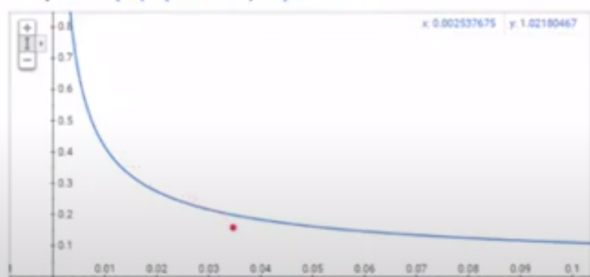
образование

Embeddings: word2vec

Subsampling frequent words.

w_i is the word, $z(w_i)$ is the fraction of this word in the whole text

Graph for $(\sqrt{x/0.001} + 1) \cdot 0.001/x$



$P(w_i)$ is the probability of keeping the word:

$$P(w_i) = \left(\sqrt{\frac{z(w_i)}{0.001}} + 1 \right) \cdot \frac{0.001}{z(w_i)}$$

Source: <http://mccormickml.com/2017/01/11/word2vec-tutorial-part-2-negative-sampling/>

Embeddings: negative sampling

Negative Sampling idea: only few words error is computed. All other words have zero error, so no updates by the backprop mechanism.

More frequent words are selected to be negative samples more often. The probability for selecting a word is just its weight divided by the sum of weights for all words.

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=0}^n (f(w_j)^{3/4})}$$

53

одинаково часто на равно вероятными словах и реже на равно вероятными

из равной вероятности расстояния класса

Continuous BOW (CBOW)

$$p(w_i | w_{i-h}, \dots, w_{i+h})$$

Predict center word from (bag of) context words

- Predicting one word each time
- Relatively fast

Skip-gram

$$p(w_{i-h}, \dots, w_{i+h} | w_i)$$

Predict context ("outside") words (position independent) given center word

- Predicting context by one word
- Much slower
- Better with infrequent words