

1) tokenization - разбиение последовательности (предложение/текст) на слова, пары слов/подслова.

2) удаление стоп-слов (слишком частые/бесполезные слова)

3) нормализация

а) стемины - обрезаем слова до корня

б) лемматизация - приводим слова в единую нормальную форму

замечка:

- пунктуация - отдельные токены
- если слова составные не важно, но всё приводим к нижнему регистру

Porter stemmer

- Published in 1979
- Base starting option

Lancaster stemmer

- Published in 1990
- The most aggressive
- Easy adding of your own rules

Snowball stemmer (Porter 2)

- Based on Porter
- More aggressive
- Most popular option now

Bag of words (BoW)

представим текст → каждому слову присваиваем индекс → получаем числовой вектор

the dog is on the table

0 0 1 1 0 1 1 1
are cat dog is now on table the

векторное представление последовательности

+

-

- очень просто 😊

- не учитывается порядок слов в предложении
- размерность векторов большая и разреженная
- векторы не нормализованы
- одни и те же слова могут иметь разные формы и окончания

замечка: мы можем учитывать локальный порядок слов, используя n-граммы.

проблема: их слишком много

- решение:
- удаляем слишком частые/редкие
 - удаляем 2-граммы с артиклями/союзными

идея: добавляем устойчивые словосочетания как отдельные токены

TF-IDF

вопрос: как определить важность слов для текста в BoW?

решение: попарно TF-IDF

term frequency - inverse document frequency

$$tf(t, d) = \#t \text{ в } d$$

- чем больше tf , тем важнее слово для документа
- чем больше idf , тем важнее слово t , тем сильнее дискриминируется важность слова в целом

$$idf(t) = \log \frac{N}{\#d, \text{ в которых есть } t}$$

замечка: если слово встречается во всех документах, то $idf(t) = \log \frac{N}{N} = \log 1 = 0$
слово будет иметь нулевую важность

идея: можно использовать $tf-idf$ как веса для усреднения, когда есть представления слов, а лексикон получил представление документа

давать хотя бы как-то использовать контекст

способ: (контекст = документ/предложение/п-грамма)

$PMI \rightarrow SVD \rightarrow$ векторные представления слов
(#слов на #слов) (#слов на K)

- Better solution: 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)**

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

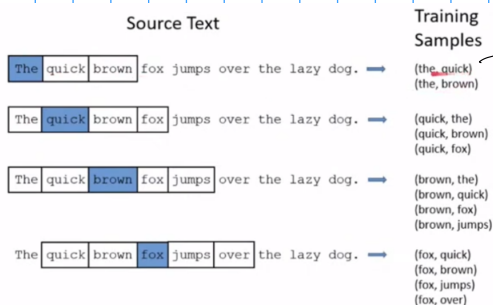
Word 2 Vec

идея метода: - давать по слову предсказывать контекст или слово по контексту

- слова, встречающиеся в похожих контекстах, близкие по значению и эмбедингам

выпрямляем:

- 1) Берём окно размера k и итерируемся им по тексту



похож Skip-gram

$O(1)$, потому что когда подаем 1 слово, то мы подаем строку с одной 1, то есть при умножении на $m \times n$ (матрица слов), мы берём 1 столбец.
Получается словарь

- 2) составляем парки для обучения зависимости от порядка (CBOW или Skip-gram)

а) CBOW - continuous bag of words
предсказываем слово по контексту

б) Skip-gram

предсказываем контекст по слову

Continuous BOW (CBOW)

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

Predict center word from (bag of) context words

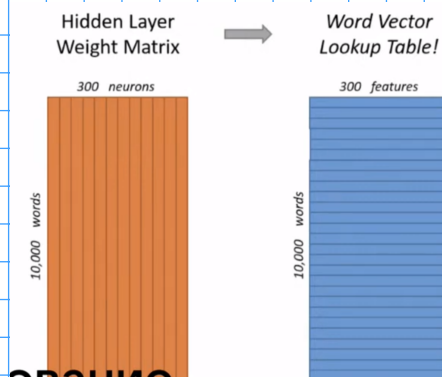
- Predicting one word each time
- Relatively fast

Skip-gram

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

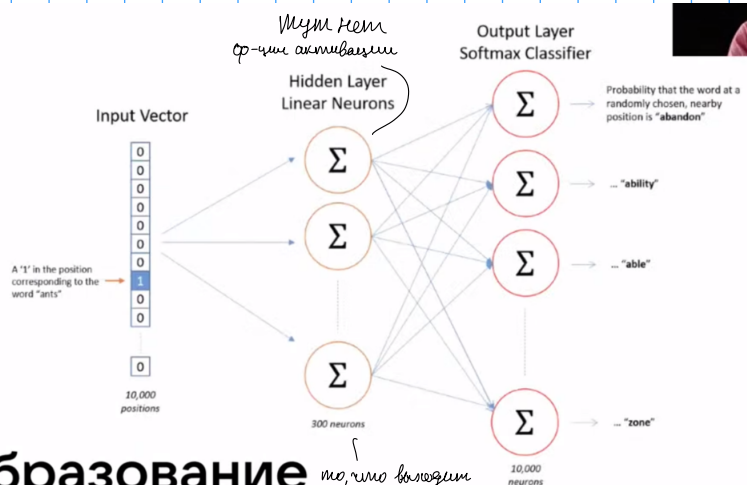
Predict context ("outside") words (position independent) given center word

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



ование

- 3) обучаем 1-слойную нейронку на классификацию



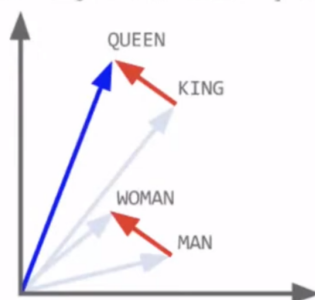
образование

то, что выводит со слова через классификацию считаем эмбедингом

- W2V можно строить

- св-во попарности эмб.

So king - man + woman = queen



интегральные моменты при обучении W2V

проблема: основная проблема - дисбаланс "классов" и перенос токенов в ток, на чьи обучаемся

решение:

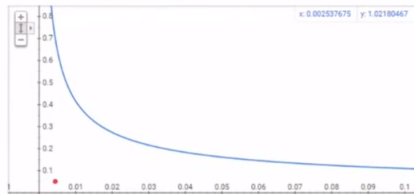
а) перенос токенов

Subsampling - оставляем часть токенов в зависимости от того как часто встречается ток

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/>

б) дисбаланс классов

нормализуем семанты - Токены не все равны, а только частоты
и чем чаще класс встречается, тем чаще он
показывается как 0 в других случаях