

Векторные представления слов: word2vec

Скачков Николай Андреевич
ВМК МГУ

2017/14/02

Как мы умеем представлять слова из словаря Ω в виде объектов?

Как мы умеем представлять слова из словаря Ω в виде объектов?

One-hot-encoding:

$$w \longrightarrow v_w : \quad v_w^i = [\#w = i]$$

У этого способа кодирования очень много недостатков!

Как мы умеем представлять слова из словаря Ω в виде объектов?

One-hot-encoding:

$$w \longrightarrow v_w : \quad v_w^i = [\#w = i]$$

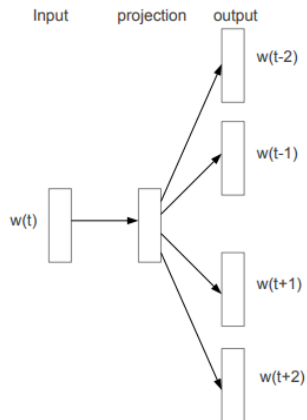
У этого способа кодирования очень много недостатков!
Хотим получить:

- представление слов в пространстве H невысокой размерности
- в данном пространстве H близкие по смыслу слова находятся рядом
- пространство H – линейное:

$$\text{vec}(\text{Madrid}) - \text{vec}(\text{Spain}) + \text{vec}(\text{France}) = \text{vec}(\text{Paris})$$

- **Идея:** Близкие слова имеют схожие контексты.
- Построим модель, которая будет хорошо предсказывать слова по контексту.

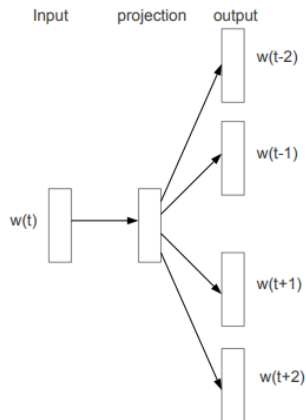
Source Text	Training Samples						
<table><tr><td>The</td><td>quick</td><td>brown</td></tr></table> fox jumps over the lazy dog. ➡	The	quick	brown	(the, quick) (the, brown)			
The	quick	brown					
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td></tr></table> jumps over the lazy dog. ➡	The	quick	brown	fox	(quick, the) (quick, brown) (quick, fox)		
The	quick	brown	fox				
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td></tr></table> over the lazy dog. ➡	The	quick	brown	fox	jumps	(brown, the) (brown, quick) (brown, fox) (brown, jumps)	
The	quick	brown	fox	jumps			
<table><tr><td>The</td><td>quick</td><td>brown</td><td>fox</td><td>jumps</td><td>over</td></tr></table> the lazy dog. ➡	The	quick	brown	fox	jumps	over	(fox, quick) (fox, brown) (fox, jumps) (fox, over)
The	quick	brown	fox	jumps	over		



- Параметры W , W' :

$v_w = W_w$ – input vector.

$v'_w = W'_w{}^T$ – output vector.

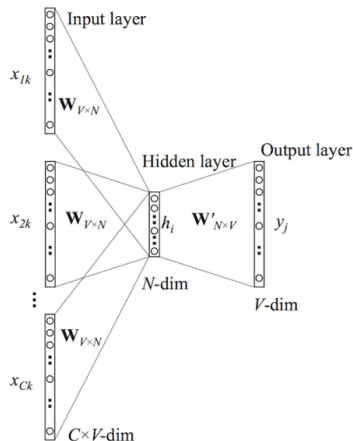


- Параметры W, W' :
 $v_w = W_w$ – input vector.
 $v'_w = W'^T_w$ – output vector.
- Подсчет условной вероятности (softmax):

$$p(w_o | w_i) = \frac{\exp(v'_{w_o}{}^T v_{w_i})}{\sum_{w \in \Omega} \exp(v'_w{}^T v_{w_i})}$$

- Функция потерь:

$$L(w_t) = - \sum_{|j|=1}^k \log p(w_{t+j} | w_t)$$



- Параметры W, W' :
 $v_w = W_w$ – input vector.
 $v'_w = W'_w{}^T$ – output vector.
- Подсчет условной вероятности (softmax):

$$h = \frac{1}{2k-2} \sum_{|j|=1}^k v_{w_j}$$

$$p(w_o|h) = \frac{\exp(v'_{w_o}{}^T h)}{\sum_{w \in \Omega} \exp(v'_w{}^T h)}$$

- Функция потерь:

$$L(w_t) = - \sum_{|j|=1}^k \log p(w_t|h)$$

Skip-gram:

- + Хорошо предсказывает контексты по редким словам, так как модель предсказывает по ним их соседей.
- Обучается в разы дольше, чем CBOW.

CBOW:

- + Быстрое обучение
- Редкие слова почти не влияют на обучение модели и выучиваются значительно хуже.

Model	Vector Dimensionality	Training words	Accuracy [%]			Training time [days]
			Semantic	Syntactic	Total	
3 epoch CBOW	300	783M	15.5	53.1	36.1	1
3 epoch Skip-gram	300	783M	50.0	55.9	53.3	3

- Количество параметров:

$$|\Omega| = 10000, \quad |h| = 300 \quad \longrightarrow \quad |W| = 3 \cdot 10^9$$

- Вычисление градиента softmax происходит за $O(|\Omega|)$

- Количество параметров:

$$|\Omega| = 10000, \quad |h| = 300 \quad \longrightarrow \quad |W| = 3 \cdot 10^9$$

- Вычисление градиента softmax происходит за $O(|\Omega|)$

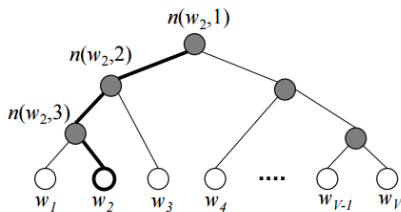
Необходимо оптимизировать вычисление функции потерь!

Для этого изменим ее:

- 1 Иерархический softmax
- 2 Negative Sampling

Над словами строится бинарное дерево. Вероятность i -го определяется как вероятность попасть в i -ый лист дерева из корня:

Над словами строится бинарное дерево. Вероятность i -го определяется как вероятность попасть в i -ый лист дерева из корня:



$$p(w|w_l) = \prod_i^{L(w)-1} \sigma([n(w, j+1) = ch(n(w, j))] \cdot v'_{n(w, j)}{}^T v_l),$$

где $v'_{n(w, j)}$ – дополнительные параметры модели.

Теперь функция потерь вычисляется за $O(\log |\Omega|)$

- Наша функция потерь поощряет одно значение быть 1, а все остальные – 0.
- Изменим функцию потерь так, чтобы градиент считался только по тому значению, которое ожидается быть близким к 1 и $q \in [5, 20]$ другим значениям.

$$\log p(w_O | w_I) \leftarrow \log \sigma(v'_{w_O}{}^T v_{w_I}) + \sum_{i=1}^q \mathbb{E}_{w_i \sim P(w)} \left[\sigma(-v'_{w_i}{}^T v_{w_I}) \right]$$

$P(w)$ – считаем, используя счётчики в степени 3/4

- Частотные слова не несут в себе большой смысловой нагрузки, но при этом значительно увеличивают размер выборки.
- Будем удалять слова из обучающей выборки с вероятностью: $p(w) = 1 - \sqrt{\frac{t}{\text{freq}(w)}}$, $t = 10^{-5}$

- Частотные слова не несут в себе большой смысловой нагрузки, но при этом значительно увеличивают размер выборки.
- Будем удалять слова из обучающей выборки с вероятностью: $p(w) = 1 - \sqrt{\frac{t}{\text{freq}(w)}}$, $t = 10^{-5}$

Method	Dimensionality	No subsampling [%]	10^{-5} subsampling [%]
NEG-5	300	24	27
NEG-15	300	27	42
HS-Huffman	300	19	47

- Negative Sampling работает быстрее, чем Hierarchical Softmax, но хуже работает с менее частотными словами.
- Стандартный выбор: skip-gram + negative sampling + окно (5-10) + Subsampling + ($|h| = 300$)

